

BERT & Family Eat Word Salad: Experiments with Text Understanding

Ashim Gupta, Giorgi Kvernadze, Vivek Srikumar

University of Utah
{ashim, giorgi, svivek}@cs.utah.edu

Abstract

In this paper, we study the response of large models from the BERT family to incoherent inputs that should confuse any model that claims to understand natural language. We define simple heuristics to construct such examples. Our experiments show that state-of-the-art models consistently fail to recognize them as ill-formed, and instead produce high confidence predictions on them. As a consequence of this phenomenon, models trained on sentences with randomly permuted word order perform close to state-of-the-art models. To alleviate these issues, we show that if models are explicitly trained to recognize invalid inputs, they can be robust to such attacks without a drop in performance.

Introduction

The BERT family of models (Devlin et al. 2019; Liu et al. 2019, and others) form the backbone of today’s NLP systems. At the time of writing, all eleven systems deemed to outperform humans in the GLUE benchmark suite (Wang et al. 2018) belong to this family. Do these models understand language? Recent work suggests otherwise. For example, Bender and Koller (2020) point out that models trained to mimic linguistic form (i.e., language models) may be deficient in understanding the meaning conveyed by language.

In this paper, we show that such models struggle even with the form of language by demonstrating that they force meaning onto token sequences devoid of any. For instance, consider the natural language inference (NLI) example in fig. 1. A RoBERTa-based model that scores $\sim 89\%$ on the Multi-NLI dataset (Williams, Nangia, and Bowman 2018) identifies that the premise entails the hypothesis. However, when the words in the hypothesis are sorted alphabetically (thereby rendering the sequence meaningless), the model still makes the same prediction with high confidence. Indeed, across Multi-NLI, when the hypotheses are sorted alphabetically, the model retains the same prediction in 79% of the cases, with a surprisingly high average confidence of $\sim 95\%$! We argue that a reliable model should not be insensitive to such a drastic change in word order.

We study the response of large neural models to *destructive transformations*: perturbations of inputs that ren-

Premise	In reviewing this history, it’s important to make some crucial distinctions.
Original Hypothesis	Making certain distinctions is imperative in looking back on the past. ENTAILMENT <i>with probability 0.99</i>
Sorted Hypothesis	back certain distinctions imperative in is looking making on past the . ENTAILMENT <i>with probability 0.97</i>

Figure 1: An Example for Natural Language Inference from the MNLi (-m) validation set. For the original premise and hypothesis, RoBERTa fine-tuned model makes the correct prediction that the premise entails the hypothesis. Alphabetically sorting the hypothesis makes it meaningless, but the model still retains the prediction with high confidence.

der them meaningless. Figure 1 shows an example. We define several such transformations, all of which erase meaning from the input text and produce token sequences that are not natural language (i.e., word salad).

We characterize the response of models to such transformations using two metrics: its ability to predict valid labels for invalid inputs, and its confidence on these predictions. Via experiments on three tasks from the GLUE benchmark, we show that the labels predicted by state-of-the-art models for destructively transformed inputs bear high agreement with the original ones. Moreover, the models are highly confident in these predictions. We also find that models trained on meaningless examples perform comparably to the original model on unperturbed examples, despite never having encountered any well-formed training examples. Specifically, models trained on meaningless sentences constructed by permuting the word order perform almost as well as the state-of-the-art models. These observations suggest that, far from actually understanding natural language, today’s state-of-the-art models have trouble even *recognizing* it.

Finally, we evaluate strategies to mitigate these weaknesses using regularization that makes models less confident in their predictions, or by allowing models to reject inputs.

In summary, our contributions are¹:

¹Our code is available at <https://github.com/utahnlp/word-salad>

Dataset	Transform	Input	Prediction
Natural Language Inference MNL1	Original Shuffled PBSMT-E	P: As with other types of internal controls, this is a cycle of activity, not an exercise with a defined beginning and end.	
		H: There is no clear beginning and end, it’s a continuous cycle.	Ent (99.48%)
		H ₁ : , beginning end no there clear ’s continuous is a it and cycle .	Ent (99.60%)
		H ₂ : The relationship of this is not a thing in the beginning .	Ent (94.82%)
Paraphrase Detection QQP	Original Repeat CopySort	Q1: How do I find out what operating system I have on my Macbook?	
		Q2: How do I find out what operating system I have?	Yes (99.53%)
		Q2: out out i find out what out out i find?	Yes (99.98%)
		Q2: ? do find have how i i macbook my on operating out system what	Yes (98.52%)
Sentiment Analysis SST-2	Original Sort Drop	A by-the-numbers effort that won’t do much to enhance the franchise.	-ve (99.96%)
		a by-the-numbers do effort enhance franchise much n’t that the to wo.	-ve (99.92%)
		a-n won do to franchise.	-ve (99.96%)

Table 1: We generate invalid token sequences using destructive transformations that render the inputs meaningless. A fine-tuned RoBERTa model assigns a high probability (in parenthesis) to the same label as the original example. For NLI, the model chooses between *entail*, *contradict*, and *neutral*. For sentiment analysis, possible labels are -ve or +ve. For paraphrase detection, model answers if the two texts are paraphrases of each other (Yes or No). The appendix contains more such examples.

1. We define the notion of destructive input transformations to test the ability of text understanding models at processing word salad. We introduce nine such transformation functions that can be used by practitioners for diagnostic purposes without requiring additional annotated data.
2. We show via experiments that today’s best models force meaning upon invalid inputs; i.e., they are not using the right kind of information to arrive at their predictions.
3. We show that simple mitigation strategies can teach models to recognize and reject invalid inputs.

Tasks and Datasets

Our goal is to demonstrate that state-of-the-art models based on the BERT family do not differentiate between valid and invalid inputs, and that this phenomenon is ubiquitous. To illustrate this, we focus on three tasks (table 1), which also serve as running examples.

Natural language inference (NLI) is the task of determining if a premise entails, contradicts, or is unrelated to a hypothesis. We use the MNLI (Williams, Nangia, and Bowman 2018) and SNLI (Bowman et al. 2015) datasets.

Paraphrase detection involves deciding if two sentences are paraphrases of each other. For this task, we use the Microsoft Research Paraphrase Corpus (MRPC, Dolan and Brockett 2005), and Quora Question Pair (QQP) dataset².

Sentiment classification requires predicting whether a sentence has a positive or negative sentiment. We use the Stanford Sentiment Treebank (SST-2, Socher et al. 2013).

Destructive Transformations

There has been a growing interest in studying input perturbations (e.g., Ebrahimi et al. 2018; Alzantot et al. 2018; Wallace et al. 2019; Jin et al. 2020; Ren et al. 2019; Garg et al.

²<https://www.quora.com/q/quoradata/First-Quora-Dataset-Release-Question-Pairs>

2020). Given a model for a task, some input perturbations preserve labels. For example, a true paraphrase of a sentence should not change its sentiment. Certain other perturbations force labels to change in controlled ways. For example, a negated hypothesis should change the label in the NLI task.

In this paper, we focus on a new class of perturbations—*destructive transformations*—which render inputs invalid. Because any informative signal in the input is erased, the transformed examples should not have *any* correct label.³

For example, in the NLI task, a hypothesis whose words are shuffled is (with high probability) not a valid English sentence. The transformed hypothesis cannot contain information to support an *entail* or a *contradict* decision. Moreover, it is not a sentence that is unrelated to the premise—it is not a sentence at all! Indeed, the transformation creates an example that lies outside the scope of the NLI task.

Yet, the premise and transformed hypothesis fit the interface of the problem. This gives us our key observation: While NLP models are typically trained to work on well-formed sentential inputs, they accept any sequence of strings as inputs. Of course, not all token sequences are meaningful sentences; we argue that we can gain insights about models by studying their response to meaningless input sequences (i.e., *invalid inputs*) that are carefully crafted from valid inputs. One advantage of this protocol of using transformations is that we do not need to collect any new data and can use original training and validation sets.

