

Adversarially Constructed Evaluation Sets Are More Challenging, but May Not Be Fair

Anonymous ACL submission

Abstract

Large language models increasingly saturate existing task benchmarks, in some cases outperforming humans, leaving little headroom with which to measure further progress. Adversarial dataset creation, which builds datasets using examples that a target system outputs incorrect predictions for, has been proposed as a strategy to construct more challenging datasets, avoiding the more serious challenge of building more precise benchmarks by conventional means. In this work, we study the impact of applying three common approaches for adversarial dataset creation: (1) filtering out easy examples (AFLite), (2) perturbing examples (TextFooler), and (3) model-in-the-loop data collection (ANLI and AdversarialQA), across 18 different adversary models. We find that all three methods can produce more challenging datasets, with stronger adversary models lowering the performance of evaluated models more. However, the resulting ranking of the evaluated models can also be unstable and highly sensitive to the choice of adversary model. Moreover, we find that AFLite oversamples examples with low annotator agreement, meaning that model comparisons hinge on the examples that are most contentious for humans. We recommend that researchers tread carefully when using adversarial methods for building evaluation datasets.

1 Introduction

Large-scale language models have attained strong performance across a variety of language understanding tasks, including question-answering, natural language inference (NLI), coreference resolution and paraphrase identification. Standard benchmarking tasks such as SQuAD (Rajpurkar et al., 2016; Lee et al., 2020) and multi-task benchmarks such as GLUE (Wang et al., 2018) and SuperGLUE (Wang et al., 2019) have seen models attain scores higher than humans. This has left little headroom

with which to measure further improvements in models and progress in NLP.

Prior work such as Le Bras et al. (2020) and Nie et al. (2020a) have proposed to construct more challenging datasets adversarially, either by only selecting examples that a given model predicts incorrectly, or by constructing new examples to deliberately stump a model. Both of approaches aim to raise the difficulty of task datasets by leveraging highly capable models (known as the *adversary model*) to assist with example selection or creation. However, one potential issue is that an adversarially constructed dataset that targets a specific model may bias the resulting data, creating datasets that are unduly challenging for one class of models but not others (Bowman and Dahl, 2021).

In contrast to other work focused on adversarial dataset creation for training (Wallace et al., 2021) or training and evaluation data (Le Bras et al., 2020; Nie et al., 2020b), we focus solely on evaluation data, and whether the choice of adversary model can introduce unwanted biases into an evaluation dataset. Ideally, an adversarially created dataset should be more difficult for all models, regardless of the choice of the adversary. In this work, we investigate three different approaches to create more challenging task evaluation datasets using adversary models: (1) *adversarial filtering*, which filters out examples from a static dataset that are identified to be easy for a given adversary model, (2) *adversarial perturbation*, wherein examples are modified to reduce performance of an adversary model, and (3) *model-in-the-loop adversarial data collection*, where human annotators interactively create examples that stump an adversary model.

For adversarial *filtering*, we study AFLite (Sakaguchi et al., 2020; Le Bras et al., 2020), an algorithm that identifies challenging subsets of an existing dataset. For adversarial *perturbation*, we evaluate TextFooler (Jin et al., 2019), a popular method for adversarially perturbing examples to

083 lower performance on a target model via word sub-
084 stitution. We apply both AFLite and TextFooler
085 in experiments across four English-language NLP
086 datasets and 18 different models to study the in-
087 teraction between the choice of adversary model
088 and the evaluation and ranking of systems on the
089 resulting dataset. For adversarial data collection,
090 we evaluate a range of models against two adversar-
091 ially collected datasets: ANLI (Nie et al., 2020a)
092 and AdversarialQA (Bartolo et al., 2020).

093 We find that all three classes of methods do re-
094 sult in more challenging evaluation datasets, but
095 with some notable drawbacks. For both adversarial
096 filtering and adversarial perturbation, the general
097 outcome is to lower performance across the board,
098 with stronger adversary models leading to more
099 challenging subsets of examples. However, with
100 both methods, the relative order of model perfor-
101 mance is not preserved, with large random varia-
102 tion in model ranks as stronger adversaries are used.
103 Performance on the resulting datasets is also much
104 worse if the evaluated and adversary models are de-
105 rived from same pretrained model, which can lead
106 to the difficulty of the adversarially constructed
107 dataset being overstated. Furthermore, adversar-
108 ial filtering tends to oversamples examples with
109 low annotator agreement, which means that the
110 selected examples are often contentious even for
111 human annotators. On the other hand, TextFooler
112 perturbations can introduce errors and lead to label
113 flips, and the distortion in resulting model rankings
114 is generally larger for TextFooler than for AFLite.
115 Jointly, these suggest that using adversarially fil-
116 tered or perturbed datasets naively for benchmark-
117 ing models is problematic.

118 We find that adversarially *collected* datasets
119 ANLI (Nie et al., 2020a) and AdversarialQA (Bar-
120 tolo et al., 2020) are also more challenging for all
121 models while also showing signs of disproportion-
122 ately disadvantaging the adversary model. How-
123 ever, with only a small number of such datasets
124 available, it is difficult to draw strong conclusions
125 about the overall efficacy or potential drawbacks of
126 the approach. Importantly, unlike for adversarial fil-
127 tering and perturbation, we cannot easily swap out
128 the adversary model for analysis, as the adversarial
129 data collection procedure can involve hundreds of
130 hours of human labor per adversary.

131 In all three cases, our findings do not outright
132 preclude the viability of adversarial dataset creation
133 for evaluation purposes, but we urge researchers

134 to keep these issues in mind when evaluating or
135 comparing models based on adversarial datasets.

2 Related Work 136

137 AFLite is an adversarial filtering algorithm pro-
138 posed by Sakaguchi et al. (2020), which also in-
139 troduced Winogrande, an adversarial Winograd
140 Schema Challenge dataset. Le Bras et al. (2020)
141 provided theoretical and empirical justification for
142 AFLite, showing that models trained on AFLite-
143 filtered data generalize better to out-of-domain
144 datasets. Other datasets constructed using adversar-
145 ial filtering include SWAG (Zellers et al., 2018) and
146 HellaSwag (Zellers et al., 2019), two adversarially
147 filtered commonsense multiple-choice datasets.

148 Given the over-parameterization of deep neu-
149 ral networks, adversarial perturbations (Goodfel-
150 low et al., 2015) have been identified as a par-
151 ticular weakness of these models. Within NLP,
152 most adversarial attacks tend to occur at the token
153 or word rather than continuous embedding level.
154 TextFooler (Jin et al., 2019) is a popular adver-
155 sarial textual perturbation method that swaps out
156 words for synonyms while attempting to retain the
157 semantic content and grammaticality of the origi-
158 nal example. TextBugger (Li et al., 2019) works
159 in a similar manner by searching for replacement
160 words in the neighborhood within a context-aware
161 embedding space, as well as using character-level
162 edits. BERT-ATTACK (Li et al., 2020) and BAE
163 (Garg and Ramakrishnan, 2020) use BERT’s mask
164 language modeling capability to substitute tokens.

165 An alternative approach is to collect data using
166 a model in the loop, where human example-writers
167 are given immediate feedback on whether an ad-
168 versary model is able to correctly answer their ex-
169 ample, and are incentivized to write examples on
170 which the models fail. ANLI (Nie et al., 2020b)
171 is an adversarial NLI dataset with multiple rounds
172 of data collection. Williams et al. (2020) further
173 provide fine-grained analysis of examples arising
174 from ANLI’s creation procedure. Bartolo et al.
175 (2020) introduce AdversarialQA, an adversarial
176 question-answering dataset. Kiela et al. (2021) fur-
177 ther extend this approach, building a platform for
178 continuous human-and-model-in-the-loop data cre-
179 ation. Using adversarially collected data as training
180 data has been shown to lead to better performance
181 on other adversarial datasets, but worse on out-
182 of-domain datasets (Kaushik et al., 2021; Bowman
183 et al., 2020). However, models trained on adversari-

ally collected data through many successive rounds have been shown to attain better performance (Wallace et al., 2021). In this work, we choose instead to focus exclusively on using adversarial examples as evaluation data.

In concurrent work, Adversarial Glue (Wang et al., 2021) applying a range of textual adversarial attacks to a subset of GLUE tasks to build a new language understanding benchmark. Importantly, they find that many adversarial attacks are prone to generating invalid examples, and perform careful manual filtering of the resulting examples.

