

Causal Augmentation for Causal Sentence Classification

Anonymous ACL submission

Abstract

Scarcity of corpora with annotated causal texts can lead to poor robustness when training state-of-the-art language models for causal sentence classification. In particular, we find that these models misclassify on augmented sentences that have been negated or strengthened in terms of their causal meaning. This is worrying because minor linguistic changes in causal sentences can have disparate meanings. To resolve such issues, we propose to generate counterfactual causal sentences by creating contrast sets (Gardner et al., 2020). However, we notice an important finding that simply introducing edits is not sufficient to train models with counterfactuals. We thus introduce heuristics, like sentence shortening or multiplying key causal terms, to emphasize semantically important keywords to the model. We demonstrate these findings on different training setups and across two out-of-domain corpora. Our proposed mixture of augmented edits consistently achieves improved performance compared to baseline across two models and both within and out of corpus’ domain, suggesting our proposed augmentation also helps the model generalize.

1 Introduction

Causality is an important concept for knowledge discovery as it conveys the idea of cause and effect. In the simplest sense, a causal relation exists between entities A and B through the statement “A causes B” or “B is caused by A”. In recent years, causal relation extraction from text has garnered large interests in Natural Language Processing (NLP) (Asghar, 2016; Xu et al., 2020).

Causal sentence classification (CSC) is the task of identifying sentences that contain causality information. Figure 1 demonstrates examples where similar claims are categorized by their causal strengths. CSC is challenging because the syntax of causality varies in context. Thus, it is difficult

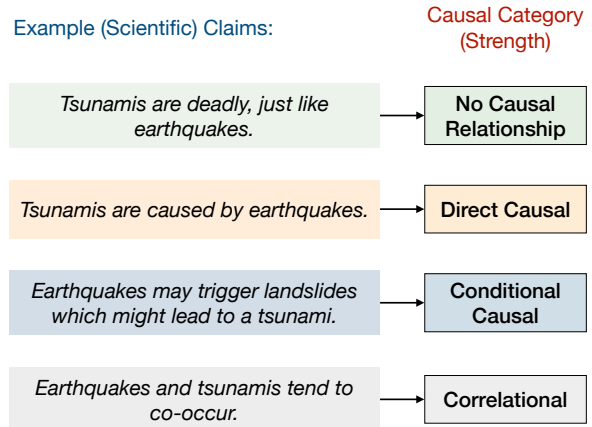


Figure 1: Causal sentence classification classifies textual claims into various categories of causal strengths.

to exhaustively capture causal expressions, especially for implicit occurrences (Asghar, 2016). Additionally, negations and the absence of causality complicate automatic causality identification tasks (Heindorf et al., 2020).

Furthermore, there is a lack of good quality CSC datasets (Asghar, 2016; Xu et al., 2020). Most NLP datasets typically treat causal relation extraction as a subtask of relation extraction, where “Cause-Effect” is one of the many relation labels. However, causality is a complex relation best learned using dedicated causal relation datasets. Dedicated causal relation corpora that exist are mostly small in size (< 5000 sentences), except for AltLex (Hidey and McKeown, 2016) that has over 40000 sentences. Datasets also tend to label causal relations in an overly simplistic binary level (as ‘causal’ or ‘not causal’). Only some works classify text by causal strengths (Girju and Moldovan, 2002; Yu et al., 2019; Sumner et al., 2014).

Data augmentation is a natural avenue for handling small-sized datasets. Augments created must be meaningful to explain representation gaps in the current datasets. In causality, both the causal direction and strength matter. As such, models need

to be sensitive towards negations and semantics of words to avoid misclassification. For example, in Figure 1, two similar sentences “*Tsunamis are caused by earthquakes.*” and “*Earthquakes may trigger ... a tsunami.*” differ in causal strengths. Therefore, we propose artificially constructing meaningful counterfactuals that would reflect the model’s decision boundaries. We do so by applying rule-based schemes that negate causal relations or strengthen conditionally causal sentences. Additionally, we explore heuristical edits on CSC performance.

We find that state-of-the-art (SOTA) language models, such as BERT (Devlin et al., 2019) with MLP or SVM classifiers, achieve improvements in classification performance when trained with our created counterfactuals. In addition, our evaluation on cross-domain datasets shows that training on augmented datasets (original plus edits) improves model generalization to out-of-domain (OOD) contexts. This is consistent with findings from (Kaushik et al., 2020a,b) in sentiment analysis and natural language inference contexts. In summary, we make the following contributions:

1. We propose causal negation and strengthening schemes based on dependency and part-of-speech (POS) tags to augment causal sentences. To our knowledge, we are the first to study the effects of counterfactual augmentation in the context of causal claims.
2. We show that current SOTA models are not robust to minimally perturbed sentences that differ in causal direction and strength.
3. We observe that simple heuristical edits on these counterfactuals make models more effective for low resource CSC with limited number of causal and conditional causal sentences.
4. We show that a mixture of counterfactuals improves performance in the trained domain and also generalize better to OOD corpora such as SCITE (Li et al., 2021) and AltLex (Hidey and McKeown, 2016).

2 Related Works

2.1 Causal Sentence Classification

Although causality is an important concept for knowledge discovery, benchmarking datasets and standardization of labeling rules have been limited,

prohibiting empirical comparisons across methodologies (Asghar, 2016; Xu et al., 2020). Most NLP benchmarking datasets define causal relations as just one out of many class labels (e.g. Part-Whole) (Jurgens et al., 2012; Gábor et al., 2018; Caselli and Vossen, 2017; Mirza et al., 2014; Mirza and Tonelli, 2016). Others, focused on causal relations, define such relations as a binary label (Li et al., 2021; Hidey and McKeown, 2016). However, causality may not always occur at extremes in real-life statements, and correlation can get confused for causation (Buhse et al., 2018). As such, instead of using a binary model of causality, a better way is to classify varying “strengths” of causal relations in sentences. In fact, a seven-point scheme¹ was proposed by Sumner et al. (2014) to categorize causal statements from health-related news and academic press releases. Subsequently, Yu et al. (2019) adapted this for scientific texts into a four-level system. In this work, we adopt the four-level causality labeled corpus and classification model by Yu et al. (2019)².

There is also an often observed issue that NLP systems that perform well on task datasets do not generalize to “real-life scenarios”, thereby misleading and overstating the accuracies and usefulness of their models. Ensuring model generalizability to other domains can be challenging. For example, Ramesh et al. (2012) showed discourse triggers are different between the biomedical and general domains. In recent years, more focus has been placed in the field to ensure sufficient data representativeness and transferability of results onto OOD settings. In this work, we will also evaluate the generalizability of our models to classify causal sentences from other domains.

2.2 Counterfactuals in NLP

Counterfactual generation is a popular strategy for NLP researchers to test and improve model robustness via adversarial learning and attacks (Morris et al., 2020; Mahler et al., 2017) or for mitigating bias (Kaushik et al., 2020a; Maudslay et al., 2019).

Gardner et al. (2020) proposed using counterfactuals to fill local theoretical gaps in a model’s

¹The seven levels of causal strengths are (1) no statement, (2) explicit statement of no relation, (3) correlational, (4) ambiguous (i.e. relationship is present but the direction and level is ambiguous), (5) conditional causal, (6) can cause, and (7) unconditionally causal.

²We were unable to work on Sumner et al.’s dataset as it was not publicly available and had very limited samples per class label.