Let us formally define destructive transformations. Consider a task with input $x \in X$ and an oracle function f that maps inputs to labels $y \in Y$. A destructive transformation $\pi : X \rightarrow X$ is a function that operates on x to produce transformed inputs $x' = \pi(x)$ such that $f(x')$ is *undefined*. That is, none of the labels (i.e., the set Y) can apply to x' .

³We refer to such changes to inputs as transformations, instead of perturbations (as in adversarial perturbations), to highlight the fact that change in the inputs need not be small.

Name	Description
Sort	Sort the input tokens
Reverse	Reverse the token sequence
Shuffle	Randomly shuffle tokens
CopySort	Copy one of the input texts and then sort it to create the second text. (Only applicable when the input is a pair of texts)

Table 2: Lexical-overlap based transformations

Destructive transformations can be chained: if π_1 and π_2 are destructive transformations for a given input x , then $\pi_1(\pi_2(x))$ is also a destructive transformation. We can think of such chaining as combining different diagnostic heuristics. For tasks whose input is a pair of texts, $x = (x_1, x_2)$, transforming either or both the components should destroy any meaningful signal in the input x . For example, in the NLI task, given an input premise and hypothesis, destroying either of them renders the example invalid. For tasks with a pair of texts as input, for our analyses, we only transform one of the inputs, although there is no such requirement.

Next, let us look at different classes of transformations.

Lexical Overlap-based Transformations

These transformation operators preserve the bag-of-words representation of the original input but change the word order. They are designed to diagnose the sensitivity of models to the order of tokens in inputs. Table 2 shows the four lexical-overlap based transformations we define here.

We ensure that `Shuffle` sufficiently changes the input by repeatedly shuffling till no bigram from original input is retained. The `CopySort` operation only applies to tasks that have multiple input texts such as NLI and paraphrase detection. As an example, for the NLI task, given a premise-hypothesis pair, it creates a transformed pair whose hypothesis is the alphabetically sorted premise.

Gradient-based Transformations

These transformations seek to study the impact of removing, repeating, and replacing tokens. To decide which tokens to replace, they score input tokens in proportion to their relative contribution to the output. One (admittedly inefficient) way to compute token importance is to calculate the change in output probability when it is removed. Recent work (e.g., Ebrahimi et al. 2018; Feng et al. 2018) suggests that a gradient-based method is a good enough approximation and is much more efficient. We adopt this strategy here.

Given a trained neural model \mathcal{M} , and the task loss function \mathcal{L} , the change in the loss for the i^{th} input token is approximated by the dot product of its token embedding \mathbf{t}_i and the gradient of the loss propagated back to the input layer $\nabla_{\mathbf{t}_i, \mathcal{M}} \mathcal{L}$. That is, the i^{th} token is scored by $\mathbf{t}_i^T \nabla_{\mathbf{t}_i, \mathcal{M}} \mathcal{L}$.

These token scores approximate the relative importance of a token; a higher score denotes a more important token. We use the tokens in the *bottom* $r\%$ as per their score—the least important tokens—to define our gradient-based trans-

Name	Description
Drop	Drop the least important tokens.
Repeat	Replace the least important tokens with one of the most important ones.
Replace	Replace the least important tokens with random tokens from the vocabulary
CopyOne	Copy the most important token from one text as the sole token in the other. (Only applicable when the input is a pair of texts)

Table 3: Gradient-based transformations

formations. We use $r = 50\%$. Table 3 summarizes the transformations that use importance ranking of the tokens.

Statistical transformation: PBSMT

Recent analyses on the NLI task have shown that neural models rely excessively on shallow heuristics (Gururangan et al. 2018; Poliak et al. 2018; McCoy, Pavlick, and Linzen 2019). In particular, Gururangan et al. (2018) showed that annotation artifacts lead to certain words being highly correlated with certain inference classes. For example, in the SNLI data, words such as *animal*, *outdoors* are spuriously correlated with the *entail* label.

Inspired by this observation, we design a transformation scheme that creates invalid examples, and yet exhibit such statistical correlations. We employ a traditional phrase-based statistical machine translation (PBSMT) system to generate examples that use phrasal co-occurrence statistics.

For each label in the task, we train a separate sequence generator that uses co-occurrence statistics for that label. For example, for the NLI task, we have three separate generators, one for each label. Suppose we have a premise-hypothesis pair that is labeled as *entail*. We destroy it using the premise as input to a PBSMT system that is trained only on the entailment pairs in the training set. We use the Moses SMT toolkit (Koehn et al. 2007) for our experiments.

Why should a system trained to generate a sentence that has a certain label (e.g., an entailment) be a destructive transformer? To see this, note that unlike standard machine translation, we use very limited data for training. Moreover, the language models employed (Heafield 2011) are also trained only on examples of one class. As a result, we found that the produced examples are non-grammatical, and often, out of context. The hypothesis H_2' in table 1, generated using PBSMT-E (i.e., PBSMT for entailments), is one such sequence. We refer to this transformation as PBSMT.

Are the Transformations Destructive?

To ascertain whether our nine transformations render sentences invalid, we asked crowd workers on Amazon’s Mechanical Turk to classify each instance as *valid* or *invalid*, the latter category is defined as an example that is incomprehensible, and is therefore meaningless.

We sampled 100 invalid sentences generated by each transformation (900 in total) and an equal number of sentences from the original (*un-transformed*) validation sets.

Transformation	% Invalid
Un-transformed	7.83
Sort	94.07
Reverse	95.59
Shuffle	94.20
CopySort	95.42
Avg. Lexical	94.82
Replace	91.21
Repeat	100.00
Drop	85.79
CopyOne	100.00
Avg. Gradient	94.25
PBSMT	79.92

Table 4: Results of crowdsourcing experiments where annotators are asked to label sentences as meaningful or not. We aggregate sentence labels from three workers.

For each sentence, we collect validity judgments from three crowd workers and use the majority label. Table 4 shows the percent of sentences marked as invalid; we see that all the transformations make their inputs incomprehensible.

Measuring Responses to Invalid Inputs

We will now define two metrics to quantify model behavior for invalid inputs. Invalid inputs, by definition, are devoid of information about the label. Consequently, they *do not* have a correct label. If the transformations are truly destructive, a *reliable* model will pick one of the labels at random and would do so with low confidence. That is, a reliable model should exhibit the following behavior: a) the agreement between original predictions and predictions on their transformed invalid variants should be random, and, b) predictions for invalid examples should be uncertain. These expected behaviors motivate the following two metrics.

Agreement is the % of examples whose prediction remains same after applying a destructive transformation. A model with agreement closer to random handles invalid examples better. For the operators designed for tasks with a pair of inputs, namely `CopySort` and `CopyOne`, tokens from one of the inputs are copied into another. In such cases, measuring agreement with original input is not useful. Instead, we measure agreement with a default label. For the NLI task, the default label is *entail*, because neural models tend to predict entailment when there is a lexical overlap between the premise and hypothesis (McCoy, Pavlick, and Linzen 2019). Following the same intuition, for paraphrase detection, the default is to predict that the pair is a paraphrase.

Confidence is defined as the average probability of the predicted label. We want this number to be closer to $\frac{1}{N}$, where N is the number of classes.⁴

⁴We could alternatively define confidence using the entropy of the output distribution. In our experiments, we found that confidence, as defined here, and entropy reveal the same insights.

Dataset	Accuracy	Confidence
SNLI	90.87	98.38
MNLI	87.31	98.27
QQP	90.70	98.89
MRPC	89.46	98.40
SST-2	94.04	99.75

Table 5: **Baseline Performance**, Accuracy and Average Confidence for RoBERTa-*base* on validation sets. For MNLI, we used MNLI-matched for experiments.

