

POLYJUICE: Automated, General-purpose Counterfactual Generation

Tongshuang Wu¹

¹University of Washington
wtshuang@cs.uw.edu

Marco Tulio Ribeiro²

²Microsoft Research
marcotcr@gmail.com

Jeffrey Heer¹

³Allen Institute for Artificial Intelligence
{jheer,weld}@cs.uw.edu

Daniel S. Weld^{1,3}

Abstract

Counterfactual examples have been shown to be useful for many applications, including calibrating, evaluating, and explaining model decision boundaries. However, previous methods for generating such counterfactual examples have been tightly tailored to a specific application, used a limited range of linguistic patterns, or are hard to scale. We propose to disentangle counterfactual generation from its use cases, *i.e.*, gather general-purpose counterfactuals first, and then select them for specific applications. We frame the automated counterfactual generation as text generation, and fine-tune GPT-2 into a generator, POLYJUICE, which produces fluent and diverse counterfactuals. Our method also allows control over where perturbations happen and what they do. We show POLYJUICE supports multiple use cases: by generating diverse counterfactuals for humans to label, POLYJUICE helps produce high-quality datasets for model training and evaluation, requiring 40% less human effort. When used to generate explanations, POLYJUICE helps augment feature attribution methods to reveal models’ erroneous behaviors.

1 Introduction

Counterfactual reasoning — mentally simulating what *would have happened* if conditions were different — is a common tool for making causality assessments (Kahneman and Tversky, 1981), which in turn are crucial for model training, evaluation, and explanation (Miller, 2019). For example, in Figure 1, “It is great for kids” is perturbed into multiple variations, and each variation brings a unique insight by simulating what would have happened if the sentence was different.

Applications of counterfactual reasoning to Natural Language Processing (NLP) generally specify the relationship $x \rightarrow \hat{x}$, and then create \hat{x} according to the relationship. As a result, generation methods

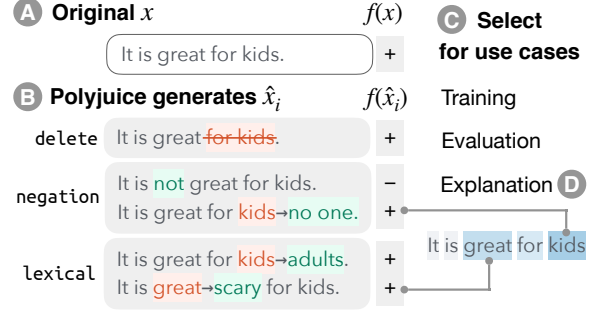


Figure 1: Overview: given an original sentiment analysis instance x (A), POLYJUICE generates (B) a large number of \hat{x} , which are then (C) independently selected for downstream use cases. For example, in (D), we select \hat{x} as counterfactual explanations, based on whether they complement the feature attributions: though both “great” and “kids” are deemed important, the selected \hat{x} show that perturbing them may not affect the prediction $f(x) = f(\hat{x}) = \text{positive}$, revealing model errors.

are heavily constrained by the application, each having their own limitations. For example, a minimal edit of a sentence x that results in a different label is useful for model training and evaluation. Such \hat{x} are usually gathered with human efforts, *i.e.*, human annotators manually create counterfactuals (Gardner et al., 2020) or perturbation functions that generate counterfactuals (Wu et al., 2019a). These are costly to generate — taking 4-5 minutes per counterfactual (Kaushik et al., 2020) — and may miss important patterns due to their reliance on human intuition (*e.g.*, humans may cover **great** \rightarrow **not great**, but can easily miss **kids** \rightarrow **no one** in Figure 1B). Adversarial examples are a different form of counterfactual reasoning: x and \hat{x} have different model predictions *despite* being minimally edited and semantically equivalent — the latter limiting most perturbations to be automated word replacements or other forms of paraphrasing (Iyyer et al., 2018; Ribeiro et al., 2018b).

However, we observe that generation methods do not have to be isolated, as various applications

share similar requirements on $x \rightarrow \hat{x}$ (e.g., minimal edits). In fact, the exhaustive search of automated methods can cover the training and evaluation more comprehensively, and non-paraphrasing changes like *add negation* are valuable adversarials for tasks like named entity recognition.

In this work, we formalize the task of *automatic counterfactual generation*, which disentangles the generation from the application (§3). First, in generation, given an input x , we produce a set of counterfactuals $\hat{X} = \{\hat{x}_1, \hat{x}_2, \dots\}$ with reasonable but *application-agnostic* relationships $x \rightarrow \hat{x}_i$ (Figure 1B). We require the counterfactuals to be *fluent*, *diverse*, and *close to x*. Afterwards, we use *application-specific* selection methods to find subsets of \hat{x} that are most effective for the applications of interest (Figure 1C). We frame the generation step as text generation, and finetune GPT-2 (Radford et al., 2019) into a generator called POLYJUICE using datasets of (x, \hat{x}) pairs. We also allow for targeted counterfactuals, by specifying where the perturbation occurs in the sentence (Donahue et al., 2020) and using control codes such as *negation* or *delete* (Figure 1B).

We propose simple yet effective selection strategies, and show that a single POLYJUICE model can support multiple downstream applications: First, it can *facilitate counterfactual training and evaluation* (§4). By asking humans to *label* a set of *unique* counterfactuals, we produce high-quality contrast sets (Gardner et al., 2020) at a fraction of the annotation effort compared to creating them from scratch (Kaushik et al., 2020). We similarly produce training data that improves model generalization in three classification tasks.

Second, POLYJUICE helps produce *black-box counterfactual explanations* (§5). In a user study, expert users only did slightly better than random (accuracy: $55 \pm 6\%$) at predicting what a model would do on POLYJUICE counterfactuals that are selected to *highlight mistaken model predictions*, as in Figure 1D. In fact, the explanations found more bugs within spots where users considered inspected after viewing feature weights and manually performing counterfactual analysis. This indicates that such counterfactuals can complement existing explanation and analysis methods, echoing evidence from social science (Miller, 2019). In summary, we:

1. Formalize the general-purpose counterfactual generation task. By *separating the generation from the use cases*, we generate fluent and di-

verse counterfactuals that bypass application-specific constraints.

2. Finetune a generator called POLYJUICE, by collecting paired sentences and enhancing controls with infilling structures and control codes—the control is *the backbone of* various downstream applications. The model is at <https://huggingface.co/uw-hai/polyjuice>.
3. Apply POLYJUICE to *model training, evaluation, and explanation*, using various selection methods (which we will open-source). POLYJUICE helps collect high-quality training and evaluation data with 40% less annotation effort, and find model bugs that on top of feature attribution explanations and counterfactual analysis.

2 Related Work

The applications of counterfactuals. Counterfactuals in NLP are most broadly used for model training, evaluation, and explanation. In training, they usually augment the training data to improve model robustness (Garg et al., 2019; Wu et al., 2019b; Wei and Zou, 2019; Kumar et al., 2020) and generalizability. Kaushik et al. (2020) and Teney et al. (2020) showed that counterfactual augmentations help mitigate the weights of spurious features. Evaluation focuses on similar aspects as training, mostly through adversarial attacks (Song et al., 2020), contrast sets (Kaushik et al., 2020) and challenge sets (Geiger et al., 2019; Liu et al., 2019a). We show the POLYJUICE generated counterfactuals are also useful in such cases.

Counterfactuals also naturally support model explanations, as “explanations are sought in response to particular counterfactual cases or foils” (Miller, 2019). Popular feature importance attribution methods (Lundberg and Lee, 2017; Ribeiro et al., 2016) all retrieve token importance through masking, which can be viewed as a form of (incomplete) counterfactual. Some work also explores directly presenting simple counterfactual examples (Ross et al., 2020; Vig et al., 2020; Kang et al., 2020), yet they lose the overview in feature attributions. We explore the selection of such examples, and combine them with existing feature attribution methods.

Counterfactual generation. As in §B, existing generation methods are usually for particular applications. Most automated ones are for adversarial examples. They enumerate through candidate word replacements (Alzantot et al., 2018a; Garg and Ramakrishnan, 2020) or paraphrases (Iyyer

et al., 2018; Malandrakis et al., 2019) that can maintain semantic meanings while exposing model errors. The exhaustive search is hard for people to scale to (Ribeiro et al., 2018c). On the other hand, those evaluating and improving model decision boundaries rely on manual efforts (Ribeiro et al., 2020), which cover more diverse patterns, but at a higher cost. Moreover, though human annotators also share strategies like negation, subject-object swapping (Kaushik et al., 2020; Gardner et al., 2020), the strategies have yet to be captured for targeted use. Template-based methods (McCoy et al., 2019; Nie et al., 2019) are more systematic and cost-effective, but usually are only applicable to a small subset of data (Li et al., 2020a). We unify the generation for different use cases, which bridges the scalability and diversity. We only distinguish the generated ones through posterior constraints (Morris et al., 2020; Alzantot et al., 2018b).

Some concurrent work is similar to ours, using language models to generate counterfactual sentences. GYC (Madaan et al., 2021) adjusts controlled text generation methods and modifies x in the latent space, whereas MiCE (Ross et al., 2020) finetunes T5. Both target at generating counterfactuals that *flip the class label*, which can be useful, but neither approach emphasizes generality, with GYC focusing on evaluation and MiCE aimed at explanation. We instead finetune POLYJUICE on general-purpose datasets, design POLYJUICE’s controls to be task-agnostic, and use them for broadening the coverage on a large variety of perturbation patterns. As we demonstrate later in §4 and §5, such controls additionally support application that require specifications on where and how to perturb.

3 General-purpose Counterfactuals

Definition. Given an instance $x \in \mathbf{X}$, a generator g should produce a set of counterfactuals $\hat{\mathbf{X}} = \{\hat{x}_1, \hat{x}_2, \dots\}$ with varying relationships $x \rightarrow \hat{x}_i$ (referred as $r(\hat{x}_i)$ for simplicity). For example, *great* \rightarrow *not great*, *kids* \rightarrow *no one* in Figure 1B are both instances of the *negation* relationship. The changes also involve a *label flipping* relationship, with both counterfactuals becoming negative sentences after the perturbation.

As illustrated in §3.1, each \hat{x} should (1) be *close* to x and (2) maintain *fluency*. Together, they should form (3) a *diverse* $\hat{\mathbf{X}}$. The three desiderata is applicable *regardless of* the use cases, while the relationship $r(\hat{x})$ would *vary with* use cases. They

inspire our modeling in §3.2, including the training prompt design, datasets, and filtering strategies.