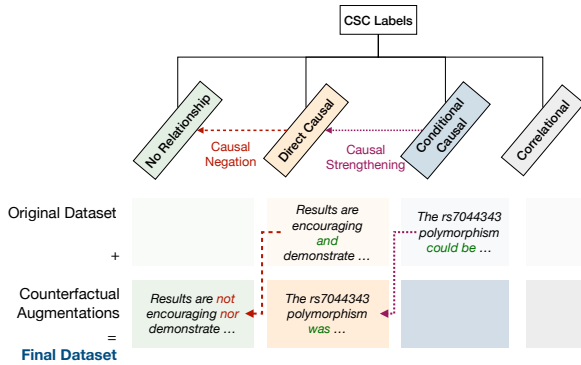


Figure 2: Strategies to generate counterfactual examples for CSC.

decision boundary. They relied on expert judgments to generate similar but meaningfully different sentences. They showed that across a variety of tasks (e.g. reading comprehension, sentiment analysis, visual reasoning) and input-output formats (e.g. classification, span extraction, structured prediction), SOTA models struggle on contrast sets compared to original test sets. In our work, we generate counterfactuals meaningful for CSC, such as moving sentences across labels when we perform Negation (*causal* \rightarrow *no relationship*) and Strengthening (*conditional causal* \rightarrow *causal*) strategies. While Gardner et al. (2020) have noted that it is challenging to come up with automated construction of contrast sets and proposed authors to manually perturb statements, we provide an automatic rule-based schema to negate and strengthen causal statements at scale.

Kaushik et al. (2020a) manually revised documents that would correspond to a counterfactual target label for sentiment analysis and natural language inference tasks. They showed that training with similar quantities of augmented data compared to the original improves generalization ability to OOD datasets. In this paper, we have also found that counterfactuals can help to improve model generalizability for CSC. Again, our linguistics-based augments do not rely on human intervention.

3 Methodology

3.1 Task Details

The CSC task involves classifying a span of text with a causal label based on its intended meaning. We use the PubMed-based corpus CSci³, provided by Yu et al. (2019) comprising of 3061 sentences

³<https://github.com/junwang4/causal-language-use-in-science>

annotated with 4 different levels of causal relation: *no relationship* (c_0), *causal* (c_1), *conditional causal* (c_2), and *correlational* (c_3).

3.2 Counterfactual Generation

In a low-resource setting, our proposal is that researchers should create counterfactuals that push causal sentences across labels so as to improve the robustness of their models. Figure 2 demonstrates the two main strategies to generate counterfactual examples for CSC, namely (1) Causal Negation and (2) Causal Strengthening. We discuss these strategies next⁴.

3.2.1 Causal Negation

In this strategy, we negate the direction of causal statements from *causal* (c_1) to *no relationship* (c_0).

After obtaining POS tags and root words based on dependency trees⁵, we performed negations around the root word. Our coding schema (Algorithm 1 in the Appendix) inserts negative words like ‘no’, ‘not’, ‘nor’ or ‘did not’ to flip the meaning of the sentence. 12 negation linguistic templates were used. Successfully negated sentences are termed as ‘Edit’ sentences. If no matching template was found, the sentence was skipped. Of the 493 original (causal) sentences from the CSci corpus, 384 sentences had available negations.

To improve text flow, we used antonyms to replace negated edits where applicable. We do so by searching for antonyms of the original root word based on WordNet (Miller, 1995) and termed successful antonym edits as ‘Edit-Alt’. To ensure similar tense was used, we detected the original word’s tense and applied the same tense onto the antonym word using the Pattern package (De Smedt and Daelemans, 2012). An example ‘Edit’ and ‘Edit-Alt’ sentence is shown in Table 1.

To select between ‘Edit’ and ‘Edit-Alt’ versions, we calculated the Levenshtein edit distance of the original word versus the antonym. We select ‘Edit-Alt’ only if the edit distance is less than or equal to 30% of the length of the longer word, rounded to the nearest integer. This allows us to keep conversions like ‘able \rightarrow unable’ for more natural word flow, but discard bolder and more drastic changes like ‘safe \rightarrow dangerous’ and ‘had \rightarrow refused’ that

⁴Our edit schemes, model pipeline, datasets and supplementary materials can be found on Github at <https://xxx.xxx.xx> (Also uploaded under Software)

⁵We used NLTK (Wagner, 2010) to obtain POS tags in PennTreeBank format and spaCy (Honnibal et al., 2020) for dependency tree extraction.

Conversion	Edit Type	Sentence
Negation	Original	TyG is effective to identify individuals at risk for NAFLD.
	Regular (Edit)	TyG is not effective to identify individuals at risk for NAFLD.
	Regular (Edit-Alt)	TyG is ineffective to identify individuals at risk for NAFLD.
	Shorten	TyG is ineffective
	Multiples	is ineffective is ineffective is ineffective
Strengthen	Original	Moreover, TT genotype may reduce the risk of CAD in diabetic patients.
	Regular (Edit)	Moreover, TT genotype will reduce the risk of CAD in diabetic patients.

Table 1: Examples of counterfactual causal sentence augments. *Notes.* Interventions are highlighted in green. Causal Strengthening can also have Shorten and Multiples edits but is excluded due to space constrains.

were either suggesting causality in the opposite direction (rather than *no relationship*) or outright wrong. Finally, after dropping duplicates, we obtained 381 sentences that represent non-causality.

We were able to apply 11 out of the 12 linguistic templates to generate causal negation for the sentences in CSci. Most edits fall into the category where we negate the root verb or adjective of the sentence. Appendix Table A1 shows one randomly sampled example per available negation method when applied onto the CSci corpus. With respect to this table, Appendix Section A.1 briefly discusses the grammatical sanity of these sentences. We inspected these randomly sampled counterfactuals to verify that sentence flows are natural and desirable.

3.2.2 Causal Strengthening

We also increased the strength of causal statements from *conditional causal* (c_2) to *causal* (c_1) by exploiting modal words. Similar to Negation, we first obtain the POS tags and dependency trees for each sentence.

Algorithm 2 in the Appendix outlines the rule-based pseudo-code. In general, the 5 linguistic templates created converts modals based on the dictionary: $\{‘could’, ‘should’, ‘would’\} \rightarrow ‘would’$ and $\{‘can’, ‘may’, ‘might’, ‘will’\} \rightarrow ‘will’$. When modals interact with verbs with lemma *‘be’*, we replace *‘modal+be’* with *‘was’* instead to convey certainty in causal meaning. For special cases when modal terms interact with *‘have’* which forms conditional perfect tense, we convert them into simple past tense by replacing *‘modal+have’* with *‘had’*. When a modal is followed by an adverb (E.g. *“can possibly”*), the adverb is removed to avoid any deviation of the causal meaning from certainty.

Table A2 shows a randomly sampled example per causal strengthening method when applied onto the CSci corpus. Of the 213 available sentences, we successfully augmented 174 of them.

3.3 Dataset Processing

Duplicates exists in the original CSci corpus and arise when we append the edits with the original sentences. De-duplication based on priority rules discussed in Appendix Section A.2 was applied.

Our augmentations would increase the sample size for particular class labels. To combat this, we randomly selected sentences such that the original class distribution is maintained. Our main analysis focuses on randomly sampled datasets to eliminate the concern that the improved performance might result from increased data size or advantageous train set distribution.⁶ However, note that the final dataset size is always slightly smaller than the original baseline due to the de-duplication step. After random sampling, the distribution thus slightly differs. The final sample counts across class labels per augmented dataset are summarized in Appendix Table A5.

3.4 Further Heuristics

Later in results Section 4.4.1, we observe that simple edits which highlight the main counterfactual phrase to the model helps improve performance. Although these heuristics result in non-grammatical sentences, we believe these edits explicitly emphasize augmented keywords for the model to learn the local syntactic changes better. Since we still train the model with the original sentences (in fact, the majority), the model will not memorise on only non-grammatical examples.

An example sentence is detailed in Table 1 with the two augmentation variations as follows:

- **Shorten:** We reduce the sentence length based on target/root word to cover a minimally interpretable phrase based on depen-

⁶We want to show that any improvements in our scores are due to increased variations of examples per class label. These variations must be meaningful for any improvement in scores.