Transform	MNLI	SNLI	QQP	MRPC	SST2
Sort	79.1	82.6	88.3	81.1	83.3
Reverse	76.9	75.1	86.8	77.9	82.5
Shuffle	79.4	81.1	88.4	80.4	84.8
CopySort	90.5	81.3	93.5	96.8	–
Avg. Lex.	82.4	80.1	89.3	84.1	83.5
Replace	63.0	51.9	69.9	56.6	78.1
Repeat	49.7	68.5	77.1	68.1	81.3
Drop	69.4	72.7	80.4	76.7	82.5
CopyOne	80.4	83.7	98.9	100	–
Avg. Grad.	65.6	69.2	81.6	75.4	80.6
PBSMT	57.0	65.6	72.5	–	75.2
Random	33.3	33.3	50.0	50.0	50.0

Table 6: Agreement scores between predictions from transformed validation set and original validation set. The closer the numbers are to random better the model behavior is. ‘–’ means the transformation is not defined for that dataset. We do not use PBSMT for MRPC as it is a much smaller dataset.

Experiments

For our primary set of experiments, we use RoBERTa (Liu et al. 2019) as a representative of the BERT family, whose fine-tuned versions have tended to outperform their BERT counterparts on the GLUE tasks (Dodge et al. 2020). We use the *base* variant of RoBERTa that is fine-tuned for three epochs across all our experiments, using hyperparameters suggested by the original paper. These models constitute our baseline. Table 5 shows the accuracy and average confidence of the baseline on the original validation sets.

Results and Observations

We apply the destructive transformation functions described earlier to each task’s validation set. To account for the randomness in the `Shuffle` transformation, its results are averaged across five random seeds. For `PBSMT` on SST-2, we use the first half sentence as input and train to predict the second half of the sentence. Table 6 shows the agreement results, and table 7 shows average confidence scores.

High Agreement Scores The high agreement scores show that models retain their original predictions even when label-

	MNLI	SNLI	QQP	MRPC	SST-2
Baseline	94.63	92.46	98.78	97.77	99.13
Random	33.33	33.33	50.00	50.00	50.00

Table 7: Average Confidence over predictions from transformed validation set. We want the the numbers to be closer to random (last row). Refer appendix for full results.

bearing information is removed from examples. This is a puzzling result: the transformations render sentences meaningless to humans, but the model knows the label. How can the model make sense of these nonsensical inputs?

We argue that this behavior is not only undesirable but also brings into question the extent to which these models understand text. It is possible that, rather than understanding text, they merely learn spurious correlations in the training data. That is, models use the wrong information to arrive at the right answer.

High Confidence Predictions. Not only do models retain a large fraction of their predictions, they do so with high confidence (table 7). This behavior is also undesirable: a *reliable* model should know what it does not know, and should not fail silently. It should, therefore, exhibit be uncertain on examples that are uninformative about the label.

Research on reliability of predictions suggests that these models are poorly calibrated (Guo et al. 2017).

Specific transformations. Invalid examples constructed by lexical transformations are more effective than others, with all agreements over 80%. Examples from such transformations have high lexical overlap with the original input. Our results suggest that models do not use input token positions effectively. We need models that are more sensitive to word order; lexical transformations can be used as a guide without the need for new test sets.

We find that SNLI models have higher agreement scores than MNLI ones for both gradient and statistical correlation based invalid examples. This could mean that the former are more susceptible to gradient-based adversarial attacks. Moreover, the lower scores for PBSMT on the MNLI model shows that it relies less on these statistical clues than SNLI—corroborating an observation by Gururangan et al. (2018).

Human response to transformed inputs. Results in table 4 show that transformed sentences are invalid. We now perform another set of human experiments to determine if the invalid examples generated (by transformations) make it difficult to perform the classification task. This mimics the exact setting that all models are evaluated on by asking humans to perform classification tasks on invalid inputs. Concretely, we ask turkers to perform the NLI task on 450 destructively transformed inputs (50 for each transformation) by “reconstructing the inputs to the best of their abilities”. We found that turkers can only ‘predict’ the correct label for invalid examples in 35% of the cases as opposed to 77% for

Calibration Method	Accuracy	ECE
Baseline	87.31	0.11
Label Smoothing	86.89	0.06
Focal Loss	86.98	0.05
Temperature Scaling	87.31	0.09

Table 8: Accuracy on the original validation set and the Expected Calibration Error (ECE) on the validation set for MNLI. Accuracy with temperature scaling is the same as baseline since it is a post-training method for calibration.

	B	LS	FL	B + TS
Lexical	82.35	81.85	80.39	81.49
Gradient	60.65	59.51	60.32	59.18
PBSMT	57.02	56.30	56.49	57.04

Table 9: Average agreement for three calibration methods on MNLI. Calibration does not improve model’s response to invalid inputs. B: Baseline, LS: Label Smoothing, FL: Focal Loss, B + TS: Temperature Scaling on baseline.

original un-transformed examples. These results reinforce the message that large transformer models can make sense of meaningless examples, whereas humans are near-random.

Analysis & Discussion

Are calibrated models more reliable? Neural networks have been shown to produce poorly calibrated probabilities, resulting in high confidence even on incorrect predictions (Guo et al. 2017). Research in computer vision has shown that improving model calibration improves adversarial robustness as well as out-of-distribution detection (Hendrycks and Gimpel 2017; Thulasidasan et al. 2019; Hendrycks, Lee, and Mazeika 2019). Given the confidence scores in table 7, a natural question is: *Does improving the calibration of BERT models improve their response to invalid examples?* We answer this question by training confidence calibrated classifiers using three standard methods.

First we use *label smoothing*, in which training is done on soft labels, with loss function being a weighted average of labels and uniform probability distribution (Pereyra et al. 2017). *Focal loss* prevents the model from becoming overconfident on examples where it is already correct. Mukhoti et al. (2020) showed that focal loss improves calibration of neural models. *Temperature scaling* is a simple calibration method that scales the network’s logit values before applying the softmax (Guo et al. 2017; Desai and Durrett 2020).

We use *Expected Calibration Error* (ECE, Naeini, Cooper, and Hauskrecht 2015) to measure a model’s calibration error. Due to space constraints, we refer the reader to the original work for a formal definition. Better calibrated models have lower ECE. All three methods improve calibration of the original model; table 8 shows results on the MNLI validation data. However, table 9 shows that none of them improve model response to invalid examples.

Impact of pretraining tasks. We now investigate the impact of pre-training tasks on a model’s response to invalid examples. Both BERT and RoBERTa use a word-based masked language modeling (W-MLM) as the auto-encoding objective. BERT uses Next Sentence Prediction (NSP) as an additional pre-training task. We experiment with other BERT variants pre-trained with different tasks: ALBERT (Lan et al. 2019) uses Sentence Order Prediction (SOP), SpanBERT (Joshi et al. 2020) and BART (Lewis et al. 2020) use Span-based MLM (S-MLM) instead of W-MLM. SpanBERT additionally uses NSP, while BART uses a Sentence Shuffling (SS) pretraining objective. ELECTRA (Clark et al. 2019) uses a Replaced Token Detection (RTD) instead of an MLM objective.

These models are trained on different corpora, and use different pre-training tasks. Despite their differences, the results presented in table 10 suggest that all of these models are similar in their responses to invalid examples. These results highlight a potential weakness in our best text understanding systems.

Different Inductive Bias. All variants of BERT considered thus far are trained with one of the auto-encoding (AE) objectives and perform rather poorly. This raises a question: *Would models that explicitly inject a word order based inductive bias into the model perform better?*

To answer this question, we consider three auto-regressive (AR) models with a recurrent inductive bias, namely, ESIM-Glove (Chen et al. 2017), ESIM-ELMo, and XLNet (Yang et al. 2019). Both ESIM models are LSTM based models, while XLNet is a transformer-based model that is trained using an auto-regressive language modeling objective along with Permutation LM (P-LM). ESIM-Glove does not use any other pre-training task, while ESIM-ELMo is based on ELMo (Peters et al. 2018) which is pre-trained as a traditional auto-regressive LM.

