040

Like a Good Nearest Neighbor: Practical Content Moderation with Sentence Transformers

Anonymous ACL submission

Abstract

Modern text classification systems have impressive capabilities but are infeasible to deploy and use reliably due to their dependence on prompting and billion-parameter language models. SetFit (Tunstall et al., 2022) is a recent, practical approach that fine-tunes a Sentence Transformer under a contrastive learning paradigm and achieves similar results to more unwieldy systems. Text classification is important for addressing the problem of domain drift in detecting harmful content, which plagues all social media platforms. Here, we propose Like a Good Nearest Neighbor (LAGONN), an inexpensive modification to SetFit that requires no additional parameters or hyperparameters but modifies input with information about its nearest neighbor, for example, the label and text, in the training data, making novel data appear similar to an instance on which the model was optimized. LAGONN is effective at the task of detecting harmful content and generally improves SetFit's performance. To demonstrate LAGONN's value, we conduct a thorough study of text classification systems in the context of content moderation under four label distributions.¹

1 Introduction

Text classification is the most important tool for NLP practitioners, and there has been substantial progress in advancing the state-of-the-art, especially with the advent of large, pretrained language models (PLM) (Devlin et al., 2019). Modern research focuses on in-context learning (Brown et al., 2020), pattern exploiting training (Schick and Schütze, 2021a,b, 2022), adapter-based fine-tuning with learned label embeddings (Karimi Mahabadi et al., 2022), and parameter efficient fine-tuning (Liu et al., 2022a). These methods have achieved impressive results on the SuperGLUE (Wang et al., 2019) and RAFT (Alex et al., 2021)

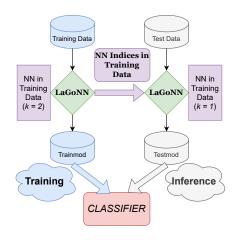


Figure 1: We embed training data, retrieve the text, gold label, and distance for each instance from its second nearest neighbor (k=2) and modify the original text with this information. Then we embed the modified training data and train a classifier. During inference, the NN from the training data is selected (k=1), the original text is modified with the text, gold label, and distance from the NN, and the classifier is called.

few-shot benchmarks, but most are difficult to use because of their reliance on billion-parameter PLMs and prompting. Constructing prompts is not trivial and may require domain expertise. 041

042

043

044

045

046

048

050

051

059

060

One exception to these cumbersome systems is SetFit. SetFit does not rely on prompting or billion-parameter PLMs, and instead fine-tunes a pretrained Sentence Transformer (ST) (Reimers and Gurevych, 2019) under a contrastive learning paradigm. SetFit has comparable performance to more unwieldy systems while being one to two orders of magnitude faster to train and run inference.

An important application of text classification is aiding or automating content moderation, which is the task of determining the appropriateness of user-generated content on the Internet (Roberts, 2017). From fake news to toxic comments to hate speech, it is difficult to browse social media without being exposed to potentially dangerous posts that may have an effect on our ability to reason (Ecker

¹Code and data: https://github.com/[REDACTED]

et al., 2022). Misinformation spreads at alarming rates (Vosoughi et al., 2018), and an ML system should be able to quickly aid human moderators. While there is work in NLP with this goal (Markov et al., 2022; Shido et al., 2022; Ye et al., 2023), a general, practical and open-sourced method that is effective across multiple domains remains an open challenge. Novel fake news topics or racial slurs emerge and change constantly. Retraining of ML-based systems is required to adapt this concept drift, but this is expensive, not only in terms of computation, but also in terms of the human effort needed to collect and label data.

061

067

072

081

084

091

100

103

105

106

107

108

SetFit's performance, speed, and low cost would make it ideal for effective content moderation, however, this type of text classification poses a challenge for even state-of-the-art approaches. For example, detecting hate speech on Twitter (Basile et al., 2019), a subtask on the RAFT few-shot benchmark, appears to be the most difficult dataset; at time of writing, it is the only task where the human baseline has not been surpassed, yet SetFit is among the top ten most performant systems.²

Here, we propose a modification to SetFit, called Like a Good Nearest Neighbor (LAGONN). LAGONN introduces no parameters or hyperparameters and instead modifies input text by retrieving information about the nearest neighbor (NN) seen during optimization (see Figure 1). Specifically, we append the label, distance, and text of the NN in the training data to a new instance and encode this modified version with an ST. By making input data appear more similar to instances seen during training, we inexpensively exploit the ST's pretrained or fine-tuned knowledge when considering a novel example. Our method can also be applied to the linear probing of an ST, requiring no expensive fine-tuning of the large embedding model. Finally, we propose a simple alteration to the SetFit training procedure, where we fine-tune the ST on a subset of the training data. This results in a more efficient and performant text classifier that can be used with LAGONN. We summarize our contributions as follows:

- 1. We propose LAGONN, an inexpensive modification to SetFit- or ST-based text classification.
- 2. We suggest an alternative training procedure

to the standard fine-tuning of SetFit, that can be used with or without LAGONN, and results in a cheaper system with similar performance to the more expensive SetFit.

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

 We perform an extensive study of LAGONN, SetFit, and standard transformer fine-tuning in the context of content moderation under different label distributions.

2 Related Work

There is not much work on using sentence embeddings as features for classification despite the pioneering work being roughly five years old (Perone et al., 2018). STs are pretrained with the objective of maximizing the distance between semantically distinct text and minimizing the distance between text that is semantically similar in feature space. They are composed of a Siamese and triplet architecture that encodes text into dense vectors which can be used as features for ML. STs were first used to encode text for classification by Piao (2021), however, the authors relied on pretrained representations.

SetFit uses a contrastive learning paradigm (Koch et al., 2015) to optimize the ST embedding model. The ST is fine-tuned with a distance-based loss function, like cosine similarity, such that examples with different labels are separated in feature space. Input text is then encoded with the fine-tuned ST and a classifier, such as logistic regression, is trained. This approach creates a strong, few-shot text classification system, transforming the ST from a sentence encoder to a topic encoder.

Most related to LAGONN is work done by Xu et al. (2021), who showed that retrieving and concatenating text from training data and external sources, such as ConceptNet (Speer et al., 2017) and the Wikitionary³ definition, can be viewed as a type of external attention that does not modify the architecture of the Transformer in question answering. Liu et al. (2022b) used PLMs, including STs, and k-NN lookup to prepend examples that are similar to a GPT-3 query sample to aid in prompt engineering for in-context learning. Wang et al. (2022) demonstrated that prepending and appending training data can benefit PLMs in the tasks of summarization, language modelling, machine translation, and question answering, using BM25 as their retrieval model for speed (Manning et al., 2008; Robertson and Zaragoza, 2009).

²https://huggingface.co/spaces/ought/raft-leaderboard (see "Tweet Eval Hate").

³https://www.wiktionary.org/

Training Data	Test Data
"I love this." [positive 0.0] (0)	"So good!" [?] (?)
"This is great!" [positive 0.5] (0)	"Just terrible!" [?] (?)
"I hate this." [negative 0.7] (1)	"Never again." [?] (?)
"This is awful!" [negative 1.2] (1)	"This rocks!" [?] (?)

LAGONN Configuration

Train Modified

LABEL	"I love this. [SEP] [positive 0.5]" (0)
TEXT	"I love this. [SEP] [positive 0.5] This is great!" (0)
BOTH	"I love this. [SEP] [positive 0.5] This is great! [SEP] [negative 0.7] I hate this." (0)
	Test Modified
LABEL	"So good! [SEP] [positive 1.5]" (?)
TEXT	"So good! [SEP] [positive 1.5] I love this." (?)
BOTH	"So good! [SEP] [positive 1.5] I love this. [SEP] [negative 2.7] This is awful!" (?)

Table 1: Toy training and test data and different LAGONN configurations considering the first training example. Train and Test Modified are altered instances that are input into the final embedding model for training and inference, respectively. The input format is "original text [SEP] [NN gold label distance] NN instance text". Input text is in quotation marks, the NN's gold label and distance from the training data are in square brackets, and the integer label is in parenthesis (see Appendix A.4 for examples of LAGONN BOTH modified text).

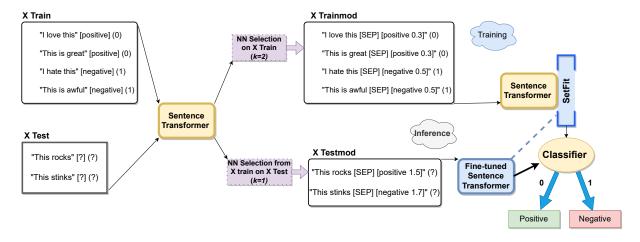


Figure 2: LAGONN LABEL uses an ST to encode training data, performs NN lookup, appends the second NN's (k=2) gold label and distance, and optionally SetFit to fine-tune the embedding model. We then embed this new instance and train a classifier. During inference, we use the embedding model to modify the test data with its NN's gold label and distance from the training data (k=1), compute the final representation, and call the classifier. Input text is in quotation marks, the NN's gold label and distance are in brackets, and the integer label is in parenthesis.

We alter the SetFit training procedure by using fewer examples to adapt the embedding model for many-shot learning. LAGONN decorates input text with its nearest neighbor's gold label, Euclidean distance, and text from the training data to exploit the ST's optimized representations. Compared to retrieval-based methods, LAGONN uses the same model for both retrieval and encoding, which can be fine-tuned via SetFit. We only retrieve information from the training data for text classification.

3 Like a Good Nearest Neighbor

Xu et al. (2021) formulate a type of external attention, where textual information is retrieved from multiple sources and added to text input to give the model stronger reasoning ability without altering the internal architecture. Inspired by this approach, LAGONN exploits pretrained and finetuned knowledge through external attention, but the information we retrieve comes only from data used during optimization. We consider an embedding function, f, that is called on both training and test

data, $f(X_{train})$ and $f(X_{test})$. Considering its success and speed on realistic, few-shot data and our goal of practical content moderation, we choose an ST that can be fine-tuned with SetFit as our embedding function.

179

180

181

185

187

191

192

193

195

196

198

199

204

206

207

209

210

213

214

215

216

217

218

219

Encoding training data and nearest neighbors LAGONN first uses a pretrained Sentence Transformer to embed training text in feature space, $f(X_{train})$. We perform NN lookup with scikitlearn (Buitinck et al., 2013) on the resulting embeddings and query the second closest NN (k=2). We do not use the NN because it is the example itself.

Nearest neighbor information We extract text from the second nearest neighbor and use it to decorate the original example. We experimented with different text that LAGONN could use. The first configuration we consider is the gold label and Euclidean distance of the NN, which we call LA-BEL. We then considered the gold label, distance, and the text of the NN, which we refer to as TEXT. Finally, we tried the same format as TEXT but for all possible labels, which we call BOTH (see Table 1 and Figure 2).⁴ Information from the second NN is appended to the text following a separator token to indicate this instance is composed of multiple sequences. While the BOTH and TEXT configurations are arguably the most interesting, we find LABEL to result in the most performant version of LAGONN, and this is the version about which we report results.

Training LAGONN encodes the modified training data and optionally fine-tunes the embedding model via SetFit, $f(X_{trainmod})$. After fine-tuning, we train a classifier $CLF(f(X_{trainmod}))$, like logistic regression.

Inference LAGONN uses information from the nearest neighbor in the training data to modify input text. We compute the embeddings on the test data, $f(X_{test})$, and query the NN lookup, selecting the NN (k=1) in the training data and extracting information from the training text. LAGONN then decorates the input instance with information from the NN in the training data. Finally, we encode the modified data with the embedding model and call the classifier, $CLF(f(X_{testmod}))$.

Intuition As f is the same function, we hypothesize that LAGONN's modifications will make

a novel instance more semantically similar to its NNs in the training data. The resulting representation should be more akin to an instance on which the embedding model and classifier were optimized. Our method also leverages both distance-based (NN lookup) and probabilistic algorithms (logistic regression) for its final prediction.

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

259

260

261

262

263

264

265

4 Experiments

4.1 Data and label distributions

In our experiments, we study LAGONN's performance on four binary and one ternary classification dataset related to the task of content moderation. Each dataset is composed of a training, validation, and test split.

Here, we provide a summary of the five datasets we studied. LIAR was created from Politifact⁵ for fake news detection and is composed of the data fields context, speaker, and statement, which are labeled with varying levels of truthfulness (Wang, 2017). We used a collapsed version of this dataset where a statement can only be true or false. We did not use speaker, but did use context and statement, separated by a separator token. Quora Insincere Questions⁶ is composed of neutral and toxic questions, where the author is not asking in good faith. Hate Speech Offensive⁷ has three labels and is composed of tweets that can contain either neutral text, offensive language, or hate speech (Davidson et al., 2017). Amazon Counterfactual⁸ contains sentences from product reviews, and the labels can be "factual" or "counterfactual" (O'Neill et al., 2021). "Counterfactual" indicates that the customer said something that cannot be true. Finally, Toxic Conversations⁹ is a dataset of comments where the author wrote a comment with unintended bias 10 (see Table 2).

We study our system by simulating growing training data over ten discrete steps sampled under four different label distributions: extreme, imbalanced, moderate, and balanced (see Table 3). On

⁴LAGONN requires a mapping from the label to the text the label represents, for example, 0 – positive and 1 – negative.

⁵https://www.politifact.com/

⁶https://www.kaggle.com/c/

 $^{{\}tt quora-insincere-questions-classification}$

⁷https://huggingface.co/datasets/hate_speech_
offensive

[%]https://huggingface.co/datasets/SetFit/ amazon_counterfactual_en

⁹https://huggingface.co/datasets/SetFit/toxic_ conversations

¹⁰https://www.kaggle.com/c/
jigsaw-unintended-bias-in-toxicity-classification/
overview

Dataset (and Detection Task)	Number of Labels
LIAR (Fake News)	2
Insincere Questions (Toxicity)	2
Hate Speech Offensive	3
Amazon Counterfactual (English)	2
Toxic Conversations	2

Table 2: Summary of datasets and number of labels. We provide the type of task in parenthesis in unclear cases.

each step we add 100 examples (100 on the first, 200 on the second, etc.) from the training split sampled under one of the four ratios. 11 On each step, we train our method with the sampled data and evaluate on the test split. Considering growing training data has two benefits: 1) We can simulate a streaming data scenario, where new data is labeled and added for training and 2) We can investigate each method's sensitivity to the number of training examples. We sampled over five seeds, reporting the mean and standard deviation.

Regime	Binary	Ternary
Extreme	0: 98% 1: 2%	0: 95%, 1: 2%, 2: 3%
Imbalanced	0: 90% 1: 10%	0: 80%, 1: 5%, 2: 15%
Moderate	0: 75% 1: 25%	0: 65%, 1: 10%, 2: 25%
Balanced	0: 50% 1: 50%	0: 33%, 1: 33%, 2: 33%

Table 3: Label distributions for sampling training data. 0 represents neutral while 1 and 2 represent different types of undesirable text.

4.2 Baselines

We compare LAGONN against a number of strong baselines, detailed below. We used default hyperparameters in all cases unless stated otherwise.

RoBERTa RoBERTa-base is a pretrained language model (Liu et al., 2019) that we fine-tuned with the transformers library (Wolf et al., 2020). We select two versions of RoBERTa-base: an expensive version, where we perform standard fine-tuning on each step (RoBERTa_{full}) and a cheaper version, where we freeze the model body after step one and update the classification head on subsequent steps (RoBERTa_{freeze}). We set the learning rate to $1e^{-5}$, train for a maximum of 70 epochs, and use early stopping, selecting the best model after training. We consider RoBERTa_{full} an upper bound as it has the most trainable parameters and requires the most time to train of all our methods.

Linear probe We perform linear probing of a pretrained Sentence Transformer by fitting logistic regression with default hyperparameters on the training embeddings on each step. We choose this baseline because LAGONN can be applied as a modification in this scenario. We select MPNET (Song et al., 2020) as the ST, for SetFit, and for LAGONN.¹² We refer to this method as Probe.

Logistic regression Here, we perform standard fine-tuning with SetFit on the first step, and then on subsequent steps, freeze the embedding model and retrain only the classification head. We choose this baseline as LAGONN also uses logistic regression as its final classifier and refer to this method as Log Reg.

k-nearest neighbors Similar to the above baseline, we fine-tune the embedding model via SetFit, but swap out the classification head for a kNN classifier, where k=3. We select this baseline as LAGONN also relies on an NN lookup. k=3 was chosen during our development stage as it yielded the strongest performance. We refer to this method as kNN.

SetFit For this baseline we perform standard fine-tuning with SetFit on each step. On the first step, this method is equivalent to Log Reg.

LAGONN cheap This method modifies data via LAGONN before fitting a logistic regression classifier. Even without adapting the embedding model, as the training data grow, modifications made to the test data may change. We refit the classification head on each step and refer to this method as LAGONN $_{cheap}$, which is comparable to Probe.

LAGONN On the first step, we use LAGONN to modify our training data and then perform standard fine-tuning with SetFit. On subsequent steps, we freeze the embedding model and use it to modify our data. We fit logistic regression on each step and refer to this method as LAGONN. It is comparable to Log Reg.

LAGONN expensive This version is identical to LAGONN, except we fine-tune the embedding model on each step. We refer to this method as LAGONN $_{exp}$ and it is comparable to SetFit. On the first step, this method is equivalent to LAGONN.

¹¹For Hate Speech Offensive, 0 and 2 denote undesirable text and 1 denotes neither.

¹²https://huggingface.co/sentence-transformers/
paraphrase-mpnet-base-v2

Method	a et	InsincereQs	1 oth		a et	AmazonCF	$a \circ th$	
Extreme	1^{st}	5^{th}	10^{th}	Average	1^{st}	5^{th}	10^{th}	Average
$RoBERTa_{full}$	$19.9_{8.4}$	$30.9_{7.9}$	$42.0_{7.4}$	$33.5_{6.7}$	$21.8_{6.6}$	$63.9_{10.2}$	$72.3_{3.0}$	$59.6_{16.8}$
SetFit	$24.1_{6.3}$	$29.2_{6.7}$	$36.7_{7.3}$	$31.7_{3.4}$	$22.3_{8.8}$	$64.2_{3.3}$	$68.6_{4.6}$	$56.8_{14.9}$
$LaGoNN_{exp}$	$30.7_{8.9}$	$37.6_{6.1}$	$39.0_{6.1}$	$36.1_{2.3}$	26.1 _{17.5}	68.4 _{4.4}	74.9 _{2.9}	63.2 _{16.7}
$RoBERTa_{freeze}$	19.9 _{8.4}	$34.1_{5.4}$	$37.9_{5.9}$	$32.5_{5.5}$	21.8 _{6.6}	$41.0_{12.7}$	$51.3_{10.7}$	40.6 _{8.9}
kNN	$6.8_{0.42}$	$15.9_{3.4}$	$16.9_{4.3}$	$14.4_{3.0}$	$10.3_{0.2}$	$15.3_{4.2}$	$18.4_{3.7}$	$15.6_{2.4}$
Log Reg	$24.1_{6.3}$	$31.7_{4.9}$	$36.1_{5.4}$	$31.8_{3.6}$	$22.3_{8.8}$	$32.4_{11.5}$	$42.3_{8.8}$	$34.5_{5.9}$
LaGoNN	$30.7_{8.9}$	$39.3_{4.9}$	$41.2_{4.7}$	$38.4_{3.0}$	26.1 _{17.5}	$31.1_{19.4}$	$33.0_{19.1}$	$30.9_{2.3}$
Probe	24.3 _{8.4}	39.8 _{5.6}	44.8 _{4.2}	$38.3_{6.2}$	24.29.0	46.3 _{4.4}	$54.6_{2.0}$	45.1 _{10.3}
$LaGoNN_{cheap}$	$23.6_{7.8}$	40.7 _{5.9}	45.3 _{4.4}	38.6 _{6.6}	$20.1_{6.9}$	$38.3_{4.9}$	$47.8_{3.4}$	$38.2_{9.5}$
Balanced								
$RoBERTa_{full}$	$47.1_{4.2}$	$52.1_{3.6}$	$55.7_{2.6}$	$52.5_{2.9}$	$73.6_{2.1}$	$78.6_{3.9}$	$82.4_{1.1}$	$78.9_{2.2}$
SetFit	$43.5_{4.2}$	$47.1_{4.6}$	$48.5_{3.9}$	$48.0_{1.7}$	$73.8_{4.4}$	$69.8_{4.0}$	$64.1_{4.6}$	$69.6_{3.6}$
$LaGoNN_{exp}$	$42.8_{5.3}$	$47.6_{2.9}$	$47.0_{1.7}$	$46.2_{2.0}$	76.0 _{3.0}	$73.4_{2.6}$	$72.3_{2.9}$	$72.5_{3.4}$
$RoBERTa_{freeze}$	$47.1_{4.2}$	$52.1_{0.4}$	$53.3_{1.7}$	$51.5_{2.1}$	$73.6_{2.1}$	$76.8_{1.6}$	$77.9_{1.0}$	$76.5_{1.3}$
kNN	$22.3_{2.3}$	$30.2_{2.3}$	$30.9_{1.8}$	$29.5_{2.5}$	$41.7_{3.4}$	$57.9_{3.3}$	$58.3_{3.3}$	$56.8_{5.1}$
Log Reg	$43.5_{4.2}$	$53.8_{2.2}$	$55.5_{1.6}$	$52.8_{3.5}$	$73.8_{4.4}$	$79.2_{1.9}$	$80.1_{1.0}$	$78.6_{1.8}$
LaGoNN	$42.8_{5.3}$	$54.1_{2.9}$	$56.3_{1.3}$	$53.4_{3.7}$	76.0 _{3.0}	80.1 _{2.0}	$81.4_{1.1}$	79.8 _{1.4}
Probe	$47.5_{1.6}$	$52.4_{1.7}$	$55.3_{1.1}$	$52.2_{2.5}$	$52.4_{3.4}$	$64.7_{2.5}$	$67.5_{0.4}$	$63.4_{4.4}$
$LaGoNN_{cheap}$	49.3 _{2.6}	54.4 _{1.4}	57.6 _{0.7}	54.2 _{2.7}	$48.1_{3.4}$	$62.0_{2.0}$	$65.3_{0.8}$	$60.5_{5.0}$

Table 4: Average performance (average precision \times 100) on Insincere Questions and Amazon Counterfactual. The first, fifth, and tenth step are followed by the average over all ten steps. The average gives insight into the overall strongest performer by aggregating all steps. We group methods with a comparable number of trainable parameters together. The extreme label distribution results are followed by balanced (see Appendix A.2 for additional results).