3 Experimental Setup

Models The focus of our investigation is how the filtered dataset changes based on the choice of the adversary model. We consider a diverse set of pre-trained Transformer models: BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), ALBERT (Lan et al., 2020), XLM-R (Conneau et al., 2020), ELECTRA (Clark et al., 2020), MiniBERTa (Zhang et al., 2021), BART (Lewis et al., 2020), and DeBERTa-V2 and DeBERTa-V3 (He et al., 2021).

Tasks We consider four task datasets for our AFLite experiments. MNLI (Williams et al., 2018) and SNLI (Bowman et al., 2015), Cosmos QA (Huang et al., 2019) and SocialiQA (Sap et al., 2019). We chose these tasks based on several criteria: having a large enough training set to be suitable for AFLite, being in a format suitable for AFLite (i.e. classification), and no model-adversarial procedure already having been applied in the creation of the dataset. We use MNLI and SNLI for TextFooler experiments, as TextFooler was designed with NLI tasks in mind.

Fine-Tuning For all models, we execute two separate fine-tuning setups. First, we perform full fine-tuning on the training set, across 3 random restarts. Second, we perform fine-tuning on a smaller held-out subset of training examples¹—these will serve as the representations $\Phi(X)$ for AFLite, as well as the adversary models for TextFooler. We also repeat this subsampling across 3 random seeds, performing fine-tuning for each one. All of the AFLite and TextFooler results are averaged across the 3 fine-tuning and 3 AFLite runs. Refer to Appendix E for more details. All models were trained

¹This follows the AFLite protocol of training a set of weaker classifiers to learn representations that effectively serve as the adversary models.

using `jiant` (Phang et al., 2020), which is built on Transformers (Wolf et al., 2020) and PyTorch (Paszke et al., 2019).

Model	MNLI	SNLI	Cosmos	SIQA
MiniBERTa-S-1M	60.2	73.4	41.6	42.4
MiniBERTa-B-1B	79.3	87.2	55.0	57.3
BERT-Base	82.7	89.5	57.8	59.8
XLM-R-Base	81.2	87.4	59.3	63.1
BART-Base	84.6	89.8	63.4	65.2
BERT-Large	85.5	91.0	61.9	65.5
ALBERT-Large	86.3	89.9	62.3	68.5
RoBERTa-Base	86.1	91.1	67.1	69.6
ALBERT-XLarge	87.2	91.6	70.9	71.2
XLM-R-Large	88.3	90.8	70.6	72.5
ELECTRA-Base	87.4	91.5	69.9	73.4
BART-Large	89.1	91.2	76.7	77.3
DeBERTa-V3-Base	89.8	92.6	74.4	77.7
RoBERTa-Large	89.6	91.8	78.5	77.4
ELECTRA-Large	90.3	92.7	83.2	79.7
DeBERTa-V2-Large	90.5	92.7	85.5	79.1
DeBERTa-V2-XLarge	90.2	92.7	87.0	78.1
DeBERTa-V3-Large	90.8	93.1	87.6	81.2

Table 1: Accuracy (%) of fully fine-tuned models on full validation sets. Models are sorted in order of average performance across all four tasks.

Table 1 shows the performance of fully fine-tuned models on the validation set of each task. In this and subsequent visualizations, we sort the models based on the average full fine-tuned performance on the four tasks, from weakest to strongest.

4 Adversarially Filtering Evaluation Sets

AFLite (Sakaguchi et al., 2020; Le Bras et al., 2020) is an adversarial filtering algorithm that iteratively removes “easy” examples from a dataset. To apply AFLite, given a dataset $D = (X, Y)$ of inputs X and labels Y , we first compute a learned representation $\Phi(x)$ for each example based on the adversary model. In each iteration, we sample multiple random subsets of the remaining data, fit weak classifiers on the subsets and compute predictions on the held-out examples. If an example is predicted correctly by more than a threshold τ of weak classifiers, it is removed from the dataset. This procedure is repeated until the number of examples removed in an iteration falls below a set threshold, resulting in a reduced dataset. Details can be found in the original work (Le Bras et al., 2020).

Sakaguchi et al. (2020) and Le Bras et al. (2020) apply AFLite before applying train/validation/test splits. However, because we are interested in the impact of the adversarial filtering on evaluation datasets,² we do not want to use evaluation exam-

²In our experiments, we use the validation set of each task as the evaluation set.

260 ples to train the weak classifiers or influence the
261 filtering procedure. Hence, we tweak the AFLite al-
262 gorithm to separately filter out evaluation examples.
263 We accomplish this by running the standard AFLite
264 on the training examples, but in each round, we use
265 the same weak classifiers and removal criteria to
266 filter out “easy” evaluation examples. We show our
267 modified AFLite in Algorithm 1 in the Appendix.

268 4.1 Results on AFLite Across Adversary and 269 Fine-tuned Models

270 We apply AFLite using all 18 models as $\Phi(X)$, and
271 evaluate against each fully fine-tuned model. The
272 statistics of the examples filtered per model and
273 task can be found in Section B in the Appendix.
274 Figure 1 shows the results of fine-tuned models on
275 validation sets filtered via AFLite using different
276 adversary models.³

277 Overall, using AFLite with stronger adversary
278 models leads to lower performance across all fine-
279 tuned models, across all four tasks. Using a suffi-
280 ciently strong adversary model for filtering pushes
281 the performance of all tuned models to only slightly
282 above chance: For instance, while most mod-
283 els score 80-90% on the unfiltered MNLI valida-
284 tion set, filtering using AFLite with DeBERTa-V3-
285 Large results in no model scoring better than 45%.

286 We also observe a mild pattern of the weakest
287 models performing slightly better as stronger ad-
288 versaries are used in MNLI, SNLI, and SocialQA.
289 One explanation is that weaker models rely on eas-
290 ily learned heuristics (McCoy et al., 2019), and the
291 weak classifiers in AFLite select examples that go
292 against these heuristics, which weaker models sub-
293 sequently perform poorly on. In contrast, stronger
294 adversaries may filter out these examples.

295 4.1.1 Impact on Model Comparison

296 Evaluation datasets are often used to compare mod-
297 els, so we analyze the impact of adversarial filter-
298 ing on the resulting sorting order of model perfor-
299 mance. For each adversary model, we evaluate the
300 fine-tuned models on the AFLite filtered dataset
301 and sort the models by performance. In Figure 2,
302 we show the ranked performance of models using
303 different adversaries for a subset of the largest ad-
304 versary models. We find that the sorting order of
305 models is generally not consistent across adversary
306 models. This is the case even if we ignore cases
307 where the fine-tuned and adversary models share

308 the same pretrained model. For MNLI and SNLI,
309 evaluating on the datasets filtered by stronger ad-
310 versaries appears to greatly distort the relative ranking
311 of models. For Cosmos QA and SocialQA, we
312 observe that even when filtering with stronger ad-
313 versaries, stronger models still tend to rank better
314 than weaker models, but the ranking order is still
315 not consistent across adversaries.

316 One interpretation of this result is that adversar-
317 ial filtering may not give us evaluation data that
318 is reliable for benchmarking and comparing mod-
319 els. An alternative interpretation is that as stronger
320 adversary models are used, a larger proportion of
321 remaining examples are challenging and therefore
322 models are more likely to perform at chance on
323 them. As such, we ought to expect stronger adver-
324 saries will lead to more randomness in the model
325 rankings. In the extreme, if the weak classifiers
326 in AFLite are as capable as the best-performing
327 model, all models should perform at chance on
328 the remaining examples. While performance on
329 the strongest adversarially filtered datasets is still
330 above chance for most models, we see that in
331 MNLI and SNLI, all models converge to a small
332 range of performance (35%–45%), meaning that a
333 small variation in the number of correctly predicted
334 examples can lead to a large change in model rank.
335 This can lead to a distorted ranking of models.

336 We might also be concerned that the impact of
337 adversarial filtering if the fine-tuned and adversary
338 models are based on the same pretrained model. To
339 measure this, we compute the rank of each model
340 when no filtering is applied, and show how much
341 the rank changes when filtering using the same pre-
342 trained model. However, as we show in Figure 10,
343 the impact of filtering with the same pretrained
344 model is disproportionately large, with all models
345 except the weakest ones—which by definition can-
346 not fall in rank—falling several positions in relative
347 rankings. This implies that adversarial filtering for
348 evaluation sets is very sensitive to the choice of
349 model, and the resulting dataset is unfairly chal-
350 lenging if the adversary and evaluated models are
351 based on the same pretrained model.