The results are shown in table 10. Again, the results are similar to models trained with auto-encoding objective. Surprisingly, even a strong recurrent inductive bias is unable to make the models sensitive to the order of words in their inputs: all the AR models have high agreement scores (over 75%) on lexical overlap-based transformations. We refer the reader to the appendix for more results.

Bigger is not always better. While larger BERT-like models show better performance (Devlin et al. 2019; Raffel et al. 2020), we find that same does *not* hold for their response on invalid examples. Table 10 shows that larger BERT models (Large vs Base) do not improve response to invalid examples (recall that smaller agreement scores are better). We see that both BERT variants outperform the RoBERTa counterparts; BERT-*base* provides over 4.5% improvement over RoBERTa-*base* in terms of agreement on invalid examples.

Small vs. large perturbations. Previous work on adversarial robustness (Alzantot et al. 2018; Jin et al. 2020; Ebrahimi et al. 2018) suggests that robustness of the model to small input perturbations is desirable, meaning that a

Class	Pretraining	Model	Agreement
AE	W-MLM + NSP	BERT-B	67.1
		BERT-L	69.0
	W-MLM	RoBERTa-B	71.7
		RoBERTa-L	73.5
	W-MLM + SOP	ALBERT-B	67.6
AE	S-MLM + NSP	SpanBERT-B	67.6
	S/W-MLM + SS	BART-B	70.0
	RTD	ELECTRA-B	68.8
AR	P-LM	XLNet-B	70.1
	LM	ESIM- ELMo	75.8
	-	ESIM- Glove	73.5

Table 10: Agreement score of different models on MNLI. **B** refers to the base variant, **L** refers to the large one. **AE** refers to models pretrained with auto-encoding objective, **AR** refers to auto-regressive models. Refer text for full key.

	MNLI	SNLI	QQP	MRPC	SST-2
Shuffled	84.56	89.44	92.15	84.80	91.97
Original	87.31	90.70	94.04	89.46	94.04

Table 11: Training on only invalid examples generated from `Shuffle`, evaluation is on original validation data.

model’s prediction should not change for small perturbations in the input. However, excessive invariance to large input perturbations is undesirable (Jacobsen et al. 2019). Our focus is not on small input changes, rather large ones that destroy useful signals (i.e., destructive transformations). The three types of transformations we discuss in this work achieve this in different ways. We argue that language understanding systems should not only provide robustness against small perturbations (adversarial robustness) but also recognize and reject large perturbations (studied in this work).

Are models learning spurious correlations? The results presented in this work raise an important question: *Why does this undesirable model behavior occur in all models, irrespective of the pretraining tasks, and is even seen in models with a recurrent inductive bias?* We hypothesize that this behavior occurs because these large models learn spurious correlations present in the training datasets, studied previously by Gururangan et al. (2018); Min et al. (2020). A simple experiment substantiates this claim. So far, we trained models on valid data and evaluated them on both valid and invalid examples. We now flip this setting: we train on *invalid* inputs generated by a transformation and evaluate on well-formed examples from the validation sets.

Table 11 presents accuracies on the original validation examples for five datasets. We observe that models trained only on shuffled examples perform nearly as well (within 97% for MNLI) as the ones trained on valid examples (second row)! These observations demonstrate that our models do not use the right kind of evidence for their predictions, or at least, the

	B + Th	Ent + Th	B + Invalid
MNLI	83.40/ 57.95	86.11 / 89.22	85.44/ 97.10
SNLI	88.01/ 54.68	89.65/ 93.54	90.88 / 98.41
QQP	90.25 / 29.20	90.29 / 88.72	90.08 / 95.24
MRPC	89.46 / 36.82	88.24 / 99.43	88.73 / 99.78
SST-2	90.37 / 35.79	92.66 / 95.41	92.78 / 96.35

Table 12: **Comparison of mitigation strategies.** First number in each cell is accuracy on original validation set. Second number is the % of examples correctly classified as invalid. Test set for invalid contains examples generated with all nine transformation functions. **B** refers to the baseline model.

kind of evidence a human would use. This result should raise further concerns about whether we have made real progress on language understanding.

Mitigation Strategies

We evaluate three mitigation strategies to alleviate the problem of high certainty on invalid inputs. The goal is to give models the capability to recognize invalid inputs. Two strategies augment the training data with invalid examples. All three introduce new hyperparameters, which are tuned on a validation set constructed by sampling 10% of the training set. The final models are then trained on the full training set.

Entropic Regularization The central problem that we have is that the models have high certainty on invalid inputs. To directly alleviate the issue, we can explicitly train them to be less certain on invalid inputs. We do so by augmenting the loss function with a new term. Let D be the original dataset and D' be its complementary invalid dataset. The new training objective is then defined as

$$L(\text{model}) = L_D(\text{model}) + \lambda H_{D'}(\text{model}) \quad (1)$$

where L_D is the standard cross-entropy loss, and $H_{D'}$ denotes the entropy of model probabilities over invalid examples. The hyperparameter λ weighs the relative impact of the two terms. Feng et al. (2018) used similar entropic regularization to improve interpretability of neural models.

We initialize the model with fine-tuned weights from the baseline model and train for three more epochs with the new training objective. The appendix provides further details.

Thresholding Model Probabilities From our results in table 7, we observe that although models are confident on invalid examples, their confidence is higher on valid ones. Following this, we experiment with a straightforward approach that thresholds output probabilities to tell valid and invalid examples apart. We used temperature scaling to ensure that the classifier probabilities are calibrated.

This approach is parameterized by a threshold θ : if the maximum probability of the classifier’s output is below θ , we deem the input invalid. We used grid search in the range $[\frac{1}{N}, 1.0]$ to find the best performing θ on a separate validation set. Here, N represents the number of labels.

Invalid as an extra class (B + Invalid) Since one of the goals is to be able to recognize invalid inputs, we can explicitly introduce a new class, *invalid*, to our classification task. The training objective for this new $N + 1$ class classification task remains the same, i.e. cross-entropy loss.

Results

With entropic regularization, we observe a significant drop in agreement scores on invalid examples. Indeed, the agreement scores on invalid examples decrease to an average of 35% after regularization. We also notice a significant increase in uncertainty on invalid examples. However, we see that in some cases, accuracies on the original validation set drop by over 1%, suggesting a trade-off between accuracy on valid examples and reliable response on invalid examples.

A well-behaved model should maintain high accuracy on valid data and also reject invalid inputs. To compare the three mitigation strategies on an equal footing, we measure accuracy on the original validation data and the percentage of invalid examples correctly identified. Since entropic regularization does not explicitly recognize invalid inputs, we apply the thresholding strategy to it (**Ent + Th**).

Table 12 compares the three methods. We see that simple thresholding (**B + Th**) does not work well and having the model learn from invalid examples is beneficial. It appears that, out of the three methods, training with an extra *invalid* label best maintains the balance between accuracy and invalid input detection.

We also studied if mitigating one kind of transformation helps against others. Using (B + Invalid) we train on one transformation and test on the rest. We found that mitigation can be transferable. The appendix provides detailed results.

Final Words

The main message of this paper is that today’s state-of-the-art models for text understanding have difficulties in telling the difference between valid and invalid text. This observation is congruent with several other recent lines of work that highlight the deficiencies of today’s text understanding systems. For example, Feng et al. (2018) construct irregular examples by successively removing words without affecting a neural model’s predictions. Adversarial attacks on NLP models (e.g. Jin et al. 2020) expose their vulnerabilities; for example, Wallace et al. (2019) offer an illustrative example where a model’s prediction can be arbitrarily changed.

Statistical models need not always get well-formed inputs. Consequently, when models are deployed, they should be guarded against invalid inputs, not just in NLP, but also beyond (e.g., Liang, Li, and Srikant 2018). Krishna et al. (2020) showed that it is possible to steal BERT-like models by using their predictions on meaningless inputs. Our work can be seen as highlighting why this might be possible: model predictions are largely not affected even if we destructively transform inputs. Our work shows simple mitigation strategies that could become a part of the standard modeling workflow.