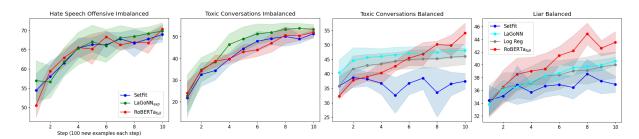


Figure 3: Average performance in the imbalanced and balanced regimes relative to comparable methods. We include $RoBERTa_{full}$ results for reference. The metric is macro-F1 for Hate Speech Offensive, average precision elsewhere.

5 Results

Table 4 and Figure 3 show our results. In the cases of the extreme and imbalanced regimes, Set-Fit's performance steadily increases with the number of training examples. As the label distribution shifts to the balanced regime, however, Set-Fit's performance quickly saturates or even degrades as the number of training examples grows. LAGONN, RoBERTa $_{full}$, and Log Reg, other finetuned PLM classifiers, do not exhibit this behavior. LAGONN $_{exp}$, being based on SetFit, exhibits a similar trend, but the performance degradation is mitigated; on the 10^{th} step of Amazon Counterfac-

tual in Table 4 SetFit's performance decreased by 9.7, while LAGONN $_{exp}$ only fell by 3.7.

LAGONN and LAGONN $_{exp}$ generally outperform Log Reg and SetFit, respectively, often resulting in a more stable model, as reflected in the standard deviation. We find that LAGONN and LAGONN $_{exp}$ exhibit stronger predictive power with fewer examples than RoBERTa $_{full}$ despite having fewer trainable parameters. For example, on the first step of Insincere Questions under the extreme setting, LAGONN's performance is more than 10 points higher.

LAGONN_{cheap} outperforms all other methods

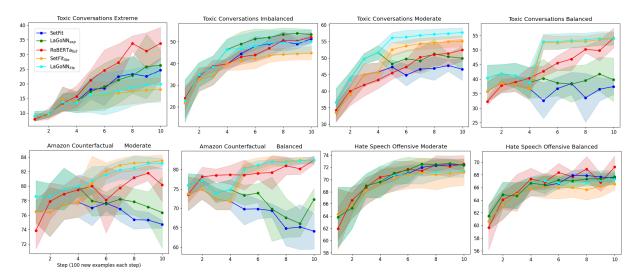


Figure 4: Average performance for all sampling regimes on Toxic Conversations and the moderate and balanced regimes for Amazon Counterfactual and Hate Speech Offensive. More expensive models, such as $LAGONN_{exp}$, SetFit, and $RoBERTa_{full}$ perform best when the label distribution is imbalanced. As the distribution becomes more balanced, inexpensive models, such as $LAGONN_{lite}$, show similar or improved performance. The metric is macro-F1 for Hate Speech Offensive, average precision elsewhere (see Appendix A.3 for additional results).

on the Insincere Questions dataset for all balance regimes, despite being the third fastest (see Table 5) and having the second fewest trainable parameters. We attribute this result to the fact that this dataset is composed of questions from Quora¹³ and our ST backbone was pretrained on similar data. This intuition is supported by Probe, the cheapest method, which despite having the fewest trainable parameters, shows comparable performance.

371

374

376

377

378

379

390

391

395

5.1 SetFit for efficient many-shot learning

Respectively comparing SetFit to Log Reg and LAGONN $_{exp}$ to LAGONN suggests that fine-tuning the ST embedding model on moderate or balanced data hurts model performance as the number of training samples grows. We therefore hypothesize that randomly sampling a subset of training data to fine-tune the encoder, freezing, embedding the remaining data, and training the classifier will result in a stronger model.

To test our hypothesis, we add two models to our experimental setup: SetFit $_{lite}$ and LaGoNN $_{lite}$. SetFit $_{lite}$ and LaGoNN $_{lite}$ are respectively equivalent to SetFit and LaGoNN $_{exp}$, except after the fourth step (400 samples), we freeze the encoder and only retrain the classifier on subsequent steps, similar to Log Reg and LaGoNN.

Figure 4 shows our results with these two new models. As expected, in the cases of extreme and imbalanced distributions, LAGONN_{exp}, SetFit,

and RoBERTa $_{exp}$, are the strongest performers on Toxic Conversations. We note very different results for both LAGONN lite and SetFit lite compared to $LaGoNN_{exp}$ and SetFit on Toxic Conversations and Amazon Counterfactual under the moderate and balanced label distributions. As their expensive counterparts start to plateau or degrade on the fourth step, the predictive power of these two new models dramatically increases, showing improved or comparable performance to RoBERTa_{full}, despite being optimized on less data; for example, LAGONN_{lite} reaches an average precision of approximately 55 after being optimized on only 500 examples. RoBERTa_{full} does not exhibit similar performance until the tenth step. Finally, we point out that LAGONN-based methods generally provide a performance boost for SetFit-based classification.

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

5.2 LAGONN's computational expense

LAGONN is more computationally expensive than Sentence Transformer- or SetFit-based text classification. LAGONN introduces additional inference with the encoder, NN-lookup, and string modification. As the computational complexity of transformers increases with sequence length (Vaswani et al., 2017), additional expense is created when LAGONN appends textual information before inference with the ST. In Table 5, we provide a speed comparison between Probe, Log Reg, SetFit, and LAGONN classification computed on the same

¹³https://www.quora.com/

Method	Time in seconds
Probe	22.9
$LaGoNN_{cheap}$	44.2
Log Reg	42.9
LaGoNN	63.4
SetFit	207.3
$LaGoNN_{exp}$	238.0
$RoBERTa_{full}$	446.9

Table 5: Speed comparison between LAGONN and comparable methods. Time includes training each method on 1,000 examples and performing inference on 51,000 examples.

hardware.¹⁴ On average, LAGONN introduced 24.2 additional seconds of computation compared to its relative counterpart.

6 Discussion

Modern research has achieved impressive results on a variety of text classification tasks and with limited training data. SetFit is one such example and can be used practically, but based on our results, the task of text classification for content moderation presents a challenge even for state-of-the-art approaches. It is imperative that we develop reliable methods that can be feasibly and quickly applied. These methods should be as inexpensive as possible such that we can re-tune them for novel forms of hate speech, toxicity, and fake news.

Our results suggest that $LAGONN_{exp}$ or SetFit, relatively expensive techniques, can detect harmful content when dealing with imbalanced label distributions, as is common with realistic datasets. This finding is intuitive from the perspective that less common instances are more difficult to learn and require more effort. The exception to this would be our examination of Insincere Questions, where $LAGONN_{cheap}$ excelled. This highlights the fact that we can inexpensively extract pretrained knowledge if PLMs are chosen with care for related tasks.

Standard fine-tuning with SetFit does not help performance on more balanced datasets that are not few-shot. SetFit was developed for few-shot learning, but we have observed that it should not be applied "out of the box" to balanced, non-few-shot data. This can be detrimental to performance and has a direct effect on our approach. However, we have observed that LAGONN can stabilize Set-

Fit's predictions and reduce its performance drop. Figures 3 and 4 show that when the label distribution is moderate or balanced (see Table 3), Set-Fit plateaus, yet less expensive systems, such as LAGONN, continue to learn. We believe this is due to SetFit's fine-tuning objective, which optimizes a Sentence Transformer using cosine similarity loss to separate examples belonging to different labels in feature space by assuming independence between labels. This may be too strong an assumption as we optimize with more examples, which is counter-intuitive for data-hungry transformers. RoBERTa $_{full}$, optimized with cross-entropy loss, generally showed improved performance as we added training data.

When dealing with balanced data, it is sufficient to fine-tune the Sentence Transformer via SetFit with 50 to 100 examples per label, while 150 to 200 instances appear to be sufficient when the training data are moderately balanced. The encoder can then be frozen and all available data embedded to train a classifier. This improves performance and is more efficient than full-model fine-tuning. LAGONN is directly applicable to this case, boosting the performance of SetFit_{lite} without introducing trainable parameters. In this setup, all models fine-tuned on Hate Speech Offensive exhibited similar, upward-trending learning curves, but we note the speed of LAGONN relative to RoBERTa_{full} or SetFit (see Figure 4 and Table 5).

7 Conclusion

We have proposed LAGONN, a simple and inexpensive modification to Sentence Transformer- or SetFit-based text classification. LAGONN does not introduce any trainable parameters or new hyperparameters, but typically improves SetFit's performance. To demonstrate the merit of LAGONN, we examined text classification systems in the context of content moderation under four label distributions on five datasets and with growing training data. To our knowledge, this is the first work to examine SetFit in this way. When the training labels are imbalanced, expensive systems, such as LAGONN $_{exp}$ are performant. However, when the distribution is balanced, standard fine-tuning with SetFit can actually hurt model performance. We have therefore proposed an alternative fine-tuning procedure to which LAGONN can be easily utilized, resulting in a powerful, but inexpensive system capable of detecting harmful content.

¹⁴We used a 40 GB NVIDIA A100 Tensor Core GPU.

8 Limitations

510

511

512 513

514

515

516

517

518

519

520

522

525

526

529

530

534

537

539

541

544

545

551

553

555

556

In the current work, we have only considered text data, but social media content can of course consist of text, images, and videos. As LAGONN depends only on an embedding model, an obvious extension to our approach would be examining the modifications we suggest, but on multimodal data. This is an interesting direction that we leave for future research. We have also considered English data, but harmful content can appear in any language. The authors demonstrated that SetFit is performant on multilingual data, the only necessary modification being the underlying pretrained ST. We therefore suspect that LAGONN would behave similarly on non-English data, but this is not something we have tested ourselves. In order to examine our system's performance under different label-balance distributions, we restricted ourselves to binary and ternary text classification tasks, and LAGONN therefore remains untested when there are more than three labels. We did not study our method when there are fewer than 100 examples, and investigating LAGONN in a few-shot learning setting is fascinating topic for future study. Finally, we note that our system could be misused to detect undesirable content that is not necessarily harmful. For example, a social media website could detect and silence users who complain about the platform. This is not our intended use case, but could result from any classifier, and potential misuse is an unfortunate drawback of all technology.

9 Ethics Statement

It is our sincere goal that our work contributes to the social good in multiple ways. We first hope to have furthered research on text classification that can be feasibly applied to combat undesirable content, such as misinformation, on the Internet, which could potentially cause someone harm. To this end, we have tried to describe our approach as accurately as possible and released our code and data, such that our work is transparent and can be easily reproduced and expanded upon. We hope that we have also created a useful but efficient system which reduces the need to expend energy in the form expensive computation. For example, LAGONN does not rely on billion-parameter language models that demand thousand-dollar GPUs to use. LAGONN makes use of GPUs no more than SetFit, despite being more computationally expensive. We have additionally proposed a simple method to make

SetFit, an already relatively inexpensive method, even more efficient.

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

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

References

Neel Alex, Eli Lifland, Lewis Tunstall, Abhishek Thakur, Pegah Maham, C. Jess Riedel, Emmie Hine, Carolyn Ashurst, Paul Sedille, Alexis Carlier, Michael Noetel, and Andreas Stuhlmüller. 2021.
 RAFT: A real-world few-shot text classification benchmark. In Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2).

Valerio Basile, Cristina Bosco, Elisabetta Fersini, Debora Nozza, Viviana Patti, Francisco Manuel Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti. 2019. SemEval-2019 task 5: Multilingual detection of hate speech against immigrants and women in Twitter. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 54–63, Minneapolis, Minnesota, USA. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems, volume 33, pages 1877–1901. Curran Associates, Inc.

Lars Buitinck, Gilles Louppe, Mathieu Blondel, Fabian Pedregosa, Andreas Mueller, Olivier Grisel, Vlad Niculae, Peter Prettenhofer, Alexandre Gramfort, Jaques Grobler, Robert Layton, Jake VanderPlas, Arnaud Joly, Brian Holt, and Gaël Varoquaux. 2013. API design for machine learning software: experiences from the scikit-learn project. In ECML PKDD Workshop: Languages for Data Mining and Machine Learning, pages 108–122.

Thomas Davidson, Dana Warmsley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In *Proceedings of the 11th International AAAI Conference on Web and Social Media*, ICWSM '17, pages 512–515.

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.

Ullrich K. H. Ecker, Stephan Lewandowsky, John Cook, Philipp Schmid, Lisa K. Fazio, Nadia Brashier, Panayiota Kendeou, Emily K. Vraga, and Michelle A. Amazeen. 2022. The psychological drivers of misinformation belief and its resistance to correction. *Nature Reviews Psychology*, 1(1):13–29.

Rabeeh Karimi Mahabadi, Luke Zettlemoyer, James Henderson, Lambert Mathias, Marzieh Saeidi, Veselin Stoyanov, and Majid Yazdani. 2022. Promptfree and efficient few-shot learning with language models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 3638–3652, Dublin, Ireland. Association for Computational Linguistics.

Gregory Koch, Richard Zemel, Ruslan Salakhutdinov, et al. 2015. Siamese neural networks for one-shot image recognition. In *ICML Deep Learning Workshop*, volume 2, page 0. Lille.

Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin Raffel. 2022a. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *arXiv* preprint arXiv:2205.05638.

Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022b. What makes good in-context examples for GPT-3? In Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures, pages 100–114, Dublin, Ireland and Online. 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. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. *Introduction to Information Retrieval*. Cambridge University Press, USA.

Todor Markov, Chong Zhang, Sandhini Agarwal, Tyna Eloundou, Teddy Lee, Steven Adler, Angela Jiang, and Lilian Weng. 2022. A holistic approach to undesired content detection. *arXiv preprint arXiv*:2208.03274.

James O'Neill, Polina Rozenshtein, Ryuichi Kiryo, Motoko Kubota, and Danushka Bollegala. 2021. I wish I would have loved this one, but I didn't – a multilingual dataset for counterfactual detection in product review. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7092–7108, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Christian S. Perone, Roberto Pereira Silveira, and Thomas S. Paula. 2018. Evaluation of sentence embeddings in downstream and linguistic probing tasks. *arXiv preprint arXiv:1806.06259*.

Guangyuan Piao. 2021. Scholarly text classification with sentence bert and entity embeddings. In *Trends and Applications in Knowledge Discovery and Data Mining*, pages 79–87, Cham. Springer International Publishing.

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 and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.

Sarah T. Roberts. 2017. *Content Moderation*, pages 1–4. Springer International Publishing, Cham.

Stephen Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: Bm25 and beyond. *Found. Trends Inf. Retr.*, 3(4):333–389.

Timo Schick and Hinrich Schütze. 2021a. Exploiting cloze-questions for few-shot text classification and natural language inference. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 255–269, Online. Association for Computational Linguistics.

Timo Schick and Hinrich Schütze. 2021b. It's not just size that matters: Small language models are also few-shot learners. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2339–2352, Online. Association for Computational Linguistics.

Timo Schick and Hinrich Schütze. 2022. True few-shot learning with Prompts—A real-world perspective. *Transactions of the Association for Computational Linguistics*, 10:716–731.

Yusuke Shido, Hsien-Chi Liu, and Keisuke Umezawa. 2022. Textual content moderation in C2C market-place. In *Proceedings of the Fifth Workshop on e-Commerce and NLP (ECNLP 5)*, pages 58–62, Dublin, Ireland. Association for Computational Linguistics.

Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. Mpnet: Masked and permuted pre-training for language understanding. In *Advances in Neural Information Processing Systems*, volume 33, pages 16857–16867. Curran Associates, Inc.

Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1).

Lewis Tunstall, Nils Reimers, Unso Eun Seo Jo, Luke Bates, Daniel Korat, Moshe Wasserblat, and Oren Pereg. 2022. Efficient few-shot learning without prompts. arXiv preprint arXiv:2209.11055.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. The spread of true and false news online. *Science*, 359(6380):1146–1151.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Shuohang Wang, Yichong Xu, Yuwei Fang, Yang Liu, Siqi Sun, Ruochen Xu, Chenguang Zhu, and Michael Zeng. 2022. Training data is more valuable than you think: A simple and effective method by retrieving from training data. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3170–3179, Dublin, Ireland. Association for Computational Linguistics.

William Yang Wang. 2017. "Liar, liar pants on fire": A new benchmark dataset for fake news detection. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 422–426, Vancouver, Canada. 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.

Yichong Xu, Chenguang Zhu, Shuohang Wang, Siqi Sun, Hao Cheng, Xiaodong Liu, Jianfeng Gao, Pengcheng He, Michael Zeng, and Xuedong Huang. 2021. Human parity on commonsenseqa: Augmenting self-attention with external attention. *arXiv* preprint arXiv:2112.03254, abs/2112.03254.

Meng Ye, Karan Sikka, Katherine Atwell, Sabit Hassan, Ajay Divakaran, and Malihe Alikhani. 2023. Multilingual content moderation: A case study on reddit. *arXiv preprint arXiv:2302.09618*.

A Appendix

A.1 Observations about LAGONN

Our original goal was to construct a system that did not need to be updated after step one and could simply perform inference on subsequent steps, an active learning setup. While the performance of this version of LAGONN did not degrade, it also did not appear to learn anything and we found it necessary to update parameters on each step. We additionally tried fine-tuning the embedding model via SetFit first before modifying data, however, this hurt performance in all cases. We include this information for transparency and because we find it interesting.

A.2 Additional results for initial experiments

Here we provide additional results from our initial experimental setup that, due to space limitations, could not be included in the main text. We note that a version of LAGONN outperforms or has the same performance of all methods, including our upper bound RoBERTa_{full}, on 54% of all displayed results, and is the best performer relative to Sentence Transformer-based methods on 72%. This excludes LAGONN_{cheap}. This method showed strong performance on the Insincere Questions dataset, but hurts performance in other cases. In cases, when SetFit-based methods do outperform our system, the performances are comparable, yet they can be quite dramatic when LAGONN-based methods are the strongest. Below, we report the mean average precision $\times 100$ for all methods over five seeds with the standard deviation, except in the case of Hate Speech Offensive, where the evaluation metric is the macro-F1. Each table shows the results for a given dataset and a given label-balance distribution on the first, fifth, and tenth step followed by the average for all ten steps. The Liar dataset seems to be the most difficult for all methods. This is expected because it likely does not include enough context to determine the truth of a statement.