3.1 Desiderata

Individual counterfactuals. A counterfactual \hat{x} should be (1) *close* to x , preferably only involving the minimal changes necessary for establishing a certain effect while leaving the rest of the instance intact (Pearl, 2018), allowing users to make causality assessments from $r(\hat{x})$. It has long been observed that humans strongly favor counterfactuals that are closer to the original instance (Kahneman and Tversky, 1981). Following common practice in NLP research, we estimate *closeness* using the semantic and syntactic distance between x and \hat{x} (Morris et al., 2020; Madaan et al., 2021).

However, \hat{x} also needs to be *fluent*, *i.e.*, grammatically correct (Morris et al., 2020) and semantically meaningful (*e.g.*, “It’s scary for water” is grammatically correct but not meaningful.) This is estimating the likelihood of a sentence occurring in reality; as Kahneman and Tversky (1981) stated, the counterfactual scenario should be one that could have easily happened, without rare assumptions or coincidences (Kahneman and Tversky, 1981).

Sets of counterfactuals. There can be infinite numbers of “what-ifs” (Pearl, 2018; Kahneman and Tversky, 1981), especially as we allow \hat{x} to deviate from x . Since our goal is to produce a general-purpose set of counterfactuals, we expect the set to be (3) *diverse* in terms of relationships between x and \hat{x} . However, when the targeted relationship is unclear, an approximation could be the similarity of counterfactuals to *each other* (using *e.g.*, self-BLEU (Hu et al., 2017)).

3.2 Modeling through Text Generation

We frame counterfactual generation as a text generation task using language models (LMs). As we demonstrate in an intrinsic evaluation (§A.3), large pre-trained LMs like GPT-2 (Radford et al., 2019) are capable of generating *fluent* text, and are more flexible than word substitutions (Garg et al., 2019) and templates (Ribeiro et al., 2018c; Wu et al., 2019a), allowing for increased *diversity*. We detail how we finetune GPT-2 to generate \hat{x} that are *close* to an original instance x (rather than arbitrary text). Importantly, we present a novel method for controlling the $r(\hat{x})$ relationship, which is essential for supporting a wide range of applications, *e.g.*, counterfactual analysis where we group \hat{x} based on how they are perturbed both semantically and

Control code	Definitions and POLYJUICE-generated Examples	Training datasets
negation	A dog is not embraced by the woman.	(Kaushik et al., 2020)
quantifier	A dog is \rightarrow Three dogs are embraced by the woman.	(Gardner et al., 2020)
shuffle	To move (or swap) key phrases or entities around the sentence. A dog \rightarrow woman is embraced by the woman \rightarrow dog .	(Zhang et al., 2019b)
lexical	To change just one word or noun chunks without breaking the POS tags. A dog is embraced \rightarrow attacked by the woman.	(Sakaguchi et al., 2020)
resemantic	To replace short phrases or clauses without affecting the parsing tree. A dog is embraced by the woman \rightarrow wrapped in a blanket .	(Wieting and Gimpel, 2018)
insert	To add constraints without affecting the parsing structure of other parts. A dog is embraced by the little woman.	(McCoy et al., 2019)
delete	To remove constraints without affecting the parsing structure of other parts. A dog is embraced by the woman .	(McCoy et al., 2019)
restructure	To alter the dependency tree structure, e.g., changing from passive to positive. A dog is embraced by \rightarrow hugging the woman.	(Wieting and Gimpel, 2018)

Table 1: A list of control codes used for semantically driving the counterfactual generation in POLYJUICE, their examples, and the representative training datasets for the corresponding patterns. More examples are in §E.

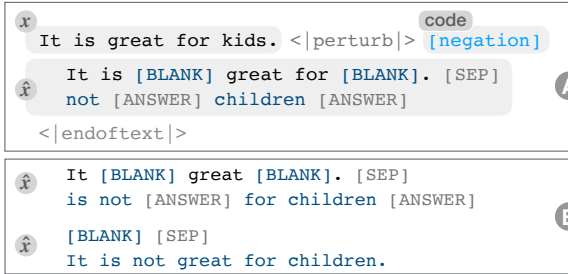


Figure 2: Forming the training text. (A) Given a pair (x, \hat{x}) , we extract the control codes and convert \hat{x} to the infilling structure by blanking the changed tokens, and then concatenate x , the code, and the blanked \hat{x} with special tokens. (B) We get multiple training texts per pair, by blanking \hat{x} subtrees that contain the change, or the entire sentence. At generation time, POLYJUICE accepts prompts that just include x , or with the code and the BLANK placement.

syntactically, counterfactual explanations on the salient features, etc.

Controlling counterfactual generation. To achieve control over $r(\hat{x})$, we need to be able to control both where in x the perturbation happens, and the type of change. To address the former, we extend the Infilling by Language Modeling (ILM) framework (Donahue et al., 2020), such that x is always part of the prompt, and \hat{x} contains [BLANK] tokens where perturbations are to be applied (Figure 2). ILM allows for perturbations of any length (additions and deletions) beyond single word substitutions, e.g., Figure 2B.

To control the type of perturbations, we additionally condition the generation on special control codes (Raffel et al., 2020a; Dathathri et al., 2020),

e.g., negation in Figure 2A. Inspired by prior work categorizing manually created counterfactuals (Kaushik et al., 2020; Gardner et al., 2020), we design 8 codes (Table 1) that distinguish lexical, syntactic, and semantic perturbations. We verify through an ablation study in §A.3.1 that training POLYJUICE with control codes greatly improves the success rate when generating counterfactuals that have these properties (by $29\% \pm 18\%$).

As Figure 2 indicates, we can control the counterfactual generation conditioning exclusively on x (in which case POLYJUICE selects the control code and the location of [BLANK] tokens), conditioning on x and the control code, or specifying both the code and where the perturbations happen.

Finetuning dataset. We rely on six existing datasets of paired sentences, each of which contains counterfactuals with a subset of $r(\hat{x})$ relationships of interest (see Table 1). To increase diversity, we additionally find naturally occurring pairs in non-paired datasets like SQAuD (Rajpurkar et al., 2016), using heuristics on the editing distance. A combination of these datasets yields a finetuned model that can produce *diverse* counterfactuals. For each paired (x, \hat{x}) , we create various training prompts as in Figure 2. We vary the location of [BLANK]s according to the example’s parse tree, and compute the primary code based on linguistic features such as part-of-speech tags or dependency trees. We use the pairs as negative samples if the change is too extensive. We form 657,144 training prompts from 191,415 sentence pairs. More dataset details and code statistics are in §A.

Filtering for fluency. Certain combinations of control codes and [BLANK] tokens cause POLYJUICE to generate ungrammatical or nonsensical counterfactuals. Similar to the language model constraints used by Morris et al. (2020), we score both x and \hat{x} using unfinetuned GPT-2, and filter out \hat{x} if its log-probability (either on the full sentence or on the perturbed chunks) decreases more than 10 points relative to x . We verify the impact of the filtering through human evaluation, as part of the labeling task in §4. As shown in §4.2, it removes a large portion of unrealistic counterfactuals.

Using counterfactuals. POLYJUICE is general-purpose and task-agnostic, and generates a variety of counterfactuals given one x . The most natural way to use POLYJUICE is to generate $\hat{\mathbf{X}}$ without constraints, and then select the counterfactuals in $\hat{\mathbf{X}}$ according to what is most appropriate for the task at hand (we provide examples of selection strategies for labeling and explanations in §4 and §5). However, users can also make use of the various control mechanisms, *e.g.*, generating counterfactuals by perturbing only specific subtrees of interest (an example is given in §4.1).

4 POLYJUICE for Training & Evaluation

Manually created counterfactuals are useful for evaluating and improving models (Gardner et al., 2020; Kaushik et al., 2020), but *creating* variations is much more difficult than *validating* them (Ribeiro et al., 2018c).

Here, we ask crowdworkers to only label counterfactuals (§4.2) that are automatically generated by POLYJUICE (§4.1). We explore whether the labeled \hat{x} can support counterfactual evaluation (§4.3) and training (§4.4), with experiments on three datasets (more in §C.1): (1) *Sentiment Analysis* with Stanford Sentiment Treebank (SST-2) (Socher et al., 2013), (2) *Natural Language Inference (NLI)* with SNLI (Bowman et al., 2015), and (3) *Duplicate Question Detection (QQP)* (Wang et al., 2018). As §4.5 summarizes, the approach can serve the purposes at a much lower cost.

4.1 Selection for Labeling

Starting with a set of \mathbf{X} , our goal is to create a set of counterfactuals \mathbf{L} for the crowdworker to label. For each $x_i \in \mathbf{X}$, we generate a large set of candidate counterfactuals $\hat{\mathbf{X}}_i$, and sample three of $\hat{x} \in \hat{\mathbf{X}}_i$ to \mathbf{L} , such that crowdworkers label the same numbers of variations for each x . To ensure we label the

most effective counterfactuals, we *generate* $\hat{\mathbf{X}}$ ’s that target at models’ known blind spots, and then *select* the unique ones to form \mathbf{L} , as detailed below.

Targeted counterfactuals. Longpre et al. (2020) found that typical data augmentations are likely to bring redundant benefits as pre-training. They suggested that the method should focus on where current models fail, which we follow in §4.4. We adjust Chen et al. (2019)’s intuitive data slicing functions to find the x with certain patterns of interest. Then, we simply blank the corresponding patterns so POLYJUICE can perform relevant changes. For example, to highlight the impact of prepositions (“His surfboard is **beneath** \rightarrow **lying on** him”), we first filter examples that have prepositions, and generate blanked prompts like “[resemantic] His surfboard is [BLANK] him.”

Prioritize unique \hat{x} . To cover more variations around local decision boundaries, we select unique counterfactuals through submodular optimization. As mentioned in §3, the selection is based on relationship $r(\hat{x})$. In this case, it involves a set of components: { the base x , tokens removed from x , added to \hat{x} , the affected parsing tree structure, and the corresponding control code. } We define a distance function on the relationships $D_L(r(\hat{x}_i), r(\hat{x}_j))$, which is a weighted combination of the distances between each individual components. Using $D_L[\cdot]$, we greedily select \hat{x} whose $r(\hat{x})$ is the least similar to those already in \mathbf{L} . For example, if “a **man** \rightarrow **woman** walks” is already in \mathbf{L} , we would penalize “the **man** \rightarrow **woman** is dancing” (same changed text) or “that dog is with the **man** \rightarrow **girl**” (similar lexical change), but prioritize “**two** people talk.”