352 4.2 Label Agreement

353 To investigate the kinds of examples being iden-
354 tified as challenging via AFLite, we use the per-
355 annotator labels of the MNLI and SNLI datasets.
356 In the original data creation procedure, each
357 validation-set example is annotated by 5 crowd-

³We present the same information in heatmaps in Figure 7 in the Appendix.

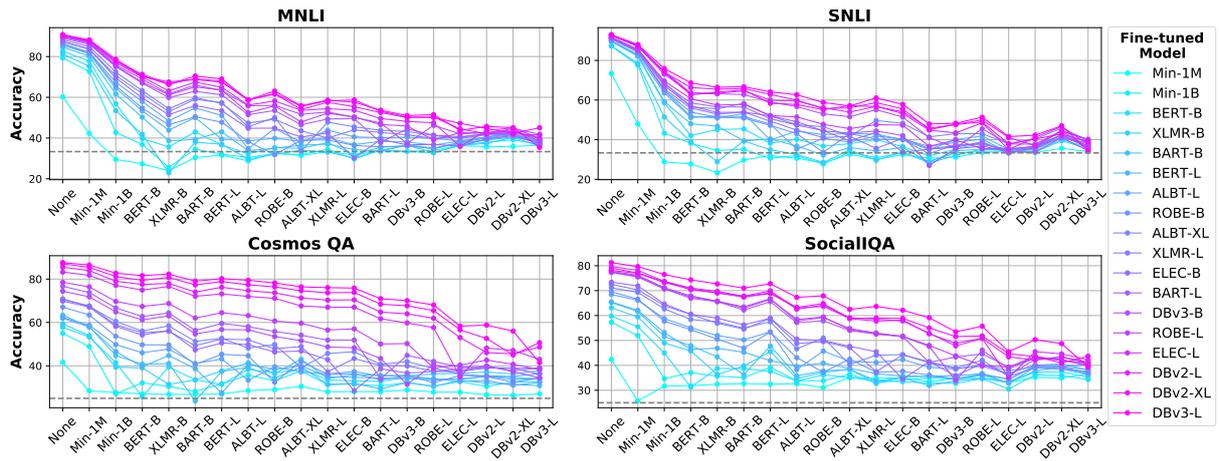


Figure 1: Performance of fine-tuned models on validation sets filtered via AFLite using adversary models. ‘None’ indicates the full unfiltered validation set. The dotted line indicates performance at chance for each task. Filtering with stronger adversary models leads to lower performance on the filtered dataset.

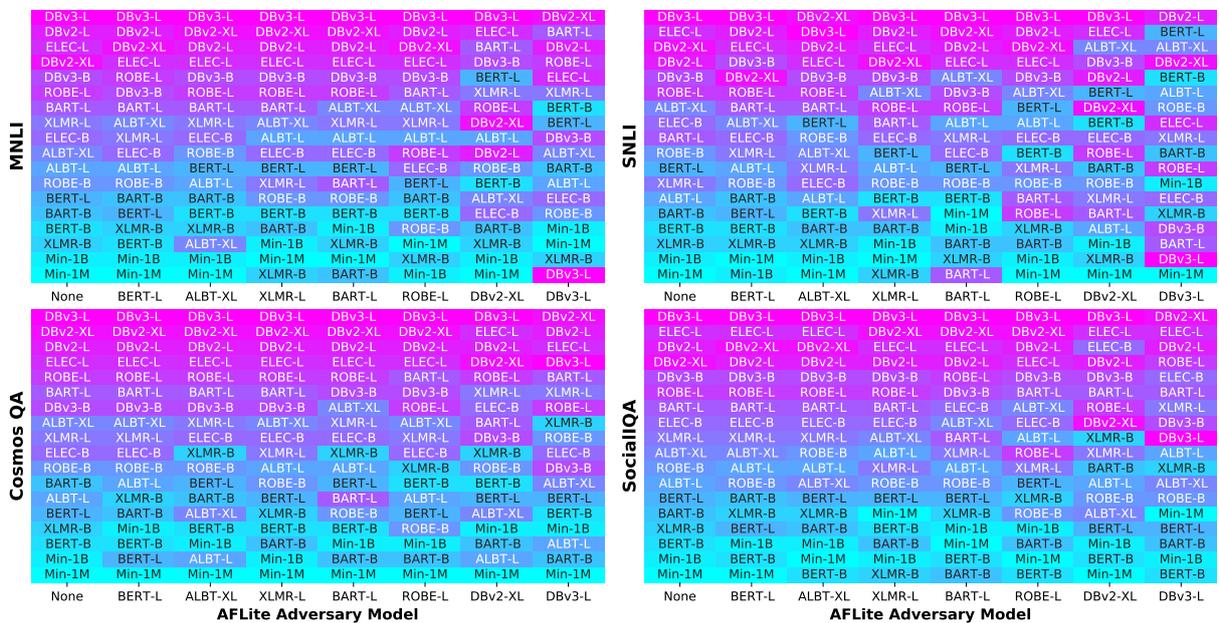


Figure 2: Ranked performance of fine-tuned models on validation sets filtered via AFLite using adversary models. We depict here a subset of the larger adversary models. For each AF Selected dataset, we sort models by their performance (Figure 1) from best (top) to worst (bottom). ‘None’ indicates the full validation set with no filtering applied. We find that the sorting order of model performance is not consistent across adversary models.

workers, and candidate examples are only accepted if at least 3 out of 5 crowdworkers agree on the label. We show in Figure 3 the average annotator agreement in the AFLite-selected examples across adversary models. For comparison, we also show the agreement rate among examples eliminated in the very first round of the AFLite procedure.

We observe a clear pattern across both datasets that filtering with stronger adversary models selects for examples with lower annotator agreement. Combined with our results above on lower model

performance on filtered datasets, we take this as good evidence that the AFLite procedure indeed selects for the most challenging examples. It is unclear if these examples are challenging because they are genuinely difficult, where humans can easily make mistakes on them, genuinely ambiguous, or simply mislabeled. Conversely, we see that the first-pass filtered examples have consistently high annotator agreement, and that this rate does not vary across strength of the adversary models.

Oversampling low-agreement examples is not

Adversary Model	MNLI		SNLI	
	Filtered In First Iteration	AF Selected	Filtered In First Iteration	AF Selected
None		88.5%		88.1%
Min-1M	90.4%	87.2%	89.8%	85.8%
Min-1B	90.8%	83.4%	89.8%	81.2%
BERT-B	91.0%	81.3%	89.7%	79.2%
XLMR-B	90.9%	79.0%	89.9%	78.0%
BART-B	90.9%	80.1%	89.9%	79.2%
BERT-L	90.9%	79.8%	89.6%	77.8%
ALBT-L	90.8%	77.7%	89.8%	77.6%
ROBE-B	90.9%	78.0%	89.7%	76.2%
ALBT-XL	90.6%	75.9%	89.6%	76.9%
XLMR-L	90.8%	77.1%	89.7%	77.1%
ELEC-B	90.8%	77.2%	89.7%	76.6%
BART-L	90.8%	75.8%	89.7%	73.9%
DBv3-B	90.7%	75.5%	89.7%	74.6%
ROBE-L	90.8%	75.4%	89.7%	75.0%
ELEC-L	90.7%	73.5%	89.6%	72.1%
DBv2-L	90.9%	73.9%	89.8%	72.5%
DBv2-XL	90.8%	74.0%	89.7%	73.9%
DBv3-L	90.6%	73.2%	89.6%	73.2%

Figure 3: Label agreement among the adversarially filtered datasets from human annotators. *AF Selected* indicates examples that are not filtered out. *None* indicates no filtering applied i.e. agreement over the full validation set. Label agreement for the AF-selected datasets falls as better adversary models are used, indicating that AFLite may be selecting for the examples with the most ambiguity or labeling noise.

necessarily a bad thing if evaluated appropriately. Pavlick and Kwiatkowski (2019) and Nie et al. (2020c) show that there can be genuine disagreement between annotators over the example label, and argue that we should go beyond optimizing for model accuracy and instead train models to predict the full distribution of human judgements. The current format of scoring models on simple accuracy is inadequate to evaluate on low-agreement examples, as the distribution of labels is reduced to a single label based on majority vote. Hence, if AFLite selects for low-agreement examples, one potential improvement would be to adjust the evaluation format to reflect annotator disagreement over labels.