Method <i>Imbalanced</i>	1^{st}	$\begin{array}{c} \textbf{Insincere-Questions} \\ 5^{th} \end{array}$	10^{th}	Average	Method Extreme	1^{st}	Toxic Conversations 5^{th}	10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	39.8 _{5.5} 43.7 _{2.7} 44.5 _{4.5}	53.1 _{4.6} 52.2 _{1.9} 52.7 _{2.4}	55.7 _{1.2} 53.8 _{0.9} 55.4 _{2.0}	50.6 _{4.4} 51.4 _{2.9} 51.8 _{3.0}	$\begin{array}{c} {\rm RoBERTa}_{full} \\ {\rm SetFit} \\ {\rm LaGoNN}_{exp} \end{array}$	$7.9_{0.5}$ $8.8_{1.2}$ $8.9_{1.7}$	$21.2_{3.7} \\ 18.1_{3.4} \\ 17.4_{6.6}$	33.8 _{5.5} 24.7 _{4.1} 26.4 _{5.2}	$21.9_{9.3} \\ 17.6_{5.5} \\ 17.9_{6.0}$
RoBERTa _{freeze} kNN Log Reg LAGONN	39.8 _{5.5} 23.9 _{2.2} 43.7 _{2.7} 44.5 _{4.5}	$44.1_{3.6} \\ 30.3_{3.0} \\ 47.6_{1.6} \\ 48.1_{2.2}$	$46.3_{2.4} \\ 31.6_{2.4} \\ 50.1_{2.1} \\ 50.3_{1.7}$	$44.0_{2.0} 30.0_{2.1} 47.6_{1.8} 48.1_{1.9}$	RoBERTa _{freeze} kNN Log Reg LAGONN	$7.9_{0.5}$ $7.9_{0.0}$ $8.8_{1.2}$ $8.9_{1.7}$	$12.8_{2.4} \\ 8.7_{0.4} \\ 13.1_{2.5} \\ 13.8_{3.9}$	$19.1_{3.2} \\ 8.7_{0.2} \\ 16.3_{3.0} \\ 17.1_{4.8}$	$13.5_{3.5} 8.5_{0.3} 13.0_{2.6} 13.4_{2.6}$
Probe LAGONN _{cheap}	$40.4_{4.2} \\ 40.8_{4.3}$	$49.4_{2.3} \\ 51.1_{2.4}$	$52.3_{1.7} $ $54.5_{1.4}$	49.0 _{3.3} 50.4 _{4.0}	Probe LAGONN _{cheap}	13.1 _{2.8} 11.3 _{2.2}	24.6 _{2.6} 21.7 _{2.7}	$30.1_{2.1} \\ 27.4_{2.3}$	23.9 _{5.6} 21.3 _{5.3}

Table 6 Table 10

Method <i>Moderate</i>	1^{st}	$\begin{array}{c} \textbf{Insincere Questions} \\ 5^{th} \end{array}$	10^{th}	Average
RoBERTa _{full} SetFit LAGONN _{exp}	48.1 _{2.3} 48.9 _{1.7} 49.8 _{1.6}	54.7 _{1.9} 53.9 _{0.7} 52.2 _{1.9}	57.5 _{1.5} 54.2 _{1.5} 53.2 _{3.3}	$53.9_{2.9}$ $52.3_{1.6}$ $52.0_{1.4}$
RoBERTa _{freeze} kNN Log Reg LAGONN	$48.1_{2.3} \\ 28.0_{2.4} \\ 48.9_{1.7} \\ \textbf{49.8}_{1.6}$	$50.2_{2.2} \ 33.9_{2.8} \ 53.6_{1.9} \ 54.4_{1.3}$	$52.0_{1.4} \\ 33.6_{2.0} \\ 55.8_{1.7} \\ 56.9_{0.5}$	$50.2_{1.4}$ $33.5_{1.9}$ $53.3_{2.2}$ $54.2_{2.2}$
Probe LAGONN _{cheap}	$45.7_{2.1} \\ 45.7_{2.2}$	$52.3_{1.8}$ $54.4_{1.6}$	$54.4_{1.1}$ $56.4_{0.6}$	$51.4_{2.5}$ $53.2_{3.2}$

Method <i>Imbalanced</i>	1^{st}	Toxic Conversations 5^{th}	10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	24.1 _{5.6} 21.8 _{6.6} 22.7 _{9.8}	$43.1_{3.4} \\ 44.5_{4.1} \\ 49.1_{5.6}$	52.1 _{2.5} 51.4 _{1.9} 53.4 _{2.3}	$42.4_{8.2} 42.1_{9.3} 45.6_{9.8}$
RoBERTa _{freeze} kNN Log Reg LaGONN	$\begin{array}{c} \textbf{24.1}_{5.6} \\ 11.5_{2.5} \\ 21.8_{6.6} \\ 22.7_{9.8} \end{array}$	$31.2_{4.4}$ $14.7_{4.0}$ $26.7_{5.3}$ $27.6_{8.9}$	$34.0_{4.0}$ $15.3_{3.2}$ $30.2_{4.0}$ $30.3_{8.7}$	$30.5_{3.1} \\ 14.6_{1.1} \\ 26.6_{2.7} \\ 27.4_{2.4}$
Probe LAGONN _{cheap}	$23.3_{2.7} \\ 20.5_{3.2}$	$33.0_{2.8}$ $31.1_{3.2}$	$37.1_{1.8}$ $35.6_{1.8}$	$32.5_{4.2} \\ 30.5_{4.6}$

Table 7 Table 11

Method Imbalanced	1^{st}	Amazon Counterfactual 5^{th}	10^{th}	Average
RoBERTa _{full} SetFit	$68.2_{4.5}$ $72.0_{2.1}$	81.0 _{1.7} 78.4 _{2.8}	82.2 _{1.0} 78.8 _{1.2}	$79.2_{3.9}$ $78.0_{2.1}$
LAGONN _{exp}	74.3 _{3.8}	$80.1_{1.4}$	$79.0_{1.6}$	79.5 _{1.9}
$RoBERTa_{freeze}$	$68.2_{4.5}$	$75.0_{2.2}$	$77.0_{2.4}$	$74.2_{2.6}$
kNN	$51.0_{4.1}$	$60.0_{3.1}$	$61.3_{2.1}$	$59.7_{3.0}$
Log Reg	$72.0_{2.1}$	$74.4_{2.3}$	$76.7_{1.8}$	$74.8_{1.4}$
LaGoNN	74.3 _{3.8}	$76.1_{3.6}$	$77.3_{3.2}$	$76.1_{1.0}$
Probe	$46.6_{2.8}$	$60.3_{1.4}$	$64.2_{1.2}$	$59.2_{5.2}$
$LaGonn_{cheap}$	$38.2_{3.2}$	$55.3_{1.8}$	$61.0_{1.2}$	$54.4_{6.7}$

Method <i>Moderate</i>	1^{st}	Toxic Conversations 5^{th}	10^{th}	Average
RoBERTa _{full} SetFit LAGONN _{exp}	34.2 _{3.4} 33.6 _{2.9} 36.6 _{4.2}	45.5 _{1.9} 47.2 _{2.2} 48.2 _{2.7}	52.4 _{3.3} 46.6 _{3.3} 49.9 _{3.7}	45.7 _{5.6} 44.3 _{4.3} 48.0 _{4.4}
RoBERTa _{freeze} kNN Log Reg LAGONN	34.2 _{3.4} 19.4 _{1.9} 33.6 _{2.9} 36.6 _{4.2}	$38.4_{2.1}$ $21.5_{3.4}$ $39.2_{2.9}$ $42.7_{3.7}$	$\begin{array}{c} 39.5_{1.8} \\ 22.4_{2.9} \\ 41.6_{2.7} \\ 45.0_{3.5} \end{array}$	$38.0_{1.5}$ $21.6_{0.8}$ $38.6_{2.4}$ $42.0_{2.5}$
Probe LaGoNN _{cheap}	$29.0_{2.7} \\ 26.1_{2.7}$	$36.1_{1.2} \ 34.3_{1.3}$	$39.1_{1.5}$ $37.5_{1.8}$	$35.5_{3.3}$ $33.6_{3.6}$

Table 8 Table 12

Method Moderate	1^{st}		10^{th}	Average
$RoBERTa_{full}$	$73.9_{2.5}$	80.01.0	$80.1_{2.3}$	$79.1_{2.1}$
SetFit	$76.5_{1.6}$	$77.0_{2.4}$	$74.7_{0.5}$	$76.5_{1.0}$
$LaGonn_{\it exp}$	$78.6_{2.2}$	$78.0_{2.1}$	$76.3_{4.9}$	$78.2_{1.0}$
$RoBERTa_{freeze}$	73.9 _{2.5}	76.6 _{1.4}	$78.5_{0.7}$	76.4 _{1.7}
kNN	$54.5_{3.1}$	$64.2_{1.9}$	$66.6_{1.3}$	$64.7_{3.5}$
Log Reg	$76.5_{1.6}$	$80.6_{0.5}$	$81.2_{0.3}$	$80.0_{1.4}$
LaGoNN	$78.6_{2.2}$	81.2 _{1.4}	$81.6_{1.1}$	$80.8_{0.9}$
Probe	$52.3_{2.0}$	64.1 _{1.8}	$67.2_{1.4}$	63.14.3
$LaGoNN_{cheap}$	$47.3_{3.4}$	$60.7_{1.5}$	$65.2_{1.4}$	$59.5_{5.2}$

Method Balanced	1^{st}	Toxic Conversations 5^{th}	10^{th}	Average
RoBERTa $_{full}$ SetFit LaGoNN $_{exp}$	$32.3_{1.1}$ $35.7_{3.4}$ 40.4 _{4.4}	$42.7_{1.8} \\ 32.6_{6.2} \\ 40.2_{6.6}$	54.1 _{3.4} 37.4 _{2.7} 39.8 _{7.5}	$43.8_{6.3} \\ 36.5_{1.9} \\ 40.0_{1.2}$
RoBERTa _{freeze} kNN Log Reg LAGONN	$32.3_{1.1} \\ 17.4_{0.8} \\ 35.7_{3.4} \\ \textbf{40.4}_{4.4}$	$39.2_{1.5}$ $23.7_{2.6}$ $44.5_{2.9}$ $46.6_{2.7}$	$41.0_{0.6} 24.3_{2.7} 46.1_{2.8} 48.1_{2.2}$	$38.5_{2.4}$ $23.1_{2.0}$ $43.6_{2.9}$ 46.1 _{2.2}
Probe LAGONN _{cheap}	$29.5_{2.4} \\ 26.8_{2.7}$	$35.9_{0.9} \ 34.5_{1.3}$	$40.2_{0.9} \\ 38.5_{0.8}$	$36.1_{3.5} \\ 34.4_{3.7}$

Table 9 Table 13

Method Extreme	1^{st}		10^{th}	Average
$\begin{array}{c} {\rm RoBERTa}_{full} \\ {\rm SetFit} \\ {\rm LaGoNN}_{exp} \end{array}$	$30.2_{1.4}$ $30.3_{0.8}$ $30.3_{0.7}$	43.5 _{2.5} 44.0 _{1.3} 40.7 _{2.9}	51.2 _{2.2} 51.1 _{2.0} 49.1 _{4.4}	44.3 _{7.4} 43.8 _{6.5} 42.2 _{6.2}
RoBERTa _{freeze} kNN Log Reg LAGONN	$30.2_{1.4} \\ 31.5_{1.2} \\ 30.3_{0.8} \\ 30.3_{0.7}$	$33.5_{3.1}$ $35.9_{2.7}$ $38.4_{2.5}$ $35.7_{2.6}$	$\begin{array}{c} 34.4_{3.4} \\ 37.4_{2.0} \\ 41.1_{1.5} \\ 39.1_{2.4} \end{array}$	$33.1_{1.4}$ $35.8_{1.7}$ $37.8_{3.3}$ $35.6_{2.7}$
Probe LAGONN _{cheap}	$29.0_{0.2} \\ 29.0_{0.1}$	$34.7_{1.5} \ 36.9_{1.8}$	$40.1_{2.1} \\ 40.5_{2.1}$	$35.1_{3.8} \ 36.2_{3.7}$

_	_				
	г	1_1	_	-1	_ /

Method Extreme	1^{st}	$\begin{array}{c} \textbf{Liar} \\ 5^{th} \end{array}$	10^{th}	Average
RoBERTa _{full} SetFit	32.0 _{2.7} 31.2 _{3.8}	34.7 _{2.9} 30.4 _{3.1}	$35.1_{4.3}$ $31.8_{2.9}$	$33.7_{1.0}$ $31.5_{0.7}$
$LaGoNN_{exp}$	$30.6_{4.7}$	$30.3_{2.0}$	$31.3_{2.0}$	$31.1_{0.6}$
RoBERTa $_{freeze}$ k NN Log Reg	$32.0_{2.7} 27.0_{0.5} 31.2_{3.8}$	$32.8_{4.5}$ $27.3_{0.8}$ $33.7_{5.1}$	34.2 _{5.0} 27.9 _{0.8} 35.7 _{5.1}	33.2 _{0.7} 27.4 _{0.3} 34.3 _{1.6}
LaGoNN	$30.6_{4.7}$	$32.0_{4.6}$	$33.7_{5.4}$	$32.6_{0.9}$
Probe LAGONN _{cheap}	$30.7_{2.0} \ 30.7_{2.0}$	$30.6_{3.9} \ 30.5_{3.8}$	$31.7_{2.9}$ $31.4_{2.6}$	$31.1_{0.4}$ $31.0_{0.4}$

Table 18

Method Imbalanced	1^{st}		10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	50.6 _{3.0} 54.4 _{4.3} 57.0 _{5.2}	$65.2_{3.9} \\ 66.3_{1.8} \\ 67.0_{4.4}$	70.3 _{1.2} 68.9 _{2.0} 69.8 _{2.1}	64.2 _{5.3} 64.3 _{4.5} 64.9 _{4.6}
RoBERTa _{freeze} kNN Log Reg LAGONN	$50.6_{3.0} \\ 55.6_{4.8} \\ 54.4_{4.3} \\ \textbf{57.0}_{5.2}$	$54.1_{1.6} 57.3_{2.3} 57.0_{3.9} 58.2_{4.1}$	55.3 _{2.3} 58.8 _{3.6} 58.2 _{3.8} 58.3 _{3.4}	$54.1_{1.3} 57.4_{1.1} 57.2_{1.1} 58.3_{0.6}$
Probe LAGONN $_{cheap}$	$46.5_{2.2} \\ 47.1_{1.3}$	$57.8_{1.7} $ $56.5_{2.2}$	$60.3_{1.2} \\ 59.5_{2.5}$	$56.5_{4.5} $ $55.6_{3.8}$

Table 15

Method Moderate	1^{st}		10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	61.9 _{3.4} 64.3 _{4.2} 63.8 _{4.9}	$70.8_{1.0} $ $70.6_{2.4} $ 71.0 _{2.1}	72.5 _{1.4} 72.4 _{0.5} 72.3 _{1.0}	69.9 _{3.2} 69.8 _{2.8} 70.0 _{3.0}
$\begin{array}{c} {\rm RoBERTa}_{freeze} \\ k{\rm NN} \\ {\rm Log~Reg} \\ {\rm LaGoNN} \end{array}$	61.9 _{3.4} 64.3 _{4.0} 64.3 _{4.2} 63.8 _{4.9}	$63.2_{4.1} \\ 63.3_{2.9} \\ 67.3_{3.2} \\ 65.0_{5.3}$	$64.1_{4.5} \\ 63.9_{2.5} \\ 67.6_{2.3} \\ 66.7_{5.9}$	$63.2_{0.6} \\ 63.7_{0.4} \\ 66.9_{1.1} \\ 65.3_{0.9}$
Probe LAGONN _{cheap}	$55.6_{1.7}$ $56.0_{3.6}$	$63.8_{0.8} \\ 62.2_{1.4}$	$66.1_{0.3} \\ 66.0_{0.9}$	$63.2_{3.0} \\ 62.3_{2.9}$

Table 16

Method <i>Imbalanced</i>	1^{st}	$\begin{array}{c} \textbf{Liar} \\ 5^{th} \end{array}$	10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	31.4 _{3.2} 32.3 _{4.5} 32.3 _{4.6}	$35.8_{2.6}$ $35.9_{3.1}$ $35.7_{3.4}$	40.0 _{4.3} 36.4 _{2.2} 36.5 _{2.3}	$36.2_{2.4}$ $35.2_{1.1}$ $35.7_{1.4}$
RoBERTa _{freeze} kNN Log Reg LAGONN	31.4 _{3.2} 27.0 _{0.2} 32.3 _{4.5} 32.3 _{4.6}	34.1 _{2.6} 28.5 _{1.0} 36.5 _{3.1} 34.9 _{2.2}	$\begin{array}{c} 35.6_{3.2} \\ 29.0_{1.0} \\ 38.5_{3.4} \\ 36.9_{2.5} \end{array}$	34.0 _{1.4} 28.7 _{0.7} 36.3 _{2.0} 35.3 _{1.4}
Probe LAGONN _{cheap}	$30.7_{3.0} \ 30.4_{3.0}$	$32.8_{1.8} \\ 32.9_{1.8}$	$35.0_{1.6}$ $35.4_{1.7}$	$33.5_{1.5} \ 33.5_{1.7}$

Table 19

Method Balanced	1^{st}	Hate Speech Offensive 5^{th}	10^{th}	Avaraga
RoBERTa _{full} SetFit LAGONN _{exp}	59.7 _{3.5} 60.7 _{1.3} 61.5 _{1.7}	66.9 _{1.2} 66.3 _{1.6} 66.4 _{1.4}	69.2 _{1.8} 67.5 _{0.9} 67.7 _{0.9}	Average 66.4 _{2.7} 65.9 _{2.2} 66.1 _{1.8}
$\begin{array}{c} \text{RoBERTa}_{freeze} \\ k \text{NN} \end{array}$	59.7 _{3.5} 60.7 _{1.3}	60.4 _{2.7} 59.6 _{2.8}	63.1 _{2.3} 59.5 _{2.5}	61.0 _{1.3} 59.5 _{0.5}
Log Reg LAGONN	60.7 _{1.3} 61.5 _{1.7}	62.5 _{0.7} 62.8 _{1.5}	63.4 _{1.0} 64.2 _{1.0}	62.3 _{1.0} 63.0 _{0.9}
Probe LAGONN _{cheap}	$54.9_{1.4}$ $54.2_{2.3}$	$58.5_{0.9} $ $58.6_{0.6}$	$60.9_{0.4} \\ 60.6_{0.5}$	$58.7_{1.7}$ $58.5_{1.8}$

Table 17

Method <i>Moderate</i>	1^{st}	$\begin{array}{c} \textbf{Liar} \\ 5^{th} \end{array}$	10^{th}	Average
RoBERTa $_{full}$	$33.9_{3.1}$	$38.4_{2.7}$	43.9 _{2.2}	39.5 _{3.0}
SetFit	$33.0_{2.6}$	$37.2_{1.8}$	$38.7_{1.5}$	$37.4_{1.6}$
$LaGonn_{exp}$	34.1 _{3.4}	38.7 _{2.3}	$39.0_{1.8}$	$37.8_{1.5}$
$RoBERTa_{freeze}$	$33.9_{3.1}$	$35.3_{2.6}$	$36.8_{2.2}$	$35.4_{1.0}$
kNN	$29.2_{0.8}$	$29.7_{1.5}$	$30.0_{0.6}$	$29.8_{0.3}$
Log Reg	$33.0_{2.6}$	$37.2_{3.9}$	$39.4_{3.5}$	$37.0_{1.8}$
LaGoNN	34.1 _{3.4}	$37.0_{3.1}$	$38.6_{3.0}$	$36.8_{1.3}$
Probe	$31.6_{1.1}$	$34.7_{2.5}$	$37.0_{2.5}$	$34.9_{1.7}$
$LaGonn_{cheap}$	$31.4_{0.9}$	$35.3_{2.3}$	$37.6_{2.0}$	$35.3_{1.9}$