Ethics Statement

Our work points to a major shortcoming of the BERT family of models: they have a hard time recognizing ill-formed inputs. This observation may be used to construct targeted attacks on trained models, especially publicly available ones. We call for a broader awareness of such vulnerabilities among NLP practitioners, and recommend that NLP models should be actively equipped with mitigations against such attacks.

Acknowledgments

We thank members of the Utah NLP group for their valuable insights, and reviewers for their helpful feedback. This material is based upon work supported by NSF under grants #1801446 (SATC) and #1822877 (Cyberlearning).

References

- Alzantot, M.; Sharma, Y.; Elgohary, A.; Ho, B.-J.; Srivastava, M.; and Chang, K.-W. 2018. Generating Natural Language Adversarial Examples. In *EMNLP*.
- Bender, E. M.; and Koller, A. 2020. Climbing towards NLU: On Meaning, Form, and Understanding in the Age of Data. In *ACL*.
- Bowman, S.; Angeli, G.; Potts, C.; and Manning, C. D. 2015. A Large Annotated Corpus for Learning Natural Language Inference. In *EMNLP*.
- Chen, Q.; Zhu, X.; Ling, Z.-H.; Wei, S.; Jiang, H.; and Inkpen, D. 2017. Enhanced LSTM for Natural Language Inference. In *ACL*.
- Clark, K.; Luong, M.-T.; Le, Q. V.; and Manning, C. D. 2019. ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators. In *ICLR*.
- Desai, S.; and Durrett, G. 2020. Calibration of Pre-trained Transformers. In *EMNLP (1)*.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL*.
- Dodge, J.; Ilharco, G.; Schwartz, R.; Farhadi, A.; Hajishirzi, H.; and Smith, N. 2020. Fine-tuning Pretrained Language Models: Weight Initializations, Data Orders, and Early Stopping. *arXiv preprint arXiv:2002.06305*.
- Dolan, W. B.; and Brockett, C. 2005. Automatically Constructing a Corpus of Sentential Paraphrases. In *IWP@IJCNLP*.
- Ebrahimi, J.; Rao, A.; Lowd, D.; and Dou, D. 2018. HotFlip: White-Box Adversarial Examples for Text Classification. In *ACL*.
- Feng, S.; Wallace, E.; Grissom II, A.; Iyyer, M.; Rodriguez, P.; and Boyd-Graber, J. 2018. Pathologies of Neural Models Make Interpretations Difficult. In *EMNLP*.
- Garg, S.; Kumar, A.; Goel, V.; and Liang, Y. 2020. Can Adversarial Weight Perturbations Inject Neural Backdoors. In *CIKM*.
- Guo, C.; Pleiss, G.; Sun, Y.; and Weinberger, K. Q. 2017. On Calibration of Modern Neural Networks. In *ICML*.
- Gururangan, S.; Swayamdipta, S.; Levy, O.; Schwartz, R.; Bowman, S.; and Smith, N. A. 2018. Annotation Artifacts in Natural Language Inference Data. In *NAACL*.
- Heafield, K. 2011. KenLM: Faster and smaller language model queries. In *WMT@EMNLP*.
- Hendrycks, D.; and Gimpel, K. 2017. A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks. In *ICLR*.
- Hendrycks, D.; Lee, K.; and Mazeika, M. 2019. Using Pre-Training Can Improve Model Robustness and Uncertainty. In *ICML*.
- Jacobsen, J.; Behrmann, J.; Zemel, R. S.; and Bethge, M. 2019. Excessive Invariance Causes Adversarial Vulnerability. In *ICLR*.
- Jin, D.; Jin, Z.; Zhou, J. T.; and Szolovits, P. 2020. Is BERT Really Robust? A Strong Baseline for Natural Language Attack on Text Classification and Entailment. In *AAAI*.
- Joshi, M.; Chen, D.; Liu, Y.; Weld, D. S.; Zettlemoyer, L.; and Levy, O. 2020. SpanBERT: Improving Pre-training by Representing and Predicting Spans. *TACL* 8.
- Koehn, P.; Hoang, H.; Birch, A.; Callison-Burch, C.; Federico, M.; Bertoldi, N.; Cowan, B.; Shen, W.; Moran, C.; Zens, R.; et al. 2007. Moses: Open source Toolkit for Statistical Machine Translation. In *ACL (Demos)*.
- Koehn, P.; Och, F. J.; and Marcu, D. 2003. Statistical Phrase-based Translation. In *NAACL*.
- Krishna, K.; Tomar, G. S.; Parikh, A. P.; Papernot, N.; and Iyyer, M. 2020. Thieves on Sesame Street! Model Extraction of BERT-based APIs. In *ICLR*.
- Lan, Z.; Chen, M.; Goodman, S.; Gimpel, K.; Sharma, P.; and Soricut, R. 2019. ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. In *ICLR*.
- Lewis, M.; Liu, Y.; Goyal, N.; Ghazvininejad, M.; Mohamed, A.; Levy, O.; Stoyanov, V.; and Zettlemoyer, L. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *ACL*.
- Liang, S.; Li, Y.; and Srikant, R. 2018. Enhancing the Reliability of Out-of-distribution Image Detection in Neural Networks. In *ICLR*.
- Lin, T.-Y.; Goyal, P.; Girshick, R.; He, K.; and Dollár, P. 2017. Focal Loss for Dense Object Detection. In *ICCV*.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv preprint arXiv:1907.11692*.
- McCoy, T.; Pavlick, E.; and Linzen, T. 2019. Right for the Wrong Reasons: Diagnosing Syntactic Heuristics in Natural Language Inference. In *ACL*.
- Min, J.; McCoy, R. T.; Das, D.; Pitler, E.; and Linzen, T. 2020. Syntactic Data Augmentation Increases Robustness to Inference Heuristics. In *ACL*.

Mukhoti, J.; Kulharia, V.; Sanyal, A.; Golodetz, S.; Torr, P. H. S.; and Dokania, P. K. 2020. Calibrating Deep Neural Networks using Focal Loss. In *NeurIPS*.

Naeini, M. P.; Cooper, G.; and Hauskrecht, M. 2015. Obtaining well Calibrated Probabilities using Bayesian Binning. In *AAAI*.

Pereyra, G.; Tucker, G.; Chorowski, J.; Kaiser, L.; and Hinton, G. E. 2017. Regularizing Neural Networks by Penalizing Confident Output Distributions. In *ICLR (Workshop)*.

Peters, M.; Neumann, M.; Iyyer, M.; Gardner, M.; Clark, C.; Lee, K.; and Zettlemoyer, L. 2018. Deep Contextualized Word Representations. In *NAACL*.

Poliak, A.; Naradowsky, J.; Haldar, A.; Rudinger, R.; and Van Durme, B. 2018. Hypothesis Only Baselines in Natural Language Inference. In *NAACL*.

Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research* 21(140): 1–67.

Ren, S.; Deng, Y.; He, K.; and Che, W. 2019. Generating Natural Language Adversarial Examples through Probability Weighted Word Saliency. In *ACL*.

Socher, R.; Perelygin, A.; Wu, J.; Chuang, J.; Manning, C. D.; Ng, A.; and Potts, C. 2013. Recursive Deep Models for Semantic Compositionality over a Sentiment Treebank. In *EMNLP*.

Thulasidasan, S.; Chennupati, G.; Bilmes, J. A.; Bhattacharya, T.; and Michalak, S. 2019. On Mixup Training: Improved Calibration and Predictive Uncertainty for Deep Neural Networks. In *NeurIPS*.

Wallace, E.; Feng, S.; Kandpal, N.; Gardner, M.; and Singh, S. 2019. Universal Adversarial Triggers for Attacking and Analyzing NLP. In *EMNLP-IJCNLP (1)*.

Wang, A.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. 2018. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. In *BlackboxNLP@EMNLP*.

Williams, A.; Nangia, N.; and Bowman, S. 2018. A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference. In *NAACL*.