311 dency parser. The final sentence might not be
312 a consecutive slice from the original.

- 313 • **Multiples:** We define a phrase as one
314 word before and after the target/root word.
315 That is, we define $PhraseLength = 3$.
316 Phrases are then duplicated by a multiple of
317 $OriginalSentenceLength/PhraseLength$
318 rounded to the nearest integer. This ensures
319 that the final sentence is up to as long as the
320 original length. Note that in the ‘*edit-alt*’
321 example of Table 1, ‘is ineffective’ represents
322 ‘is not effective’. Thus, although the actual
323 phrase length is 2, the intended meaning
324 is based off the latter phrase that had a
325 length of 3. Hence, we maintained a fixed
326 $PhraseLength$ for all sentences.

327 3.5 Out-of-domain Testing

328 We train our models on the CSci corpus and con-
329 duct testing on SCITE (Li et al., 2021)⁷ and AltLex
330 (Hidey and McKeown, 2016)⁸ corpora to show
331 that inclusion of meaning counterfactuals during
332 model training aids in OOD applications. While the
333 CSci corpus is constructed from scientific-based
334 PubMed sentences, the SCITE corpus contains gen-
335 eral sentences extended from the SemEval 2010
336 task 8 dataset. AltLex consists of sentences from
337 English Wikipedia. AltLex was built for causal
338 relation identification, and therefore, has multiple
339 entries per sentence based on different entities and
340 relations. We revised the format of the corpus such
341 that if a sentence has any one causal relation, the
342 sentence is considered causal. Additionally, be-
343 cause SCITE and AltLex labels are binary, we cre-
344 ated two measures of accuracy. The first, ‘Acc’,
345 considers only exact class labels (*no relationship*
346 (c_0) and *causal* (c_1)) (i.e. predicting other labels
347 are considered wrong). The second, ‘Acc_{Group}’,
348 calculates accuracy after grouping [*no relationship*,
349 *correlational*] into *no relationship* (c_0) and [*causal*,
350 *conditional causal*] into *causal* (c_1) to align with
351 the binary labels. In total, we test on 4439 sen-
352 tences from SCITE and 37677 sentences from Al-
353 tLex.

354 3.6 Modeling

355 In each setting, we train and validate using K=5
356 folds, with 5 epochs per fold. For loss, we use the

standard cross-entropy loss for multi-class classi-
357 fication. For OOD testing, we take the majority
358 prediction from the five trained models of the five
359 folds. We explore the results with the following
360 two models: 361

362 3.6.1 BERT+MLP (MLP)

363 We replicate the best performing model on the CSci
364 corpus (Yu et al., 2019) which is a BioBERT (Lee
365 et al., 2020) plus multi-layer perceptron (MLP)
366 pipeline. The default architecture was: BioBERT
367 embeddings were fed into a single MLP layer that
368 served as the classifier.

369 3.6.2 BERT+MLP+SVM (SVM)

370 Instead of applying LinearSVM based off unigrams
371 and bigrams like the original authors (Yu et al.,
372 2019), we believe a fairer comparison would be
373 to use BERT embeddings as inputs into an SVM
374 model. To allow for representation updates, for
375 each sentence (s), the BioBERT encoder is applied.
376 The BERT output (z) runs through two MLP layers
377 (MLP_1 and MLP_2) to predict class labels. The
378 second layer is ultimately is discarded, and we take
379 the hidden representation (r) as fixed inputs into
380 the SVM classifier after all epochs. The equations
381 below outlines our pipeline,

$$382 z = BERT(s), \quad z \in \mathbb{R}^{h_1} \quad (1)$$

$$383 r = MLP_1(z), \quad r \in \mathbb{R}^{h_2} \quad (2)$$

$$384 o = MLP_2(r), \quad o \in \mathbb{R}^c \quad (3)$$

$$385 p = SVM(r), \quad p \in \mathbb{R}^1, \quad (4)$$

386 where, p represents the final predicted label, and
387 $h_1 = 768$, $h_2 = 24$, and $c = 4$.

388 4 Results & Discussion

389 4.1 Baseline

390 Table 2 reports our performance on the CSci cor-
391 pus. With the MLP baseline model, we were un-
392 able to replicate the reported scores by Yu et al.
393 (2019) of 90.1% accuracy and 88.1% macro F-
394 score. We achieved slightly lower scores of 89.15%
395 and 87.01% respectively. For SVM, our proposed
396 implementation using updated BERT embeddings
397 with a detached head is superior over Yu et al.
398 (2019)’s unigram and bigrams method as we ob-
399 serve significant improvements of accuracy from
400 77.2% to 88.86% and macro F-score from 72.2%
401 to 86.95%.

402 The inclusion of a mixture of edits (Negation*Shorten with Strengthen*Regular) returns the
403

⁷<https://github.com/Das-Boot/scite>

⁸<https://github.com/chridey/AltLex>

Conversion	Edit Type	MLP				SVM			
		F1	Acc	F1 _{Orig}	Acc _{Orig}	F1	Acc	F1 _{Orig}	Acc _{Orig}
Yu et al. (2019)		88.10	90.10	88.10	90.10	72.20	77.20	72.20	77.20
Ours (Base)		87.01	89.15	87.01	89.15	86.95	88.86	86.95	88.86
Negation	Regular	-1.55	-1.92	-0.19	-0.95	-2.33	-1.99	-1.18	-1.28
Negation	Shorten	+1.06	+0.89	+0.57	-0.04	+0.95	+1.19	+0.38	+0.18
Negation	Multiples	+1.46	+1.45	+0.93	+0.49	+1.14	+1.28	+0.60	+0.32
Strengthen	Regular	+1.75	+1.14	+0.80	+0.84	+0.73	+0.49	-0.28	+0.20
Strengthen	Shorten	+1.08	+0.91	+0.16	+0.62	+0.86	+1.08	-0.24	+0.71
Strengthen	Multiples	+0.98	+0.98	-0.05	+0.57	+0.62	+0.82	-0.50	+0.38
Both	Shorten, Regular	+2.80	+2.33	+1.73	+1.35	+1.45	+1.38	+0.14	+0.19

Table 2: Performance metrics on CSci corpus. *Notes.* BioBERT models trained on variations of CSci corpus (Original plus edits), with edits matching existing labels and randomly sampled to match base class distribution. Results are for test set when trained and predicted over 5-folds. Macro F-score (F1) and accuracy (Acc) are in percentages. Columns with lower subscript “Orig” are calculated for original sentences only (i.e. performance for edits is ignored). Rows below “Ours (Base)” report relative changes to it. The best performance per column is **bolded**. Precision and Recall scores are available in Appendix Tables A7 and A8.

best performance across all metrics: Accuracy improves by 1.35% our MLP baseline, achieving Acc_{Orig} of 90.60%⁹. Notice that we find improvements of accuracy and F-score beyond the original reported scores, even though our replicated scores were lower. The SVM model also demonstrates that the inclusion of edits improves performance.

Intuitively, we are exposing the model to more sentence types of the real world. We are also specifically choosing sentences near the boundaries of the labels (i.e. with minor edits, sentences’ labels can change). Therefore, the model is able to learn better in the CSC task. Interestingly, we noticed that shorten or multiples edits improved performance for negated edits, seemingly more than regular edits itself. Section 4.4.1 expands on this finding.

4.2 Robustness on Edits

Table 3 highlights how current SOTA models are not robust to minimally altered sentences that changes in causal direction and strength.