4.2 Labeling Procedure & Efficiency

Procedure. We crowd label the counterfactuals on Amazon Mechanical Turk. For each round, the annotator is given the original x and its groundtruth as references¹, and is asked to label three counterfactuals by (1) the class label and (2) fluency (“*likely written by a native English speaker*”). We carefully remove noisy workers using hidden *gold rounds*, as well as filters on label distributions and completion time. We also remove noisy labels through majority votes. More details are in §C.2.

Fluency. Crowdworkers rated most counterfactuals to be fluent: 75% for SST-2, 70% for QQP, and 82% for SNLI. One of the authors also manu-

¹For QQP and NLI, we only perturbed *duplicate* and *entailment* examples, as others are significantly harder to flip (called *asymmetric counterfactuals* (Garg et al., 2019)).

Task	Dev.	Orig. set	Contrast set	Consistency
<i>Sentiment</i>	94.3	93.8	84.9 (-8.9)	76.1
<i>NLI</i>	86.5	91.6	72.3 (-19.3)	56.4
<i>QQP</i>	91.7	87.5	75.3 (-12.2)	61.1

Table 2: Counterfactuals as contrasts sets, revealing model insufficiency. The table contains accuracies on the development set, the original x (*Orig. set*), the contrast sets, and the consistency between them.

ally labeled 600 counterfactuals of 120 instances, and arrived at similar fluency: The rate of unfiltered counterfactuals was 61%, which increased to 78% after filtering, showing that the filtering is effective.

Efficiency. Labeling three counterfactuals of a given example is reasonably easy, as (1) annotators are better at *verifying* counterfactuals than manually *generating* them (Ribeiro et al., 2018c), and (2) annotators only need to focus on the reference example and the corresponding perturbed phrases, rather than re-parsing the full instance for each label they submit (Khashabi et al., 2020). As such, the median time for labeling one round (three \hat{x}) is 30 seconds. Even in the most extreme (and very rare) case where we only keep 100 labeled \hat{x} that change the groundtruth after 500 rounds (*Sentiment* in §4.3), the average human annotation time per \hat{x} is 2.5 minutes, which is still 40% more efficient than manual creation: Kaushik et al. (2020) reported that workers spent roughly 5 minutes to revise an IMDB review and 4 minutes a sentence (for *NLI*), even prior to additional filtering and validation.

4.3 Experiment: Evaluation

Do POLYJUICE counterfactuals serve as a useful test set for uncovering model errors? Gardner et al. (2020) created *contrast sets*, i.e., minimally edited instances for inspecting model deficiencies, which we try to replicate. Following the definition, we filter the counterfactuals to only keep those whose groundtruth label is different from x ’s, resulting in contrast sets with sizes 100–300.

We test finetuned BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019b), and DistilBERT (Sanh et al., 2019) models opensourced on Huggingface (Wolf et al., 2020), and report the best performing models on the validation set, which happen to be the RoBERTa model across tasks.² We report model accuracies on the full validation sets, the original examples for collecting counterfactuals, and the contrast sets. We also report consis-

²huggingface.co/{roberta-large-mnli, textattack/roberta-base-SST-2, ji-xin/roberta_base-QQP-two_stage}

tency, i.e., cases where the model predicts both the originals and the counterfactuals correctly. Table 2 shows that all state-of-the-art models perform significantly worse on our contrast sets, and the performance decreases for a similar amount as in (Gardner et al., 2020). Results on other unreported models are consistent. In other words, POLYJUICE counterfactuals can reveal models’ limited capabilities as the original contrast sets.

4.4 Experiment: Training

Do POLYJUICE counterfactuals help yield better model generalization? We label counterfactuals for data augmentation. Because counterfactuals that maintain the groundtruth can still improve model stability, we keep all the valid \hat{x} for one x , as long as at least one of them flips the groundtruth.

We finetune roberta-base models in HuggingFace. For each augmented model (aug), we include m counterfactuals, as well as n examples from the original dataset (the original x of each \hat{x} will always be included). Their baselines comp replace the m counterfactuals with another m original examples. Similar to Khashabi et al. (2020), this baseline helps evaluate whether labeling *counterfactuals* is more effective than *non-counterfactual* data. The reported performances are averaged across multiple data samples and random seeds (More in §C.3).

Results. As shown in Table 3–5, compared to adding the same amount of original data, *the counterfactuals improve models’ generalization* on out-of-domain datasets, challenge and contrast sets, as well as CheckList tests, while maintaining the in-domain accuracy. Importantly, in both *NLI* and *QQP*, *just adding $m/n < 10\%$ data is sufficient to boost the performance*.

However, targeted augmentation is necessary. While randomly generated counterfactuals are beneficial for *Sentiment*, similar to Huang et al. (2020), they are not effective for either *NLI* or *QQP* (numbers omitted for space.) Instead, we apply the slicing strategy in §4.1 to prioritize counterfactuals related to known error cases, e.g., prepositions in DNC (Kim et al., 2019) for *NLI*, and entity orders in CheckList for *QQP*. The improved aug in Table 4 and 5 indicates that such methods may help avoid repetitive augmentations.

The counterfactuals usually improve the model without hurting its counterpart. In *NLI*, DNC includes pairs of probing examples, one from the original MNLI and one manually and minimally

n	m	model	SST-2	Senti140	SemEval	Amzbook	Yelp	IMDB	IMDB-Cont.	IMDB-CDA
4,000	2,000	comp	92.9 \pm 0.2	88.9 \pm 0.3	84.8 \pm 0.5	85.1 \pm 0.4	90.0 \pm 0.3	90.8 \pm 0.5	92.2 \pm 0.6	86.5 \pm 0.2
4,000	2,000	aug	92.7 \pm 0.2	90.7 \pm 0.4	86.4 \pm 0.1	85.6 \pm 0.8	90.1 \pm 0.0	90.6 \pm 0.3	94.0 \pm 0.3	89.7 \pm 0.5

Table 3: *Sentiment* model performances. aug maintains the in-domain and out-of-domain accuracies on reviews (SST-2, Amzbook, Yelp, IMDb Movie Review (Ni et al., 2019; Asghar, 2016; Maas et al., 2011)), but improves on Twitter data (Senti140 and SemEval 2017 (Go et al., 2009; Rosenthal et al., 2017)), likely because their distributions are less similar to SST-2 than the reviews. The model also improves on the contrast sets (IMDb-Contrast and IMDb-CAD (Gardner et al., 2020; Kaushik et al., 2020)).

n	m	model	SNLI	MNLI-m	MNLI-mm	SNLI-CDA	break	DNC	stress	diagnostic
20,000	1,574	comp	85.7 \pm 0.4	86.1 \pm 0.2	86.6 \pm 0.2	72.8 \pm 0.3	86.4 \pm 1.5	54.5 \pm 0.6	65.1 \pm 0.6	56.0 \pm 0.8
20,000	1,574	aug	85.3 \pm 0.3	86.0 \pm 0.1	86.4 \pm 0.0	73.6 \pm 0.2	89.1 \pm 1.2	57.7 \pm 0.3	65.1 \pm 0.2	57.5 \pm 0.5

Table 4: *NLI* model performances. aug performs better than comp on DNC (Kim et al., 2019), which is the target of the augmentation. It also improves on other contrast/challenge sets (Naik et al., 2018; Glockner et al., 2018; Wang et al., 2018), while maintaining the in-domain and out-of-domain (Williams et al., 2018) accuracies.

TESTNAME	Δ fail%
Order does not matter for symmetric relations	-18.4%
Order does not matter for comparison	-26.5%
Order does matter for asymmetric relations	-14.5%
How can I become $\{more\ X \neq less\ X\}$	-30.7%
How can I become $\{more\ X = less\ antonym(X)\}$	28.0%
How can I become $\{X \neq not\ X\}$	-10.4%
How can I become $\{X \neq not\ antonym(X)\}$	-5.5%

Table 5: Sample *QQP* CheckList tests (Ribeiro et al., 2020), with Δ fail% denoting the failure rate change from comp to aug. With $n = 20,000$ and $m = 1,911$, aug reduced the failure rates on 11 tests (out of the 27 where comp failed as defined in §C.4), while only increasing it for 2. The model improves consistently on most related cases, but possibly overfits on more/less.

perturbed. The improvement on it shows that aug did not overfit to the augmented pattern. Similarly, the gain on CheckList tests are mostly consistent in *QQP* (e.g., all tests related to entity ordering were improved), except for the more/less contrast in Table 5. Future work should further strategize the ranking to more equally cover competing patterns.

4.5 Takeaways

Our experiments show that POLYJUICE supports both counterfactual evaluation and training. More importantly, it achieves these at a noticeable lower cost, not only through shorter crowdlabelling time per instance, but also through targeted perturbation: data slicing and blanking help allocate the labeling efforts to best compensate the existing data, rather than wasting them on repetitive patterns.

5 POLYJUICE for Explanations

Counterfactual explanations have been elusive in NLP, despite evidence from social science research

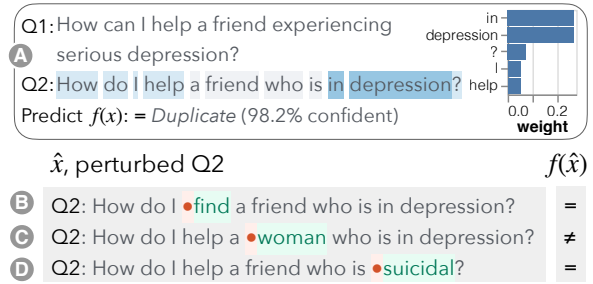


Figure 3: (A) A *QQP* instance with its prediction³ showing 98.2% confidence that the two sentences are duplicative (=), as well as the SHAP feature weights. Counterfactual explanations complement SHAP with concrete, readable examples, and alert abnormalities missed by SHAP. e.g., (C) depicts a surprising flipped prediction (≠).

which indicates that they may be more intuitive than feature attribution methods (Miller, 2019). Compared to the opaque feature weights in Figure 3A, Figure 3B more intuitively shows that the BERT *QQP* model considers “help” trivial: *help* \rightarrow *find* in Q2 does not change the prediction.

Here, we first select counterfactuals to serve as local explanations in §5.1, and evaluate their effectiveness in §5.2. We discuss the use of \hat{x} beyond local explanations in §5.3.