Wolf, T.; Debut, L.; Sanh, V.; Chaumond, J.; Delangue, C.; Moi, A.; Cistac, P.; Rault, T.; Louf, R.; Funtowicz, M.; Davison, J.; Shleifer, S.; von Platen, P.; Ma, C.; Jernite, Y.; Plu, J.; Xu, C.; Le Scao, T.; Gugger, S.; Drame, M.; Lhoest, Q.; and Rush, A. 2020. Transformers: State-of-the-Art Natural Language Processing. In *EMNLP (Demos)*.

Yang, Z.; Dai, Z.; Yang, Y.; Carbonell, J.; Salakhutdinov, R. R.; and Le, Q. V. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *NeurIPS*.

More Examples

Table 13 shows more invalid examples generated using our transformations.

Additional Technical Details and Observations

The following list contains some of the technical details and observations not mentioned in the paper:

1. For all the lexical overlap transformations, we ensure that the punctuation at the end of sentence remains at the end.
2. We found that some of the datasets had a formatting problem. For example, one example in QQP validation set only had one input. We exclude such examples from our data.
3. We found that for PBSMT models created for NLI task, agreement scores were much higher for classes: entailment, and contradiction, than neutral class.

Additional Details for Crowd-Sourcing Experiments

We collected crowd annotations for ascertaining the (in)validity of examples generated by our transformation functions. Here, we discuss the instructions we provided to annotators along with costs. We also briefly describe the quality control procedure we used.

Crowd Workers. We collected annotations from three Amazon’s MTurk workers. We used the majority label among them for our analysis. They were asked to classify sentence as invalid only *if the text is not natural and text is incomprehensible*.

Every worker was paid \$0.05 for each Human Intelligence Task (HIT), and every HIT consisted of 5 examples.

To ensure that annotation quality, in every HIT, we insert an example which was already annotated (by graduate students). We reject all the annotations from a HIT if worker provides a wrong answer to that example. In all, we only rejected 3% of the total hits. Because we rejected some annotations, the increments are not in 0.33% multiples.

A total of 1800 sentences were annotated in the first human experiment. This had equal number of valid and invalid sentences. This was to ensure that valid and invalid sentences have equal representation and the estimate we get from crowd sourced experiment is un-biased.

Human response to classification on invalid examples

For ascertaining the validity of generated examples from our transformations, first we performed a crowd-sourcing experiment where humans were asked to rate sentences as valid or invalid. We describe an additional set of experiments where humans are asked to perform classification on invalid examples, a setting identical to the one our models work with. We perform this experiment to understand the contrast between humans and models w.r.t their response when presented with invalid examples.

Specifically, for every transformation we study, we sample 50 invalid examples generated (450 total) and ask turkers to perform classification task on these examples. Additionally, we have 50 clean examples to measure their performance on valid examples. For every example, we collect annotations from three annotators and use the majority label for our analyses. Since, invalid examples are expected to be incomprehensible, we provide an additional instruction to the turkers to perform the task by *”reconstructing the inputs to the best of their abilities”*. Rest of the instructions are taken from the original task instructions.

We perform this analyses on two tasks: NLI (MNLI), and Sentiment Analysis (SST2). We use NLI as a task representative of tasks that require a pair of inputs, and use SST2 as representative of tasks that require only one input. In the main paper, we have only presented these results for MNLI due to space constraints. Here we describe results on both tasks.

Results of crowd-sourced experiments are presented in table 14. We note that accuracy for humans on invalid examples is much lower than their accuracy on valid (un-transformed) examples. This shows that our proposed transformations in fact, do destroy the label determining information in the inputs. On invalid examples, for MNLI the humans perform with 35% accuracy, while on SST2 the accuracy is around 49%. Both of these numbers are near-random: 33% for MNLI, 50% for SST2.

Additional Experimental Results

RoBERTa is highly confident: Detailed results

In results section of the main paper, we show that fine-tuned RoBERTa model achieves a high agreement score with high confidence on invalid examples generated using our proposed transformations. We provided average confidence values for the model on each dataset. Here we provide detailed results for all datasets on all transformations. These results are shown in table 15.

Because MRPC has less than 4000 training examples, we could not generate PBSMT based invalid samples.

Comparing models from the BERT-family

We now provide detailed results for four models from the BERT-family on MNLI. We aim to compare these models for their response on invalid inputs generated using our proposed transformation functions. Agreement scores are provided in table 16. Confidence scores are provided in table 17

Mitigation results

To augment the training data with invalid examples, we sample 50% examples from training set and for each of them, generate all nine types of invalid examples as described in section on Destructive Transformations. These examples are then augmented to the original training set to get our final augmented training set. We initialize the model with the fine-tuned weights from the baseline model and then train it further for three more epochs with the loss function of eq. (1).

Dataset	Transform	Input	Prediction
Natural Language Inference - MNL	Original	P: I feel that you probably underestimate the danger, and therefore warn you again that I can promise you no protection. H: I warn you again, that I can promise you no protection, as I feel that you probably underestimate the danger.	Ent (99.65%)
		H': , , again as can danger feel i i i no probably promise protection that that the underestimate warn you you you .	Ent (98.38%)
	Original Reversed	P: However, the specific approaches to executing those principles tended to differ among the various sectors. H: Specific approaches to each principle is the same in each sector. H': sector each in same the is principle each to approaches specific .	Con (99.96%) Con (99.94%)
		P: Nash showed up for an MIT New Year's Eve party clad only in a diaper. H: Nash had too many nasty pictures on Instagram. H': too had had too many many inst inst..	Neu (99.84%) Neu (99.84%)
	Original PBSMT-C	P: The number of steps built down into the interior means that it is unsuitable for the infirm or those with heart problems. H: The interior is well suited for those with cardiac issues. H': there was no way to the is unsuitable for the park and infirm .	Con (99.43 %) Con (99.71 %)
		Original Repeat	Q1: How long can I keep a BigMac in my fridge before eating it? Q2: How long do refried beans last in the fridge after you open the can? Q2': how longDonald cluesried beans International in the fridge111 dessert open the Mass squ
Original CopyOne	Q1: How do I post a question in quora? Q2: How can I ask my question on Quora? Q2': ora		Yes (99.97 %) Yes (99.98 %)
	Original Shuffle		Q1: What is the best programming language to learn first and why? Q2: What programming language is best (easiest) to learn first? Q2': (first programming best language learn is easiest what ?) to
Sentiment Analysis - SST-2	Original	Directed in a paint-by-numbers manner .	-ve (99.95 %)
	Repeat	directed in a a-by-n in inby	-ve (98.45 %)
	Original	old-form moviemaking at its best .	+ve (99.98 %)
	Replace	CBS scare 1966 moviem GS at;200b; NZ add	+ve (99.98 %)

Table 13: More Examples for invalid examples generated using destructive transformations that render the inputs meaningless. A fine-tuned RoBERTa model makes the same predictions with very high probability (in parenthesis). PBSMT-C uses a generation model to learn statistical correlations for class Contradict. For NLI, the model chooses between *entail* (Ent), *contradict* (Con), and *neutral* (Neu). For sentiment analysis, possible labels are -ve or +ve. For paraphrase detection, model answers if the two texts are paraphrases of each other (Yes or No). Notice that in some of the examples, tokens are just the sub-words from the RoBERTa tokenizer.

Dataset	Valid	Invalid
MNLI	77%	35%
SST-2	83%	49%

Table 14: Human Accuracy for Valid and Invalid examples on two datasets. Invalid examples are generated from our proposed transformations, while valid examples come from the original validation sets.

After applying entropic regularization, we observe a significant drop in model's agreement scores on invalid examples (refer fig. 2). Indeed, the agreement scores on invalid examples have decreased to an average of 35% after regularization. Moreover, owing to the regularization, we notice a significant increase in uncertainty on nonsensical examples. However, we also notice that in some cases accuracies on original validation set decrease by more than 1%. For instance, for MNLI, accuracy decreased from 87.31 to 86.11.