5 Adversarially Perturbed Datasets

An alternative approach to creating more challenging datasets is to modify examples to lower the performance of a given adversary model. TextFooler (Jin et al., 2019) is a popular algorithm for adversarially perturbing examples, which involves swapping out words in the input for synonyms, with additional constraints aimed at retaining grammaticality and semantic content of the original example. We refer to the original work for full details of the TextFooler algorithm.

We apply the TextFooler methods to the SNLI and MNLI datasets, following the setup of the original work, where only the hypotheses are perturbed. For full comparability of results, we use the same models for $\Phi(X)$ in AFLite as the adversaries for

TextFooler, and evaluate using the models fine-tuned on the full training set. We also use the same hyperparameters for the TextFooler algorithm as in the original work.

Overall, we find that adversarial perturbations impact model evaluation in a similar manner to adversarial filtering. TextFooler adversarial perturbations generally lower model performance across the board, with stronger adversaries leading to lower average performance across all models. However, we also observe a strong trend of models performing much worse on data perturbed adversarially to the same pretrained model, as seen by the distinct dip in performance along the diagonals of Figure 12. Moreover, using stronger adversaries also distorts the order of model performance, as seen in Figure 4, much more so than with adversarial filtering. For instance, using DeBERTa-V3-Large as the adversary, ALBERT-XLarge is the best performing model; however, with ALBERT-XLarge as the adversary, DeBERTa-V3-Large is in turn the best performing model. The inconsistent impact on model performance across adversary models makes the naive usage of adversarial perturbations to create more challenging evaluation datasets problematic.

One major concern is whether the TextFooler-perturbed examples still retain the same semantic content and, more importantly, the labels of the original examples. Jin et al. evaluated a 100 examples from the validation set and found that the perturbations somewhat distorted the grammaticality and the meaning of the SNLI examples with a BERT adversary. We also observe a fall in quality of the examples in our experiments, with label flips arising from modifying key words, which greatly hurts the validity of performance obtained on these examples. We show a number of examples perturbed with different adversary models in (Table 4). Without any additional checks, using automatically perturbed examples for evaluating models can lead to extremely misleading results. Works such as AdversarialGLUE (Wang et al., 2021) perform an additional step of human validation and filtering of model-perturbed examples, which are necessary to ensure reliable evaluation data.

6 Model-in-the-Loop Adversarially Collected Datasets

In model-in-the-loop adversarial data collection, human crowdworkers are tasked with writing examples that a given adversary model will incorrectly

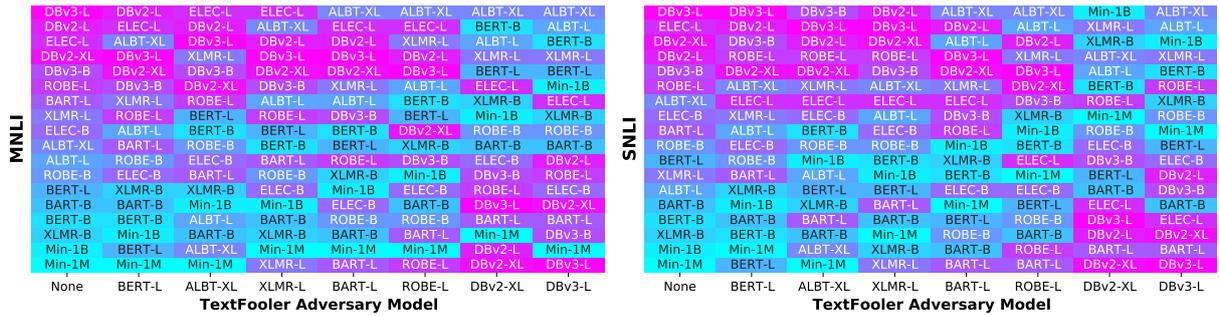


Figure 4: Ranked performance of fine-tuned models on validation sets perturbed via TextFooler using adversary models. For each perturbed dataset, we sort models by their performance (Figure 1) from best (top) to worst (bottom). ‘None’ indicates the full validation set with no filtering applied. We find that the sorting order of model performance is not consistent across adversary models, with dramatic reversals in many cases.

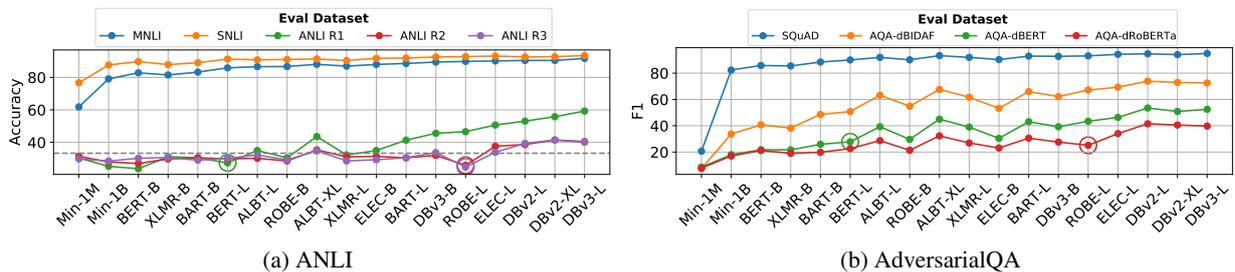


Figure 5: Measuring the performance of models on adversarially collected datasets. Exact Match scores for AdversarialQA are shown in Figure 13 in the Appendix. For each adversarially created dataset, the corresponding base adversary model used in model-in-the-loop data creation is circled in the corresponding color for that dataset. Performance at chance on ANLI is shown with a dotted line. While adversarial dataset creation appears to create datasets that are slightly harder for the adversary model compared to other models, the resulting datasets are harder across the board for all models, with stronger models still performing relatively better.

label. We consider two established model-in-the-loop adversarially collected datasets. ANLI (Nie et al., 2020b) is an NLI dataset adversarially collected through three iterative rounds, where the data for each round is written to be adversarial to models trained on a combination of MNLI, SNLI, and data from previous rounds. BERT-Large is used as the adversary model for round 1 of data collection, while RoBERTa-Large is used for rounds 2 and 3. AdversarialQA (Bartolo et al., 2020), is an adversarial question-answering dataset in the format of SQuAD 1.1 (Rajpurkar et al., 2016). Unlike ANLI, it consists of separately collected examples based on three adversary models: BiDAF (Seo et al., 2017), BERT-Large, and RoBERTa-Large.

While both datasets come with training, validation and test data splits, we conduct our analysis on the validation data. For both datasets, we fine-tune models on the conventional training data for each task,⁴ before evaluating on both the standard and adversarial validation datasets.

⁴MNLI and SNLI for ANLI, and SQuAD 1.1 for AdversarialQA.

Figure 5 shows results on both model-in-the-loop datasets. For ANLI, about half of the models perform at chance for ANLI R1, whereas the stronger models perform significantly above chance. On the other hand, for ANLI R2 and R3, most models perform at chance except for the largest DeBERTa models. These results show that the ANLI data-generating procedure leads to examples that are harder for all models. However, we also observe that for ANLI R2 and R3, the performance of the adversary model, RoBERTa-large, is markedly below chance. This supports our observation above that while adversarial dataset creation can lower performance across the board, it still tends to hurt the adversary model more than others.

We see similar results for AdversarialQA, with models performing poorer as the datasets are generated with stronger adversaries. Unlike for ANLI, models do significantly better than chance on the adversarial datasets, with almost all models staying above 20% F1 and 10% EM.

Compared to our more extensive experiments on adversarial filtering, there are fewer datasets

collected using different adversary models, given the financial cost and manual writing needed to obtain examples. Hence we cannot draw strong conclusions about the efficacy of adversarial data collection for evaluation data from the current set of results. Moreover, the adversaries used in ANLI and AdversarialQA are not among the strongest models we used in our adversarial filtering experiments, where we saw the greatest distortion in the ranking of models. However, we do find that adversarial data collection leads to harder examples with stronger adversary models. As more work is done on adversarially collecting datasets and building benchmarks based on them (Kielar et al., 2021), we recommend that researchers pay close attention to the impact of the choice of adversary model and evaluate across a range of different models.

7 Discussion

We highlight that this work has not investigated the nature of the adversarial examples outside of the impact on model performance and annotator agreement. Works such as Williams et al. (2020) will be important for understanding exactly what examples are considered adversarial and why they are challenging to different models.