To conduct the experiment, we randomly split the available negated edits (n=381) by half, keeping 191 negated sentences for training and the remaining 190 for testing. The 190 original sentences that corresponds to the negated test set were removed from the original CSci corpus to avoid exposing models to highly similar sentences during training¹⁰. Models trained with this base train set danger-

⁹The full original set achieved 90.33% accuracy if we were to include the subset that is dropped out due to random sampling. We predict the labels for this dropped-out subset like an OOD dataset, i.e. taken across 5-folds after training completes.

¹⁰In experiments not shown, the models trained on the full

ously predicted 157 out of 190 test sentences in the opposite direction as *causal* instead of *no relationship*. A shockingly dismal test accuracy of 12.63% was attained at best, and prediction counts are available in Appendix Table A6.

Our finding surfaces the problem that the models are likely memorizing on key causal terms instead of understanding sentence structure and flow. Therefore, they were unable to discern the negation involved. Inclusion of counterfactual examples help to fill this representation gap. We created augmented sets by combining the base train set with the 191 negated train sentences for retraining. Once we exposed the models to these negated examples during training, the same models could predict the right label with up to 73.68% accuracy.

We also tested the models’ efficacy on strengthened sentences converted from *conditional causal* to *causal*. Once counterfactual examples were included in the train set, improvements on test accuracy was obtained to a significant, but smaller, extent of +13.79% improvement at best.

4.3 Improving Generalization

In Table 4, we show that inclusion of edits during training also helps to improve generalization in cross-domain applications. Although our train dataset is an academic and scientific-based text represented by a BioBERT language model, we

original CSci corpus almost certainly wrongly predicts the 190 negated sentences as *causal*. To prove our point that models are memorizing causal terms, we removed the overlapping sentences to eliminate the possibility of the models memorizing similar sentences in train and test set instead.

Conversion	n	MLP	SVM
Original	190	12.63	10.53
Negation	190	+61.05	+62.63
Original	87	77.01	73.56
Strengthen	87	+11.49	+13.79

Table 3: Accuracy (in percentage) of BioBERT models trained on a subset of CSci corpus and predicted on a fully augmented difference set. *Notes.* The best performance per section per column is **bolded**.

show that when we apply the same model to the general-based SCITE and Wikipedia-based AltLex corpora, inclusion of edits improved classification performance. We were unable to find improvements in generalisation for MLP model on SCITE dataset, which could be due to our limited edit schemes. However, for AltLex, there are consistent improvements for almost all types of edits across both models. The mixture of edits (Negation*Shorten with Strengthen*Regular) again reports the best generalisation outcomes by showing improvements in accuracy (up to +0.94%) across all models and datasets, except a negligibly small reduction (-0.02%) for MLP model when tested on SCITE.

4.4 Ablations

4.4.1 Need for Heuristical Edits

Earlier in Table 2, we noted that models exposed to Negation*Regular edits are unable to learn the boundaries effectively: Acc_{Orig} fell by 0.95% and 1.28% for the MLP and SVM models respectively from our baselines. However, when we perform simple heuristics like Shorten, accuracy could improve to -0.04% (negligible reduction) and +0.18% for MLP and SVM respectively.

We study the net change in classification counts per model per label in Table 5 to explore this phenomenon. Given class labels i and j predicted by a model and our baseline respectively, we report the model’s $NetChange_i = Right_i - Wrong_i = \sum_{j \neq i} n_{(i=true)j} - \sum_{i \neq j} n_{i(j=true)}$, where $i, j = c_0, c_1, c_2, c_3$ and n refers to the number of observations. $Right_i$ ($Wrong_i$) is the number of observations where a model predicts correctly (wrongly) for class label i but baseline predicts wrongly (correctly). When either MLP or SVM model is trained with the augmented Negation*Regular dataset, the model becomes confused and predicts poorly for *causal* (c_0) and *no relationship* (c_0) classes. Once

the edits were presented in the shortened form, this situation improves. This short exploration points us to believe exposing sample-curated features is needed in our low-resource setting. Highlighting the model to the short spans of (non-)causality helps point out the exact borders we want the models to become sensitive to.

Interestingly, we observe improvements in classification for labels we did not edit (c_3) in the majority of settings. This highlights the possibility that exposing models to minimally perturbed sentences around label boundaries could improve comprehension beyond the introduced edits.

4.4.2 Capturing Causal Strengths

By capitalizing on CSci’s labels, our methodology allows us to expose causal strengths in SCITE and AltLex corpora beyond the original binary labels. For SCITE, the baseline MLP model originally labeled five sentences as *conditional causal*. Applying the model trained with Strengthen*Regular edits, we observed that four remained as *conditional causal* (c_2) while one of the sentence¹¹ correctly switched label to *causal* (c_1). For the baseline SVM model, seven sentences were tagged as c_2 , of which four remained, and the same one as MLP’s converted to c_1 . One¹² correctly switched to *no relationship* (c_0) as labeled, while the last sentence¹³ converted to *correlational* (c_3), which is surprising because we did not edit any sentences to or from class c_3 . Unfortunately, the authors of SCITE tagged this sentence as causal, which means this is considered to be mislabeled. However, the sentence contains signals like ‘corresponds to’, which should be correlational, not causal.

We believe the model might be picking up on what makes something *conditional causal* versus all other labels (not just comparing to the one we edit to). Our short qualitative analysis again supports the earlier quantitative study that exposing models to meaningfully augmented sentences across labels could improve classification even for other uninvolved labels.

¹¹“In the present recession, which has been triggered by a collapse in land prices, land-value taxation would reverse the collapse - not by re-inflating a temporary speculative bubble, but by inducing investment in infrastructure that permanently enhances the utility of the land.”

¹²“The glass tealight holder appears to float inside the metal spiral as it spins in the gentle breeze.”

¹³“The increase of the signal might correspond to formation of the high-density excitons, while the reduction of the signal originates from the relaxation.”

Conversion	Edit Type	SCITE				AltLex			
		MLP		SVM		MLP		SVM	
		Acc	Acc _{Group}	Acc	Acc _{Group}	Acc	Acc _{Group}	Acc	Acc _{Group}
Ours (Base)		86.28	85.83	85.04	84.50	85.57	84.64	85.91	84.68
Negation	Regular	-1.46	-1.67	-0.36	-0.41	-0.22	-0.44	+0.18	+0.41
Negation	Shorten	-0.20	-0.27	+0.02	+0.02	+0.61	+0.54	+0.74	+1.05
Negation	Multiples	-0.18	-0.16	-0.38	-0.38	+0.89	+0.95	+1.19	+1.58
Strengthen	Regular	-0.27	-0.14	+1.01	+1.10	+0.51	+0.69	+0.54	+0.84
Strengthen	Shorten	-3.40	-3.36	-0.11	-0.05	+0.30	+0.37	+0.99	+1.38
Strengthen	Multiples	-1.31	-1.28	-0.90	-0.90	+0.88	+0.99	+0.07	+0.29
Both	Shorten, Regular	-0.02	-0.05	+0.79	+0.63	+0.94	+0.84	+0.31	+0.41

Table 4: Performance metrics of on OOD datasets. *Notes.* BioBERT models trained on variations of CSci corpus (Original plus edits), with edits matching existing labels and randomly sampled to match base class distribution. Results are for out-of-domain SCITE and AltLex corpus taking mode class predicted over 5-folds. Accuracies are reported in percentages. Columns ‘Acc’ considers only exact class labels, while ‘Acc_{Group}’ calculates accuracy after converting the four labels to form the binary labels. Rows below “Ours (Base)” report relative changes to it. The best performance per column is **bolded**.