5.1 Selection: Abnormalities as Explanations

Because complete explanation is overwhelming, Miller (2019) concluded that people usually only expect explanations on contrastive cases (“foils”) that they consider *abnormal*. We similarly select abnormal counterfactuals as explanations, i.e., counterfactuals where the expected and the actual

³From BERT *QQP* model: <https://huggingface.co/textattack/bert-base-uncased-QQP>

changes in the prediction do not match. In other words, the relationship $r(\hat{x})$ used in the selection is the change of model prediction.

Given a model f , we define the actual prediction change as $D_f(\hat{x}, x) = |f_p(\hat{x}) - f_p(x)|$, where $f_p(x)$ denotes the prediction probability of f on x . Meanwhile, the *human expectation* on the prediction change is $\mathbf{H}[D_f(\hat{x}, x)]$. The mismatch between the reality and the expectation is then:

$$\Delta D_f(\hat{x}, x) = D_f(\hat{x}, x) - \mathbf{H}[D_f(\hat{x}, x)]$$

We select two *expectation violation* counterfactuals, *i.e.*, (1) large prediction change from trivial perturbations (small expectations): $\hat{x}_L = \arg \max_{\hat{x}} \Delta D_f(\hat{x}, x)$, or (2) unchanged predictions from large changes: $\hat{x}_U = \arg \max_{\hat{x}} -\Delta D_f(\hat{x}, x)$.

$\mathbf{H}[\cdot]$ can take various forms. Though we can use the cosine distance in the latent space (Reimers and Gurevych, 2019) to select standalone explanations, a more effective use case is to *complement* existing feature attribution methods like SHAP (Lundberg and Lee, 2017) or LIME (Ribeiro et al., 2016). These methods provide overviews of feature weights, which are not available in pure counterfactual explanations. However, they estimate weights by *masking* words, and therefore cannot reflect how models would react to replacements or additions.

To highlight this usually overlooked nuance, we define $\mathbf{H}[\cdot]$ using the importance (weights) of the perturbed tokens in x , estimated by SHAP. Abnormalities selected this way highlight the missed pieces, and better calibrate users’ trust in the predictor. Figure 3C shows such a \hat{x}_L : **friend** \rightarrow **woman** changes the prediction from *Duplicate* to *Non-Duplicate*, even though “friend” is trivial according to SHAP; Whereas changing the important “in depression” in Figure 3D still results in *Duplicate* (\hat{x}_U). Detailed distance functions are in §D.1.

5.2 Evaluation on Local Explanations

We conduct a user study to verify whether our local counterfactual explanations can complement SHAP. The study takes the form of counterfactual simulation (Hase and Bansal, 2020), with participants predicting a model’s behavior on \hat{x} . Intuitively, the more they simulate incorrectly, the more information they grasp *if we show the counterfactuals*.

Procedure. We recruited 13 graduate students who have intermediate NLP knowledge and have experience using model explanations, and asked them to simulate the aforementioned *QQP* model

for 20 rounds. In each round, the participants were given a base example with the model’s prediction, as well as the SHAP weights, highlighted in the text and with a bar chart (Figure 3A). Moreover, they could create up to ten counterfactuals on their own, and query the model predictions on them. More interactions with the predictor usually result in better mental models (Miller, 2019), and we are interested in whether our counterfactuals *still add information* after such nearly unlimited predictor access (participants submitted 6.3 ± 3.2 queries per round.) Participants then simulated the model’s predictions on six given counterfactuals, two from each of the following three conditions. We concluded the study with surveys on their counterfactual creation and simulation strategies. The interface is in §D.2.

Conditions. We compare three types of counterfactuals: (1) *SHAP-c*, the POLYJUICE-generated counterfactuals, selected to complement SHAP; (2) *Random*, the randomly selected POLYJUICE counterfactuals; (3) *Human*, the human-generated counterfactuals, in which two graduate students (not participants) played with the model, and each created one \hat{x} where the prediction was incorrect and counterintuitive according to the SHAP score on x .

Results. As a within-subject study, we compared the error rate of human simulations across the three conditions. Participants were able to simulate the cases in *Random* (error rate $e = 23\% \pm 6\%$), possibly because *Random* selections contained more minor variations that aligned with the SHAP values and participants’ intuitions. They missed more *Human* ($e = 39\% \pm 11\%$) cases, and were even slightly worse on *SHAP-c* ($45\% \pm 6\%$, only slightly better than random guess).

This shows that *SHAP-c* counterfactuals are beyond participants’ learnings from feature attributions and manual counterfactual analysis, and *would still add value if they were presented*. They are also at least as effective as the *Human* ones, which is very expensive to create — each graduate student spent 1.5–2 hours on the task.

Usually, participants simulated the model incorrectly because they missed the inspection spots. For example, they repeatedly perturbed “depression” in Figure 3A, and therefore had to “guess using intuitions” when simulating Figure 3B. However, in 24% of the missed *SHAP-c* cases, participants successfully covered the related pattern,⁴

⁴At least one of their queries perturbed the same spans as

but were misled by their inspections — “labeled based on similar examples I tried,” as one subject articulated. It was hard for them to imagine the model predicting *Duplicate* on Figure 3B (**help** → **find**), when the model predicted *Non-Duplicate* on their query “How do I **help** → **play with**...?” The number dropped to 15% for the *Human* condition. In other words, SHAP-c *found more bugs within spots where humans considered inspected*.

Takeaway. POLYJUICE counterfactuals complement feature attribution methods and counterfactual analysis, as effectively as hiring a second expert for the analysis (but much cheaper). In particular, they highlight erroneous spots where humans may be misled by their own analysis.

5.3 Beyond Local Explanations

Interactive explanations: user-selected abnormalities. Automated selections focus on general abnormality, but users should be able to point towards the part *they* do not understand (Miller, 2019). The targeted perturbation seamlessly supports such interactive explanation. For example, an analyst can follow up on Figure 3C by BLANKing “friend” in Q2, and observe the model’s unstable behaviors: it predicts *Non-duplicate* when **friend** is changed to **woman**, **kid**, **professional**, but remains *Duplicate* when the noun is **man**, **student**.

Global Explanations: Recurring edits with unclear impact. *Global* explanations provide systematic understandings beyond individual instances — yet another important aspect of model understanding (Miller, 2019). We define global abnormality as *perturbation patterns whose impacts on the prediction are hard to generalize*. To locate them, we group counterfactuals based on their $r(\hat{x})$, which includes { tokens removed from x , added to \hat{x} , and the corresponding control code. } Then, we select groups with unstable prediction changes, quantified as large entropy on $f(x) \rightarrow f(\hat{x})$ for its \hat{x} . For example, one global abnormality for the NLI RoBERTa model in §4.3 is **two** → **three** in the hypothesis. Out of the 253 \hat{x} whose original $f(x) = \textit{entailment}$, 138 flipped to *contradiction*, 22 to *neutral*, yet 93 was intact, resulting in entropy $I = 0.91$. We observe that most cases flipped to *contradiction* have the explicit word “two” in the premise, whereas the prediction-intact ones suggest that the model struggles with counting:

<p>P: Two women having drinks at the bar. H: Two → Three women are at a bar. $f(x) \rightarrow f(\hat{x})$: entailment → contradiction — P: A boy and a girl gaze in a clothing store window. H: Two → Three kids are looking in a store window. $f(x) \rightarrow f(\hat{x})$: entailment → entailment</p>

6 Discussion

We create counterfactuals through *task-agnostic* generation, and *task-specific* selection. As a result, all applications have access to the same pool of counterfactuals generated by POLYJUICE, which are fluent, diverse, and close to the original sentence. With additional selection methods, POLYJUICE supports various downstream tasks, including counterfactual data augmentation, contrast set generation, as well as counterfactual explanation. We have made POLYJUICE available, and plan to opensource the implementation of the selection methods. Below, we discuss promising future work:

More balanced sentence pairs and sampling. POLYJUICE inherits some most salient contrasts pairs from existing paired datasets. For example, it is more likely to change **man** to **woman** than **child**. To further improve vocabulary diversity, we can emphasize more on *finding naturally occurring sentence pairs in non-paired dataset*. We used heuristics on text overlaps to broaden some control codes, but methods like syntactic tree editing or knowledge graph distances can also be useful.

Supporting existing counterfactual reasoning tools. Besides independent use cases, POLYJUICE may also be plugged into existing tools that involve counterfactual generation. For example, domain experts can create perturbation tests more easily if POLYJUICE serves as an additional building block for model testing (Ribeiro et al., 2020). Or, they may translate $r(\hat{x})$ into labeling functions for data programming (Ratner et al., 2017).

Human-in-the-loop counterfactual generation. Humans can add more values than placing blanks and labeling counterfactuals, *e.g.*, supply missing perturbation patterns after seeing POLYJUICE’s generations. In tasks that require context-dependent generations — difficult for the sentence-based model alone — human annotators can also initiate seed counterfactuals for POLYJUICE to extend. For example, to perturb question answering instances (Gardner et al., 2020), the human can add the compositional reasoning steps *related to* the corresponding paragraph for POLYJUICE to perturb around.

\hat{x} , and query text overlaps with the \hat{x} for over 70%.

Acknowledgements

The project was supported by ONR grant N00014-18-1-2193, NSF RAPID grant 2040196, NSF award IIS-1901386, the University of Washington WRF/Cable Professorship, and the Allen Institute for Artificial Intelligence (AI2). We gratefully thank Jim Chen, Scott Lundberg, Hao Peng, Sameer Singh, and Sitong Zhou for their helpful comments. We also appreciate the valuable input from our user study participants.