While our adversarial filtering experiments were performed on single adversary models, a possible alternative is to ensemble a diverse set of adversary models when running AFLite, or weight examples based on the AFLite example selection based on each adversary. This approach may help reduce the issue of disproportionate impact on any given adversary model’s performance, and weighting evaluation across different example subsets may also potentially reduce the unstable ranking of models. However, this would significantly increase the cost of running the algorithm, would not address the issue of oversampling low-agreement examples, and may simply create a bias in favor of models that are greatly dissimilar from the pool of adversaries (e.g. non-Transformer models).

8 Conclusion

In this work, we have investigated three different approaches to adversarially constructing more challenging evaluation datasets in extensive experiments across 18 different pretrained models.

While all three methods result in lower scores as stronger adversaries are used, our takeaways on the viability of adversarial data creation to construct

more challenging evaluation datasets are mixed.

Using a modified AFLite to adversarially filter evaluation examples, we find that there is a disproportionately large impact on the performance of fine-tuned models derived from the same pretrained model as the adversary, that the resulting ranking of models is unstable across the choice of adversary model especially as stronger adversaries are used, and that the filtering selects for examples with low annotator agreement over labels. On the other hand, the impact on model rankings is somewhat expected as a higher proportion of difficult examples remain after filtering.

Using TextFooler to perturb examples, we find even greater distortion in model rankings with stronger adversaries, and that examples can often be perturbed to the point of flipping labels, which is dire for model evaluation.

In a smaller set of experiments on adversarially collected datasets, we find hints of datasets being more challenging for the same pretrained model as the adversary, consistent with (Nie et al., 2020a), but are unable to draw stronger conclusions given the small number of adversaries used.

As the cost of using models goes down and their capabilities improve, we are likely to see more attempts to involve them in dataset creation. Models may be used adversarially as discussed above, or used to assist in writing examples via text generation models, or used in other ways, such as automatically identifying outliers or low-quality human-written examples. In any of these cases, it is possible to create an adverse and undesirable feedback loop in the data creation procedure.

While we believe that adversarial data methods can be helpful in creating more challenging evaluation benchmarks, we should take extra care to avoid the pitfalls of these approaches. Careful human validation of examples, as in (Wang et al., 2021), can help to address some pitfalls of adversarial approaches, but addressing biases at the dataset rather than example level can still be challenging. Importantly, adversarial datasets must accurately reflect the core task or capability being measured, ideally with a diverse set of examples that have good coverage of the linguistic phenomena associated with the task. As it stands, we do not yet find any free lunch for creating more challenging evaluation datasets.

References

- 602
603 Max Bartolo, Alastair Roberts, Johannes Welbl, Sebastian Riedel, and Pontus Stenetorp. 2020. [Beat the AI: Investigating adversarial human annotation for reading comprehension](#). *Transactions of the Association for Computational Linguistics*, 8:662–678.
- 608 Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.
- 614 Samuel R. Bowman and George Dahl. 2021. [What will it take to fix benchmarking in natural language understanding?](#) In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4843–4855, Online. Association for Computational Linguistics.
- 621 Samuel R. Bowman, Jennimaria Palomaki, Livio Baldini Soares, and Emily Pitler. 2020. [New protocols and negative results for textual entailment data collection](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8203–8214, Online. Association for Computational Linguistics.
- 628 Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. [ELECTRA: Pre-training text encoders as discriminators rather than generators](#). In *International Conference on Learning Representations*.
- 633 Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- 642 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- 651 Siddhant Garg and Goutham Ramakrishnan. 2020. [BAE: BERT-based adversarial examples for text classification](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6174–6181, Online. Association for Computational Linguistics.
- 657 Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. [Explaining and harnessing adversarial examples](#). Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. [Deberta: Decoding-enhanced bert with disentangled attention](#). Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. [Deberta: Decoding-enhanced bert with disentangled attention](#). In *International Conference on Learning Representations*.
- 667 Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. [Cosmos QA: Machine reading comprehension with contextual commonsense reasoning](#). 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)*, pages 2391–2401, Hong Kong, China. Association for Computational Linguistics.
- 676 Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2019. Is bert really robust? natural language attack on text classification and entailment. *arXiv preprint arXiv:1907.11932*.
- 680 Divyansh Kaushik, Douwe Kiela, Zachary C. Lipton, and Wen-tau Yih. 2021. [On the efficacy of adversarial data collection for question answering: Results from a large-scale randomized study](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6618–6633, Online. Association for Computational Linguistics.
- 689 Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina Williams. 2021. [Dynabench: Rethinking benchmarking in NLP](#). *arXiv preprint*.
- 697 Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. [ALBERT: A lite BERT for self-supervised learning of language representations](#). In *8th International Conference on Learning Representations, ICLR 2020*.
- 703 Ronan Le Bras, Swabha Swayamdipta, Chandra Bhagavatula, Rowan Zellers, Matthew Peters, Ashish Sabharwal, and Yejin Choi. 2020. [Adversarial filters of dataset biases](#). In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1078–1088. PMLR.
- 710 Gyeongbok Lee, Seung-won Hwang, and Hyunsouk Cho. 2020. [SQuAD2-CR: Semi-supervised annotation for cause and rationales for unanswerability in SQuAD 2.0](#). In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 5425–5432, Marseille, France. European Language Resources Association.

717	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880, Online. Association for Computational Linguistics.	
718		
719		
720		
721		
722		
723		
724		
725		
726	Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. 2019. Textbugger: Generating adversarial text against real-world applications . <i>Proceedings 2019 Network and Distributed System Security Symposium</i> .	
727		
728		
729		
730		
731	Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020. BERT-ATTACK: Adversarial attack against BERT using BERT . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6193–6202, Online. Association for Computational Linguistics.	
732		
733		
734		
735		
736		
737		
738	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 . <i>arXiv preprint</i> .	
739		
740		
741		
742		
743	Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3428–3448, Florence, Italy. Association for Computational Linguistics.	
744		
745		
746		
747		
748		
749		
750	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020a. Adversarial NLI: A new benchmark for natural language understanding . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4885–4901, Online. Association for Computational Linguistics.	
751		
752		
753		
754		
755		
756		
757	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020b. Adversarial NLI: A new benchmark for natural language understanding . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> . Association for Computational Linguistics.	
758		
759		
760		
761		
762		
763		
764	Yixin Nie, Xiang Zhou, and Mohit Bansal. 2020c. What can we learn from collective human opinions on natural language inference data? In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9131–9143, Online. Association for Computational Linguistics.	
765		
766		
767		
768		
769		
770		
771	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library . In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, <i>Advances in Neural Information Processing Systems 32</i> , pages 8024–8035. Curran Associates, Inc.	774 775 776 777 778 779 780 781 782 783
772		
773		
774	Ellie Pavlick and Tom Kwiatkowski. 2019. Inherent disagreements in human textual inferences . <i>Transactions of the Association for Computational Linguistics</i> , 7:677–694.	784 785 786 787
775		
776		
777		
778		
779		
780		
781		
782		
783		
784	Jason Phang, Phil Yeres, Jesse Swanson, Haokun Liu, Ian F. Tenney, Phu Mon Htut, Clara Vania, Alex Wang, and Samuel R. Bowman. 2020. jiant 2.0: A software toolkit for research on general-purpose text understanding models . http://jiant.info/ .	788 789 790 791 792 793
785		
786		
787		
788		
789		
790		
791		
792		
793		
794	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	794 795 796 797 798 799
795		
796		
797		
798		
799		
800	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Winogrande: An adversarial winograd schema challenge at scale . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pages 8732–8740.	800 801 802 803 804
801		
802		
803		
804		
805	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social IQa: Commonsense reasoning about social interactions . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4463–4473, Hong Kong, China. Association for Computational Linguistics.	805 806 807 808 809 810 811 812 813
806		
807		
808		
809		
810		
811		
812		
813		
814	Min Joon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. 2017. Bidirectional attention flow for machine comprehension . In <i>5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings</i> .	814 815 816 817 818 819
815		
816		
817		
818		
819		
820	Eric Wallace, Adina Williams, Robin Jia, and Douwe Kiela. 2021. Analyzing dynamic adversarial training data in the limit .	820 821 822
821		
822		
823	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems . In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, <i>Advances in Neural Information Processing Systems 32</i> , pages 3266–3280. Curran Associates, Inc.	823 824 825 826 827 828 829 830 831
824		
825		
826		
827		
828		
829		
830		
831		

- 832 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018.
833 [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- 840 Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and Bo Li. 2021. [Adversarial GLUE: A multi-task benchmark for robustness evaluation of language models](#). *CoRR*, abs/2111.02840.
- 845 Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. [A broad-coverage challenge corpus for sentence understanding through inference](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics.
- 853 Adina Williams, Tristan Thrush, and Douwe Kiela. 2020. [ANLizing the adversarial natural language inference dataset](#). *arXiv preprint*.
- 856 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- 868 Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. [Swag: A large-scale adversarial dataset for grounded commonsense inference](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- 873 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [HellaSwag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.
- 880 Yian Zhang, Alex Warstadt, Xiaocheng Li, and Samuel R. Bowman. 2021. [When do you need billions of words of pretraining data?](#) In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1112–1125, Online. Association for Computational Linguistics.