Conversion	Edit Type	MLP					SVM				
		c_0	c_1	c_2	c_3	Total	c_0	c_1	c_2	c_3	Total
Negation	Regular	-13	-18	+10	0	-21	-9	-21	+1	-2	-31
Negation	Shorten	-15	+9	+9	+2	+5	-4	+7	+1	+6	+10
Negation	Multiples	-5	+9	+5	+9	+18	-1	+8	+3	+3	+13
Strengthen	Regular	-7	+12	+10	+9	+24	+1	+1	+5	-4	+3
Strengthen	Shorten	-7	+11	+6	+6	+16	0	+8	+3	+5	+16
Strengthen	Multiples	-14	+13	+7	+10	+16	-12	+11	+6	+3	+8
Both	Shorten, Regular	+2	+20	+10	+9	+41	-2	+12	+1	-4	+7

Table 5: Net change in correct classification counts on CSci corpus compared to “Ours (Base)” for original. *Notes.* Recall that Negation is the conversion of $c_1 \rightarrow c_0$; Strengthen is the conversion of $c_2 \rightarrow c_1$;

4.4.3 Other Experiments

We also explored other popular methodologies that did not produce consistent and significant improvements from baseline. These include, *i*) creating more edit types (using masking, synonyms and paraphrasers), *ii*) extending to a five-way classification problem (by labelling negated edits as a new class label, different from *no relationship* (c_0)), and *iii*) experimenting on contrastive learning loss functions. Appendix Section A.3 details these experiments further for interested readers.

5 Conclusion and Future Work

We explored the task of CSC in a low-resource setting. Following recent literature, we generated counterfactual sentences via rule-based edits that change sentences’ causal direction and strength. We show that SOTA CSC models worryingly misclassifies on such augmented sentences. We demonstrate that inclusion of our edits during training can help to improve classification performance both on original and edit sentences, and within and outside

of the corpus’ domain. However, we find that simple edits, such as negation, might be insufficient to teach effective decision boundaries given limited data size. We thus propose heuristic edit schemes and find performance improvements across both training and OOD contexts too.

Moving forward, we plan to replicate our findings on more datasets to further demonstrate our augmentation scheme’s widespread applicability. Yu et al. (2020)’s recent corpus based on scientific press statements annotated with the four class labels of causality as per our set up is a promising option. Additionally, causality is hard as one has to distinguish between causal effects as factual events of real-world or at the level of “meta-causality” (Andersson et al., 2020). In our work, we did not go beyond the “correctness” of the claims. Grounding the claims to world knowledge is an important future work. Lastly, we wish to find alternative models which can learn directly from original plus regular edits without the need for heuristics.

583
584
585
586
587
588
589

590
591
592

593
594
595
596
597

598
599
600
601
602
603

604
605
606
607
608
609
610

611
612

613
614
615
616
617
618
619
620
621
622

623
624
625
626
627
628
629
630
631

632
633
634
635
636
637
638
639

References

Marta Andersson, Murathan Kurfalı, and Robert Östling. 2020. [A sentiment-annotated dataset of English causal connectives](#). In *Proceedings of the 14th Linguistic Annotation Workshop*, pages 24–33, Barcelona, Spain. Association for Computational Linguistics.

Nabiha Asghar. 2016. [Automatic extraction of causal relations from natural language texts: A comprehensive survey](#). *CoRR*, abs/1605.07895.

Susanne Buhse, Anne Christin Rahn, Merle Bock, and Ingrid Mühlhauser. 2018. Causal interpretation of correlational studies—analysis of medical news on the website of the official journal for german physicians. *PLoS One*, 13(5):e0196833.

Tommaso Caselli and Piek Vossen. 2017. [The event storyline corpus: A new benchmark for causal and temporal relation extraction](#). In *Proceedings of the Events and Stories in the News Workshop@ACL 2017, Vancouver, Canada, August 4, 2017*, pages 77–86. Association for Computational Linguistics.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. 2020. [A simple framework for contrastive learning of visual representations](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR.

Tom De Smedt and Walter Daelemans. 2012. Pattern for python. *J. Mach. Learn. Res.*, 13(1):2063–2067.

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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.

Kata Gábor, Davide Buscaldi, Anne-Kathrin Schumann, Behrang QasemiZadeh, Haifa Zargayouna, and Thierry Charnois. 2018. [Semeval-2018 task 7: Semantic relation extraction and classification in scientific papers](#). In *Proceedings of The 12th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2018, New Orleans, Louisiana, USA, June 5-6, 2018*, pages 679–688. Association for Computational Linguistics.

Matt Gardner, Yoav Artzi, Victoria Basmova, 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 *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings, EMNLP 2020, Online Event, 16-20 November 2020*, pages 1307–1323. Association for Computational Linguistics.

Roxana Girju and Dan I. Moldovan. 2002. [Text mining for causal relations](#). In *Proceedings of the Fifteenth International Florida Artificial Intelligence Research Society Conference, May 14-16, 2002, Pensacola Beach, Florida, USA*, pages 360–364. AAAI Press.

Stefan Heindorf, Yan Scholten, Henning Wachsmuth, Axel-Cyrille Ngonga Ngomo, and Martin Potthast. 2020. [Causenet: Towards a causality graph extracted from the web](#). In *CIKM ’20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*, pages 3023–3030. ACM.

Christopher Hidey and Kathy McKeown. 2016. [Identifying causal relations using parallel wikipedia articles](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics.

Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. [spaCy: Industrial-strength Natural Language Processing in Python](#).

David Jurgens, Saif Mohammad, Peter D. Turney, and Keith J. Holyoak. 2012. [Semeval-2012 task 2: Measuring degrees of relational similarity](#). In *Proceedings of the 6th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2012, Montréal, Canada, June 7-8, 2012*, pages 356–364. The Association for Computer Linguistics.

Divyansh Kaushik, Eduard H. Hovy, and Zachary Chase Lipton. 2020a. [Learning the difference that makes a difference with counterfactually-augmented data](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Divyansh Kaushik, Amrith Setlur, Eduard H. Hovy, and Zachary C. Lipton. 2020b. [Explaining the efficacy of counterfactually-augmented data](#). *CoRR*, abs/2010.02114.

Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020. [Supervised contrastive learning](#). *arXiv preprint arXiv:2004.11362*.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. [Biobert: a pre-trained biomedical language representation model for](#)