References

- Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani Srivastava, and Kai-Wei Chang. 2018a. [Generating natural language adversarial examples](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2890–2896, Brussels, Belgium. Association for Computational Linguistics.
- Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani Srivastava, and Kai-Wei Chang. 2018b. [Generating natural language adversarial examples](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2890–2896, Brussels, Belgium. Association for Computational Linguistics.
- Nabiha Asghar. 2016. Yelp dataset challenge: Review rating prediction. *arXiv preprint arXiv:1605.05362*.
- 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*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Vincent Chen, Sen Wu, Alexander J Ratner, Jen Weng, and Christopher Ré. 2019. Slice-based learning: A programming model for residual learning in critical data slices. In *Advances in neural information processing systems*, pages 9397–9407.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. [Plug and play language models: A simple approach to controlled text generation](#). In *International Conference on Learning Representations*.
- 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.
- Chris Donahue, Mina Lee, and Percy Liang. 2020. [Enabling language models to fill in the blanks](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2492–2501, Online. Association for Computational Linguistics.
- Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2018. [HotFlip: White-box adversarial examples for text classification](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 31–36, Melbourne, Australia. Association for Computational Linguistics.
- Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hannaneh Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace, Ally Zhang, and Ben Zhou. 2020. [Evaluating models’ local decision boundaries via contrast sets](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1307–1323, Online. Association for Computational Linguistics.
- Sahaj Garg, Vincent Perot, Nicole Limtiaco, Ankur Taly, Ed H Chi, and Alex Beutel. 2019. Counterfactual fairness in text classification through robustness. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 219–226.
- 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.
- Atticus Geiger, Ignacio Cases, Lauri Karttunen, and Christopher Potts. 2019. [Posing fair generalization tasks for natural language inference](#). 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 4485–4495, Hong Kong, China. Association for Computational Linguistics.
- Max Glockner, Vered Shwartz, and Yoav Goldberg. 2018. [Breaking NLI systems with sentences that require simple lexical inferences](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 650–655, Melbourne, Australia. Association for Computational Linguistics.
- Alec Go, Richa Bhayani, and Lei Huang. 2009. Twitter sentiment classification using distant supervision. *CS224N project report, Stanford*, 1(12):2009.
- Emily Goodwin, Koustuv Sinha, and Timothy J. O’Donnell. 2020. [Probing linguistic systematicity](#).

- In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1958–1969, Online. Association for Computational Linguistics.
- Peter Hase and Mohit Bansal. 2020. [Evaluating explainable AI: Which algorithmic explanations help users predict model behavior?](#) In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5540–5552, Online. Association for Computational Linguistics.
- Zhiting Hu, Zichao Yang, Xiaodan Liang, R. Salakhutdinov, and E. Xing. 2017. Toward controlled generation of text. In *ICML*.
- William Huang, Haokun Liu, and Samuel R Bowman. 2020. Counterfactually-augmented snli training data does not yield better generalization than unaugmented data. *arXiv preprint arXiv:2010.04762*.
- Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke Zettlemoyer. 2018. [Adversarial example generation with syntactically controlled paraphrase networks](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1875–1885, New Orleans, Louisiana. Association for Computational Linguistics.
- Robin Jia, Aditi Raghunathan, Kerem Göksel, and Percy Liang. 2019. [Certified robustness to adversarial word substitutions](#). 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 4129–4142, Hong Kong, China. Association for Computational Linguistics.
- Yichen Jiang and Mohit Bansal. 2019. [Avoiding reasoning shortcuts: Adversarial evaluation, training, and model development for multi-hop QA](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2726–2736, Florence, Italy. Association for Computational Linguistics.
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2020. Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 8018–8025.
- Daniel Kahneman and Amos Tversky. 1981. The simulation heuristic. Technical report, Stanford Univ CA Dept of Psychology.
- Sin-Han Kang, Hong-Gyu Jung, Dong-Ok Won, and Seong-Whan Lee. 2020. Counterfactual explanation based on gradual construction for deep networks. *arXiv preprint arXiv:2008.01897*.
- Divyansh Kaushik, Eduard Hovy, and Zachary C Lipton. 2020. Learning the difference that makes a difference with counterfactually-augmented data. *International Conference on Learning Representations (ICLR)*.
- Pride Kavumba, Naoya Inoue, Benjamin Heinzerling, Keshav Singh, Paul Reiser, and Kentaro Inui. 2019. [When choosing plausible alternatives, clever hans can be clever](#). In *Proceedings of the First Workshop on Commonsense Inference in Natural Language Processing*, pages 33–42, Hong Kong, China. Association for Computational Linguistics.
- Daniel Khashabi, Tushar Khot, and Ashish Sabharwal. 2020. [More bang for your buck: Natural perturbation for robust question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 163–170, Online. Association for Computational Linguistics.
- Najoung Kim, Roma Patel, Adam Poliak, Patrick Xia, Alex Wang, Tom McCoy, Ian Tenney, Alexis Ross, Tal Linzen, Benjamin Van Durme, Samuel R. Bowman, and Ellie Pavlick. 2019. [Probing what different NLP tasks teach machines about function word comprehension](#). In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)*, pages 235–249, Minneapolis, Minnesota. Association for Computational Linguistics.
- Varun Kumar, Ashutosh Choudhary, and Eunah Cho. 2020. [Data augmentation using pre-trained transformer models](#). In *Proceedings of the 2nd Workshop on Life-long Learning for Spoken Language Systems*, pages 18–26, Suzhou, China. Association for Computational Linguistics.
- Chuanrong Li, Lin Shengshuo, Zeyu Liu, Xinyi Wu, Xuhui Zhou, and Shane Steinert-Threlkeld. 2020a. [Linguistically-informed transformations \(LIT\): A method for automatically generating contrast sets](#). In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 126–135, Online. Association for Computational Linguistics.
- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020b. [BERT-ATTACK: Adversarial attack against BERT using BERT](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6193–6202, Online. Association for Computational Linguistics.
- Nelson F. Liu, Roy Schwartz, and Noah A. Smith. 2019a. [Inoculation by fine-tuning: A method for analyzing challenge datasets](#). 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 2171–2179, Minneapolis, Minnesota. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b.

- Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Shayne Longpre, Yu Wang, and Chris DuBois. 2020. [How effective is task-agnostic data augmentation for pretrained transformers?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4401–4411, Online. Association for Computational Linguistics.
- Scott M Lundberg and Su-In Lee. 2017. [A unified approach to interpreting model predictions](#). In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4765–4774. Curran Associates, Inc.
- Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, pages 142–150.
- Nishtha Madaan, Inkit Padhi, Naveen Panwar, and Dip-tikalyan Saha. 2021. Generate your counterfactuals: Towards controlled counterfactual generation for text. *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Nikolaos Malandrakis, Minmin Shen, Anuj Goyal, Shuyang Gao, Abhishek Sethi, and Angeliki Metallinou. 2019. [Controlled text generation for data augmentation in intelligent artificial agents](#). In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, pages 90–98, Hong Kong. Association for Computational Linguistics.
- Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. [Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3428–3448, Florence, Italy. Association for Computational Linguistics.
- Tim Miller. 2019. [Explanation in artificial intelligence: Insights from the social sciences](#). *Artificial Intelligence*, 267:1 – 38.
- John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020. [TextAttack: A framework for adversarial attacks, data augmentation, and adversarial training in NLP](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 119–126, Online. Association for Computational Linguistics.
- Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. Stress test evaluation for natural language inference. *arXiv preprint arXiv:1806.00692*.
- Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. [Justifying recommendations using distantly-labeled reviews and fine-grained aspects](#). 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 188–197, Hong Kong, China. Association for Computational Linguistics.
- Yixin Nie, Yicheng Wang, and Mohit Bansal. 2019. Analyzing compositionality-sensitivity of nli models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6867–6874.
- Judea Pearl. 2018. Causal and counterfactual inference. *The Handbook of Rationality*, pages 1–41.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020a. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020b. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ questions for machine comprehension of text](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Yanou Ramon, David Martens, Foster Provost, and Theodoros Evgeniou. 2019. Counterfactual explanation algorithms for behavioral and textual data. *arXiv preprint arXiv:1912.01819*.
- Alexander Ratner, Stephen H Bach, Henry Ehrenberg, Jason Fries, Sen Wu, and Christopher Ré. 2017. Snorkel: Rapid training data creation with weak supervision. *Proceedings of the VLDB Endowment*, 11(3):269–282.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert: Sentence embeddings using siamese bert-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?": Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.

- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018a. Anchors: High-precision model-agnostic explanations. In *AAAI*, volume 18, pages 1527–1535.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018b. [Semantically equivalent adversarial rules for debugging NLP models](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 856–865, Melbourne, Australia. Association for Computational Linguistics.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018c. [Semantically equivalent adversarial rules for debugging NLP models](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 856–865, Melbourne, Australia. Association for Computational Linguistics.
- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. [Beyond accuracy: Behavioral testing of NLP models with CheckList](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.
- Sara Rosenthal, Noura Farra, and Preslav Nakov. 2017. Semeval-2017 task 4: Sentiment analysis in twitter. In *Proceedings of the 11th international workshop on semantic evaluation (SemEval-2017)*, pages 502–518.
- Alexis Ross, Ana Marasović, and Matthew E. Peters. 2020. [Explaining nlp models via minimal contrastive editing \(mice\)](#).
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. [Winogrande: An adversarial winograd schema challenge at scale](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8732–8740.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *ArXiv*, abs/1910.01108.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Liwei Song, X. Yu, H. Peng, and Karthik Narasimhan. 2020. Universal adversarial attacks with natural triggers for text classification. *ArXiv*, abs/2005.00174.
- Samson Tan, Shafiq Joty, Min-Yen Kan, and Richard Socher. 2020. [It’s morphin’ time! Combating linguistic discrimination with inflectional perturbations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2920–2935, Online. Association for Computational Linguistics.
- Damien Teney, Ehsan Abbasnejad, and Anton Hengel. 2020. [Learning What Makes a Difference from Counterfactual Examples and Gradient Supervision](#), pages 580–599.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. Investigating gender bias in language models using causal mediation analysis. *Advances in Neural Information Processing Systems*, 33.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. [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.
- Jason Wei and Kai Zou. 2019. [EDA: Easy data augmentation techniques for boosting performance on text classification tasks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6382–6388, Hong Kong, China. Association for Computational Linguistics.
- John Wieting and Kevin Gimpel. 2018. [ParaNMT-50M: Pushing the limits of paraphrastic sentence embeddings with millions of machine translations](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 451–462, Melbourne, Australia. Association for Computational Linguistics.
- 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, New Orleans, Louisiana. Association for Computational Linguistics.
- 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.

- Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2019a. [Errudite: Scalable, reproducible, and testable error analysis](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 747–763, Florence, Italy. Association for Computational Linguistics.
- Xing Wu, Shangwen Lv, Liangjun Zang, Jizhong Han, and Songlin Hu. 2019b. Conditional bert contextual augmentation. In *International Conference on Computational Science*, pages 84–95. Springer.
- Zichao Yang, Zhiting Hu, Chris Dyer, Eric P Xing, and Taylor Berg-Kirkpatrick. 2018. Unsupervised text style transfer using language models as discriminators. In *Advances in Neural Information Processing Systems*, pages 7287–7298.
- Huangzhao Zhang, Hao Zhou, Ning Miao, and Lei Li. 2019a. [Generating fluent adversarial examples for natural languages](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5564–5569, Florence, Italy. Association for Computational Linguistics.
- Kaizhong Zhang and Dennis Shasha. 1989. Simple fast algorithms for the editing distance between trees and related problems. *SIAM journal on computing*, 18(6):1245–1262.
- Yuan Zhang, Jason Baldridge, and Luheng He. 2019b. [PAWS: Paraphrase adversaries from word scrambling](#). 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 1298–1308, Minneapolis, Minnesota. Association for Computational Linguistics.
- Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texygen: A benchmarking platform for text generation models. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 1097–1100.