Table 20

Method		Liar		
Balanced	1^{st}	5^{th}	10^{th}	Average
$RoBERTa_{full}$	$33.8_{2.1}$	$39.4_{2.4}$	43.5 _{1.7}	40.2 _{3.2}
SetFit	$34.4_{2.3}$	$36.7_{1.7}$	$37.0_{1.3}$	$36.5_{1.1}$
$LaGoNN_{exp}$	$33.8_{1.8}$	$34.2_{2.7}$	$37.2_{1.9}$	$36.2_{1.4}$
$RoBERTa_{freeze}$	$33.8_{2.1}$	$36.6_{1.6}$	$38.6_{1.5}$	$36.7_{1.5}$
kNN	$30.1_{0.4}$	$31.3_{2.1}$	$30.6_{1.1}$	$30.9_{0.4}$
Log Reg	34.4 _{2.3}	$38.3_{2.5}$	$40.0_{2.0}$	$37.9_{1.6}$
LaGoNN	$33.8_{1.8}$	$38.3_{1.3}$	$40.6_{0.6}$	$38.1_{2.0}$
Probe	$32.1_{1.9}$	$35.2_{1.4}$	$37.2_{2.5}$	$35.2_{1.7}$
$LaGonn_{cheap}$	$31.9_{1.9}$	$36.0_{1.0}$	$37.5_{2.5}$	$35.7_{1.8}$

Method		Insincere Questions		
Extreme	1^{st}	5^{th}	10^{th}	Average
$RoBERTa_{full}$	$19.9_{8.4}$	$30.9_{7.9}$	$42.0_{7.4}$	$33.5_{6.7}$
SetFit	$24.1_{6.3}$	$29.2_{6.7}$	$36.7_{7.3}$	$31.7_{3.4}$
$LaGoNN_{exp}$	30.7 _{8.9}	$37.6_{6.1}$	$39.0_{6.1}$	$36.1_{2.3}$
SetFit _{lite}	24.16.3	38.1 _{6.3}	41.1 _{6.5}	35.6 _{5.5}
${\rm LaGoNN}_{lite}$	$30.7_{8.9}$	41.8 _{8.3}	$43.4_{8.5}$	$39.3_{4.4}$
RoBERTa _{freeze}	19.98.4	34.1 _{5.4}	$37.9_{5.2}$	$32.5_{5.4}$
kNN	$6.8_{0.4}$	$15.9_{3.4}$	$16.9_{4.3}$	$14.4_{3.0}$
Log Reg	$24.1_{6.3}$	$31.7_{4.9}$	$36.1_{5.4}$	$31.8_{3.6}$
LaGoNN	$30.7_{8.9}$	$39.3_{4.9}$	$41.2_{4.7}$	$38.4_{3.0}$
Probe	24.3 _{8.4}	39.8 _{5.6}	44.84.2	38.3 _{6.2}
LAGONN _{chean}	$23.6_{7.8}$	$40.7_{5.9}$	$45.3_{4.4}$	$38.6_{6.6}$

Table 22

Table 21

A.3 Additional results for second experiment

822

823

824

825

827

829

832

833

835

837

838

839

841

842

844

845

846

847

Here we provide additional results from our second set of experiments that, due to space limitations, could not be included in the main text. We note that a version of LAGONN outperforms or has the same performance of all methods, including our upper bound RoBERTa_{full}, on 60% of all displayed results, and is the best performer relative to Sentence Transformer-based methods on 65%. This excludes LAGONN_{cheap}. This method showed strong performance on the Insincere Questions dataset, but hurts performance in other cases. In cases when SetFit-based methods do outperform our system, the performances are comparable, yet they can be quite different when LAGONN-based methods are the strongest. Below, we report the mean average precision $\times 100$ for all methods over five seeds with the standard deviation, except in the case of Hate Speech Offensive, where the evaluation metric is the macro-F1. Each table shows the results for given dataset and a given label-balance distribution on the first, fifth, and tenth step followed by the average for all ten steps. Liar appears to be the most difficult dataset for all methods. This is expected because it likely does not include enough context to determine the truth of a statement.

Method Imbalanced	1^{st}	$\begin{array}{c} \textbf{Insincere Questions} \\ 5^{th} \end{array}$	10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	39.8 _{5.5} 43.7 _{2.7} 44.5 _{4.5}	$53.1_{4.6} $ $52.2_{1.9} $ $52.7_{2.4} $	$55.7_{1.2}$ $53.8_{0.9}$ $55.4_{2.0}$	$50.6_{4.4}$ $51.4_{2.9}$ $51.8_{3.0}$
SetFit _{lite} LAGONN _{lite}	43.7 _{2.7} 44.5 _{4.5}	52.9 _{2.6} 53.5 _{2.7}	55.8 _{1.8} 55.9 _{2.4}	52.2 _{3.4} 52.6 _{3.5}
RoBERTa _{freeze} kNN Log Reg LAGONN	39.8 _{5.5} 23.9 _{2.2} 43.7 _{2.7} 44.5 _{4.5}	$44.1_{3.6} \\ 30.3_{3.0} \\ 47.6_{1.6} \\ 48.1_{2.2}$	$46.3_{2.4} \\ 31.6_{2.4} \\ 50.1_{2.1} \\ 50.3_{1.7}$	$44.0_{2.0} \\ 30.0_{2.1} \\ 47.6_{1.8} \\ 48.1_{1.9}$
Probe LAGONN _{cheap}	$40.4_{4.2} \\ 40.8_{4.3}$	$49.4_{2.3} \\ 51.1_{2.4}$	$52.3_{1.7} $ $54.5_{1.4}$	$49.0_{3.3} \\ 50.4_{4.0}$

Table 23

Method <i>Moderate</i>	1^{st}	$\begin{array}{c} \textbf{Insincere Questions} \\ 5^{th} \end{array}$	10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	48.1 _{2.3} 48.9 _{1.7} 49.8 _{1.6}	$54.7_{1.9} 53.9_{0.7} 52.2_{1.9}$	$57.5_{1.5} 54.2_{1.5} 53.2_{3.3}$	$53.9_{2.9}$ $52.3_{1.6}$ $52.0_{1.4}$
SetFit _{lite} LAGONN _{lite}	48.9 _{1.7} 49.8 _{1.6}	56.5 _{1.4} 56.1 _{2.8}	58.7 _{0.6} 58.3 _{1.5}	55.0 _{3.5} 54.6 _{3.5}
RoBERTa _{freeze} kNN Log Reg LAGONN	$48.1_{2.3} \\ 28.0_{2.4} \\ 48.9_{1.7} \\ \textbf{49.8}_{1.6}$	$50.2_{2.2}$ $33.9_{2.8}$ $53.6_{1.9}$ $54.4_{1.3}$	$52.0_{1.4} \\ 33.6_{2.0} \\ 55.8_{1.7} \\ 56.9_{0.5}$	$50.2_{1.4} \\ 33.5_{1.9} \\ 53.3_{2.2} \\ 54.2_{2.2}$
Probe LAGONN _{cheap}	$45.7_{2.1} \\ 45.7_{2.2}$	$52.3_{1.8}$ $54.4_{1.6}$	$54.4_{1.1}$ $56.4_{0.6}$	$51.4_{2.5}$ $53.2_{3.2}$

Table 24

Method Balanced	1^{st}	$\begin{array}{c} \textbf{Insincere Questions} \\ 5^{th} \end{array}$	10^{th}	Average
RoBERTa _{full} SetFit LAGONN _{exp}	$47.1_{4.2} 43.5_{4.2} 42.8_{5.3}$	$52.1_{3.6} 47.1_{4.6} 47.6_{2.9}$	$55.7_{2.6} \\ 48.5_{3.9} \\ 47.0_{1.7}$	$52.5_{2.9} 48.0_{1.7} 46.2_{2.0}$
SetFit _{lite} LAGONN _{lite}	$43.5_{4.2} \\ 42.8_{5.3}$	54.6 _{2.4} 53.5 _{3.7}	59.6 _{0.9} 58.6 _{2.5}	$53.6_{5.8}$ $52.2_{6.4}$
RoBERTa _{freeze} kNN Log Reg LAGONN	$47.1_{4.2} 22.3_{2.3} 43.5_{4.2} 42.8_{5.3}$	$52.1_{0.4}$ $30.2_{2.3}$ $53.8_{2.2}$ $54.1_{2.9}$	$53.3_{1.1} \\ 30.9_{1.8} \\ 55.5_{1.6} \\ 56.3_{1.3}$	$51.5_{2.1}$ $29.5_{2.5}$ $52.8_{3.5}$ $53.4_{3.7}$
Probe LAGONN _{cheap}	47.5 _{1.6} 49.3 _{2.6}	$52.4_{1.7}$ $54.4_{1.4}$	$55.3_{1.1}$ $57.6_{0.7}$	52.2 _{2.5} 54.2 _{2.7}

Method Balanced	1^{st}		10^{th}	Average
$RoBERTa_{full}$	$73.6_{2.1}$	$78.6_{3.9}$	$82.4_{1.1}$	$78.9_{2.2}$
SetFit	$73.8_{4.4}$	$69.8_{4.0}$	$64.1_{4.6}$	$69.6_{3.6}$
$LaGonn_{exp}$	76.0 $_{3.0}$	$73.4_{2.6}$	$72.3_{2.9}$	$72.5_{3.4}$
SetFit _{lite}	73.84.4	80.4 _{1.8}	82.4 _{0.8}	$78.3_{4.3}$
$LaGonn_{lite}$	76.0 _{3.0}	$80.0_{1.3}$	$82.5_{0.9}$	$79.2_{3.2}$
$RoBERTa_{freeze}$	$73.6_{2.1}$	$76.8_{1.6}$	$77.9_{1.0}$	$76.5_{1.3}$
kNN	$41.7_{3.4}$	$57.9_{3.3}$	$58.3_{3.3}$	$56.8_{5.1}$
Log Reg	$73.8_{4.4}$	$79.2_{1.9}$	$80.1_{1.0}$	$78.6_{1.8}$
LaGoNN	76.0 $_{3.0}$	$80.1_{2.0}$	$81.4_{1.1}$	79.8 _{1.4}
Probe	52.43.4	64.7 _{2.5}	$67.5_{0.4}$	63.44.4
$LaGonn_{cheap}$	$48.1_{3.4}$	$62.0_{2.0}$	$65.3_{0.8}$	$60.5_{5.0}$

Table 29

Table 25

Method Extreme	1^{st}	$\begin{array}{c} \textbf{Amazon Counterfactual} \\ 5^{th} \end{array}$	10^{th}	Average
RoBERTa _{full} SetFit LAGONN _{exp}	21.8 _{6.6} 22.3 _{8.8} 26.1 _{17.5}	63.9 _{10.2} 64.2 _{3.3} 68.4 _{4.4}	72.3 _{3.0} 68.6 _{4.6} 74.9 _{2.9}	59.6 _{16.8} 56.8 _{14.9} 63.2 _{16.7}
$SetFit_{lite}$ $LaGoNN_{lite}$	22.3 _{8.8} 26.1 _{17.5}	$62.4_{5.1} \\ 68.3_{4.3}$	67.5 _{5.2} 68.9 _{4.3}	56.5 _{14.7} 60.6 _{15.1}
RoBERTa _{freeze} kNN Log Reg LAGONN	21.8 _{6.6} 10.3 _{0.2} 22.3 _{8.8} 26.1 _{17.5}	$41.0_{12.7} \\ 15.3_{4.2} \\ 32.4_{11.5} \\ 31.1_{19.4}$	$51.3_{10.7} \\ 18.4_{3.7} \\ 42.3_{8.8} \\ 33.0_{19.1}$	$40.6_{8.9} \\ 15.6_{2.4} \\ 34.5_{5.9} \\ 30.9_{2.3}$
Probe LAGONN _{cheap}	$24.2_{9.0}$ $20.1_{6.9}$	$46.3_{4.4} \\ 38.3_{4.9}$	54.6 _{2.0} 47.8 _{3.4}	$45.1_{10.3} \\ 38.2_{9.5}$

Table 26

Method Imbalanced	1^{st}		10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	68.2 _{4.5} 72.0 _{2.1} 74.3 _{3.8}	81.0 _{1.7} 78.4 _{2.8} 80.1 _{1.4}	82.2 _{1.0} 78.8 _{1.2} 79.0 _{1.6}	$79.2_{3.9} 78.0_{2.1} 79.5_{1.9}$
$SetFit_{lite}$ $LAGONN_{lite}$	72.0 _{2.1} 74.3 _{3.8}	$79.1_{1.4} \\ 79.2_{1.7}$	$81.6_{1.3} \\ 81.9_{1.1}$	$79.1_{2.7}$ 80.2 _{2.2}
RoBERTa _{freeze} kNN Log Reg LAGONN	$68.2_{4.5} \\ 51.0_{4.1} \\ 72.0_{2.1} \\ \textbf{74.3}_{3.8}$	$75.0_{2.2} \\ 60.0_{3.1} \\ 74.4_{2.3} \\ 76.1_{3.6}$	$77.0_{2.4} \\ 61.3_{2.1} \\ 76.7_{1.8} \\ 77.3_{3.2}$	$74.2_{2.6} 59.7_{3.0} 74.8_{1.4} 76.1_{1.0}$
Probe LAGONN _{cheap}	$46.6_{2.8} \\ 38.2_{3.2}$	$60.3_{1.4} \\ 55.3_{1.8}$	$64.2_{1.2} \\ 61.0_{1.2}$	$59.2_{5.2}$ $54.4_{6.7}$

Toxic Conversations 5^{th} Method 10^{th} Average Extreme $\overline{21.2}_{3.7}$ $\overline{33.8}_{5.5}$ $\overline{21.9}_{9.3}$ $RoBERTa_{full}$ $7.9_{0.5}$ SetFit $18.1_{3.4}$ $24.7_{4.1}$ $8.8_{1.2}$ $17.6_{5.5}$ $26.4_{5.2}$ $LaGonn_{exp}$ $17.4_{6.6}$ $8.9_{1.7}$ $17.9_{6.0}$ $SetFit_{lite}$ $\overline{15.9}_{4.8}$ $18.0_{3.9}$ $14.9_{3.2}$ $8.8_{1.2}$ $LaGoNN_{lite}$ $8.9_{1.7}$ $16.1_{5.9}$ $19.8_{6.0}$ $15.5_{3.7}$ ${\bf RoBERTa}_{freeze}$ $12.8_{2.4}$ $19.1_{3.2}$ $7.9_{0.5}$ $13.5_{3.5}$ kNN $8.7_{0.4}$ $7.9_{0.0}$ $8.7_{0.2}$ $8.5_{0.3}$ $8.8_{1.2}$ $16.3_{3.0}$ Log Reg $13.1_{2.5}$ $13.0_{2.6}$ LaGoNN $8.9_{1.7}$ $13.8_{3.9}$ $17.1_{4.8}$ $13.4_{2.6}$ **24.6**_{2.6} 21.7_{2.7} Probe $13.1_{2.8}$ $30.1_{2.1}$ $23.9_{5.6}$ $27.4_{2.3}$ $LaGoNN_{\it cheap}$ $21.3_{5.3}$ $11.3_{2.2}$

Table 30

Table 27

Method Moderate	1^{st}		10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	73.9 _{2.5} 76.5 _{1.6} 78.6 _{2.2}	$80.0_{1.0} 77.0_{2.4} 78.0_{2.1}$	$80.1_{2.3} \\ 74.7_{0.5} \\ 76.3_{4.9}$	$79.1_{2.1} 76.5_{1.0} 78.2_{1.0}$
SetFit _{lite} LAGONN _{lite}	76.5 _{1.6} 78.6 _{2.2}	$80.4_{3.8} 80.8_{1.9}$	83.5 _{0.8} 83.1 _{0.7}	80.3 _{2.8} 81.0 _{1.7}
RoBERTa _{freeze} kNN Log Reg LAGONN	$73.9_{2.5} \\ 54.5_{3.1} \\ 76.5_{1.6} \\ \textbf{78.6}_{2.2}$	76.6 _{1.4} 64.2 _{1.9} 80.6 _{0.5} 81.2 _{1.4}	$78.5_{0.7} \\ 66.6_{1.3} \\ 81.2_{0.3} \\ 81.6_{1.1}$	$76.4_{1.7} \\ 64.7_{3.5} \\ 80.0_{1.4} \\ 80.8_{0.9}$
Probe LAGONN _{cheap}	$52.3_{2.0} \\ 47.3_{3.4}$	$64.1_{1.8} \\ 60.7_{1.5}$	$67.2_{1.4} \\ 65.2_{1.4}$	$63.1_{4.3}$ $59.5_{5.2}$

Method Imbalanced	1^{st}	Toxic Conversations 5^{th}	10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	24.1 _{5.6} 21.8 _{6.6} 22.7 _{9.8}	43.1 _{3.4} 44.5 _{4.1} 49.1 _{5.6}	52.1 _{2.5} 51.4 _{1.9} 53.4 _{2.3}	42.4 _{8.2} 42.1 _{9.3} 45.6 _{9.8}
SetFit _{lite} LAGONN _{lite}	$21.8_{6.6} \\ 22.7_{9.8}$	$41.4_{4.4} \\ 47.0_{6.3}$	$44.8_{3.1} \\ 50.2_{5.4}$	$39.0_{7.0} $ $43.7_{8.6}$
RoBERTa _{freeze} kNN Log Reg LAGONN	24.1 _{5.6} 11.5 _{2.5} 21.8 _{6.6} 22.7 _{9.8}	$31.2_{4.4} \\ 14.7_{4.0} \\ 26.7_{5.3} \\ 27.6_{8.9}$	$34.0_{4.0} \\ 15.3_{3.2} \\ 30.2_{4.0} \\ 30.3_{8.7}$	$30.5_{3.1} \\ 14.6_{1.1} \\ 26.6_{2.7} \\ 27.4_{2.4}$
Probe LAGONN _{cheap}	$23.3_{2.7} \\ 20.5_{3.2}$	$33.0_{2.8}$ $31.1_{3.2}$	$37.1_{1.8}$ $35.6_{1.8}$	$32.5_{4.2} \ 30.5_{4.6}$

Table 28 Table 31

Method <i>Moderate</i>	1^{st}	Toxic Conversations 5^{th}	10^{th}	Average	Method Imbalanced
RoBERTa _{full} SetFit LAGONN _{exp}	34.2 _{3.4} 33.6 _{2.9} 36.6 _{4.2}	45.51.9 47.22.2 48.22.7	52.4 _{3.3} 46.6 _{3.3} 49.9 _{3.7}	45.7 _{5.6} 44.3 _{4.3} 48.0 _{4.4}	RoBERTa _{full} SetFit LAGONN _{exp}
SetFit _{lite} LAGONN _{lite}	33.6 _{2.9} 36.6 _{4.2}	52.6 _{2.0} 56.1 _{1.5}	55.1 _{1.6} 57.7 _{1.4}	48.8 _{7.3} 52.3 _{6.8}	$\begin{array}{c} \textbf{SetFit}_{lite} \\ \textbf{LAGONN}_{lite} \end{array}$
RoBERTa _{freeze} kNN Log Reg LAGONN	34.2 _{3.4} 19.4 _{1.9} 33.6 _{2.9} 36.6 _{4.2}	$38.4_{2.1}$ $21.5_{3.4}$ $39.2_{2.9}$ $42.7_{3.7}$	$39.5_{1.8} 22.4_{2.9} 41.6_{2.7} 45.0_{3.5}$	$38.0_{1.5}$ $21.6_{0.8}$ $38.6_{2.4}$ $42.0_{2.5}$	RoBERTa $_{free}$ k NN Log Reg LAGONN
Probe	29.0 _{2.7}	36.1 _{1.2} 34.3 _{1.2}	39.1 _{1.5} 37.5 _{1.8}	35.5 _{3.3}	Probe LAGONN _{chee}