696	biomedical text mining. <i>Bioinform.</i> , 36(4):1234–1240.	754
697		755
698	Zhaoning Li, Qi Li, Xiaotian Zou, and Jiangtao Ren. 2021. Causality extraction based on self-attentive bilstm-crf with transferred embeddings. <i>Neurocomputing</i> , 423:207–219.	756
699		757
700		758
701		
702	Taylor Mahler, Willy Cheung, Micha Elsner, David King, Marie-Catherine de Marneffe, Cory Shain, Symon Stevens-Guille, and Michael White. 2017. Breaking NLP: Using morphosyntax, semantics, pragmatics and world knowledge to fool sentiment analysis systems. In <i>Proceedings of the First Workshop on Building Linguistically Generalizable NLP Systems</i> , pages 33–39, Copenhagen, Denmark. Association for Computational Linguistics.	759
703		760
704		761
705		762
706		
707		
708		
709		
710		
711	Rowan Hall Maudslay, Hila Gonen, Ryan Cotterell, and Simone Teufel. 2019. It’s all in the name: Mitigating gender bias with name-based counterfactual data substitution. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 5266–5274. Association for Computational Linguistics.	763
712		764
713		765
714		766
715		767
716		768
717		769
718		770
719		
720		
721	George A. Miller. 1995. Wordnet: A lexical database for english. <i>Commun. ACM</i> , 38(11):39–41.	771
722		772
723		773
724	Paramita Mirza, Rachele Sprugnoli, Sara Tonelli, and Manuela Speranza. 2014. Annotating causality in the TempEval-3 corpus. In <i>Proceedings of the EACL 2014 Workshop on Computational Approaches to Causality in Language (CAtoCL)</i> , pages 10–19, Gothenburg, Sweden. Association for Computational Linguistics.	774
725		775
726		776
727		
728		
729		
730	Paramita Mirza and Sara Tonelli. 2016. CATENA: causal and temporal relation extraction from natural language texts. In <i>COLING 2016, 26th International Conference on Computational Linguistics, Proceedings of the Conference: Technical Papers, December 11-16, 2016, Osaka, Japan</i> , pages 64–75. ACL.	777
731		778
732		779
733		780
734		781
735		782
736		783
737	John X. 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 <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2020 - Demos, Online, November 16-20, 2020</i> , pages 119–126. Association for Computational Linguistics.	784
738		785
739		786
740		787
741		788
742		789
743		790
744		791
745		
746	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library. In <i>Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada</i> , pages 8024–8035.	792
747		793
748		794
749		795
750		796
751		797
752		798
753		
	Balaji Polepalli Ramesh, Rashmi Prasad, Tim Miller, Brian Harrington, and Hong Yu. 2012. Automatic discourse connective detection in biomedical text. <i>J. Am. Medical Informatics Assoc.</i> , 19(5):800–808.	
	Petroc Sumner, Solveiga Vivian-Griffiths, Jacky Boivin, Andy Williams, Christos A Venetis, Aimée Davies, Jack Ogden, Leanne Whelan, Bethan Hughes, Bethan Dalton, Fred Boy, and Christopher D Chambers. 2014. The association between exaggeration in health related science news and academic press releases: retrospective observational study. <i>BMJ</i> , 349.	
	Wiebke Wagner. 2010. Steven bird, ewan klein and edward looper: Natural language processing with python, analyzing text with the natural language toolkit - o’reilly media, beijing, 2009, ISBN 978-0-596-51649-9. <i>Lang. Resour. Evaluation</i> , 44(4):421–424.	
	Jinghang Xu, Wanli Zuo, Shining Liang, and Xianglin Zuo. 2020. A review of dataset and labeling methods for causality extraction. In <i>Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020</i> , pages 1519–1531. International Committee on Computational Linguistics.	
	Bei Yu, Yingya Li, and Jun Wang. 2019. Detecting causal language use in science findings. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 4663–4673. Association for Computational Linguistics.	
	Bei Yu, Jun Wang, Lu Guo, and Yingya Li. 2020. Measuring correlation-to-causation exaggeration in press releases. In <i>Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020</i> , pages 4860–4872. International Committee on Computational Linguistics.	

A Appendix

A.1 Negation Examples

Appendix Table A1 shows one randomly sampled example per available negation method when applied onto the CSci corpus. As shown, most examples fall into ‘VB_3.1’, ‘VB_5.1’, ‘JJ_1.3’ and ‘VB_1.2’ types, for which the templates in Algorithm 1 work well for¹⁴. For rarer method types, like ‘VB_2.1’, the templates seem to work poorly. Further investigation shows that the error arose from the POS tagging step: “Both” was tagged as a VB but should have been a DT or CC, for which, we have no template for at the moment, so the example would have been correctly skipped. As for ‘VB_4.1’, the negated example is unnatural but not grammatically wrong.

A.2 De-duplication

After appending original sentences with edits, we conduct de-duplication. Appendix Table A3 shows problematic duplicates that had differing labels. The original CSci corpus contained 7 duplicate sentences instances which were removed. 6 of them were exact duplicates (same label, same sentence), while the last 1 (sentence S/N 1) was duplicated with different labels (c_0 and c_2). We manually changed this to keep only the c_0 label. The total data size thus reduces from $n=3061$ to $n=3054$. We also take this chance to highlight concerns that some sentences in CSci were labelled contrary to how we understood them.

Subsequent duplicates were handled via rule-based removal. The motivation was to ensure identical sentences do not have different labels which adds noise to our training. Our assumption is that if an edit was performed but remained identical to the original, the original must have been mislabelled sentence. We note that our rule-based de-duplication cannot accommodate multi-label cases, as there was one sentence (S/N 4) that correctly reflected both c_0 and c_1 labels in different parts of the sentence, but due to de-duplication, we only kept the c_0 label.

A.3 Other Experiments

Other experiments conducted but did not produce significant improvements are mentioned here.

¹⁴We highlight the main POS tags used and mentioned: VB (verbs, e.g. ‘eating’), JJ (adjective, e.g. ‘big’), IN (preposition or subordinating conjunction, e.g. ‘by’), DT (determiner, e.g. ‘he’), CC (coordinating conjunction, e.g. ‘and’), MD (modal, ‘should’).

Other Edit Types Three were explored:

- **Mask:** Based on POS, all nouns are replaced by the token “[MASK]”.
- **Synonyms:** Using WordNet synonyms, we skip common words¹⁵ and randomly replace up to 5 words. Synonyms match tense and plurality of original words using Pattern package, which is imperfect.
- **T5Para:** We run the sentence through a pre-trained T5-paraphraser model¹⁶.

Appendix Table A4 shows an example sentence with the above edits for the same causal sentence of Table 1. With the SVM model, only Strengthen*Synonyms appended with original increased accuracy on CSci by 1.01% while Strengthen*T5Para increased accuracy by 0.39%. However, these findings could not be replicated across to the MLP model nor for Negation.

Extending to a Five-way Classification In our main set up, we focused on edits that matched the original labels and are randomly sampled such that the unified train set matches base class distribution for fairer comparison to baseline. Current negations are labelled *no relationship* (c_0). However, to the extent that we believe negated causal statements deserve a class of their own, we also explore the event when negations are labelled with a new level *not causal* (c_4) instead. Based on the set up for Table 3, we obtained even higher improvements in accuracy of +70.53% and +74.74% for the MLP and SVM model respectively. This could be due to the clearer distinction of a *not causal* sentence structure compared to if we were to combine them with other *no relationship* statements. When we extended the MLP and SVM model to work with such a five-way classification set up, we did observe improvements in Acc_{Orig} for shorten, multiples and synonyms versions of edits. However, because we cannot truly balance the dataset (random sampling does not apply here because we have a whole new class), we cannot be certain if the improvements were due to the larger dataset or the model picking up on the boundaries. Furthermore, the improvements did not generalize on our OOD set ups.

¹⁵We do not try to find synonyms for common words with these POS types: ‘DT’, ‘IN’, ‘EX’, ‘CC’, ‘MD’, ‘WP’, ‘WD’, ‘WR’, ‘UH’, ‘RP’, ‘SY’, ‘PO’

¹⁶https://huggingface.co/ramsrigouthamg/t5_paraphraser

888 **Other Training Setups** In addition to standard
889 cross-entropy based supervised learning, we also
890 explored contrastive learning schemes. In particu-
891 lar, we trained with Supervised Contrastive Loss
892 (SupCon) (Khosla et al., 2020; Chen et al., 2020)
893 and Triplet Margin Loss (Paszke et al., 2019). In
894 the contrastive setup, we introduced counterfac-
895 tuals as the negative examples for each anchor
896 sentence. For positive samples, we used shorten,
897 synonyms, and T5Para augmentation strategies on
898 the original anchor sentence. However, the re-
899 sults did not provide performance improvements
900 in either CSci or OOD datasets, which highlights
901 the challenge of building a generalized scheme of
902 counterfactual generations. Exploring avenues in
903 contrastive-learning remains a critical future work.