Algorithm 1: AFLite for Evaluation Data

Input: training dataset $D_T = (X_T, Y_T)$, **evaluation dataset** $D_V = (X_V, Y_V)$, pre-computed representation $(\Phi(X_T), \Phi(X_V))$, model family \mathcal{M} , target dataset size n , number of random partitions m , training set size $t < n$, slice size $k \leq n$, early-stopping threshold τ

Output: **Filtering history of evaluation examples** H , **remaining evaluation examples** R

```
 $S = D_T$   
 $R = D_V$   
while  $|S| > n$  do  
  // Filtering phase  
  forall  $i \in S$  do  
    Initialize multiset of out-of-sample training predictions  $E_T(i)$ ;  
    forall  $i \in R$  do  
      Initialize multiset of out-of-sample evaluation predictions  $E_V(i)$ ;  
    for iteration  $j : 1..m$  do  
      Randomly partition  $S$  into  $(T_j, S \setminus T_j)$  s.t.  $|S \setminus T_j| = t$ ;  
      Train a classifier  $\mathcal{L} \in \mathcal{M}$  on  $\{(\Phi(x), y) | (x, y) \in S \setminus T_j\}$ ;  
      forall  $i = (x, y) \in T_j$  do  
        Add the prediction  $\mathcal{L}(\Phi(x))$  to  $E_T(i)$ ;  
      forall  $i = (x, y) \in R$  do  
        Add the prediction  $\mathcal{L}(\Phi(x))$  to  $E_V(i)$ ;  
      forall  $i = (x, y) \in S$  do  
        Compute the predictability score  $\hat{p}(i) = |\{\hat{y} \in E_T(i) \text{ s.t. } \hat{y} = y\}| / |E_T(i)|$ ;  
      forall  $i = (x, y) \in R$  do  
        Compute the predictability score  $\hat{p}(i) = |\{\hat{y} \in E_V(i) \text{ s.t. } \hat{y} = y\}| / |E_V(i)|$ ;  
      Select up to  $k$  instances  $S'$  in  $S$  with the highest predictability scores subject to  $\hat{p}(i) \geq \tau$ ;  
       $S = S \setminus S'$ ;  
      Select all instances  $R'$  in  $R$  where  $\hat{p}(i) \geq \tau$ ;  
       $R = R \setminus R'$ ;  
      Append  $R'$  to  $H$ ;  
      if  $|S'| < k$  then  
        break;  
return  $H, R$ 
```

A Modified AFLite

Algorithm 1, shows the modified AFLite algorithm, where the original algorithm applied to training examples is shown in black, and the additional lines applied to the evaluation examples are highlighted in red.

$\Phi(X)$ is the CLS or <S> embeddings of corresponding adversary model, fine-tuned on a separate held-out training set for the task (10% of the training data, following AFLite).

B AFLite Filtering Statistics

Figure 6 shows the breakdown of filtered examples when applying AFLite with different models. Each example in the validation set falls into one of three categories: examples filtered out on the

first iteration of AFLite, examples filtered in all subsequent iterations, and examples remaining after applying AFLite (AF Selected). In most cases, more than half the validation datasets are filtered out within the first iteration, meaning that these examples were largely correctly predicted by a set of weak classifiers using the learned representations of partially tuned adversary models. Moreover, the stronger the adversary model, the more examples tend to be removed in the first iteration. Subsequent filtering iterations remove comparatively much fewer examples.

Among the AF Selected examples for Cosmos QA and SocialIQA, we see a trend that the stronger the adversary model, the fewer examples remain after AFLite. We do not see the same pattern in MNLI and SNLI, where the number of AF Selected examples does not vary consistently across strength of models. We note that Cosmos QA and SocialIQA use different AFLite hyperparameters from MNLI and SNLI because of the difference in datasets sizes (Table 3).

C Additional Results

Figure 8 shows the same information as Figure 1, with fine-tuned models on the X-axis and adversary models shown in different curves. Figure 7 shows the same information in a heatmap. Figure 11 shows the average agreement across adversarially filtered datasets, including the agreement among subsequent iterations of AFLite. Figure 13 shows exact-match scores on the AdversarialQA datasets.

D Models

Table 2 shows additional details for each of the pretrained models used in our experiments.

E Fine-Tuning Details

For full fine-tuning, we fine-tune for 3 epochs for MNLI and SNLI, and 5 epochs for Cosmos QA and SocialIQA. For fine-tuning weak classifiers for $\Phi(x)$, we subsample 10% of the training examples for MNLI and SNLI, and 5000 examples for Cosmos QA and SocialIQA, fixing the subsamples across all models. We repeat the subsampling procedure three times. In both fine-tuning setups, we hold out 500 examples from the training set for early stopping. These training examples are held out for both full fine-tuning as well as the AFLite

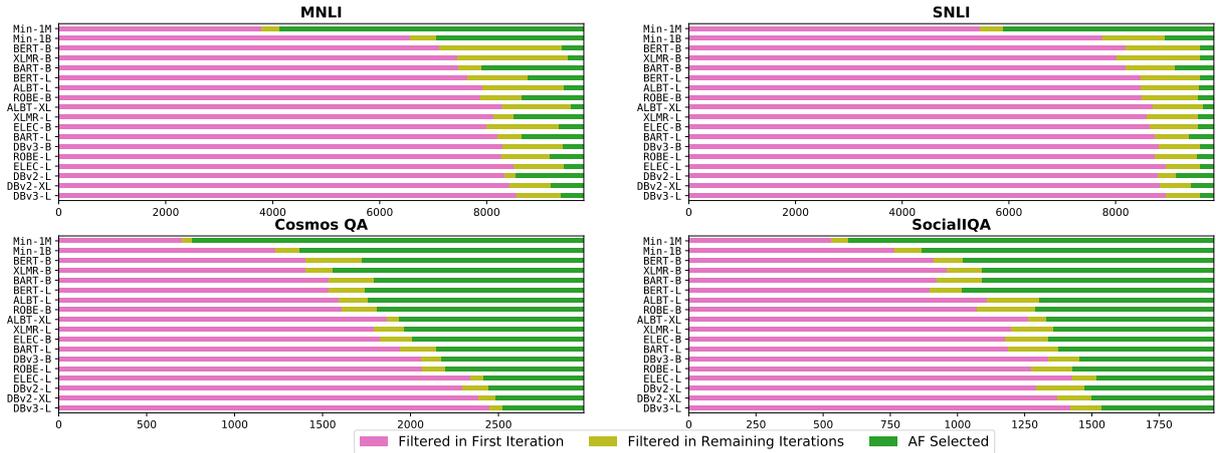


Figure 6: Statistics of AFLite-filtered datasets. We apply Algorithm 1 to the validation set of each task across adversary models, and average across three random seeds. *AF Selected* indicates examples that are included in the final evaluation set. For most models, majority of the examples are filtered out within the first iteration of AFLite.