Method		Hate Speech Offensive		
Imbalanced	1^{st}	5^{th}	10^{th}	Average
$RoBERTa_{full}$	$50.6_{3.0}$	$65.2_{3.9}$	70.3 _{1.2}	$64.2_{5.3}$
SetFit	$54.4_{4.3}$	$66.3_{1.8}$	$68.9_{2.0}$	$64.3_{4.5}$
$LaGoNN_{exp}$	$57.0_{5.2}$	$67.0_{4.4}$	$69.8_{2.1}$	$64.9_{4.6}$
SetFit _{lite}	54.44.3	65.5 _{3.0}	$65.9_{3.5}$	63.53.9
$LaGoNN_{\it lite}$	$57.0_{5.2}$	66.6 _{2.6}	$66.6_{1.9}$	$64.3_{4.1}$
$RoBERTa_{freeze}$	50.63.0	54.1 _{1.6}	$55.3_{2.3}$	54.1 _{1.3}
kNN	$55.6_{4.8}$	$57.3_{2.3}$	$58.8_{3.6}$	$57.4_{1.1}$
Log Reg	$54.4_{4.3}$	$57.0_{3.9}$	$58.2_{3.8}$	$57.2_{1.1}$
LaGoNN	$57.0_{5.2}$	$58.2_{4.1}$	$58.3_{3.4}$	$58.3_{0.6}$
Probe	$46.5_{2.2}$	57.8 _{1.7}	$60.3_{1.2}$	$56.5_{4.5}$
$LaGonn_{cheap}$	$47.1_{1.3}$	$56.5_{2.2}$	$59.5_{2.5}$	$55.6_{3.8}$

Table 32 Table 35

Method Balanced	1^{st}		10^{th}	Average
	$32.3_{1.1}$ $35.7_{3.4}$ 40.4 _{4.4}	$42.7_{1.8} \\ 32.6_{6.2} \\ 40.2_{6.6}$	$54.1_{3.4} \\ 37.4_{2.7} \\ 39.8_{7.5}$	$43.8_{6.3} \\ 36.5_{1.9} \\ 40.0_{1.2}$
SetFit _{lite} LAGONN _{lite}	35.7 _{3.4} 40.4 _{4.4}	52.7 _{2.5} 52.9 _{2.6}	53.9 _{2.2} 54.0 _{2.3}	46.8 _{7.8} 48.3 _{6.4}
RoBERTa _{freeze} kNN Log Reg LAGONN	$\begin{array}{c} 32.3_{1.1} \\ 17.4_{0.8} \\ 35.7_{3.4} \\ \textbf{40.4}_{4.4} \end{array}$	$\begin{array}{c} 39.2_{1.5} \\ 23.7_{2.6} \\ 44.5_{2.9} \\ 46.6_{2.7} \end{array}$	$41.0_{0.6} 24.3_{2.7} 46.1_{2.8} 48.1_{2.2}$	$38.5_{2.4} \\ 23.1_{2.0} \\ 43.6_{2.9} \\ 46.1_{2.2}$
Probe LAGONN _{cheap}	$29.5_{2.4} \\ 26.8_{2.7}$	$35.9_{0.9} \ 34.5_{1.3}$	$40.2_{0.9} \\ 38.5_{0.8}$	$36.1_{3.5} \ 34.4_{3.7}$

Method Moderate	1^{st}		10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	61.9 _{3.4} 64.3 _{4.2} 63.8 _{4.9}	$70.8_{1.0} \ 70.6_{2.4} \ extbf{71.0}_{2.1}$	72.5 _{1.4} 72.4 _{0.5} 72.3 _{1.0}	69.9 _{3.2} 69.8 _{2.8} 70.0 _{3.0}
$SetFit_{lite}$ $LaGONN_{lite}$	$64.3_{4.2} \\ 63.8_{4.9}$	$70.3_{2.2} $ $70.7_{1.4}$	$71.2_{2.1} \\ 71.4_{1.0}$	$69.3_{2.3} \\ 69.4_{2.5}$
$\begin{array}{c} {\rm RoBERTa}_{freeze} \\ k{\rm NN} \\ {\rm Log~Reg} \\ {\rm LaGoNN} \end{array}$	61.9 _{3.4} 64.3 _{4.0} 64.3 _{4.2} 63.8 _{4.9}	$63.2_{4.1} 63.3_{2.9} 67.3_{3.2} 65.0_{5.3}$	$64.1_{4.5} \\ 63.9_{2.5} \\ 67.6_{2.3} \\ 66.7_{5.9}$	$63.2_{0.6} \\ 63.7_{0.4} \\ 66.9_{1.1} \\ 65.3_{0.9}$
Probe LAGONN _{cheap}	$55.6_{1.7} \\ 56.0_{3.6}$	$63.8_{0.8} \\ 62.2_{1.4}$	66.1 _{0.3} 66.0 _{0.9}	63.2 _{3.0} 62.3 _{2.9}

Table 33 Table 36

Method Extreme	1^{st}		10^{th}	Average
RoBERTa _{full} SetFit	$30.2_{1.4}$ $30.3_{0.8}$	43.5 _{2.5} 44.0 _{1.3}	51.2 _{2.2} 51.1 _{2.0}	44.3 _{7.4} 43.8 _{6.5}
$LaGoNN_{exp}$	$30.3_{0.7}$	$40.7_{2.9}$	$49.1_{4.4}$	$42.2_{6.2}$
SetFit _{lite} LAGONN _{lite}	$30.3_{0.8} \ 30.3_{0.7}$	$43.4_{2.5} \\ 40.9_{3.4}$	$45.5_{3.4} \\ 41.5_{4.8}$	41.6 _{4.6} 39.1 _{3.6}
RoBERTa _{freeze} kNN Log Reg LAGONN	$30.2_{1.4} \\ 31.5_{1.2} \\ 30.3_{0.8} \\ 30.3_{0.7}$	$\begin{array}{c} 33.5_{3.1} \\ 35.9_{2.7} \\ 38.4_{2.5} \\ 35.7_{2.6} \end{array}$	$34.4_{3.4} \\ 37.4_{2.0} \\ 41.1_{1.5} \\ 39.1_{2.4}$	$\begin{array}{c} 33.1_{1.4} \\ 35.8_{1.7} \\ 37.8_{3.3} \\ 35.6_{2.7} \end{array}$
Probe LAGONN _{cheap}	$29.0_{0.2} \\ 29.0_{0.1}$	$34.7_{1.5} \ 36.9_{1.8}$	$40.1_{2.1} \\ 40.5_{2.1}$	$35.1_{3.8}$ $36.2_{3.7}$

Method Balanced	1^{st}	Hate Speech Offensive 5^{th}	10^{th}	Average
RoBERTa _{full} SetFit	$59.7_{3.5}$ $60.7_{1.3}$	66.9 _{1.2} 66.3 _{1.6}	69.2 _{1.8} 67.5 _{0.9}	66.4 _{2.7} 65.9 _{2.2}
$LaGoNN_{exp}$	61.5 _{1.7}	$66.4_{1.4}$	$67.7_{0.9}$	$66.1_{1.8}$
SetFit _{lite} LAGONN _{lite}	60.7 _{1.3} 61.5 _{1.7}	$66.3_{2.0}$ $67.1_{1.1}$	66.5 _{0.9} 67.3 _{0.8}	$65.1_{1.7} \\ 66.0_{1.7}$
$\frac{\text{RoBERTa}_{freeze}}{k \text{NN}}$	59.7 _{3.5} 60.7 _{1.3}	$60.4_{2.7}$ $59.6_{2.8}$	$63.1_{2.3}$ $59.5_{2.5}$	$61.0_{1.3}$ $59.5_{0.5}$
Log Reg LaGoNN	60.7 _{1.3} 61.5 _{1.7}	$62.5_{0.7} \\ 62.8_{1.5}$	$63.4_{1.0} \\ 64.2_{1.0}$	$62.3_{1.0} \\ 63.0_{0.9}$
Probe LAGONN _{cheap}	$54.9_{1.4} \\ 54.2_{2.3}$	$58.5_{0.9}$ $58.6_{0.6}$	60.9 _{0.4} 60.6 _{0.5}	58.7 _{1.7} 58.5 _{1.8}

Table 34 Table 37

Method Extreme	1^{st}	$egin{aligned} \mathbf{Liar} \ 5^{th} \end{aligned}$	10^{th}	Average	Method Balanced	1^{st}	$egin{array}{c} \mathbf{Liar} \ 5^{th} \end{array}$	10^{th}	Average
RoBERTa _{full} SetFit LAGONN _{exp}	32.0 _{2.7} 31.2 _{3.8} 30.6 _{4.7}	34.7 _{2.9} 30.4 _{3.1} 30.3 _{2.0}	$\begin{array}{c} 35.1_{4.3} \\ 31.8_{2.9} \\ 31.3_{2.0} \end{array}$	33.71.0 31.50.7 31.10.6	RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	33.8 _{2.1} 34.4 _{2.3} 33.8 _{1.8}	39.4 _{2.4} 36.7 _{1.7} 34.2 _{2.7}	43.5 _{1.7} 37.0 _{1.3} 37.2 _{1.9}	40.2 _{3.2} 36.5 _{1.1} 36.2 _{1.4}
SetFit $_{lite}$ LAGONN $_{lite}$	$31.2_{3.8} \ 30.6_{4.7}$	$32.7_{3.8} \ 31.8_{3.9}$	$33.5_{4.2} \\ 32.4_{2.7}$	$32.7_{0.8} \ 31.6_{0.6}$	$SetFit_{lite} \ LaGoNN_{lite}$	$34.4_{2.3}$ $33.8_{1.8}$	$38.7_{2.3}$ $37.6_{2.0}$	$40.3_{2.8} \\ 39.4_{2.8}$	$38.0_{2.1} \ 37.2_{1.9}$
RoBERTa _{freeze} kNN Log Reg LAGONN	32.0 _{2.7} 27.0 _{0.5} 31.2 _{3.8} 30.6 _{4.7}	$32.8_{4.5} 27.3_{0.8} 33.7_{5.1} 32.0_{4.6}$	34.25.0 27.90.8 27.5.1 33.75.4	33.2 _{0.7} 27.4 _{0.3} 34.3 _{1.6} 32.6 _{0.9}	RoBERTa _{freeze} kNN Log Reg LAGONN	33.8 _{2.1} 30.1 _{0.4} 34.4 _{2.3} 33.8 _{1.8}	$36.6_{1.6} \\ 31.3_{2.1} \\ 38.3_{2.5} \\ 38.3_{1.3}$	$38.6_{1.5} \\ 30.6_{1.1} \\ 40.0_{2.0} \\ 40.6_{0.6}$	$36.7_{1.5} 30.9_{0.4} 37.9_{1.6} 38.1_{2.0}$
Probe LAGONN _{cheap}	$30.7_{2.0} \ 30.7_{2.0}$	$30.6_{3.9} \\ 30.5_{3.8}$	$31.7_{2.9} \ 31.4_{2.6}$	$31.1_{0.4} \\ 31.0_{0.4}$	Probe LAGONN _{cheap}	$32.1_{1.9} \ 31.9_{1.9}$	$35.2_{1.4} \\ 36.0_{1.0}$	$37.2_{2.5} \\ 37.5_{2.5}$	35.21.7 35.71.8

Table 38 Table 41

Method <i>Imbalanced</i>	1^{st}	$\begin{array}{c} \textbf{Liar} \\ 5^{th} \end{array}$	10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	31.4 _{3.2} 32.3 _{4.5} 32.3 _{4.6}	$35.8_{2.6} \\ 35.9_{3.1} \\ 35.7_{3.4}$	40.0 _{4.3} 36.4 _{2.2} 36.5 _{2.3}	$36.2_{2.4}$ $35.2_{1.1}$ $35.7_{1.4}$
SetFit _{lite} LaGoNN _{lite}	32.3 _{4.5} 32.3 _{4.6}	$35.6_{2.7} \\ 35.2_{2.4}$	$37.4_{2.6} \\ 36.6_{2.7}$	$35.8_{1.6} \ 35.5_{1.3}$
RoBERTa _{freeze} kNN Log Reg LAGONN	31.4 _{3.2} 27.0 _{0.2} 32.3 _{4.5} 32.3 _{4.6}	34.1 _{2.6} 28.5 _{1.0} 36.5 _{3.1} 34.9 _{2.2}	$\begin{array}{c} 35.6_{3.2} \\ 29.0_{1.0} \\ 38.5_{3.4} \\ 36.9_{2.5} \end{array}$	34.0 _{1.4} 28.7 _{0.7} 36.3 _{2.0} 35.3 _{1.4}
Probe LAGONN _{cheap}	$30.7_{3.0} \ 30.4_{3.0}$	$32.8_{1.8}$ $32.9_{1.8}$	$35.0_{1.6} \\ 35.4_{1.7}$	$33.5_{1.5} \ 33.5_{1.7}$

Table 39

Method <i>Moderate</i>	1^{st}	$\begin{array}{c} \textbf{Liar} \\ 5^{th} \end{array}$	10^{th}	Average
RoBERTa $_{full}$ SetFit LAGONN $_{exp}$	33.9 _{3.1} 33.0 _{2.6} 34.1 _{3.4}	38.4 _{2.7} 37.2 _{1.8} 38.7 _{2.3}	43.9 _{2.2} 38.7 _{1.5} 39.0 _{1.8}	39.5 _{3.0} 37.4 _{1.6} 37.8 _{1.5}
$\frac{\text{SetFit}_{lite}}{\text{LaGoNN}_{lite}}$	33.0 _{2.6} 34.1 _{3.4}	$38.5_{1.3}$ $38.4_{2.0}$	$40.4_{2.0} \\ 39.6_{1.5}$	$38.2_{2.1} \\ 37.9_{1.6}$
RoBERTa _{freeze} kNN Log Reg LAGONN	33.9 _{3.1} 29.2 _{0.8} 33.0 _{2.6} 34.1 _{3.4}	$\begin{array}{c} 35.3_{2.6} \\ 29.7_{1.5} \\ 37.2_{3.9} \\ 37.0_{3.1} \end{array}$	$\begin{array}{c} 36.8_{2.2} \\ 30.0_{0.6} \\ 39.4_{3.5} \\ 38.6_{3.0} \end{array}$	$35.4_{1.0} \\ 29.8_{0.3} \\ 37.0_{1.8} \\ 36.8_{1.3}$
Probe LAGONN _{cheap}	$31.6_{1.1}$ $31.4_{0.9}$	$34.7_{2.5}$ $35.3_{2.3}$	$37.0_{2.5} \\ 37.6_{2.0}$	$34.9_{1.7} \ 35.3_{1.9}$

Table 40

A.4 Examples of LAGONN modified text

WARNING: Some of the examples below are of an offensive nature. Please view with caution.

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

In this section, we provide examples of how $LaGoNN_{\it exp}$ modifies test text from the datasets we studied under the BOTH configuration. We choose this configuration because the information it appends from a NN in the training data to a test instance encapsulates both the LABEL and the TEXT configuration. LAGONN $_{exp}$ was trained under a balanced distribution and five examples per label were chosen randomly on the first, fifth, and tenth step to demonstrate how the same test instance might be decorated with different training examples as the training data grow. We recognize that some the images below are difficult to see and have made the .csv files available with our code and data files. Note that MPNET's separator token is </s>, not [SEP].

Test Modified	Gold Label
What rapper still relevant and popular today has the best rhyme schemes? (-) chinicre question 3.59947132110997? What would be a good nickname for Trump, Donald Dumbck, and President Spaniovich? (-)s- 'valid question 4.12427423845215> What are after class 12 courses in commerce stream to choose from 1 how completed my class 12 (exceeded 69) and all mits obtaines from the int on to do job).	valid question
Which books do you suggest to someone who get a free time and will help him stay motivated? <valid 3.9509353637695312="" question=""> What are the best online courses to learn data science? <instincere 4.300448417663574="" question=""> What are the more steps in Career Oriented Education?</instincere></valid>	valid question
How will by left if amence talks badly about Kunt? (y) Knisinere question 3.5003605308331735 Why are the Uit government and the media (especially the BBC and the Guardian) demonsing ordinary British people, manipulating buzz words like &Ceath-right4GE, &Ceatislamaphobia4CE, &Cearcist4GE to suppress [espitate activates and Whitelian grooming agrang? / yo-valid question 5.5000375308379335 How for Investin exist Yession?	valid question
Why is quite HYPP inherited */p' > valid question 4.055334999032715 \(c) a vol. quest	valid question
Itsour do the Valente Stewarts leather packets schlever their quality during the manufacturing process? y/o valid question 37911384084801914+ how are the Lancaster leather sofas manufactured? y/o valinancere question 4.3559441366467285) am an experienced programmer and in my high school my teacher tries to make me use printion or late 37.05 must me, puptions just a language for beginners, thereby making it not for me." 19 set and 10.01 do asynthing more to all of all on synthym soft and the	valid question
is Ariana Grande really as mean and bitchy as she seems? <insincere 3.572277549928655="" question=""> Why is Alia Bhatt so dumb? <valid 3.924571990966797="" question=""> Do you agree with Congressman Steve King's comments on immigrant children in detention centers?</valid></insincere>	insincere ques
Do you go you know that allers are real and all those satellities we send up in space work as a sort of tracking device for them so in a few years it will be to late for Earth? **(ye)* - cinimener question 3 5094439029993004 have you noticed how conservatives are capturing the English language and modifying the identification of policits would **?(e)* - will expect so 3.0094439029993004 have you noticed how conservatives are capturing the English language and modifying the identification of the source	insincere ques
is t politically incorrect to any female privilege, but it is a more accurate term to say, white female privilege? -(2) similar cequates on 3.2328300996219379 by year the UK government and the media (especially the 850 and the Guardian) demonstrage ordinary british people, manipulating buzz words like [decah-inplicate]. Generalized this purpose is perfectly a contract or incorrect to immigrant ordinary british indicates the contract ordinary ordinary british people, manipulating buzz words like [decah-inplicate]. Generalized to suppose in the contract ordinary ordinary british people, manipulating buzz words like [decah-inplicate].	insincere ques
On Mothershift Day, it it reasonable to reflect there is some truth in the unfashonable notion than women are more driven by emotion and men more driven by reason? 4/50 - instructors question 3,49204158782599. Why we the UK government and the media (especially the 8Ex and the Guardian) demonstraing ordinary british possible, immigrately called an instruction of contrage at Muslim grooming garps? 4/50 - valid question 3,7174020818278750 Do you agree with Congressman Sees we first groomers or immigrate childer in determinent centering contents or immigrate childer in determinent centering centering contents or immigrate childer in determinent centering ce	insincere ques
If the U.S. president is a role model, bit acceptable for children to say "go fuck yourself," and use the word "pusy" to describe women? (-5) relications a role model, in a receptable for children to say "go fuck yourself," and use the word "pusy" to describe women? (-5) relications are used to a relication of the pusy of the pu	insincere que

Figure 5: Insincere Questions, step 1.