Transform	MNLI	SNLI	QQP	MRPC	SST2
Random	33.33	33.33	50.00	50.00	50.00
Sort	94.86	93.17	98.76	98.19	99.33
Reverse	94.15	89.15	98.85	97.65	99.25
Shuffle	94.94	92.53	98.77	97.67	99.17
CopySort	97.96	92.48	98.16	99.55	-
Avg. Lex.	94.48	91.83	98.63	98.27	99.25
Replace	94.81	93.35	99.56	97.62	99.21
Repeat	90.46	91.29	99.31	94.07	99.26
Drop	94.01	93.96	99.20	97.31	99.32
CopyOne	94.94	92.69	99.86	99.76	-
Avg. Grad.	93.55	92.82	99.49	97.19	99.26
PBSMT	95.05	93.52	96.56	-	98.39

Table 15: Average Confidence over predictions from transformed validation set. The closer the numbers are to random, better the model. Cells marked ‘-’ indicate that the transformation is not defined for that dataset.

Clearly there is some trade-off between accuracy on clean examples and reliable response on invalid examples.

Training Details

We will release the code to reproduce all over experiments after acceptance.

Mitigation Strategies

We used three mitigation strategies in our paper: Thresholding original model probabilities (B + Th), Thresholding Entropic Regularized Model’s probabilities (Ent + Th), and training with an extra invalid class (IC). Both (Ent + Th) and IC need invalid examples at training time. For these two methods, we augmented the training data with invalid examples generated by applying all nine types of destructive transformations to randomly selected set of 50% examples from training set. These examples are then combined with the original training set to get our final training set. We initialize the model with the fine-tuned weights from the baseline model and then train it further for two more epochs.

Since all our experiments are on validation datasets, we cannot use those for searching hyperparameters. To find best hyperparameters, we sample 10 % examples from training set and use this set for validation. After finding best hyperparameters, the models are trained on full training sets.

Finding best thresholds We used two metrics for comparing the three mitigation strategies – (1) Accuracy on original validation set (Acc.), (2) Percentage of invalid examples identified correctly (%Invalid). For second metric, we create a set consisting of equal number of invalid examples generated by each transformation. As noted in the section on mitigation strategies in our main paper, there is a trade-off between Acc. and %Invalid. This is because these models might incorrectly classify examples from original (clean)

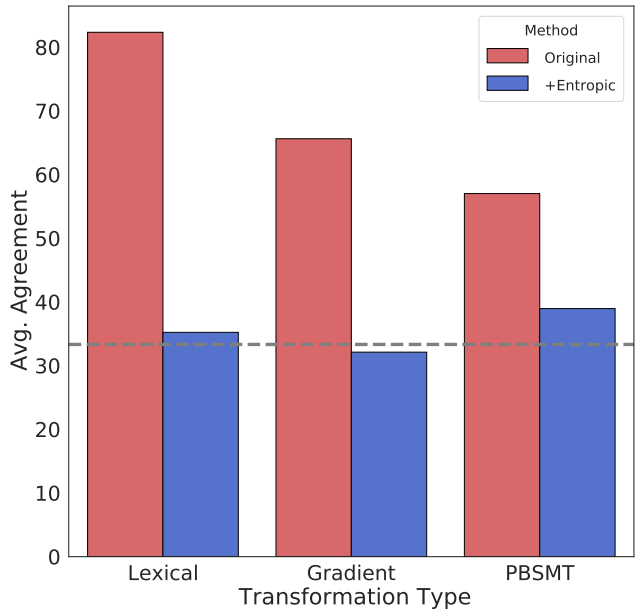


Figure 2: Results after applying entropic regularization. We found that not only did the agreement scores come down to an average of 35%, average confidence also decreased significantly. Gray line in the figure indicates agreement scores for random predictions.

validation set as invalid. For thresholding both baseline and entropic model, we have to find a balance between the two metrics, which leads to decrease in accuracy on original validation set.

We follow a simple procedure to find these thresholds using the validation set sampled from training set. We perform a grid search in range $[\frac{1}{N}, 1.0]$ with step size of 0.001, where N is the number of classes for the task. After this, we select all the thresholds that provide accuracy on (clean) validation set within a certain tolerance value (3%) of the original accuracy. Out of all these thresholds, the threshold that provides best % Invalid is selected. For example, if accuracy of the Entropic model on MNLI is 87.31, we calculate two metrics for all thresholds between $[\frac{1}{3}, 1.0]$. We then select all thresholds for which clean accuracy is within 3% of 87.31 (-84.31). From these values, the threshold that provides best detection of invalid examples is selected for analysis.

Regularization Parameter for Entropic Regularization

Entropic regularization introduces a new parameter λ that provides a trade-off between regularization and original cross entropy loss (refer section on Entropic Regularization in main paper). We found that large values of λ tend to decrease accuracy on original validation set. We fix $\lambda = 0.1$ for all our experiments. We did try other values in $\{0.01, 0.1, 0.3, 0.5, 1.0, 5.0\}$.

Training a PBSMT Model As described in main paper, we train separate sequence generation models for each la-

Transformation	BERT-base	BERT-large	RoBERTa-base	RoBERTa-large
Sort	76.97	78.54	79.19	81.20
Reverse	74.32	74.19	76.91	75.54
Shuffle	76.94	78.22	79.35	81.93
CopySort	81.49	88.86	90.48	90.48
Avg. Lexical	80.48	79.95	82.35	82.29
Replace	56.50	59.82	63.01	64.19
Repeat	45.23	49.22	49.71	49.98
Drop	61.80	64.10	69.43	71.45
CopyOne	74.58	71.22	80.41	89.75
Avg. Gradient	59.37	61.09	65.63	68.84
PBSMT	56.15	57.37	57.02	57.16
Random	33.33	33.33	33.33	33.33

Table 16: Results comparing average agreement among model sizes. Closer the numbers are to random, better is the model. ‘-’ means the transformation is not defined for that dataset.

Transformation	BERT-base	BERT-large	RoBERTa-base	RoBERTa-large
Sort	94.26	94.88	94.86	94.87
Reverse	93.45	93.8	94.15	94.39
Shuffle	94.31	94.97	94.94	94.03
CopySort	94.31	94.36	97.96	96.56
Avg. Lexical	94.08	94.50	95.48	94.96
Replace	89.66	91.34	94.81	94.68
Repeat	92.65	78.97	90.46	92.99
Drop	79.19	93.39	94.01	94.71
CopyOne	98.49	98.63	94.94	91.45
Avg. Gradient	89.99	90.58	93.55	93.46
PBSMT	94.59	95.15	95.05	95.11

Table 17: Results comparing confidence scores among model sizes. Closer the numbers are to random, better is the model. ‘-’ means the transformation is not defined for that dataset.

bel. We found the examples to be non-grammatical as well as completely out of context. With default PBSMT parameters, the model mostly duplicated inputs at its outputs, while adjusting distortion coefficient and language model coefficient achieved our desirable results.

Issues with training on Gradient based Transformations

When mitigating with invalid examples generated using gradient based methods, we noticed that sometimes if the number of gradient based invalid examples were large, the model would stop learning and start predicting the same class for every example. To overcome this behavior, simply using less number of gradient based examples worked. For mitigation, we tried using 50%, or 40%, or 30% of training examples to add to the augmented training set.

Transferability of Mitigation Strategies

We study if mitigation against one transformation helps against others using the strategy of training with the *in-*

valid label. We train models with invalid examples of one transformation and evaluate on invalid examples of all other transformations. The results in fig. 3 suggest that mitigation against transformations within a class type are transferable. For instance, training with any of the *Sort*, *Shuffle*, *Reverse* mitigates other transformations from the same group. However, to mitigate *Drop* and *PBSMT*, we need to train on invalid examples of these types. Table 18 shows all of the results.