904 **A.4 Reproducibility Checklist**

905 We include additional details about our main exper-
906 iment not highlighted in other parts of the paper.

- 907 • **Computing Infrastructure:** Tesla V100
908 SXM2 32 GB
- 909 • **MLP Hyperparameters:** “atten-
910 tion_probs_dropout_prob”: 0.1, “hidden_act”:
911 “gelu”, “hidden_dropout_prob”: 0.1, “hid-
912 den_size”: 768, “initializer_range”: 0.02,
913 “intermediate_size”: 3072, “layer_norm_eps”:
914 1e-12, “max_position_embeddings”:
915 512, “num_attention_heads”: 12,
916 “num_hidden_layers”: 12, “type_vocab_size”:
917 2, “vocab_size”: 28996
- 918 • **SVM Hyperparameters:** kernel: “linear”,
919 “C”: 1e-2
- 920 • **Average Runtime:** For 5 epochs and 5 folds,
921 our baseline MLP model takes approximately
922 22 minutes 51 seconds to train and validate.

923 **A.5 Additional figures and tables**

Algorithm 1: NegationRules – Causal negation scheme

Input: *edit_id, text_ids, text, pos, sentid2tid, max_try=2, curr_try=0***Output:** *text, method, edit_id*

```
1 curr_try ← curr_try + 1
2 curr_pos, curr_word ← pos[edit_id], text[edit_id]
3 prev_pos, prev_word ← pos[edit_id - 1], text[edit_id - 1] if valid else None
4 next_pos, next_word ← pos[edit_id + 1], text[edit_id + 1] if valid else None
5 while curr_try ≤ max_try do
6   if curr_pos = VB then
7     if curr_word = AuxilliaryType then
8       if edit_id = max(text_ids) then
9         | Insert *not* in front of curr_word // Method 'VB_1.1'
10      else if next_word = DeterminerType then
11        | Replace next_word with *no* // Method 'VB_1.2'
12        | edit_id ← edit_id + 1
13      else if next_word = NounType then
14        | Insert *not* behind of curr_word // Method 'VB_1.3'
15      else if next_pos = VB then
16        | Insert *no* behind of curr_word // Method 'VB_1.4'
17      else if edit_id = min(text_ids) then
18        | Replace curr_word with *Not* + lowercased curr_word // Method 'VB_2.1'
19      else if prev_word = NounType then
20        | Replace curr_word with *did not* + lemma(curr_word) // Method 'VB_3.1'
21      else if edit_id = max(text_ids) then
22        | Insert *not* in front of curr_word // Method 'VB_4.1'
23      else if prev_word = AuxilliaryType next_pos = IN|TO then
24        | Insert *not* in front of curr_word // Method 'VB_5.1'
25    else if curr_pos = NN then
26      | Get head_id of head word of curr_word based on dependency tree
27      | text, method, edit_id ← NegationRules(head_id, text_ids, text, pos, sentid2tid,
28      | curr_try)
29    else if curr_pos = JJ then
30      if edit_id = max(text_ids) then
31        | Insert *not* in front of curr_word // Method 'JJ_1.1'
32      else if next_word = PositiveConjunctionType then
33        | Insert *not* in front of curr_word // Method 'JJ_1.2'
34        | Replace next_word with *nor* else
35        | Insert *not* in front of curr_word // Method 'JJ_1.3'
36    else if curr_pos = IN then
37      | Insert *not* in front of curr_word // Method 'IN_1.1'
36 Define method as method name if applicable edit occurs
37 return text, method, edit_id
```

Method	Regular (Edit)	Regular (Edit-Alt)	n
VB_1.2	Eyes with better vision at baseline had no more favorable prognosis, whereas eyes with initial macular detachment, intraoperative iatrogenic break, or heavy SO showed more unfavorable outcomes.	Eyes with better vision at baseline abstained a more favorable prognosis, whereas eyes with initial macular detachment, intraoperative iatrogenic break, or heavy SO showed more unfavorable outcomes.	35
VB_1.3	Age, female sex, BMI, non-HDL cholesterol, and polyps are not independent determinants for gallstone formation.	Age, female sex, BMI, non-HDL cholesterol, and polyps differ independent determinants for gallstone formation.	12
VB_1.4	Both general and central adiposity have no causal effects on CHD and type 2 diabetes mellitus.	Both general and central adiposity refuse causal effects on CHD and type 2 diabetes mellitus.	2
VB_2.1	Not "both a low-fat vegan diet and a diet based on ADA guidelines improved glycemic and lipid control in type 2 diabetic patients."	-	1
VB_3.1	Collectively, these findings did not indicate that energy-matched high intensity and moderate intensity exercise are effective at decreasing IHL and NAFLD risk that is not contingent upon reductions in abdominal adiposity or body mass.	Collectively, these findings contraindicate that energy-matched high intensity and moderate intensity exercise are effective at decreasing IHL and NAFLD risk that is not contingent upon reductions in abdominal adiposity or body mass.	174
VB_4.1	The benefits of exercise for reducing risk of chronic disease, including CVD, are well not known.	-	1
VB_5.1	A higher BMI and a greater prevalence of comorbidities had not driven patients to seek a more radical solution for their obesity, i.e., surgery.	A higher BMI and a greater prevalence of comorbidities had attract patients to seek a more radical solution for their obesity, i.e., surgery.	81
JJ_1.1	The effects of TRT on cardiovascular risk markers were not ambiguous.	-	6
JJ_1.2	Results are not encouraging nor demonstrate that exercise was popular and conveyed benefit to participants.	Results are discouraging and disprove that exercise was popular and conveyed benefit to participants.	15
JJ_1.3	While LSG weakens the LES immediately, it does not predictably not affect postoperative GERD symptoms; therefore, distensibility is not the only factor affecting development of postoperative GERD, confirming the multifactorial nature of post-LSG GERD.	While LSG weakens the LES immediately, it does not predictably impede postoperative GERD symptoms; therefore, distensibility is not the only factor affecting development of postoperative GERD, confirming the multifactorial nature of post-LSG GERD.	53
IN_1.1	Although further investigation of long-term and prospective studies is not needed, we identified four variables as predisposing factors for higher major amputation in diabetic patients through meta-analysis.	-	1

Table A1: Example negated causal sentences per method *Notes*. “Method” refers to negation method label as per Algorithm 1. “Regular (Edit)” refers to direct negation from this Algorithm. “Regular (Edit-Alt)” refers to alternate intervention using same negation location, but based off antonyms from WordNet, if available. Interventions, excluding lemmatisation or case-changes, are highlighted in green. “n” is the number of successful conversions applicable in CSci corpus.

Algorithm 2: StrengthenRules – Causal strengthening scheme

Input: *edit_id, text_ids, text, pos, sentid2tid, curr_try=0***Output:** *text, method, edit_id*

```
1 Initialise ModalDict
2 curr_try ← curr_try + 1
3 curr_pos, curr_word ← pos[edit_id], text[edit_id]
4 next_pos, next_word ← pos[edit_id + 1], text[edit_id + 1] if valid else None
5 nnext_pos, nnext_word ← pos[edit_id + 2], text[edit_id + 2] if valid else None
6 while curr_try ≤ max_try do
7   if lemma(next_word) = ‘be’ then
8     Replace curr_word with *was* // Method ‘MOD_1.2’
9     Replace next_word with empty string
10  else if lemma(next_word) = ‘have’ then
11    if lemma(nnext_word) = ‘be’ then
12      Replace curr_word with *was* // Method ‘MOD_3.2’
13      Replace next_word and nnext_word with empty string
14    else
15      Replace curr_word with *had* // Method ‘MOD_3.1’
16      Replace next_word with empty string
17  else if curr_pos = MD & next_pos = RB then
18    Replace curr_word with ModalDict[curr_word] // Method ‘MOD_4.1’
19    Replace next_word with empty string
20  else
21    Replace curr_word with ModalDict[curr_word] // Method ‘MOD_1.1’
22 Define method as method name if applicable edit occurs
23 return text, method, edit_id
```

Method	Regular (Edit)	n
MOD_1.1	Physical therapy in conjunction with nutritional therapy may will help prevent weakness in HSCT recipients.	98
MOD_2.1	The rs7044343 polymorphism could be was involved in regulating the production of IL-33.	42
MOD_3.1	Increased titers of cows milk antibody before anti-TG2A and celiac disease indicates that subjects with celiac disease might have had increased intestinal permeability in early life.	21
MOD_4.1	Physical rehabilitation aimed at improving exercise tolerance can possibly will improve the long-term prognosis after operations for lung cancer.	13

Table A2: Example strengthened conditional causal sentences per method. *Notes.* “Method” refers to strengthening method label as per Algorithm 2, resulting in augments as per “Regular (Edit)”. Interventions, excluding lemmatisation or case-changes, are highlighted in green. Words removed from original version are striked out and highlighted in red. “n” is the number of successful conversions applicable in CSci corpus.