Test Modified	Gold Label
What rapper still relevant and popular today has the best rhyme schemes? -instrucere question 3.873907949447832> What would be a good nickname for Trump, Donald Dumbok, and President Spankovich? -(s) - valid question 4.028938797454834> Why does Dancing with the Stars not include Bachata as one their drace style?	valid question
Which book do you suggest to omenone who get a first time and will help him stay motivated? **(s)> "wall question 3.608125897209795> What is a good degree to get at community college if you want to explore different subjects and figure out your career path? **(s)> "sinsincere question 3.8502604981399254> What is a good degree to get at community college if you want to explore different subjects and figure out your career path? **(s)> "sinsincere question 3.8502604981399254> What is a good degree to get at community college if you want to explore different subjects and figure out your career path? **(s)> "sinsincere question 3.8502604981399254> What is a good degree to get at community college if you want to explore different subjects and figure out your career path? **(s)> "sinsincere question 3.8502604981399254> What is a good degree to get at community college if you want to explore different subjects and figure out your career path? **(s)> "sinsincere question 3.8502604981399254> What is a good degree to get at community college if you want to explore different subjects and figure out your career path? **(s)> "sinsincere question 3.8502604981399254> What is a good degree to get at community college if you want to explore different subjects and figure out your career path? **(s)> "sinsincere question 3.8502604981399254> What is a good degree to get at community college if you want to explore the path of t	valid question
How will you feel if someone talks badly about Kuntt? <valid 3.535563163757324="" question=""> How do I stop feeling bad after a girl had a crush on me? <insincere 3.689171075820923="" question=""> Why Indian girls go crazy about marrying Shri. Rahul Gandhi Ji?</insincere></valid>	valid question
Why is equine HYPP inherited? <inslineere 3.6035702228546143="" question=""> Can female animals with male humans sex? <insline 3.7413032054901123="" question=""> How long do guinea pigs live for?</insline></inslineere>	valid question
House of the Valent's Exercise leather pixels as chieve their quality during the manufacturing process? /p> valid question 2,7472882270812999? Now a ret la unscater leather sofas manufactured? /p> <isnincere 3,94888477708992="" question=""> Why don't all Trump supporters buy only made in USA goods, e.g. many of them have their care of a distain/European companies, spin on pilescus where more than 70% of Herms are not made in USA, as simultaneous of control of their spin of the distain and the spin of the spin of their sp</isnincere>	valid question
is Arisana Grande really as mean and bitchy as she seems? <insincere 3.252298831939697="" question=""> Why is Alia Bhatt so dumb? How do I stop feeling bad after a girl had a crush on me?</insincere>	insincere questio
Do you gay, show that alters are real and all those statistics we send up in space work as a sort of tracking device for them so in a few years it will be too late for Earth? chaincere question 3 067336559256543> Isn't it obvious now that walking on the moon by the Americans was a hoak, because walking on this principle of the property of the statistics of the stat	insincere questio
at a politically incorrect to any firming privilege, but it is a more accurate term to say, while female privilege? > https://sev.nuist.org/linear/2018/2018/2018/2018/2018/2018/2018/2018	insincere questio
On Modified "So Day, 18 reasonable to reflect there is some truth in the unfalsionable notion than women are more driven by emotion and men more driven by reason?	f insincere questio
If the U.S. persident is a role model, it it acceptable for children to say "go fully oursel," and use the word "pussy" to describe women? — (Instruction of the pussy of t	insincere questio

Figure 6: Insincere Questions, step 5.

Test Modified	Gold Label
What rapper still relevant and popular today has the best rhyme schemes? <valid 3.7103171348571777="" question=""> What is the oldest fashion trends running yet? <instincere 3.871907949447632="" question=""> What would be a good nickname for Trump, Donald Dumbck, and President Spankovich?</instincere></valid>	valid question
Which books do you suggest to somence who get a free time and will help him stay mothwated? "vp" valid question 3.1401493178300569- How can I stay mothwated when learning something new? "cp" vinincere question 3.7235569417173293- I'm hungry and I'm too lazy too get out of bed, should I get a specification of the properties of th	valid question
How will you feel if someone talks badly about Kuntt? -kinsincere question 3.4893462657928467> Does Tamil isai Soundarajan support Vijayendra for disrespecting the Tamil Anthem? -kalid question 3.5355563163757324> How do I stop feeling bad after a girl had a crush on me?	valid question
Why is equine HYPP inherited? <valid 3.5067965984344482="" question=""> What disadvantages do animals that don't have bones face? <insincere 3.6035702228546143="" question=""> Can female animals with male humans sex?</insincere></valid>	valid question
House of the Valent Stevens leather jackets oblieve their quality during the manufacturing process? 4/3- vialid question 2,747288227081299> How are the Lancaster leather softs manufactured? 4/3- 4/insincere question 3,9087233543395996- Are Newyort cigarettes designed to selectively destroy black people's DNAP	valid question
A Arians Grander early as mean and bitchy as the seeme? C/2- X-raid question 3.183897782174879: like this girl who used to be quite rude and would run through boyfriends very fast. But now that school started again, she seems to have gotten a lot nicer throughout Summer. Is she faking her pollteness, and is it worth pursuing her? Cyc. Vinificance question 3.13586900000386979. Why 13 like hatter to do units?	insincere questio
Do you guy, know that allens are real and all those statilles we send up in space work as a sort of tracking device for them so in a few years it will be too late for Earth? - (yo - kninnerer question 3.087338559258543- inn't it obvious now that walking on the moon by the Americans was a hoas, because walking on the moon by the A	insincere questio
s it politically incorrect to say femile privilege, but it is a more accurate term to say, white female privilege? (y- <inspired 2,915,800,000,000,000,000,000,000,000,000,00<="" question="" td=""><td>insincere questio</td></inspired>	insincere questio
In MotherArt's Day, is it reasonable to reflect their is some truth in the unfashionable notion than women are more driven by emotion and men more driven by reason?	Insincere questio
The U.S. president is a role model, it is exceptable for children to say "go futu' yourset," and use the word "pussy" to describe women? -(>) -	

Figure 7: Insincere Questions, step 10.

est Modified	Gold Label
lings to the wall, desent flog around when a bag is pulled out, the mess of bags falling out is gone. <not-counterfactual 3.6492768002859289="" <="" after="" br="" hopes="" it="" it's="" keep="" shape="" that="" washing.="" will=""> "in a counterfactual 4.0123462877001959" \(\text{Will delivered this immediately invoid have given this product five time because it worked.""</not-counterfactual>	not-counterfactu
like these parts they all ow enough without being inappropriets when you sit or bend over <pre>/> counterfactual 3.0/2000009972540- "But odd in rough, the bottoms are a little too loose in the wast (37) and could have used another inch or two in the inseam (i normally take a 35\" or 36\" in jeans, depending to the brand if this height-end underswere the likes of which to you guillet get at .01, so, you might get at .01, so, you mi</pre>	not-counterfactu
te was very professional and wish all transactions I make through Amazon were this good. > <counterfactual 3.4319908519927=""> I wish I had had him as an instructor at college. > <not-counterfactual 4.054038895233154=""> I worried that it would be cheap or not fit orwhateverBut WOWI</not-counterfactual></counterfactual>	not-counterfactu
Well written with a twist iddin't appear. < non-counterfactual 3.8327973194122214> "The crossover from the characters from one novel to others keeps me interested; after all, ido hate to miss a Dee-Ann or Eggle(" appearance."" < non-counterfactual 3.8320303212402344> "\\ \ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.8327973194122214> "The crossover from the characters from one novel to others keeps me interested; after all, ido hate to miss a Dee-Ann or Eggle(" appearance."" < non-counterfactual 3.8320303212402344> "\ \ \ \ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.8320303212402344> "\ \ \ \ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.8320303212402344> "\ \ \ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.8320303212402344> "\ \ \ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.8320303212402344> "\ \ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.8320303212402344> "\ \ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.8320303212402344> "\ \ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.8320303212402344> "\ \ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.8320303212402344> "\ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.8320303212402344> "\ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.832030212402344> "\ \ \ Hold reviewed this memodiately invoid how given this growther factual 3.832030212402344> "\ \ Hold reviewed this memodiately invoid how given this growther factual 3.832030212402344> "\ \ Hold reviewed this memodiately invoid how given this growther factual 3.832030212402344> "\ \ Hold reviewed this memodiately invoid how given this growther factual 3.8320303212402344> "\ \ Hold reviewed this memodiat	not-counterfactu
loesn't feel like the quality levi's I am used to. <not-counterfactual 3.2773308753967285=""> However, the fabric is not that great, it's cheap scratchy cotton. <counterfactual 3.746659755706787=""> The blanket is nice and soft but it is white, so if it doesn't light up it isn't much use!</counterfactual></not-counterfactual>	not-counterfactu
I we had will study. I believe the enfosed hardware would have been sufficient. Counterfactual 3,438643550872805 worked that it would be cheap or not fit it	counterfactual
this ever turns into a film. Hope they do it justice! (*) contounterfactual \$5.2915239341045b* "The crossover from the characters from one novel to others keeps me interested; after all, I do hate to miss a \Dee-Ann or Eggie\" appearance."" (*)> counterfactual \$3.751143217086792> \\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\	counterfactual
I you don't want a prominent diplay this rack is too large for most bed or living rooms, it is wafer and taller than my tall Broynill wardonce style demanded by the rack is too large for most bed or living rooms, it is wafer and taller than my tall Broynill wardonce style demanded by the state is reported by the room that it is solven, and a side by the room to the solvent in the state of the	counterfactual
wish I could have seen all of the places he recommends I <counterfactual 3.5627076625823975=""> I wish I had had him as an instructor at college. <not-counterfactual 4.141315937042236=""> I worried that it would be cheap or not fit orwhateverBut WOWI</not-counterfactual></counterfactual>	counterfactual
with record replace just that small stupid opecs, since there's nothing wrong with the rest of the hose assembly <counterfactual 3.0597372093200684=""> I wish the storage compartment was a little bigger and opened up instead of ididding on and off. -counterfactual 4.048873311187744> I worned that would be cheap or not fill the counterfactual 4.048873311187744> I worned that would be cheap or not fill the counterfactual 4.048873311187744> I worned that would be cheap or not fill the counterfactual 4.048873311187744> I worned that would be cheap or not fill the counterfactual 4.048873311187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.04887331187744> I worned that would be cheap or not fill the counterfactual 4.0488733118744> I worned that would be cheap or not fill the counterfactual 4.048873311874> I worned that would be cheap or not fill the counterfactual 4.048873311874> I worned that would be cheap or not f</counterfactual>	it counterfactual

Figure 8: Amazon Counterfactual, step 1.

Test Modified	Gold Label
Clings to the wall. decein 10p around when a bag is pulsed out, the mess of lags failing out is goine. <a (37)="" (i="" 35\"="" 36\"="" a="" and="" another="" are="" being="" bottoms="" brand="" but="" could="" depending="" enough,="" have="" height="" href="https://www.chi.org/brights/br</th><th>not-counterfact</th></tr><tr><td>like these jeans they at low enough without being inappropriate when you air or bend over. ()>> counterfactual Z.85083830835033: " if="" in="" inch="" inseam="" is="" jeans,="" little="" loose="" normally="" not="" oddly="" of="" of<="" on="" or="" property="" take="" td="" the="" this="" too="" two="" used="" vasit=""><td>not-counterfact</td>	not-counterfact
te was very professional and with all transactions i make through Amazon were this good. - not-counterfactual 3.459845054229738> Had the person handling the shipping of this tembers at all concerned with the use of the product at the end of the mailing process. the slightest that of cancel on the benear half of concerned with the use of the product at the end of the mailing process. the slightest that of cancel on the benear half of concerned with the use of the product at the end of the mailing process. the slightest that of cancel on the benear half of concerned with the use of the product at the end of the mailing process. the slightest that of cancel on the benear half of the process and the process a	not-counterfact
Well written with a twist I didn't expect. <not-counterfactual 2.8573162746429443=""> Fun read Could have been a little longer with more detail.</not-counterfactual>	not-counterfact
Decent Ted like the quality levis! an used to, v/a> counterfectual 2.738877182008850 it has the same great comfortable & flittering features place the great denim texture that Lee has perfected-innoorthing and stretchy without the excessive cling-but! think it must have been designed for people who have a preter unufue of their features place and the preter unufue of their features place and the preter unufue of their features place.	not-counterfact
'we had wall study. I believe the enclosed hardware would have been sufficient. (*/p> con-counterfactual 2.683815446777344- it was a little trickly to find the center of the study using my stud finder but once i felt comfortable with the lines i had drawn, I drilled the pilot holes and boiled this thing to the wall.	counterfactual
this ever turns into a film, I hope they do it justice! <not-counterfactual 2.671574354711753=""> I read this book because of the motion picture that is coming out soon. <counterfactual 3.1458709239959717=""> Was a good story, though there could have been more to it.</counterfactual></not-counterfactual>	counterfactual
you don't want a prominent display this rack is too large for most bed or living rooms, it is voider and saler than my sall Broynill waredone style deress which was the largest piece in the room until this shoe rack. (2) excounterfactual 2.7335783489938491 buought this mount because i wanted one that would sit of the student in the student intended to the student	n counterfactual
wish I could have seen all of the places he recommends I <counterfactual 2.799947738647461=""> I wish I had had him as an instructor at college. <not-counterfactual 3.3013432025909424=""> And as the ole man isn't any version of slender it was good that he got to try on some shirts before hand.</not-counterfactual></counterfactual>	counterfactual
wish rout register but that amal study discs, since there's nothing wrong with the rest of the hose assembly. <	counterfactual

Figure 9: Amazon Counterfactual, step 5.

Test Modified	Gold Label
Clings to the wall, doesn't flow around when a bag is juilled out, the mess of bags falling out is gone >> not-counterfactual 3.12405693385957- And the dvd cases were tightly packed to ensure they ddn't move around. <-counterfactual 3.289605140686035> if I had to come up with anything regative, I would see say that the attachment don'eff't seen to skep mit he vicunic flower when not in us were to properly!	not-counterfactu
like these jeans they sit low enough without being inappropriate when you sit or bend over. <not-counterfactual 2.447404623031616=""> These shorts fit really well and look good too. <counterfactual 2.550638198852539=""> The top fits great just wish the bottoms fit too.</counterfactual></not-counterfactual>	not-counterfactu
He was very professional and with all transactions I make through Amazon were this good. 4/5 -notic counterfactual 3.299117891792178- This new speaker was just what the doctor ordered and I couldn't be more pleased. 4/5 -	not-counterfactu
Well written with a twist riddin't spect. < root-counterfactual 2.597446092950035" All bit workmankler, not up to Lord's high standard of /A Might to Remember," but well-detailed, and a story hat not many now innov." <p></p>	not-counterfactu
Doesn't feel like the quality levi's I am used to. <counterfactual 2.5612902641296387=""> I was hoping the pants would be thicker but being that it's not too expensive it's understandable. <not-counterfactual 2.572395086288452=""> But it doesn't have a lining like the last couple models I bought.</not-counterfactual></counterfactual>	not-counterfactu
If we had wall studis, I believe the enclosed hardware would have been sufficient, <	counterfactual
If this ever turns into a film, Those they do It justice (1/9 controunterfactual 2.671374354171753-) read this book because of the motion picture that is coming out soon. counterfactual 3.141676187515259- Wish this story would have been longer and turned into a book, with some gut werenching action, low/hate lovers causering scene, with a happy ending at the end.	counterfactual
If you don't want a prominent display this rack is too large for most bed or living rooms, it is wafer and faller than my all Broynill wardrobe giving dessers which was the largest piace in the room until this shoe rack. 4/2 - *Counterfactual 2.753/958484895846=1 bought this mount because I wanted one that would sit on three quisis intends on the heavy to be because my try is guide heavy and loved unknown that to use it.	counterfactual
wish I could have seen all of the places he recommends! < counterfactual 2.7999041080474854) wish I had had him as an instructor at college. <not-counterfactual)="" 3.2604622840881348="" a="" afordable="" few="" hats="" him="" i="" loosing.<="" mind="" order="" td="" to="" wanted="" wouldn't=""><td>counterfactual</td></not-counterfactual>	counterfactual
I visit I could replace just that ramal studied joice, since there's nothing worse to be assembly. (5) *Contourierfactual 2.47003547541 visit I could replace just that ramal studied joice, since there's nothing worse worse the season of the	/ counterfactual

Figure 10: Amazon Counterfactual, step 10.

est Modified	Gold Label
isomos demand that you accept their fudge asking. But in one of ure ever will-(-)/- Nont taxic 3.5918/21005249023> Sounds just wirkly, working for the state that is, So it begs the question, why work for the state if the pay is no bad versus the private sector? Seems logical to just make the writch?	not toxic
don't think amyone likes this health care bill, it sthick for everyone. So years and older are going to get hammered with higher premiums. People with empressing conditions will also see their premiums go through the roor Eventually no one will be able to afford it. They're still not addressing the reason and the outdoor. Prescription from your May war and remains appaig 2 to 3 times higher for prescription drugs than canadiscans are? They plouded is down with Democratic hist not play muscled intensive thou health risk with on the history. Here is the sad that ALAPA pot out, and in that a fact check at the end. It's worth watching alto intellectually benefit and intellectually benefit and samples and samp	not toxic
ather than call you a link (1) just post a link to HART's helicopter footage of this "most densely populated" agland. https://www.youtube.com/watch?vsopitsWWVDdw And if you've been following a little closer, you'd know about the luxury towers and boutque hotels that developers are planning on fulfilling their TOD requirement with "cash gifts" to the city instead of actual affordable TOD. And if you've been following a little closer, you'd know about the luxury towers and boutque hotels that developers are planning on fulfilling their TOD requirement with "cash gifts" to the city instead of actual affordable TOD. And if you've been following a little closer, you'd know about the luxury towers and boutque hotels that developers are planning on fulfilling their TOD requirement with "cash gifts" to the city instead of actual affordable TOD. And If you've been following a little closer, you'd know about the luxury towers and boutque hotels that developers are planning to the p	not taxic
et the "real news" from Breithant and Drudge Report? Alahahahahah lar you for real? (yo - not toxic \$5594274897711185 Sounds just swful, working for the state that is. So it begs the question, why work for the state if the pay is so bad versus the private sector? Seems logical to just make the switch? - toxic 20232066293535165 vious disagrant is a miserable existence. Fine given a investigate of the state if the pay is so bad versus the private sector? Seems logical to just make the switch?	not toxic
usd live to see one of those American Boats with all the planes on it at the bottom of the sea. <i> root toxic 3.82891183531494> Sounds just avoid, working for the state that is. So it begs the question, why work for the state if the pay is so bad versus the private sector? Seems logical to just make the switch? <i> root 2.414443175799505 vooid imaging it is a insertable existince. Being a winning per losers in on way to per towards it.</i></i>	not toxic
witf, nowhere in scripture does it say that Elizabeth was a kid when married to Zachariah and nowhere does it say that Many was a child when married to Joseph. That's a complete and utter fabrication. In fact, it says in Luke 1-12, "And the spake out with a loud voice, and said, Blessed art thou among women, and lessed is the finit of thy wom's. It says blessed art thou among WOMEN, "Cy-not back, 5,6674931532715" Sounds just workly working for the state that is. So it begs the question, why work for the state if the pay is so bad versus the private sector? Seems logical to just make the switch?	toxic
rigist Markel and all other European political leaders with have aided and abetted the ongoing invasion of Europe by the forces of the crescent moon death cult should be tried as accessories to Mr. Urban's murder. Vision 1. Special Section 1. Special S	toxic
hope you don't have kids if you see this woman's actions as acceptable. And I applaud the den for kicking the kid out. She brought unwanted negative attention upon them. However, she will, and is already likely, pay the the price for her stupid stunt. <to-toxic 3.0406124591827393=""> Calling evenybody that disagrees that you are action and another than the price for her stupid stunt. <to>toxic 3.0406124591827393> Calling evenybody that disagrees that you are action and another than the price for her stupid stunt. <to>toxic 3.0406124591827393> Calling evenybody that disagrees that you are acceptance to those who determinedly violate Scripture. In fact, he often spoke of hell.</to></to></to-toxic>	toxic
one care; what a paid liberal trolling back like you believes lunatio. 4/5 > chack 2.8411785555438995 Calling even body that disagrees with you as casts and and M-Multiple and you great you will be a consist and and M-Multiple and you great you g	toxic
It all you NOP "LEAP" manifesto types, where is your knew Naomil Glen? Her favoring adoration of Chavez and Venezuelan thuggery knows no bounds. I'm sure he's a unfully hysterical over the thought that such a pathetic dictstorship could ever be sanctioned.	s

Figure 11: Toxic Conversations, step 1.