Model	Abbreviation	Reference	Parameters	Training Objective
MiniBERTa Small 1M	Min-1M	Zhang et al. (2021)	~45M	Masked language modeling
MiniBERTa Base 1B	Min-1B	Zhang et al. (2021)	~100M	Masked language modeling
BERT-Base (cased)	BERT-B	Devlin et al. (2019)	~100M	Masked language modeling + NSP
BERT-Large (cased)	BERT-L	Devlin et al. (2019)	~340M	Masked language modeling + NSP
XLM-R-Base	XLMR-B	Conneau et al. (2020)	~100M	Masked language modeling
XLM-R-Large	XLMR-L	Conneau et al. (2020)	~340M	Masked language modeling
BART-Base	BART-B	Lewis et al. (2020)	~100M	Text infilling + Sentence permutation
BART-Large	BART-L	Lewis et al. (2020)	~340M	Text infilling + Sentence permutation
ALBERT-Large (v2)	ALB-L	Lan et al. (2020)	~18M	Masked language modeling + SOP
ALBERT-XLarge (v2)	ALB-XL	Lan et al. (2020)	~60M	Masked language modeling + SOP
RoBERTa-Base	RoBE-B	Liu et al. (2019)	~100M	Masked language modeling
RoBERTa-Large	RoBE-L	Liu et al. (2019)	~340M	Masked language modeling
ELECTRA-Base	ELEC-B	Clark et al. (2020)	~100M	Replaced token detection
ELECTRA-Large	ELEC-L	Clark et al. (2020)	~340M	Replaced token detection
DeBERTa-V2-Large (v2)	DBv2-L	He et al. (2021)	~900M	Masked language modeling
DeBERTa-V2-XLarge (v2)	DBv2-XL	He et al. (2021)	~1.5B	Masked language modeling
DeBERTa-V3 Base	DBv3-B	He et al. (2021)	~100M	Replaced token detection
DeBERTa-V3 Large	DBv3-L	He et al. (2021)	~418M	Replaced token detection

Table 2: Pretrained models used in our experiments

950 procedure. As such, validation examples never in-
 951 fluence the fine-tuning or AFLite procedures, only
 952 being used when we perform AFLite and filter our
 953 validation examples as described in Algorithm 1.

954 For DeBERTa, unlike in He et al. (2020), we do
 955 not apply SiFT during fine-tuning.

956 F AFLite Hyperparameters

957 Table 3 shows the hyperparameters for our AFLite
 958 runs.

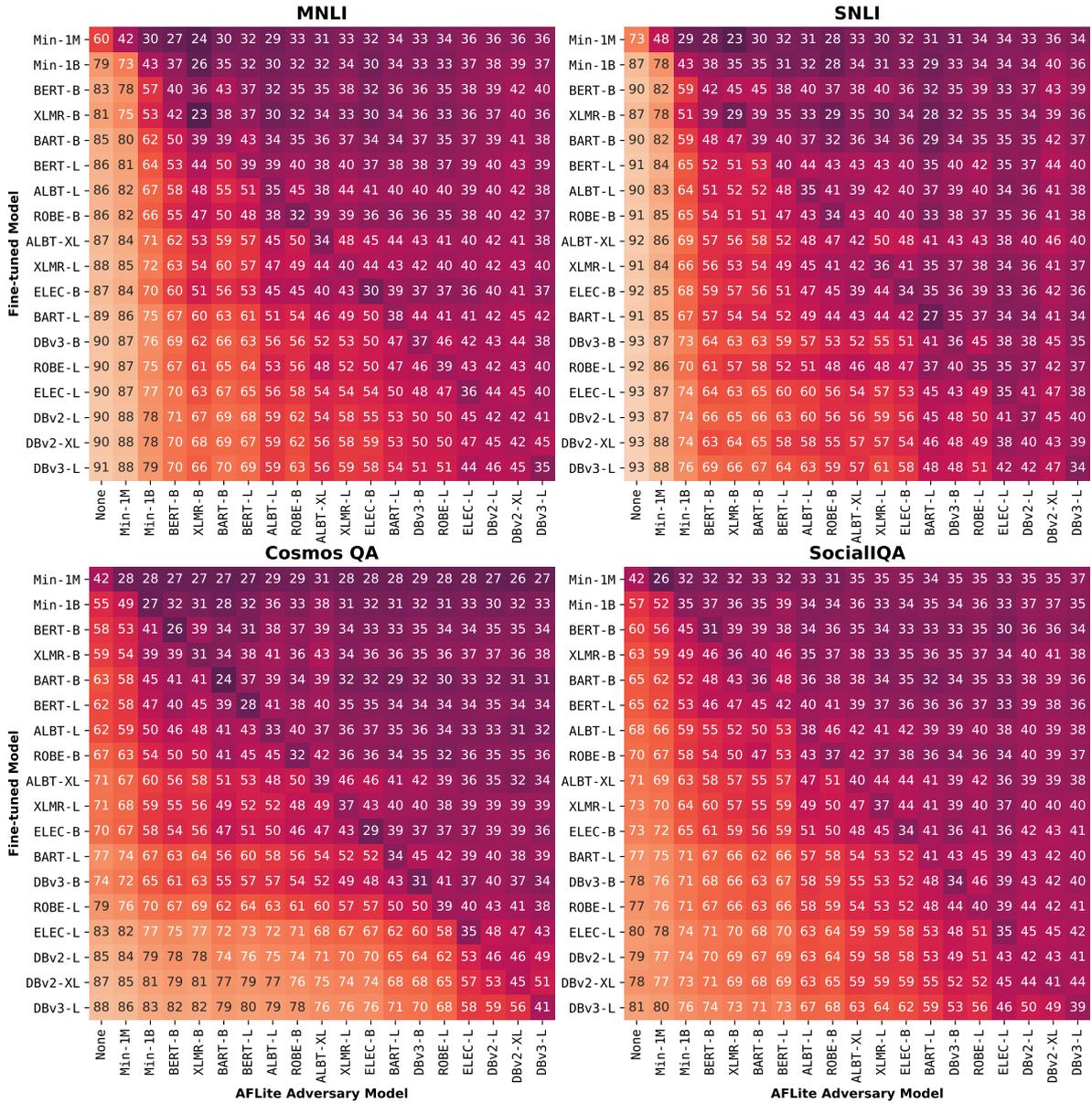


Figure 7: Performance of fine-tuned models on validation sets filtered via AFLite using adversary models. ‘None’ indicates the full validation set with no filtering applied. Filtering with stronger adversary models leads to lower performance on the filtered dataset, across all fine-tuned models. However, filtering also tends to hurt the adversary model itself more than other models on average (darker cells on the diagonal).

	MNLI	SNLI	Cosmos QA	SocialQA
m	64	64	64	64
t	50K	40K	10k	10k
k	10K	10K	500	500
τ	0.75	0.75	0.75	0.75
Taken From	Le Bras et al. (2020)	Le Bras et al. (2020)	Sakaguchi et al. (2020)	Sakaguchi et al. (2020)

Table 3: AFLite Hyperparameters

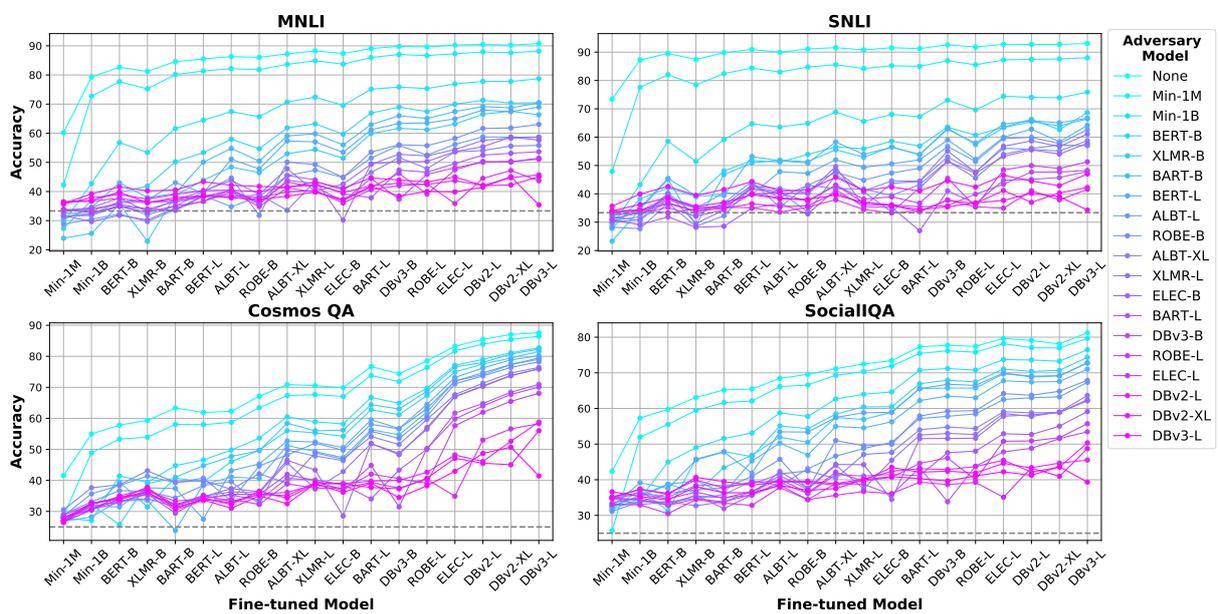


Figure 8: Performance of fine-tuned models on validation sets filtered via AFLite using adversary models. ‘None’ indicates the full validation set with no filtering applied. The dotted line indicates performance at chance for each task. Filtering with stronger adversary models leads to lower performance on the filtered dataset, across all fine-tuned models.