Dataset	negation	quantifier	leixcal	resemantic	insert	delete	restructure	shuffle	global
CAD	3,456	457	10,650	4,634	2,169	2,162	234	84	3,756
Contrast	336	436	1,607	1,291	589	586	275	149	877
HANS	50	0	0	0	3,926	3,926	494	1,602	2
ParaNMT	2,797	825	10,000	10000	6,442	6,205	5,136	1,417	10,000
PAWS	81	1,815	10,000	10000	3,630	3,403	4,551	10,000	10,000
WinoGrande	3,011	94	10,000	6,927	120	124	453	65	3184
<i>Crawled</i>	0	0	5,000	0	5,000	5,000	0	108	5,000
Total	9,731	3,627	47,257	32,852	21,876	21,406	11,143	13,425	32,819

Table 6: The datasets used for finetuning the GPT-2 generation model, and the control code distributions.

A GPT-2 as Counterfactual Generator

A.1 Training data collection

We combine the following NLP datasets for finetuning POLYJUICE. To achieve a more balanced distribution, for each dataset, we extract control codes from all the data pairs available, and randomly sample up to 10,000 instances per control code. The distribution is shown in Table 6.

Contrast set Authors of 10 existing NLP dataset each manually perturbed 100–1,000 instances in small but meaningful ways that change the gold label, so to inspect a model’s decision boundary around a local instance (Gardner et al., 2020). The perturbation patterns vary based on the tasks and the annotators, allowing us to learn diverse strategies. To make sure we can use the contrast set to evaluate the *Sentiment* model, we excluded the IMDb movie review from the training.⁵

Counterfactually-augmented data (CAD) To augment the training data, Kaushik et al. (2020) crowdsourced counterfactuals for IMDb movie review (1.7k counterfactuals on 1.7k original instances) and SNLI (6.6k counterfactuals on 1.67k original). Similar to the contrast set, CAD’s perturbation patterns vary based on the task, but can especially contribute to *negation*. We split the movie review paragraphs into paired sentences, to match the sentence length of other datasets.

WinoGrande is a large-scale dataset of 44k instances for testing common sense problems (Sakaguchi et al., 2020). The dataset contains nearly identical sentences that differ only by one trigger word (e.g., one noun), which flips the correct answer choice for certain questions. The dataset is most suitable for learning lexical exchanges.

⁵Similarly, though *QQP* would be a potentially interesting dataset for training POLYJUICE, we omitted it so *QQP* can be used in our evaluation.

ParaNMT-50M contains 50 million English-English sentential paraphrase pairs, covering various domains and styles of text, as well as different sentence structures (Wieting and Gimpel, 2018).

PAWS contains 108k paraphrasing and non-paraphrasing pairs with high lexical overlaps. Zhang et al. (2019b) created such challenging pairs through controlled word swapping and back translation. As a result, the dataset demonstrates the *shuffle* and *restructure* strategies.

HANS is a controlled evaluation dataset designed for testing decision boundaries of NLI models (McCoy et al., 2019). The dataset contains 10k pairs of premises and hypotheses created based on 10 heavily fallible syntactic templates, and therefore compensates rarer structural changes that may be missed by PAWS.

Crawled We additionally crawl naturally occurring sentence pairs from non-paired datasets like SQuAD (Rajpurkar et al., 2016) to boost some specific patterns and increase lexical diversity. We estimate *close* pairs using edit distance, and broadly included pairs as long as the editing is less than 60%. This inevitably includes cases that should not be considered counterfactuals, e.g., “how do I not be” and “how do I recover it” are incorrectly considered *negation* pairs. We only include them for the most determined patterns, i.e., *lexical*, *insert*, *delete*, and *shuffle*. Among them, we further filter the pairs using control codes (see the section below).

A.2 Training Prompts & Parameters

Given a sentence pair, we use the two sentences interchangeably as x and \hat{x} to learn the counterfactuals both ways. For a (x, \hat{x}) , we compute its primary control code based on linguistic features like part-of-speech tagging or dependency trees, and blank out the changed subtrees in \hat{x} . For example, *negation* occurs when we observe changes on

negation modifiers or specific words like “supposedly”, and `shuffle` occurs when we have overlaps between tokens deleted and added. When multiple changes occur, we label it with the primary control code, which most significantly changes the semantic meaning on the corresponding subphrase. In Figure 2A, we use the code `negation`, as `great` \rightarrow `not great` is more significant than `kids` \rightarrow `children`. If we cannot identify the control code for a (x, \hat{x}) or if the editing distance is too large, we denote it with `global` and use it as a negative training sample. Importantly, to allow flexible blanking at the generation time, we generate multiple training prompts from one (x, \hat{x}) , with different blanking strategies: (1) just the changed tokens, (2) the associated parsing structures, (3) the merged changes, and (4) the entire sentence. As a result, we form up to four unique training prompts given one (x, \hat{x}) pair (some examples are in Figure 2).

With the interchangeable orders and the blanks, we generate 657,144 training prompts from 191,415 sentence pairs. We use the data to finetune an off-the-shelf GPT-2 model from Wolf et al. (2020), but any LM can potentially be used. We finetuned the model for 3 epochs, with an initial learning rate $5e-5$, a batch size of 16 and a sequence length of 120.

A.3 Intrinsic Evaluations

A.3.1 Controllability with Ablation Studies

We finetune another GPT-2 model with training prompts that *do not* contain control codes (called POLYJUICE-*a*), and quantify the impact of the codes through an ablation study. For each control code, we compare the *control success rate* of POLYJUICE and POLYJUICE-*a* on 250 prompts (from 100 unique original sentences). For each prompt, we generate counterfactuals through beam search (beam = 10), and recompute the codes on the top three returns. We deem the control successful if at least one of the three recomputed codes matches the input (though in POLYJUICE-*a*, we only measure whether the code naturally occurs in the uncontrolled generation.) The success rate increases by $28.4\% \pm 18.2\%$ across all control codes, ranging from `quantifier` (increasing 8%, from 40.4% to 48.4%) to `insert` (64.1%, from 13.5% to 78.6%).

There are three common failure cases for the codes: (1) The dual manipulation from the control codes and the blanks can conflict, *e.g.*, “a dog is embraced by a [BLANK]” would not respond

Model	Diversity	Closeness	
	Self-BLEU \uparrow	Semantic \downarrow	Tree edit \downarrow
POLYJUICE	0.82	0.37	2.02
Masked-LM	0.66	0.27	1.89
PPLM-BoW	0.82	0.65	7.65

Table 7: The intrinsic metrics comparing POLYJUICE with Masked-LM and PPLM-BoW. \uparrow (or \downarrow) indicates whether the metric should be maximized (or minimized). As expected, POLYJUICE counterfactuals are *closer* to the original instance than PPLM-BoW, and more *diverse* than the Masked-LM ones.

to `negation`. (2) x does not have a corresponding pattern. `shuffle` is not applicable when the sentence has only one adjective or noun (*e.g.*, “the movie is good”). (3) Certain pattern is very prominent that it dominates the generation probability, *e.g.*, the model tends to perturb the quantifier “two” in “two dogs are running”, regardless of the code. In the ablation study, we filtered out prompts that fell under cases 1 and 2.

A.3.2 Closeness, Diversity, Fluency

Similar to Madaan et al. (2021), we verify that the POLYJUICE generations are fluent (“plausible” in their case) through human evaluations (§4.2). We also quantify the diversity and closeness by comparing it with baseline models.

For a given x and its generated counterfactuals $\hat{\mathbf{X}}$, we approximate *diversity* using self-BLEU (Malandrakis et al., 2019; Zhu et al., 2018) within the generated counterfactual set $\hat{\mathbf{X}}$. The higher the BLEU, the more lexically different the generated phrases are to each other (Note that this only reflects one form of $r(\hat{x})$ relationship that is task-agnostic, and is not the only possible measurement for diversity). Meanwhile, *closeness* is measured using both the average semantic and syntactic distance between x and every $\hat{x} \in \hat{\mathbf{X}}$. We compute the semantic distance using a pretrained sentence similarity model (Reimers and Gurevych, 2019), and the syntactic one using the tree edit distance (Zhang and Shasha, 1989).

We use similar baselines as Madaan et al. (2021): (1) *Masked-LM*, where we blank certain parts of x (by randomly placing up to three [MASK] tokens), and ask RoBERTa to fill in the blank. We rely on the CheckList implementation (Ribeiro et al., 2020), as it allows filling in multiple blanks at once through beam search.

(2) *PPLM-BoW* (Dathathri et al., 2020), a model that uses bag-of-words to control the generation.

Application	Strategies
Data augmentation	Lexical (Wu et al., 2019b; Wei and Zou, 2019; Kumar et al., 2020) Paraphrasing (Iyyer et al., 2018) Perturbation functions (Ratner et al., 2017)
Counterfactual data aug.	Manual (Kaushik et al., 2020) Lexical (Garg et al., 2019)
Adversarial attack	Lexical (Alzantot et al., 2018a; Garg and Ramakrishnan, 2020; Li et al., 2020b; Morris et al., 2020; Tan et al., 2020; Jin et al., 2020; Ebrahimi et al., 2018; Zhang et al., 2019a; Jia et al., 2019) Template (Jiang and Bansal, 2019) Insert (Song et al., 2020)
Contrast set	Manual (Gardner et al., 2020) Templates (Li et al., 2020a)
Challenge sets	Lexical (heuristic) (Kaushik et al., 2020; Naik et al., 2018) Paraphrasing (Kavumba et al., 2019) Templates (Geiger et al., 2019; Kaushik et al., 2020; Nie et al., 2019; McCoy et al., 2019)
Model analysis	Lexical (Garg et al., 2019) Template (Goodwin et al., 2020) Perturbation functions (Wu et al., 2019a; Bowman et al., 2015; Ribeiro et al., 2020) Controlled text generation (Madaan et al., 2021)
Explanations	Lexical (Hase and Bansal, 2020; Vig et al., 2020; Kang et al., 2020) Lexical (through masking) (Ramon et al., 2019; Ribeiro et al., 2018a) Language model (Ross et al., 2020)

Table 8: A short survey on counterfactual application papers, and their generation strategies.