	G
omos demand that you accept their fudge packing. But none of us ever will: 4/2> Qook 3.18849530810647> So you admit you would exterminate inferior humans. 4/2> 4not took 3.2964952716827393> Mark Mackinnon and the interests he work for would like us to 'get used to it', because they don't want to	
thing practicat to stop it. or in this paymont list in health care bill. It attitud for evenuous. So vears and older are point to set harmoned with histher cremiums. Peccle with previousne conditions will also see their premiums so through the roof. Eventually no one will be able to afford it. They've still not addressine the reason and the	no
non-Personal register, Way we American spring 2 to 3 three higher for prescription drugs the for a considerar set 7 may please at some with Democrates and full representations and the control of the prescription of the prescri	
Legacitions will lead of further rino deep dest.	ıd
at's if Ryan and McConnell can pass this huge sack of excrement.	no
their than call you a list, (II just post a link to HART's helicopter footage of this "most densely populated" agland.	
ttps://www.youtube.com/watch?wopissNWWOdw	
not if you're been following a titled closes, you'd know about the lawuy towers and boudge includes part developes are planning on fulfilling their TOO requirement with "cash gibt" to the city instead of actual affordable TOO. (1) o root tow's COSSSEREALISSERS DO NOT their and use about bould the read. Loss on many discount but me in their planning our prior, it is also that their planning of their planning our prior, it is also their planning our prior, it is closed as a consistency of their planning our prior, it is also their planning our prior, it is closed as a consistency of their planning our prior, it is also their planning our planning our prior, it is also the planning our planning o	
t the "real news" from breitbart and Drudge Report? Altehalhahahal Are you for real? "In one read the Dispatch one would think Trump is the most evil person on the planet."	
est only, but does, and the would be ask up by give in the Perkinson a pass of the "was art "POUL". The "but of give in the "man or give in the "	
ter the election. Something to do with their choice of "information" sources no doubt.	n
I to to see one of from a remarkan floar that which all the planes on it at the bottom of the year. (Ip- You'd, 3,800054459581377) less Regent Seens floars will never offer Mr Hammond another trip, one what or the planes, some which contains a containing to the terminals, teasures, some which do not seen that the planes is assume, some which contains the planes is assumed to the planes in the planes is assumed to the planes in the planes is assumed to the planes in the planes is assumed to the planes is assumed to the planes in the planes in the planes is assumed to the planes in the planes is assumed to the planes in the planes is assumed to the planes in the planes in the planes in the planes is assumed to the planes in the pla	
sets available on his "resulte deck" in "Tallew Off his menaining transverse was recommended from the commended and ordered and the commended and ordered and ord	
e's a total ingrated onct toxic 3.409156084060669 Now replaced by the savy EA-18Q Growler I using a preswitting Military Operating Areal Get over it!!!!!	n
for containing in incidence desired in particul desired in any section of containing in incidence and incidence an	
Rodieware then of useless grumbling, and keep your tongue from stander; because no secret word is without result, and a lying mouth destroys the soul \$45 (Wisdom 1:11)	
that is the case, then Trumpi's soul was utterly destroyed decades ago.	to
get Merkia hal of other Conspanse political landers who have a sided and on detect the copping manager of the research more death cults should be trief as accessories to No. Under 1 murder. (3) or Maria 1 2007/3827258273 http: what happens when you better the popular of the conspanse of more under the research more of the more and the second of the s	
e Libs will be happy to let this die because Monsef is now a very poor salesman given her own immigration dishonesty. That said if the election prospects sour significantly for the Libs I have no doubts that PM buts will ram through Ranked Ballot	t
ope you don't have let if you we this wemman storms as acceptable, and in aground the der for kinding the kind on the hought member degrish antendring your member degrish and and in your bear the properties of	
	to
a sild so far off the creatives had a danger to herself and to others.	
sold of fir off the concer pile is disager to herself and to others. Intelligence is a concernance is a disager to herself and to other and intelligence is a concernance in the deep and into residual extension is sold in content in the concernance is a concernance in the concernance in the concernance is a concernance in the concernance is a concernance in the concernance is a concernance in the concernance in the concernance is a concernance in the concernance in the concernance is a concernance in the concernance is a concernance in the concernance is a concernance in the concernance in the concernance is a concernance in the concernance in the concernance is a concernance in the concernance is a concernance in the concernance is a concernance in the concernance in the concernance is a concernance in the concernance in the concernance is a concernance in the concernance in the concernance in the concernance is a concernance in the	
as lists for first conserve its charger to heard and to other. It is not to be presented to the proof of the group of the deep exist into redict extensity territory one create what per post of the deep exist into redict extensity territory one create what per post of the deep exist into redict extensity territory one create what per post of the deep exist into redict extensity territory terri	sed
as also for the conserve's a stagent broad and to cross. The conserve is the conserve is a stagent broad and to cross. The conserve is a stagent broad and to cross. The conserve is a part of the conserve is a stagent broad and the conserve is a part of the conserve is a part	sed
as distant for fine consert was during the horself and to content. As a post of the post of the consert was a start of the content to content	sed

Figure 12: Toxic Conversations, step 5.

Test Modified	Gold Label
rooms demand that you accept their fudge packing. But none of us ever will- (7)>- Nonic 3.188312844570512> So you admit you would esterminate inferior humans. - Ynot toxic 3.2884285144805812> Mark Mackinnon and the interests he work for would like us to 'get used to it', because they don't want to do	
applying goodard in stopic. don't think applied like stopic. don't think applied like stopic like stopic will be able to lift, it striks for everyone. So years and older are going to get harmened with higher premiums. Repole with prevailing conditions will also see their premiums go through the conf. Eventually no one will be able to affect it. They're still not addressing the reason and the indicated and the confidence of the stopic will be able to affect it. They're still not addressing the reason and the indicated and the stopic will be able to affect it. They're still not addressing the reason and the indicated and the stopic will be able to affect it. They're still not addressing the reason and the indicated and the stopic will be able to affect it. They're still not addressing the reason and the indicated and the stopic will be able to affect it. They're still not addressing the reason and the provided and the stopic will be able to affect it. They're still not addressing the reason and the provided and the stopic will be able to affect it. They're still not addressing the reason and the provided and the stopic will be able to affect it. They're still not addressing the reason and the provided and the stopic will be able to affect it. They're still not addressing the reason and the provided and the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able to affect it. They will not address the still be able	not taxic
Into June understand Continue politicality de september per la production of the pro	•
but the bubble but over 60% of people including all those trump votes never liked or wanted observable and dont want in row, "tump has NOT bold the gips to back off you lying sack off its. he wants if gone and replaced period and the gop are doing just that together WITH trump. you will continue to turn reality into supplicity. The trump of the people including all those trump votes never liked or wanted observable want in continue to turn reality into supplicity.	not toxic
Rather than call you a liar, "Il just post a link to HART's helicopter footage of this "most densely populated" ag land.	
https://www.youtube.com/watch?huoppitRNVVVDdw	
and if you's basen following a little closure, you's innow about the favour's your war and budget in the state of the country to the city instead of actual affordable **COU.** On one took 1.90008513.12322270 is appose you just support urban synant then with prainting. Cold that close is a processing to the city instead of actual affordable **COU.** On one took 1.90008513.12322270 is appose you just support urban synant then with prainting. Cold that close is a processing to the city instead of actual affordable **COU.** On one took 1.90008513.12322270 is appose you just support urban synant then with prainting. Cold that close is a processing to the country of the country	not taxic
cat the "Yeal news" from the bitter and Drudge Report? Allahahahahaha Lay you for real? "Visit Statist 28015250555" "If one read the Dispatch one would think Trump is the most evil person on the planet." Not well, but indicate, and it would be again as his his behavior are assal fine years of 2015.	
	,
It was a gift from media that far't worried about the actual news, they are more worried about trying to influence soft heads. CNN has been screwing up a lot when it comes to Trump, same with the NYT that is now being sued for libel.	
Every stupid mistake they make gives an even larger advantage to Trump and strengthens his supporters that already believe the NSSM is biased against him and makes fence sitters begin to question what's news and what's crap. Like I've been saying, the NSSM is slitting it's own throat.	not toxic
You'd live see one of those American losts with all the planes on its the bottom of the see. <pre>//po</pre> <	
ports according to the control of th	
He was cartainly not an ideal assanger and, for one floating on a freeble, in the state of the s	not taxic
With contacts in software does it as to predict flower than as a first either mannered to Exclude the and non-thine does it as pred to day, see a coll of them mannered to Except than a complete and out of the Except than a coll of the except than a col	
by the logic of jour second paragraph, there can review be an ordinally shall secol at a, time see at by definition coor in a state of passion, which begs the question, why in this case would the Conjunes go through the touble of condemning securit immortality? This sounds like something your example of a right parties could any ordinal points (printing logical points) in a faster of passion.	toxic
Larget Namel and all other Conspany political leaders with here a lided and sheltend the originity invasion of trumps by the forces of the creasent more down. Amendood, Circles, Larget, etc. and the socialist troopes an	
Thus the Euro-socialist-bureaucrast pick the low-hanging fluit with litigious persecution of American firms which dominate because unlike their pathetic Euro-competitors, the U.S. firms are clever, hard-working, and well-capitalized.	
If the the Europeans wish to engage in this transparent financial inquisition, then the US should respond with courser litigation for Strillions against corrupt scriftings like VVI (think deset fiddlet) as well as UBS (tredit susse/HSE) (the dit superante financial inquisition, then the US should respond with courser litigation for Strillions against corrupt scriftings.	
This lay, crrupt, incompetent Expos want to by with the, the let them be fractedly included in the control of t	
Sincere question for you: if sesus and his followers celebrated Eucharist as a communal meal seated around a table, what gives Rome the right to alter this simple act of worship (perhaps "fellouship" is a better word—more suited toward love of god and neighbor), given to us by the Lord himself?	toxic
Thogs you can't have kidd if you see th in woman's action as a cosposale, and a goodwide fee for five first (e) in kidd, you have the present you make not a served you want and a served in the present and a ser	
You are completely (growned of the shortage of nursal in this country - In some case, critical shortages. And why would anyone want to be a nursal when they are disrepenced by a former case and fixederal prosecutor such as you. (In) - rong this about the impropriatement of my comment. River usely waited in special country and downers revised as an imaginary shywing, as for my compared of sharear, it was provided as an imaginary sharear of sharear and imaginary sharear	toxic
no one cares what a paid library trolling hack like you believe bursts, (1) or coincil \$3878348485644857 it shwps armuses me when a troll gets on, they like their own comments and simply assert everyone else is wrong. Never any evidence to rebut it just blind assertions. (1) or not tool: \$181038378444056. Out dioliticity to see that one commiss, a blass interest in the they take everying legs they like.	toxic
Old II you NOP "LEAF" manifacts types, where is you't here is you't here is not identify their standards with a second or of charact and innecessation through you have no bounds, i'm sure sha's andfully juparied over the thought that such a pathetic discoration could were be secretored. (In-trace) as the leaf of the second in the second	
And as for my post being "speculation? - which part - that the Liberais are the party in power, or that this involves money?	
as for me not knowing what is going on, you are correct, I am not a member of the Liberal party insider dique, as you apparently are.	toxic

Figure 13: Toxic Conversations, step 10.

Test Modified	Gold Label
Afform as yor in bloor negotations with city employees, Minkwaike Mayor Tom Barrett demanded concessions that went beyond those mandated by Gov. Scott Walkers collective bargaining law c/p a letter to members c/p - strue statement 3.833720313355912 Donald Trump as yo Libay a Makes standor (Christopher). Servers sets 60% overland prefugees to state with prevenors with one a Republicant, not not be Democrate. 4-(p-r) and interview or not implainally stated by the statement of the	true statement
Risk Scott spir All Abourd Floridals is a 100 percent privale venture. There is no state money involved: 4/>> a TV Intensive 4/>> a in Intensive 1 such support Floridation Flo	n true statement
Julie Pace asys The Obama administration is using at its ligal justification for these intrinsics (on the Islamic State), an authorization for military force statement 35 adding the present of the president himself has called for repeal of (*/p> a question to NVite House Press Secretary Join Extend (*/p> a fuertowing of the Statement 35 adding	true statement
John Kassich says We are now eighth in the nation in job creation we are No. 1 in the Milowest. (*) a news conference (*) to true statement 355900007046- Jorge Eforas asys in the lasts is year of Clancia sdministration violent crime was down in the United States. It was down in the Robe claims. (but it was so in Providence, Cry's a delates (*) and statement 40500005050000500000000000000000000000	true statement
Milke Pence says it was Itiliany Clinton who left Americans in human way in Benghal and after four Americans fell and, What difference at this point does it make? 4(p> the Republican national connection of po-strue statement 3.7440842903137207> longe Eloris asys in the last six years of Clancis administration violent; crime was down in the Freight and the region. It was down in Rhode Island. But it was up in Providence. a debate Sales tastement 3.745598958969116> Donald Trump says You will learn more about Donald Trump by going down to the Federal Elections to see the financial disclosure form than by looking at tax returns.	true statement
Rand Paul asy Of the roughly 1.5 percent of Americans undo adort have health insurance, half of them made more than \$50,000 syes. <>> in interiew on Comedy (centrals". The Daily Show" <>>> true tratement 3.7997431397100937. Bernie S says We have the highest rate of childhood goverty of any major countr on Carts. <>>> in interiew on CMM <>> fine the statement 3.99974213929914392 0 sould "Improve you have ingrished in The paul" of the statement 3.9997421392914935 0 sould "Improve you have ingrished in Statement 3.9997421392914935 0 sould "Improve you have ingrished in Statement 3.9997421392914935 0 sould "Improve you have ingrished in Statement 3.9997421392914935 0 sould "Improve you have ingrished in Statement 3.9997421392914935 0 sould "Improve you have ingrished in Statement 3.9997421392914935 0 sould "Improve you have ingrished in Statement 3.9997421392914935 0 sould "Improve you have ingrished in Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.9997421392914935 0 sould "Improve you have in the Statement 3.999742139291493 0 sould "Improve you have in the Statement 3.999742139291493 0 sould "Improve you have in the Statement 3.9997421	false statement
Banack Obamas says Stomulus tax cuts "began showing up in paychecks of 4.8 million Indiana households about three months ago." > a speech in Wakarusa, Ind. > <true 3.8199117888655089-="" <="" administration="" but="" clancis="" crime="" down="" eloras="" in="" island.="" it="" last="" lorge="" of="" providence.="" rhode="" says="" six="" statement="" states.="" the="" united="" up="" violent="" was="" years=""></true> > a debate > debate	false statement
Allen West spain if you look at the application for a security clearance, have a clearance that even the president of the United States cannot obtain because of my background. a candidate forum (1)> a candidate forum (2)> a candidate forum (3)> a candidate forum (3) (3) (3) (3) (3) (3) (4) (3) (4)	false statement
Bernië 5 aas We now work the longest hours of any people around the world. < />> 6-25 ANA interview of 2-35 Expenses 5-35 SEG95484375819105> Rush imbably as yet of many help view for the roof of hidhood poverty of any major country on Earth. a interview on CNN false statement 4.056144275183105> Rush imbably as yet of many perfut yets or for many privincience striking interview on CNN for the root of the root	false statement
Sanh Palin say Donald Trumps conversion to pro-life beliefs are shir to Justice Biefers, who said in the past that abortion was no big deal to him. A/s> an interview on CNM 4/5> "after statement 3.752/687225341797> Donald Trump say The federal government is sending refugees to states with governors who are Regulation, not to the Democrate. 4/5> an interview on LNM and personal profits of the Democrate 4/5> and personal personal profits of the Democrate 4/5> and personal pe	false statement

Figure 14: LIAR, step 1.

Test Modified	Gold Label
Afform say in labor negotiations with city employees, Miniwakee Mayor Tom Barrett demanded concessions that went beyond those mandated by Gor. Scott Walkers collective bargaining law c/p> a letter to members c/p> distor stakement 3.131746292114259 Tom Barrett says Goy. Scott Walker said no to equal pay for equal work for women. - All You A collective bargaining law c/p> a letter to members c/p> distor stakement 3.131746292114259 Tom Barrett says Goy. Scott Walker said no to equal pay for equal work for women. - All You A collective bargaining law c/p> a letter to members c/p> distor stakement 3.180082595855713> Scott Walker says if public employees dont pay more for benefits starting April 1, 2011, the equivalent is 1,500 state employee layoffs by June 30, 2011 and 10,000 to 12,000 state and local government employee layoffs in the next two years. - A remove conference	true statement
RIGK StOCK stays All Aboard Florids is 100 percent private venture. There is no ratte money involved, 479 a TV interview 470 statement 3.0582423594328834 Charle Crist says All Aboard Florids is receiving millions in Florida taxpayer dollars. 570 a fundraising email 570 a fundraising email 570 crypt examples of the statement of th	true statement
Julie Pace says The Obama administration is using as its legal justification for been entries (no the labeline State), an authorization for military brice that the president himself has called for repeal of «(p> a question to White House Press Secretary Josh Earnest «(p> true statement 2 982755860042285 Martha Raddets says The Obama administration originally wanted 01,000 troops to remain in larg — not combat troops, but military advisers, special operations forces, to watch the counterterrorism effort. «(p> comments on ABC's 'This Week' «(p> daise statement 3.246009588241577> Rick Perry says Obama has chosen to deny the vicious sub-the-mittin motivation of the attack on a loader Jewish grocery in a first. «(p) a statement in pression of the state of the press	true statement
John Kasich says We ser now eighth in the nation in job creation we are No. 1 in the Midwest 4/5/2 a network conference (yor false statement 1.50.089892073709591) Ted 597s(client days Gos. John Assich incorrectly climined Ohios economy was 58th in the nation when he took office. We were sixth in the nation in terms of economic job growth. Cyra in Interview on CNM (yor for the statement 2.08989205305405407 Tem You foundlife says if you state the population growth view of the compression of the properties of the compression of the compre	true statement
Mike Pence spit it was Hillary Clinton who left Americans in human way in Benghasi and after four Americans fell said, What difference at this point of section in make? (s) the Republican national convention (s) or true statement 2.87502/7642974554 Hillary Clinton says When terrorists silled more than 250 Americans in Lebanon under Remark Regars, the Commonst didnt make that so a startain issue. (s) or a SOUNT to make 10% or a startain issue. (s) or a SOUNT to make 10% or a startain issue and so a startain issue. (s) or a SOUNT to make 10% or a startain issue. (s) or a startain issue and so a startain issue and so a startain issue. (s) or a startain issue and so a startain issue and so a startain issue. (s) or a startain issue and so a star	true statement
Rand Paul says Of the roughly 15 percent of Americans who don't have health insurance, half of them made more than \$50,000 ayear.	false statement
Bank Obama says Stimulus tax cuts "Eagan showing up in paychecks of 4.8 million Indiana households about three months ago" «()» a speech in Wakarusa, Ind. «()» ridige statement 2,8003,213(2539945) Paul Drous says Stimulus money, funded a government board that made recommendations that would cost 37(20,00) poles and 52(3) billion in sails.« ()» or better ()» or twee statement 2,92(37)(37(35)(37(35)) and () or the favoral stimulus projects." ()» in and finders at the Tea Party Comention.	false statement
Allen West says if you look at the application for a security clearance, I have a clearance that even the president of the United States cannot obtain because of my background. > a candidate forum > false statement 3.0504826872258> Ted Cruz says One of the most troubling aspects of the Rubio-Schumer Gang of Eight bill was that it gave President Obama blanks authority to admit refugees, including syntam refugees, without mandating any background checks whatsoever. > a Republican presidential debate in Las Vegas crue statement 3.196129560470583> David Shuster says Said former U.S. Ambassador to Kenny Scion Gardon on Science Organis (President Obama Bolants).	false statement
Bernis 5 asy We now work the longest hours of any people around the world. (>> < C-SPAN interview (>> < Cue attainment 3.089577/6731079) and secretary interview (>> < Colorance Coloran	false statement
Sanh Pilli says Donald Trumps convenion to pro-ilfe beliefs are alin its Latent Biders, who said in the past that aborton was no big deal to him. (*) an interview on CNM (*) or fails estatement 3.100239932399072 Perman Cain says Said Planned Parenthoods early objective was to help kill black babies before they cannel not be word (*) or be tall at a consensate the hist has (*) or the statement 3.100239932399072 Perman Cain says Said Planned Parenthoods early objective was to help kill black babies before they cannel not be word (*) or be tall at a consensate the hist has (*) or be statement 3.100239932399072 Perman Cain says Said Planned Parenthoods early objective was to help kill black babies before they cannel not be supported by the said of the said	false statement