Premise	Original Hypothesis	Label	Perturbed Hypothesis	Issue
Adversary: BERT-Large				
In Hong Kong you can have a plate, or even a whole dinner service, hand-painted to your own design.	It's impossible to have a plate hand-painted to your own design in Hong Kong.	contradiction	It's imaginable to have a plate hand-painted to your own design in Hong Kong.	Semantics change
He hadn't seen even pictures of such things since the few silent movies run in some of the little art theaters.	He had recently seen pictures depicting those things.	contradiction	He had recently experimented pictures exposing those things.	Problematic word substitution
As a basic guide, the symbols below have been used to indicate high-season rates in Hong Kong dollars, based on double occupancy, with bath or shower.	As you can see, the symbols are of dolphins and octopuses.	neutral	As you can see, the symbols are of sharks and octopuses.	Semantics change
Adversary: DeBERTa-V3-Large				
Do you think Mrs. Inglethorp made a will leaving all her money to Miss Howard? I asked in a low voice, with some curiosity.	I yelled at the top of my lungs.	contradiction	I muttered at the superior of my lungs.	Problematic word substitution
But when the cushion is spent in a year or two, or when the next recession arrives, the disintermediating voters will find themselves playing the roles of budget analysts and tax wonks.	The cushion will likely be spent in under two years.	entailment	The cushion will likely be spent in under three years.	Semantics change
substitute my my yeah my kid'll do uh four or five hours this week for me no problem	I just can't make the time because of my job.	neutral	I just can't make the time for of my job.	Grammatical error

Table 4: TextFooler Examples

Min-1M	0	0	0	0
Min-1B	0	0	1	0
BERT-B	1	0	2	2
XLMR-B	2	1	1	1
BART-B	1	1	6	2
BERT-L	1	3	3	2
ALBT-L	2	2	4	1
ROBE-B	6	4	6	3
ALBT-XL	6	3	6	2
XLMR-L	4	2	1	2
ELEC-B	8	8	7	8
BART-L	5	9	7	2
DBv3-B	6	7	10	11
ROBE-L	4	8	2	4
ELEC-L	15	8	8	9
DBv2-L	5	2	1	5
DBv2-XL	4	4	2	4
DBv3-L	17	16	3	8
	MNLI	SNLI	Cosmos	SIQA

Figure 10: For each fine-tuned model, we compute the change in rank (1=best, 18=worst) from evaluating on the full evaluation set, and on the dataset filtered using the same pretrained model for the adversary. In almost all cases, filtering on the same pretrained model leads to a fall in ranking, indicating that the model is disproportionately affected by filtering with itself.

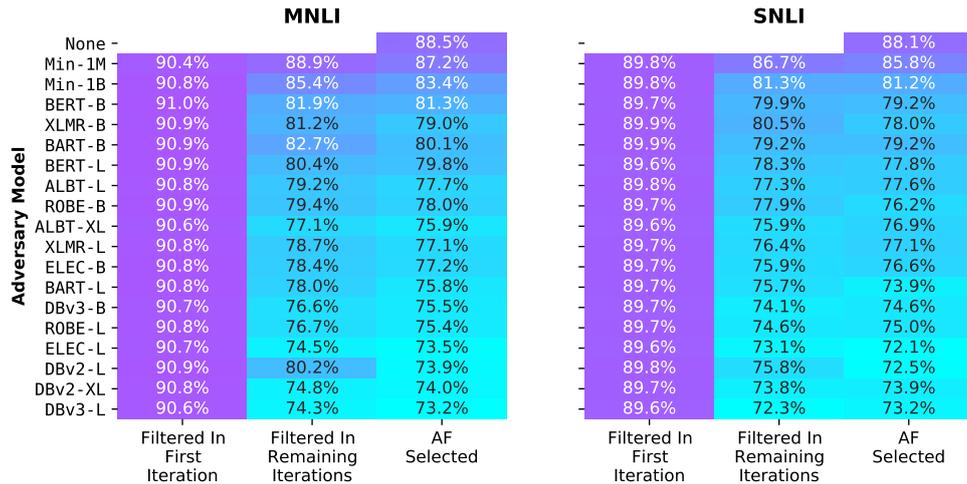


Figure 11: Label agreement among the adversarially filtered datasets from human annotators. *AF Selected* indicates examples that are not filtered out. Label agreement is very high for first pass filtered examples for all models. On the other hand, label agreement for the remainder datasets falls as better adversary models are used, indicating that AFLite may be selecting for the examples with the most ambiguity or labeling noise.

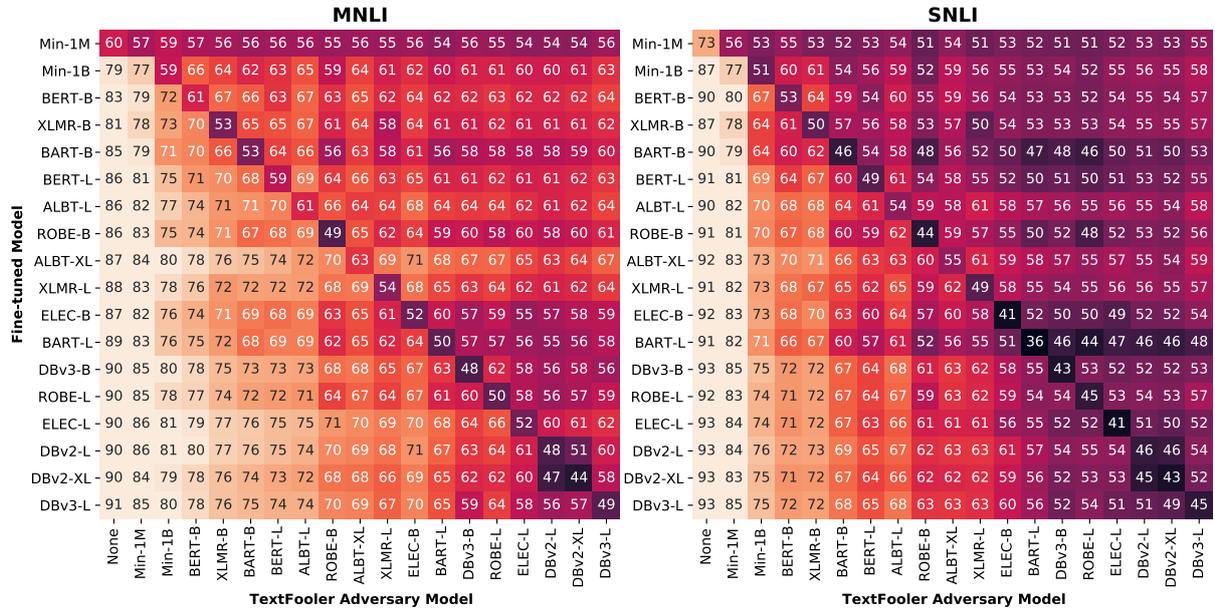


Figure 12: Performance of fine-tuned models on validation sets perturbed via TextFooler using adversary models. 'None' indicates the full validation set with no perturbation applied. Perturbing with stronger adversary models leads to lower performance on the filtered dataset, across all fine-tuned models. However, perturbing also tends to hurt the adversary model itself more than other models on average (darker cells on the diagonal).

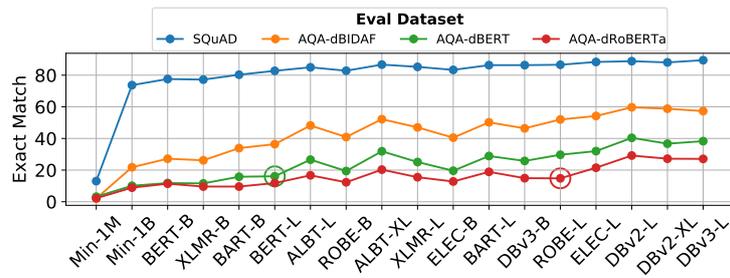


Figure 13: Measuring the performance of models on AdversarialQA. AdversarialQA models are fine-tuned on SQuAD 1.1. For each adversarially created dataset, the corresponding base adversary model used in model-in-the-loop data creation is circled in the corresponding color for that dataset.