Comparison of different models

On page 6 (table 10), we presented results of using models with different pre-training tasks and inductive biases. Detailed results are presented here in table 19. These results show that irrespective of the pre-training tasks and inductive biases, neural models perform very similarly. Particularly surprising result is that for models with recurrent inductive bias (XLNet, ESIM-*) which are trained with explicit word order based inductive bias: agreement scores are high even for lexical transformations which destroy exam-

	Sort	Reverse	Shuffle	CopySort	Replace	Repeat	Drop	CopyOne	PBSMT
Sort	99.66	99.82	99.43	99.39	41.45	77.19	29.33	89.2	49.25
Reverse	99.03	99.86	98.64	97.8	32.66	74.62	17.86	87.12	45.79
Shuffle	99.64	99.85	99.62	98.81	49.45	81.02	32.87	83.53	52.25
CopySort	75.27	77.66	67.81	99.97	1.01	10.28	0.68	29.92	0.05
Replace	50.21	87.13	49.22	3.69	99.11	45.29	19.26	36.25	65.22
Repeat	80.97	96.93	78.88	98.61	22.05	99.93	33	99.47	58.25
Drop	90.61	95.22	94.42	97.74	82.9	98.73	94.81	97.62	42.02
CopyOne	55.26	80.15	50.14	11.45	45.47	55.77	15.93	99.83	12.68
PBSMT	25.45	41.25	21.26	36.5	76.56	65.2	58.56	79.25	98.64

Table 18: Results on transfer effects of each transformation function. The labels on the left indicate the function that the model was trained on, labels on the top show the functions it was tested on. All numbers denote percentage of invalid examples detected. We note that the functions have a strong in-class transfer effect.

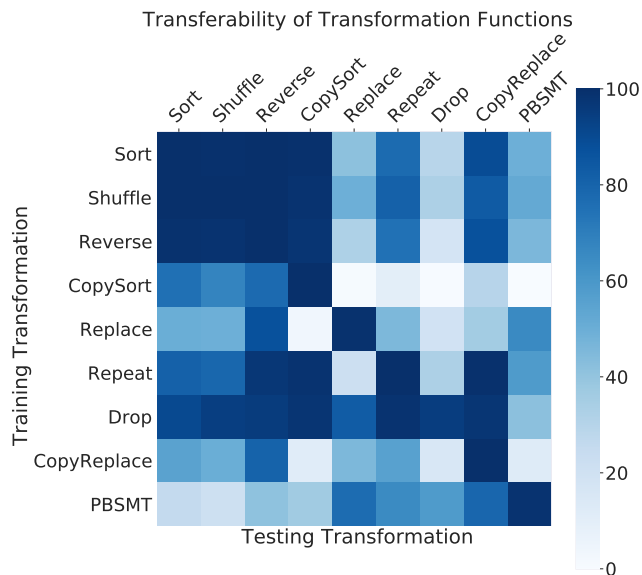


Figure 3: Transferability across transformation functions. We notice a great deal of transferability across lexical overlap transformations.

ples by changing the word order of the inputs in irregular ways.

Reproducibility

In this section, we provide details on our hyperparameter settings along with some comments on reproducibility.

Models Used As described in main paper, we performed our experiments by fine-tuning RoBERTa-*base* model (Liu et al. 2019). For our implementation, we used Huggingface’s *transformers* repository (Wolf et al. 2020). We also used RoBERTa-*large*, BERT-*base*, BERT-*large* from this repository to compare different models from BERT family.

Hyperparameter Tuning for Calibration Methods We used three methods for improving calibration of our model

– *Label Smoothing (LS)*, *Focal Loss (FL)*, and *Temperature Scaling (TS)*. Each of these methods introduce a single new hyperparameter. For Label Smoothing, this parameter (λ_{LS}) controls the smoothing of hard labels. Focal loss adds a focusing parameter γ_{FL} (Lin et al. 2017). For temperature scaling, we have an additional parameter T_{TS} that controls the amount of scaling to be applied to softmax logits (Guo et al. 2017).

Since all our experiments are on validation datasets, we cannot use those for searching hyperparameters. To find best hyperparameters, we sample 10 % examples from original training set and use this set for validation. After finding best hyperparameters, the models are trained on original, full training sets.

To find λ_{LS} for MNLI, we used grid search over $[0.1, 0.3]$ with 0.05 as step size, our best hyperparameter was $\lambda_{LS} = 0.1$. As LS needs to train a model from scratch for each value of λ_{LS} , to avoid excessive computation we used this same value across all datasets. For γ_{FL} , as suggested by Lin et al., we use grid search over $[0, 5.0]$ with step size 0.5 and find that $\gamma_{FL} = 2.0$ works best. Again to avoid excessive computation, we fix this value for all datasets. For T_{TS} , we maximize the log-likelihood of validation set (described above) using LBFSG optimizer as recommended by Guo et al.. Value for T_{TS} was 3.21. We found that while LS, FL improve calibration for all datasets, temperature scaling made calibration worse for SST-2 ($0.05 \rightarrow 0.24$).

Computing Infrastructure Used Most of our experiments (except on PBSMT) required access to GPU accelerators. We primarily ran our experiments on three machines: Nvidia Tesla V100 (16 GB VRAM), Tesla P100 (16 GB VRAM), and Nvidia TITAN X (Pascal) (12 GB VRAM).

We used the Moses SMT system for our PBSMT models (Koehn et al. 2007; Koehn, Och, and Marcu 2003). These required intensive CPU computations and were executed on a machine with an Intel(R) Xeon(R) 2.40GHz CPU with 28 cores. To speed up computation, we used 48 parallel threads.

Class	Pretraining	Model	Lexical	Gradient	PBSMT	Average
Auto-Encoding	W-MLM + NSP	BERT-B (Devlin et al. 2019)	80.48	57.01	55.02	67.13
		BERT-L	80.95	59.92	58.01	69.05
	W-MLM	RoBERTa-B (Liu et al. 2019)	82.35	65.63	57.02	71.72
		RoBERTa-L	82.90	68.17	57.42	73.52
	W-MLM + SOP	ALBERT-B (Lan et al. 2019)	78.55	57.05	56.73	66.57
Auto-Encoding	S-MLM + NSP	SpanBERT-B (Joshi et al. 2020)	81.03	57.11	55.4	67.55
	S/W-MLM + SS	BART-B (Lewis et al. 2020)	78.68	64.81	56.49	70.04
	RTD	ELECTRA-B (Clark et al. 2019)	79.36	60.95	58.11	68.81
Auto-Regressive	P-LM	XLNet-B (Yang et al. 2019)	81.33	62.13	57.42	70.14
	LM	ESIM- ELMo (Peters et al. 2018)	76.49	78.15	63.17	75.78
	-	ESIM- Glove (Chen et al. 2017)	74.19	75.31	63.41	73.49

Table 19: Agreement score of different models on MNLI. **B** refers to the base variant, **L** refers to the large one. Last four columns present agreement scores for 3 kinds of transformations discussed in this work, with last column being the average. W-MLM refers to word based Masked Language Modeling (MLM), S-MLM refers to span based MLM, NSP refers to Next Sentence Prediction, SOP refers to Sentence Order Prediction, SS refers to sentence shuffling objective, RTD refers to Replaced Token Detection. P-LM refers to permutation language modeling, LM refers to traditional auto-regressive language modeling pre-training objective used in ELMo.

Dataset and Model	Average RunTime
MNLI + Invalid	14 hr
SNLI + Invalid	18 hr
MRPC + Invalid	13 mins
QQP + Invalid	11 hr
SST-2 + Invalid	1.5 hr

Table 20: Average Training time with an additional invalid class

Average Run times We fine-tune all our models for three epochs with a batch size of 8 (discussed in next section). Average training times are presented in table 20.

Fine-tuning Details We used the RoBERTa-*base* model for most of our experiments. This model has 12 layers each with hidden size of 768 and number of attention heads equal to 12. Total number of parameters in this model is 125 million.

We also experimented with BERT-*base*, BERT-*large*, and RoBERTa-*large* for comparison, which have 110 million, 340 million, and 355 million parameters respectively. For both BERT models, we used the *uncased* variant.