We use the first two words of x as the input context (prompt), limit the length of the generation to be similar to x , and apply their default condition “positive words.” As the model generates *arbitrary text* that do not depend on x , we agree with Madaan et al. (2021) that PPLM-BoW should not satisfy the *closeness* requirement.

We randomly select 120 instances from three classification tasks, *QQP*, *NLI*, and *Sentiment* (40 per task), and generate 10 counterfactuals per x using each of the three generators. The averaged metrics are in Table 7. POLYJUICE achieves compatible diversity with PPLM-BoW (self-BLEU score), and compatible *closeness* with Masked-LM, achieving a balance between both.

Ideally, we would also like to compare POLYJUICE with our concurrent work GYC (Madaan et al., 2021) and MiCE (Ross et al., 2020). Inspired by style transfer (Yang et al., 2018) and controlled text generation, GYC performs the perturbation on the latent space of the input x . Meanwhile, MiCE uses a two-step framework to generate counterfactual explanations, with the generator being T5 (Raffel et al., 2020b) finetuned on the task-specific dataset. As mentioned in §2, both generators focus on flipping the class label of a given x . Unfortunately, both require extensive implementation or finetuning, and has yet to be opensourced.

B Survey of Perturbation Applications

Table 8 summarizes existing applications of the counterfactuals, with their generation strategies.

C Additional Details to Train & Eval §4

C.1 Tasks & Data

Sentiment Analysis (*Sentiment*) aims to determine the sentiment polarity of a given sentence (*positive* or *negative*). We select Stanford Sentiment Treebank (SST-2) (Socher et al., 2013) as the base dataset for augmentation. It contains sentences extracted from full movie reviews on Rotten Tomatoes, which is relatively more aligned with the training data for POLYJUICE. While the dataset also contains finer-grained labels on subphrases, we only use full sentences. As a result, the full training data contains 6,920 sentences.

Natural Language Inference (*NLI*) is a 3-way classification task, with inputs consisting of two sentences, a premise and a hypothesis, and the three possible labels being *entailment*, *contradiction*, and *neutral*. We augment the data based on SNLI (Bowman et al., 2015).

Duplicate Question Detection (*QQP*) analyzes whether two questions are duplicates of each other (*i.e.*, if you have the answer to one question, whether you can infer the answer for the other one.) We use QQP as the base dataset (Wang et al., 2018), a collection of question pairs from the community question-answering website Quora.

TASK DESCRIPTION

You will annotate a series of examples with two pieces of information:

1. **Natural**: Whether this sentence is likely written by a native speaker (**Valid**), or the writer doesn't speak English well, e.g., s/he makes **severe** grammar errors/the sentence is not semantically meaningful (**Invalid**, *no need to disqualify wrong spacing, short phrases or informal verbal language*).
2. **Label**: Whether two sentences/phrases S1 and S2 are captions of the same image. That is, if S1 correctly describes what's shown in an image, whether S2 can hold for the same image. (**Definitely False** / **May be True** / **Definitely True**);

For each round, you will be given a reference example:

Old S1 Two girls dancing in traditional garb .

Old S2 Two girls are dancing .

Label **Definitely True**

And you will be labeling several of its variations, with **New S2** edited. The labeling might be more intuitive if you pay attention to **what's changed**, and whether the change **affects the label in the reference example above**.

Old S1 Two girls dancing in traditional garb .

New S2 ~~Two~~ None of the girls are dancing .

Valid? Valid

Label **Definitely False**

Old S1 Two girls dancing in traditional garb .

New S2 Two girls **girls** are dancing .

Valid? Invalid

Label **Definitely True**

PROCEDURE

You will first go through a **1-round training phrase** to help you get familiar with the task. Then, you will complete **22 rounds** of labelings.

Figure 4: The instruction for the *NLI* labeling task in §4, with annotators labeling the perturbed hypotheses (*New S2*). Instructions are similar for *QQP* and *Sentiment*, except for the label definitions and the examples.

Reference Example

Old S1 Police officer with riot shield stands in front of crowd .

Old S2 A police officer stands in front of a crowd .

Label **Definitely True**

Label the following! [Review the instructions!](#)

The **green color** highlights new words added in **New S2** , compared to **Old S2** in the Reference example above. **.** indicates something is deleted.

Old S1 Police officer with riot shield stands in front of crowd .

New S2 A police officer **standing behind** a crowd .

Valid? ☐ Invalid ☒ Valid

Label ☒ Definitely False ☐ May be True ☐ Definitely True

Old S1 Police officer with riot shield stands in front of crowd .

New S2 A police officer stands **next to a truck** .

Valid? ☐ Invalid ☒ Valid

Label ☐ Definitely False ☐ May be True ☐ Definitely True

Old S1 Police officer with riot shield stands in front of crowd .

New S2 A **policeman** stands in front of a crowd .

Valid? ☐ Invalid ☒ Valid

Label ☐ Definitely False ☒ May be True ☐ Definitely True

Figure 5: A sample labeling task: The crowdworkers annotate three counterfactuals based on their validity and class label, with respect to the original instance.

C.2 MTurk Labeling Details

Procedure The study started with an introduction, in which we explained the context and tasks (Figure 4): “given a reference example, the crowdworker should annotate its counterfactual variations, based on whether the counterfactual is valid (*fluent*), and the classification task label.” To familiarize them with the task, we asked them to complete 1-2 training rounds, and explained the

expected labels. The annotator then completed 22 tasks, labeling 3 counterfactuals of a single example in each round, as in Figure 5. The 22 rounds consisted of 20 actual labeling tasks and 2 extra “gold rounds”, which included 6 unambiguous counterfactuals and known groundtruth labels in total. The gold cases later served as filters for high quality crowdworkers. As a result, each annotator contributed $20 \times 3 = 60$ labels. The median annotation time was around 15-20 minutes (14.9 for *QQP*, 16.7 for *Sentiment*, and 19.8 for *NLI*), and participants received \$2.5 on average.

Participants We recruited participants from Amazon’s Mechanical Turk (MTurk), limiting the pool to subjects from within the United States with a prior task approval rating of at least 97% and a minimum of 1,000 approved tasks.

Data quality We applied two filtering strategies: (1) *High-quality worker*. We only kept data from participants whose median labeling time per round was more than 18 seconds and correctly labeled at least 4 gold counterfactuals (out of 6), or who correctly labeled all gold ones regardless of the time spent. (2) *Majority vote labeling*. We collected two annotations per counterfactual, and only kept those

Task description

You will help us understand how people learn to **simulate the behavior of NLP models**.

DUPLICATE QUESTION DETECTION

The NLP model in question predicts whether two questions are duplicates of each other - that is, whether they should point to the same answer in a website like Quora. The labels are thus (**Non-duplicate** / **Duplicate**).

THE CONTEXT

In each round, you will see a reference example like below, with the model's prediction on it (below, the model makes a **Correct prediction**).

Old Q1 Why is the sky blue ?
Old Q2 Why is that the sky is so blue ?
Model predicts **Duplicate (98.5% confident)**
Model correct? **Correct**

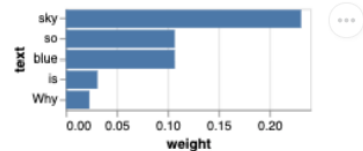
Additional Clues about the model's behavior

FEATURE IMPORTANCE

To help you understand what is driving the model's prediction, we keep **Q1** fixed and highlight on **Q2** words that are **important for the model's prediction**, in blue (we also show a bar chart with the five most important words).

These values are computed by a standard black-box explanation technique (SHAP), which **masks different groups of words in Q2** and summarizes how predictions change as a result.

Old Q1 Why is the sky blue ?
Old Q2 Why is that the sky is so blue ?
Model predicts **Duplicate (98.5% confident)**
Model correct? **Correct**



ASK THE MODEL QUESTIONS!

You can also **make small changes to Q2** and see the resulting model predictions, in order to get a better understanding of how the model behaves around the reference example. You can ask up to **10** additional questions per round to learn more about the model. For example, you may want to ask the following question:

Old Q1 Why is the sky blue ?
New Q2 Why is that the sky is so **blue dark** ?
Model predicts **Non-duplicate (99.6% confident)**

YOUR TASK

After seeing the *reference example*, *feature importances*, and *asking your own questions of the model*, we will ask you to try to **guess how the model would predict several variations of Question 2** (**New Q2**).

Beware that **the model is not perfect, so it may make mistakes** (you should try to simulate the model to the best of your ability). Below is an example of a variation of **Q2** we might ask you to label:

Old Q1 Why is the sky blue ?
New Q2 Why is that the sky is so **blue white** ?
Model will predict ☐ Non-duplicate ☒ Duplicate

As you see, the model may be incorrect, so please learn about **the model's behavior** carefully through the **Additional Clue**.

Procedure

You will first go through a **1-round training phrase** to help you get familiar with the task. Then, you will complete **20 rounds** of labelings.

Figure 6: The instruction for the explanation study in §5.2.