S/N	Sentence	Label				Conversion
		c_0	c_1	c_2	c_3	
1	None the less, both artificially sweetened beverages and fruit juice were unlikely to be healthy alternatives to sugar sweetened beverages for the prevention of type 2 diabetes.	1		1		Original
2	There was no effect on lumen volume, fibro-fatty and necrotic tissue volumes.	1	1			Negation
3	There are no indications that endogenous and exogenous gonadal hormones affect the radiation dose-response relationship.	1	1			Negation
4	In two randomized trials comparing the PCSK9 inhibitor bococizumab with placebo, bococizumab had no benefit with respect to major adverse cardiovascular events in the trial involving lower-risk patients but did have a significant benefit in the trial involving higher-risk patients.	1	1			Negation
5	Altering margin policies to follow either SSO-ASTRO or ABS guidelines would result in a modest reduction in the national re-excision rate.		1	1		Strengthen
6	Adding an allowance for accumulation of thyroidal iodine stores would produce an EAR of 72 μ g and a recommended dietary allowance of 80 μ g.		1	1		Strengthen
7	" In a randomized controlled trial of 230 infants with genetic risk factors for celiac disease, we did not find evidence that weaning to a diet of extensively hydrolyzed formula compared with cows milk-based formula would decrease the risk for celiac disease later in life.		1	1		Strengthen

Table A3: Sentences that had duplicates with differing labels. *Notes.* Rule-based de-duplication was performed, with the final label kept highlighted in green. “Conversion” refers to the augmented edit dataset that when we merge with the original, the duplicate appears. Do note that Sentence S/N 7, to us, should be labelled as *no relationship* (c_0), but was labelled as *conditional causal* (c_2) by original authors.

Conversion	Edit Type	Sentence
	Original	TyG is effective to identify individuals at risk for NAFLD.
	Regular (Edit)	TyG is not effective to identify individuals at risk for NAFLD.
	Regular (Edit-Alt)	TyG is ineffective to identify individuals at risk for NAFLD.
Negation	Shorten	TyG is ineffective.
	Multiples	is ineffective is ineffective is ineffective
	Mask	[MASK] is ineffective to identify [MASK] at [MASK] for [MASK]
	Synonyms	TyG exists inefficient to describe someone at take chances for NAFLD.
	T5Paraphraser	Ineffective for identifying individuals at risk for NAFLD.

Table A4: Extended examples of counterfactual causal sentence augments *Notes.* Interventions are highlighted in green.

Conversion	Edit Type	n_{c_0}	n_{c_1}	n_{c_2}	n_{c_3}	n
	Original (Yu et al., 2019)	1356	494	213	998	3061
Negation	Regular	1356	491	212	995	3054
Negation	Shorten	1356	491	212	995	3054
Negation	Multiples	1356	491	212	995	3054
Strengthen	Regular	1353	494	209	995	3051
Strengthen	Shorten	1353	494	209	995	3051
Strengthen	Multiples	1353	494	209	995	3051
Both	Shorten, Regular	1356	494	209	995	3054

Table A5: Number of sentences per class label after appending edits with base corpus, de-duplication and random sampling.

Conversion	True Label	c_0	c_1	c_2	c_3	Total
Negation	c_0	24	157	5	4	190
Strengthen	c_1	3	67	16	1	87

Table A6: Number of sentences predicted per class label for augmented dataset when trained on only original CSci corpus. *Notes.* Counts correspond to accuracy scores reported in Rows 1 and 3 of Table 3.

Conversion	Edit Type	P	R	F1	Acc	P_{Orig}	R_{Orig}	$F1_{Orig}$	Acc_{Orig}
Yu et al. (2019)		87.80	88.60	88.10	90.10	87.80	88.60	88.10	90.10
Ours (Base)		86.02	88.13	87.01	89.15	86.02	88.13	87.01	89.15
Negation	Regular	-1.81	-1.20	-1.55	-1.92	+0.29	-0.71	-0.19	-0.95
Negation	Shorten	+0.76	+1.45	+1.06	+0.89	+0.46	+0.78	+0.57	-0.04
Negation	Multiples	+1.47	+1.44	+1.46	+1.45	+1.05	+0.81	+0.93	+0.49
Strengthen	Regular	+1.96	+1.51	+1.75	+1.14	+0.98	+0.58	+0.80	+0.84
Strengthen	Shorten	+1.54	+0.54	+1.08	+0.91	+0.52	-0.29	+0.16	+0.62
Strengthen	Multiples	+1.51	+0.38	+0.98	+0.98	+0.53	-0.70	-0.05	+0.57
Both	Shorten, Regular	+2.98	+2.57	+2.80	+2.33	+1.90	+1.54	+1.73	+1.35

Table A7: Performance metrics of BERT+MLP on CSci corpus. *Notes.* BioBERT models trained on variations of CSci corpus (Original plus edits), with edits matching existing labels and randomly sampled to match base class distribution. Results are for test set when trained and predicted over 5-folds. Precision (P), Recall (R), macro F-score (F1) and accuracy (Acc) are reported in percentages. Columns with lower script "Orig" are calculated for base items only (i.e. performance for edits is ignored). Rows below "Ours (Base)" report relative changes to it. The best performance per column is **bolded**.

Conversion	Edit Type	P	R	F1	Acc	P_{Orig}	R_{Orig}	$F1_{Orig}$	Acc_{Orig}
Yu et al. (2019)		73.90	71.10	72.20	77.20	73.90	71.10	72.20	77.20
Ours (Base)		86.28	87.70	86.95	88.86	86.28	87.70	86.95	88.86
Negation	Regular	-2.72	-1.85	-2.33	-1.99	-0.89	-1.44	-1.18	-1.28
Negation	Shorten	+0.60	+1.36	+0.95	+1.19	+0.16	+0.67	+0.38	+0.18
Negation	Multiples	+1.18	+1.12	+1.14	+1.28	+0.68	+0.53	+0.60	+0.32
Strengthen	Regular	+0.97	+0.44	+0.73	+0.49	-0.14	-0.46	-0.28	+0.20
Strengthen	Shorten	+1.19	+0.54	+0.86	+1.08	+0.17	-0.65	-0.24	+0.71
Strengthen	Multiples	+0.92	+0.26	+0.62	+0.82	-0.21	-0.84	-0.50	+0.38
Both	Shorten, Regular	+1.25	+1.69	+1.45	+1.38	+0.00	+0.32	+0.14	+0.19

Table A8: Performance metrics of BERT+MLP+SVM on CSci corpus. *Notes.* Yu et al.'s SVM method does not use BERT inputs. Our BioBERT models are trained on variations of CSci corpus (Original plus edits), with edits matching existing labels and randomly sampled to match base class distribution. Results are for test set when trained and predicted over 5-folds. Precision (P), Recall (R), macro F-score (F1) and accuracy (Acc) are reported in percentages. Columns with lower script "Orig" are calculated for base items only (i.e. performance for edits is ignored). Rows below "Ours (Base)" report relative changes to it. The best performance per column is **bolded**.