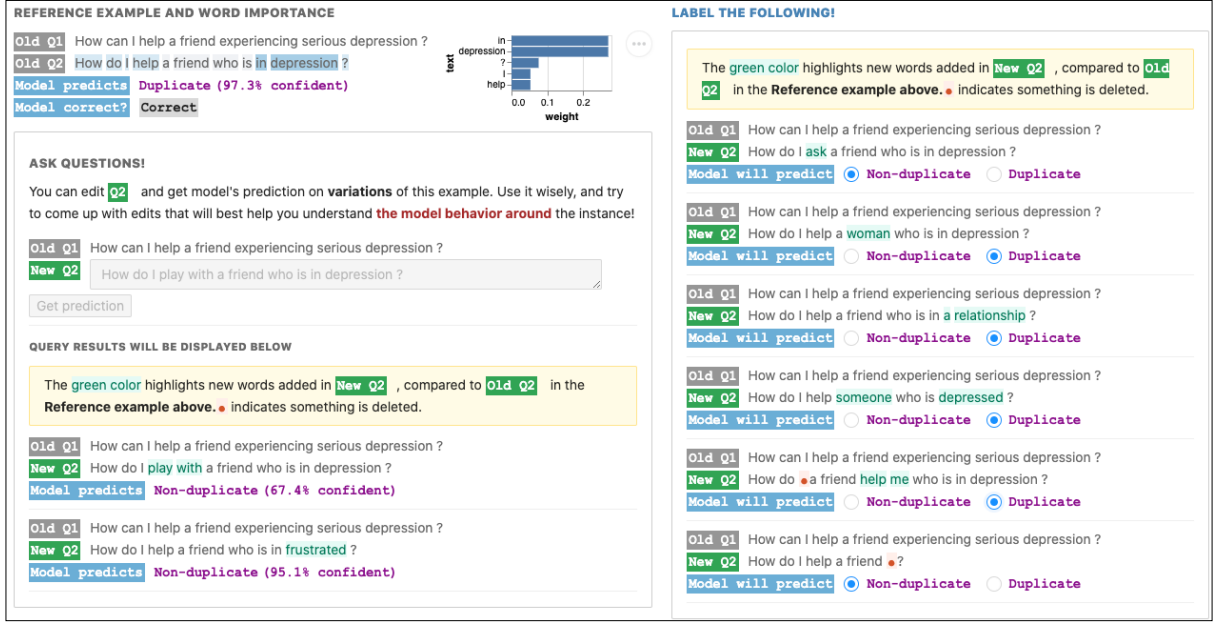


Figure 7: A sample explanation interface task for §5.2.

that at least one annotator deemed *valid/fluent*, and both annotators agreed on a particular *class label*.

Because the crowdworkers were usually noisy, when set out to collect augmentations on 1,000 original examples (thus 3,000 counterfactuals), we typically collect counterfactuals for 1,000 counterfactuals on 600 original examples. One of the authors labeled a subset of 100 counterfactuals on 100 original examples in *Sentiment*, and reached high agreement with the majority-voted results ($\kappa = 0.77$, the raw labeling agreement 88%).

C.3 Training Details for §4.4

We finetuned roberta-base models (Liu et al., 2019b) provided by HuggingFace Transformers (Wolf et al., 2020). For each given (m, n) , we created three different samples of training data, and Table 3–5 report the average and the variance of the samples. Each sample was further averaged over four random seeds. For each run, we heuristically picked initial learning rates 1e-5, 2e-5, 2e-5 for *Sentiment*, *NLI* and *QQP*, respectively, trained 20 epochs with a dropout rate of 0.1 and a batch size of 16. We selected the epoch that had the highest accuracy on the corresponding validation set (1/5 of the training data size, with the same ratio of original and counterfactual examples.)

C.4 Criteria for CheckList Tests

As a behavioral testing framework, CheckList defines multiple tests, and measures models’ linguistic capabilities using the failure rates of each test.

We define a test to have failed if the failure rate is over 20%. Because failure rates are more sensitive than the accuracy, we say a model capability is affected, if the failure rate of a test changes (increases or decreases) more than 5% (e.g., failure rates going from 20% to 21% is insignificant), and the delta accounts for 10% (e.g., failure rates decreasing 8% from 100% to 92% does not count.)

D Additional Details to Explanations §5

D.1 Selection Methods

SHAP complement as local explanation Because the SHAP weights reflect the average effect of masking a token t , instead of finding just one abnormal counterfactual, we focus on *word features that are abnormal on average*, i.e., there is a mismatch between the importance of the changed word features and the model behavior, averaged over *all counterfactuals* that have those features.

In this case, the expected change in prediction of perturbing a token t is the SHAP importance on it $\mathbf{H}[\mathbf{D}_t(t, x)] = s(t)$; For example, in Figure 3, $s(t = \text{“depression”}) = 0.276$.

The actual prediction change is the weighted average of $|f_p(x) - f_p(\hat{x})|$ for all the \hat{x} that affect t (depression \rightarrow trouble, depression \rightarrow a relationship), with the weight corresponding to the number of words modified in \hat{x} : If $e(\hat{x})$ denotes the set of edited words (replaced or deleted from x in $r(\hat{x})$), then $w(\hat{x}) = 1/|e(\hat{x})|$. Intuitively, the more words changed in \hat{x} , the less effect each word has; In Fig-

ure 3D, we regard “depression” to be responsible for half of the impact on changing **in depression** \rightarrow **suicidal**.

We group the counterfactuals based on their affected words $G_t = \{\hat{x} \mid t \in e(\hat{x})\}$. $D_f(t, x)$ then becomes:

$$D_f(t, x) = \frac{1}{|G_t| + 1} \left(s(t) + \sum_{\hat{x} \in G_t} w(t) \cdot |f_p(x) - f_p(\hat{x})| \right)$$

We add the additional SHAP weight $s(t)$ as a smoothing factor to punish dominating outlier \hat{x} . Then the gap between the expectation and the reality is (similar to §5.1):

$$\Delta D_f(t, x) = D_f(t, x) - \mathbf{H}[D_f(t, x)]$$

We first find the abnormal tokens: (1) t with small SHAP weight, but \hat{x} that change t experience large prediction change on average: $t_L = \arg \max_{t \in x} \Delta D_f(t, x)$, and (2) t with large SHAP weight, but \hat{x} with t changed usually have intact prediction: $t_U = \arg \max_{t \in x} -\Delta D_f(t, x)$.

Then, we use the most extreme cases within the groups of G_{t_L} and G_{t_U} as the concrete counterfactual explanations, based on their prediction change $|f_p(x) - f_p(\hat{x})|$, and the aggregated SHAP weights of all the changed tokens:

$$\hat{x}_L = \arg \max_{\hat{x} \in G_{t_L}} \left(|f_p(x) - f_p(\hat{x})| - \sum_{u \in r(\hat{x})} s(u) \right)$$

Global explanation To enable the grouping, we first featurize each counterfactual \hat{x} with respect to its original instance x , using (1) its control code [TAG:negation] for the example in Figure 2), (2) its remove phrases [FROM:kids], (3) its added phrases [TO:not], [TO:children], and (4) the combined template [TEMP: **kids** \rightarrow **children**]. For tokens involving multiple changes, we featurize both the primary and the combined changes, and so the example in Figure 2 also have additional features like [TAG:negation] & [TAG:lexical]. These features form the relationship $r(\hat{x})$ in use.

For each feature $h \in r(\hat{x})$, we build a group of counterfactuals that contain h (of course, these also include their corresponding x): $G_h = \{\hat{x} \mid h \in r(\hat{x})\}$. We compute the probability of the predicted label changes all the $\hat{x} \in G_h$: $Pr(y_1, y_2) = |G_h^{y_1 \rightarrow y_2}| / |G_h|$, where $G_h^{y_1 \rightarrow y_2} = \{(x, \hat{x}) \mid (x, \hat{x}) \in G_h, f(x) = y_1, f(\hat{x}) = y_2\}$. The abnormality of a feature h is represented by the entropy of the prediction

change:

$$I_h = - \sum_{y_1 \in Y, y_2 \in Y} Pr(y_1, y_2) \cdot \log Pr(y_1, y_2)$$

We find the abnormal feature with $\hat{h} = \arg \min I_h$.

D.2 User Study Details

The instruction for the user study in §5.2 is in Figure 6, and Figure 7 shows the sample interface for one round. Participants started by just seeing the reference example and the model query box on the left hand side. When they chose to start the task or after they had exhausted their ten query chances, the query box was disabled, the tasks on the right were displayed, and the participants complete the tasks.

E Additional counterfactuals

Here, we include some additional counterfactual for each classification task.

Additional examples
It sucked me in .
For those who are intrigued by politics of the '70s , the film is every bit as fascinating → flawed as it is flawed → intriguing .
So exaggerated and broad that it comes off as annoying → engaging rather than charming → annoying .
The film delivers → doesn't deliver what it promises: A look at the "wild ride" that ensues when brash young men set out to conquer the online world with laptops, cell phones and sketchy business plans.
It's a crime movie made by someone who obviously knows nothing → much about crime.
Even with → Despite all its botches, Enigma offers all → none of the pleasure of a handsome and well-made entertainment .
"Catch Me" feels almost capable of charming the masses with star power, a pop-induced score and sentimental moments that have become a Spielberg trademark.
A sentimental but entirely irresistible → unentertaining portrait of three aging sisters.
This is a movie full of grace → mistakes and, ultimately, no hope.
Its simplicity puts an exclamation point on the fact that this isn't something to be taken seriously, but it also wrecks any chance → poses an opportunity of the movie rising above similar fare.
If → Even if the film fails to fulfill → succeeds in fulfilling its own ambitious goals , it nonetheless sustains → fails to sustain interest during the long build-up of expository material .

Table 9: Additional examples for Sentiment Analysis.

Additional examples
Q1: Poor people are more generous than rich people. Why?
Q2: Is it true poor → rich people are more generous than rich → poor people?
Q2: Is it true poor people are more generous → give more to charities than rich people?
Q1: Are TripAdvisor reviews more reliable than Yelp reviews because of their review process
Q2: Are TripAdvisor reviews more → less reliable than Yelp reviews?
Q2: Are TripAdvisor reviews more reliable → unreliable than Yelp reviews?
Q2: Are TripAdvisor → Yelp reviews more reliable than Yelp → TripAdvisor reviews?
Q1: How do you describe a smell?
Q2: How will you describe a smell to → from a person?
Q2: How will you describe a sound → smell to a person?
Q1: Why are most psychopaths males and not females?and are female psychopaths different from male psychopaths?
Q2: Is there a difference between male and female psychopaths?
Q2: Is there a difference between male and female psychopaths ?
Q2: Why is there a difference between male and female psychopaths?
Q1: How can I loose 5kgs weight in a week without exercise?
Q2: How can I lose weight without → by doing exercise?
Q2: How can I lose weight without doing exercise ?
Q2: How can I → Why can't you lose weight without doing exercise?

Table 10: Additional examples for Duplicate Question Detection, with the Q1 intact and Q2 changed.

Additional examples
P: A small child stands in front of short white table.
H: A child near → sitting on a table.
P: A child sticks his head through a hole to create a picture of his head being a flower blossom
H: The child is poking his head through → to see something.
P: A young woman is playing fool.
H: The woman is old → very young and not playing any games.
P: Metal supports make a repeating X shape along the wall of the station.
H: The walls are stronger with metal → wood supports rather than wood → metal fences .
H: The walls are stronger → less sturdy with metal supports rather than wood.
P: Several gentlemen are speaking into a microphone and the man in the glasses appears to be saying something funny
H: Somebody that is shown wears → shows glasses.

Table 11: Additional examples for Natural Language Inference, with the **P**remise intact and **H**ypothesis changed.