Figure 15: LIAR, step 5

Test Modified	Gold Label
Afscme says in labor negotiations with city employees, Milwaukee Mayor Tom Barrett demanded concessions that went beyond those mandated by Gov. Scott Walkers collective bargaining law a letter to members <false 3.131746292114258="" statement=""> Tom Barrett says Gov. Scott Walker said no to equal pa</false>	v
for equal work for women. a TV ad <true 1="" <="" and="" at="" s=""> <true 1="" 3.1408446197509766="" at="" avenue="" s=""> Portland Association Teachers says Did you know that if you accepted the Districts proposal today you would have NO pay increase for 4 years? Seven years of frozen wages = Disrespect. a newsletter</true></true>	true statemen
Rick Scott say All Aboard Florids is 100 percent private wethur. There is no state money involved, (?> a Y linteriew (-)< run; et attement 3.081021266991152> Charlie Crist says All Aboard Florids is receiving millions in Florids tappyer dollars. (-)> a fundrasing email > (-)> daise statement 3.18101162109375 Corp (compressions) and the statement 3.081011626991152> Charlie Crist says All Aboard Florids is receiving millions in Florids tappyer dollars. (-)> a fundrasing email (-)> daise statement 3.18101162109375 Corp (compressions) and the statement 3.0810116291152> Charlie Crist says All Aboard Florids is receiving millions in Florids tappyer dollars. (-)> a fundrasing email (-)> daise statement 3.18101162109375 Corp (compressions) and the statement 3.081011629109375 Corp (compressions) and t	true statemen
Julie Pace says The Obama administration is using as its legal justification for these airstrikes (on the Islamic State), an authorization for military force that the president himself has called for repeal of question to White House Press Secretary Josh Earnest - the gatement 2.962770462095133> Martha Raddets says The Obama administration originally by the Companies of the Companies o	true statemen
with fission say We are now eighth in the nation in job creation, we are No. 1 in the Midwest. (%) a news conference (½) refuse statement. 2 (60)999207305911 Ted Strickland says Gov. John Assich Incorrectly claimed Othiog economy was 38th in the nation when he took office. We were sixth in the nation in zernor decomment job growth. (%) as interview on CRM (½) are instanced with season statement. 2 (80)9986240000000 John Assich history and the control of the limit o	true statemen
Wile Pence says it was Hillary Clinton who left Americans in harms way in Benghari and after four Americans fell said, What difference at this point does it make? The Republican national convention <true 2.5874825608111572="" statement=""> Hillary Clinton says When terrorists killed more than 250 Americans is absorour more Ronald Reagen, the Democrats doint make that a partisian size.</true>	true statemen
Fand Paul say Of the roughly 15 percent of Americans unde ont have health insurance, half of them made more than \$50,000 a year, <>> an interview on Conego (central). "The Daily Show" (-y) - shits statement 2800/2583778885525 Aand Paul says Over half of the young people in medical, dental and law school in the young people in medical, dental and law school in the part of the young people in medical, dental and law school in the young people in the young peop	false statemer
Barack Obama says Stimulus tax cuts "began showing up in paychecks of 4.8 million Indiana households about three months ago." > a speech in Wakarusa, Ind. > false statement 2.890828132629945> Paul Brown says Stimulus money funded a government board that made recommendations that would cost 178,000 (plos and \$28.5 billion in sales. > a tweet > to tweet The two sales of the sales of the sales of the sales receives per person by \$460. The message is the internet.	false statemer
Allen West says If you look at the application for a security destrainer, I have a clearance that even the president of the United States cannot obtain because of my background. (2) a candidate forum (2) of bits estatement 3.021407385738955 State Southerland says 92 percent of President Barack Obamss administration has never evolved outside generated. (2) or administration has never evolved outside generated. (3) or administration, and (3) are administration has never evolved outside generated. (3) or administration, and (3) are administration has never evolved outside generated. (3) or administration, and (3) are administration has never evolved outside generated. (3) or administration has never evolved outside generated (3) and (3) are administration has never evolved outside generated. (3) or administration has never evolved outside generated (3) are administration has never evolved outside generated. (3) and (3) are administration has never evolved outside generated (3) are administration and (3) are	false statemer
Jernie S asy We now work the longest hours of any people around the world. < />> \$\(\rightarrow \) \$\(\rig	false statemer
Sarah Palin says Donald Trumps conversion to pro-life beliefs are skin to Justifu Belbers, who ask in the past that abortion was no big deal to him. > interview on CNN -false statement 2.78877687452285> Donald Trump says Public support for abortion is actually going down a little bit, polls show. >	feles statement

Figure 16: LIAR, step 10

Test Modified	Gold Label
If this usy as niggah can be with her, I can too. Ti she with Chuu Bake for anyways this niggah look like a http://t.co/NVIDROPPOT (4/2-chake speech 3.880554523758336) @AtBlackface @Mn/Mooncricket @UncleUnco @Boobly-Gotood: @FuckThe-Malloops @RecisityeepoS ape? Fuck u broke as a recist (/o-chersyle las niggah can be with her, I can too. Ti she with Chuu Bake for anyways this niggah look like a http://t.co/NVIDROPPOT (4/2-chake speech 3.880554523758336) @AtBlackface @Mn/Mooncricket @UncleUnco @Boobly-Gotood @FuckThe-Malloops @RecisityeepoS ape? Fuck u broke as a recist (/o-cherther 4.20903205871582) @UNIVForum These Manusist kikes promoted Anti-Semilism by creating this film, what was their incentive? They just promote the entering.	hate speech
8#8220_@lamkrause: No need to thank me, killing niggers is it's own reward "	
Dis my pal from our graduation when we was 26 > <nether 3.52255586241455=""> @LanceTHESPOKEN Statement of the Day: Pol who RT @ChiefsatWar are even dumber than he is. The ole birds of a feather type of deal. \hatespeech 3.5796432495117188> @Brenddan_they said some bullshit. Damn ingers</nether>	hate speech
T Blassifikir. niggers are so growned http://p.co/9900/WNR9-6/schales seeds-01 326738465258317-8 @ABBlockface @Minkoncricked @Undelstood @Flucifichalilloops @Raccititeger99 ape? Fuck urboine ass noticy 50center 3.706527232714844- @VNNForum These Manist kikes promoted intelligence of the control of the	hate speech
If @Bookards. Lift futing white tash /sp- data speech 1951/070156643559 @Additional face @Minkoncrided @Unicelation @Bookshopsocod @FutinGardinal speech approach approach as present 1951/070156643559 @Additional face @Minkoncrided @Unicelation @Bookshopsocod @FutinGardinal speech approach approach as present 1951/070156643559 @Additional face @Minkoncrided @Unicelation @Bookshopsocod @FutinGardinal speech approach approach approach as present 1951/070156643559 @Additional face @Minkoncrided @Unicelation @Bookshopsocod @FutinGardinal speech approach approa	h hate speech
pmile_pay_Congrabulations, you are officially fusting retained. <pre>// ** creation 4.07/1989/80031559* RT [abs/cg_basedGod_"New gr full retain" <pre>// cy> chate speech 3.47981357744629* @Brenddan_they said some bullshit. Damn niggers <pre>// coffensive language 3.601623773574829* RBUN Imago</pre> ###200_081008664111 Light parts from from "problocks' What the full chappeared to heavy full retained" <pre>// retained // ret</pre></pre></pre>	hate speech
onna have them pussies mix up some concrete today, teach them to pose like me. I am a badass motherfucker, and I will let you be too [< (>> chart speech 3.22750224423452> @ARBitchizec @Minkooncricket @Undelunco @BobbyBotsods @FuckTheMallcops @RacistNegrod9 ape? Fuck u broke ass racist offensive imagings \$3.530353455564> BRUH imago &Be2120, @Nobeector. Them hose was rash plut in leps etain if them ("problocks: What the fuck happened to these? http://t.co/GeVidSmiQif* &BE2121, (>> reneither 3.8099786151177385> RT @Venue, Lynn: &BE2120, @plylani: look ghetto but it work truty. (Inciprion/VelPerSea88212), &BE31557, &BE	offensive languag
Immodifications us call un bestre a bitch firing passing pick is a dog that banks too much (r)> coffensive lengage 3 1484224916107179 kT gilm, Army, Bickness from out always a bitch, pometimes I sleep - shate speech 3 14039840426835744 Women who are ferminist are the ugly bitches who can't find a man for hemselves (r)> confirmed = 32107086998699999999 Garly gird use a strong it in grant in grant fill girl grant in confirmed in the first who can't find a man for hemselves (r)> confirmed = 32107086999999999999 Garly gird use a strong it grant	offensive languag
to less than 3 had bitches in my bed at a time(r) - ordinative language 3.2010555989854027. Then again i phoot bitches in the sost fording zome to update this like key my car or but my windows shate speech 3.66655732559205- bitch kill yoseff, go on to the bathroom and eat the pills bitch, all of em - intermeter 4.29817248593566-86973, SEF STALAM MASKAGE OF THE part speech - OP BED 2012-12 ** Harby (CapitaAsPLOpS)	offensive languag
TI @ The Drug Tible: many lan't a backstabbing bitch that lies and deceives me (x) - coffeesing lenguage 3,455954495568945 xf @ m. Am, Bitches I'm on a laways a bitch, sometimes I sleep - shate speech 3,505273589457178- @ vinm/Juiclous faggot I knew you weren't really my friend - cerebre I sladeof7534500597 xf @ m. Am of my file m	offensive languag
poor,Bandod,Aiphone,Binds Assex,Bxx., Elicosety Squirting pusy and fragered asthole in the full Capitate Colories by Image and State (1998) Section (April 2018) Section (Apri	offensive languag
by baby gaine pig is so cute #Addonable (s)- Other speech 3.82485881421995 @witherapgiessit y don't und us monkeys then? Under us because were better than crackers &mon, [love ft. Deep down u know were done nothin 2 u - cofferable segregation of the speech 3.8248581421995 @witherapgiessit y don't und us monkeys then? Under the because were been contained to the speech of the speech and the speech of th	neither
T @Kick_Mans Gaints Pathul_lets-Pathul_Mets-Pathul_Net	neither
Ejestoff wit could get matching burner phones and be glettoff bit for a few months - what speech a 5667954448832540 - whiteangelss84 y dot und us monkey then? U hade us because were better than crackers 8 amp; love it. Deep down u know were done nothin 2 u - roffensive language 199549432790581 - p (Filed Northand 1) with paperats buoght may be a car man People bitch about not getting what car they want when they want it, and its free8e8210.4 (*po - onether 1 8/07/20079154982104 Filed Williams) and the speech about not getting what car they want when they want it, and its free8e8210.4 (*po - onether 1 8/07/20079154982104 Filed Williams) and the speech about not getting what car they want when they want it, and its free8e8210.4 (*po - onether 1 8/07/20079154982104) and the speech about not getting what car they want when they want it, and its free8e8210.4 (*po - onether 1 8/07/20079154982104) and the speech about not getting what car they want when they want it, and its free8e8210.4 (*po - onether 1 8/07/20079154982104) and they want it is speech about not getting what car they want with a speech about not getting what car they want when they want it, and its free8e8210.4 (*po - onether 1 8/07/20079154982104) and they want it is speech about not getting what car they want it, and its free8e8210.4 (*po - onether 1 8/07/20079154982104) and they want it is speech about not getting what is s	neither
This White Iron Band plays this weekend in Fargo, ND at the Aquanium(21-). Friday(10-29-10) with Charlie Parr. The next night, Saturday> "hate speech 3.8393898010253906">88128514,88128514,88128514,812	neither
IT @distruct: @Fon/News @fjoy/T And I don't have any confidence NOWNHATSOEVER in you Barack! You're the sole reason why this country is in this 88230, INSTANCE OF CONTROL OF THE PROPRET OF THE PROP	

Figure 17: Hate Speech Offensive, step 1

est Modified	Gold Label
If this ugly assinges not be with her, I can too. Til she with Chuu Bake for anyways this niggah look like a http://t.co/IVXDRDPQOT c/u-chale speech 2.715285538183596-RT @WaavyLee: His balls suby RT @Watshing @Telsire1st: Real women do this http://t.co/IVXDRDPGDT and real faggots let em di https://t.co/IVXDRDPGDT and real faggots let em di https://t.co/IVXDR	hate speech
889220_@lamkrause: No need to thank me, killing niggers is it's own reward.888221;	
Dis my pal from our graduation when we was 25 > \nate speech 2.605855449935027> @samzbikowski some negro amigo pulled a gun on Nate and I a few weeks ago. I was STOKEDIB4128299.84128299.4128299.45> \nate State Sta	hate speech
T Glessiditic negers are so genomen http://t.co/1990/00/NIPG 6/2- chate speech 2 G7 \$500759905909- RT @WhitelOnly_1 imagers! http://t.co/hib3uisky2 6/2- cneither 2/734483385337666- Bamp; thots are wearing Uggs RT @ BigBoothyudys Lik REIs*BillBecause negros are pulling out their Timbs* 6/2- caffensive angage 2 20272944089929- RT @WhitelOnly_1 images 2 20272944089929- RT @WhitelOnly_1 images 2 2027294409929- RT @WhitelOnly_1 images 2 2027294099929- RT @WhitelOnly_1 images 2 2027294099929- RT @WhitelOnly_1 images 2 2027294099929- RT @WhitelOnly_1 images 2 2027294099- RT @WhitelOnly_1 images 2 2027294099- RT @WhitelOnly_1 images 2 2027294099- RT @WhitelOnly_1 images 2 202729409- RT @WhitelOnly_1 images 2 2027	hate speech
RT @RosieZaya1: Ur fucking white trash <a 2.422173500061035="" hrate-speech=""> @FrankielGrande fugly queer white trash <a <="" bitches="" hrate-speech="" loans="" s=""> <a 31"="" 41"="" 81"="" by="" drained="" gill-control="" gill-control<="" href="https://doi.org/10.100/journal-picked-bitches-picked-bitch</td><td>hate speech</td></tr><tr><td>Bimle, rays Congressiotions, you are distally fusing greated. (yo - thate speech 2.4547482496395559: @darfmdnana Yea you do retard. (yo - Geffenvie Inaquage 2.83156443767112 > Loll 8#8220_@tisweetz_Sitch : Coop OT @TIFRANY_PORSIONE: You little twest 8#8221; (yo - centher 2.8371688747406006 N general productions) and the state of the stat</td><td>hate speech</td></tr><tr><td>pones have them pussion mix up pome concrete today, teach them to pose like me. I am a baddes motherfucker and i will et you be too (./e/> -offense here language _758007288523140; @400z_UAN (YCM. I can't get any work done if you keep howin off you betten <./p></td><td>offensive langua</td></tr><tr><td>Broaddersoans us all ur bestre a bitch 1'm guesting first a dog that barks to munt of yo- offensive lengage 3 07/2039654751515 (all 18 8220_@15.5wetz_Bitch Ooppl 07 @TFFANY_DRSCHE: You little twist 5 88221_</td><td>offensive langua</td></tr><tr><td>to less than 3 bad bitches in my bed at a time(p> -orientive language 3.2021/304545823> Don't loss sleep bout these bitches bot they come and go 8#128076; -(p> -chate speech 3.2785214351554033> bitch kill yoself, go on to the bathroom and eat the pills bitch, all of em -(p> -creither 3.4171059131522314-
tees those away firm drained by 81" td=""><td>offensive langua</td>	offensive langua
T (F) the Companies was your to a backstabbing bitch that lies and deceives are (x) or-coffessive lenguage 2991437934179459 T (B) the Coffiffier. If the bit passes ("Only God can judge me" she's a hoe. (x) or chairs speech 3.08849452974112) RT (B corryimates I got called a flaggor for buying girl toms so now I'm gonna fucility as person in the sax of yor character 3.2884951722459 ("Inscorps of (Contractification,") (in contracting the property of the propert	offensive langua
psom, mandroid, pilybonn, Pilopad, Barra, Max, Elicisate Signating pusy and fragered asshole in the //Linch/Need/wWW c/s/> overither 1.5877738421125849+ spom, mandroid pilybonn, alped disea, More pail inter/Linch/Confordity (s/> charts peech 1.2566939382723887> fill @mitch/mancusor: Behardrinting/White pail pail inter/linch/Confordity pail inter/Linch/	offensive langua
Wy bably guines agis is so cute #46/contale <pre>// or - neather = 1,5464395959679677 Our female guines agis is pregenet & \$217281_8217281_8217281_8212825_8212825_\$21</pre>	neither
T (FIGUR) Ann Ganth PROLL, Lett-PROLL, Meth Profil, "Inniverse-PROLL METH PROLE AND PROJECT AND PROLE AND PROJECT AND PROLE AND PROLE AND PROJECT AND PROLE AND PROJECT AND PROJ	neither
Bjestoth we could get matching burner phones and be ghetto fab for a few months <offerensive 3.12552990404665="" language=""> @Zoily23 @JBillinovich we need to grow mullest together so we can get all the bitches and @HannahKubiak can hate on us <hate <="" on="" p="" us=""></hate></offerensive>	neither
This White Iron Band plays this weekend in Fargo, ND at the Aquanium(21-). Friday(10-29-10) with Charlie Parr. The next night, Saturday> relative 3.4018728733062744 Limadoo naw man RT @DipOnline Yo want in RT @HumbitonBanks: U serious bro?? lol RT @CheexModelenkinz 2-3-10am early bird special **Saturday of Saturday (Saturday Saturday (Saturday Saturday Saturda	neither
IT @dsrtvet: @foxNeus @floy? And I don't have any confidence NOWNHATSOPVER in you Barack! You're the sole reason why this country is in this8e\$230. chate speech 2.6967606\$4449463> RT @veeveeveeveeve: If I was Obama id call a press conference Samp, sit joe bidens neck on live to just 2 show these reasons are not as the property of the property	

Figure 18: Hate Speech Offensive, step 5

st Modified	Gold Label
this ugly ass niggable can be with her , I can too . Ti she with Chou Bake for anyways this niggab look like a http://t.co/N/10R0PpQT coffersive lenguage 2.65357065500805669-RT @Currently_Spitta. And if a bitch can't respect a nigga wit some paper and a fresh pair of bball shorts then she was raised terribly. > shate speech 2.7233835383939-RT @Vasay/see. His balls safty RT @Hattaling @Trelairs ist. Real women do this http://t.co/N/508aH83T682221, and real faggots let em do that . Smh cheether 3.09851055709239555 @NeonTreat @FoCBeauty	hate speech
#8220@lamkrause: No need to thank me, killing niggers is it's own reward."	
is my pal from our graduation when we was 25 - shate speech 25458652781524669 RT @The_Bhodes: “_@kim92493 @The_Bhodes: @kim92493 @patpatbush with you've been judged" it happens. #whitepowerI'll hang you nigger@#8221; wo…, - coffensive language 5644790744781496> “_@kim92493 @patpatbush with you've been judged" it happens. #whitepowerI'll hang you nigger@#8221; wo…, - coffensive language 5644790744781496> “_@kim92493 @patpatbush with you've been judged" it happens. #whitepowerI'll hang you nigger@#8221; wo…, - coffensive language 5644790744781496> “_@kim92493 @patpatbush with you've been judged" it happens. #whitepowerI'll hang you nigger@#8221; wo…,	hate speech
[# glasd/tick_ingers are to ignorant http://t.co//PpQ/cis/vick-/pthest speech 1.07/8007/80090809080 FT @Whitstork_i.i. Hispace 1 http://t.co//PpQ.sis/vick-/pthest speech 2.07/8007/80090809080 FT @Whitstork_i.i. Hispace 1 http://t.co//PpQ.sis/vick-ythest speech 2.07/8007/80090809080 FT @Whitstork_i.i. Hispace 2.07/8007/8009080 FT Whitstork_i.i. Hispace 2.07/8009080 FT Whi	hate speech
@RosieZaya1: Ur fucking white trash <hate 2.422173500061035="" speech=""> @FrankieJGrande fugly queer white trash <offensive 2.6756434440612793="" language=""> RT @Jayy_Gee96: Dumb bitches <neither 2.783188819885254=""> RT @BeardedNixon: Poont gotta be trash</neither></offensive></hate>	hate speech
mike_m/r Congrahalmons, you are officially fuching retarded	hate speech
none have them puzzies mix up some concrete today. Intent himmen to soze like me. I am a Sadass mother-ficker and i will let you be foor. (4)p- cofferein le	offensive langu
maddle-enants u.c. all or bettle a blink in guesting she's a dag that bar's too much size a September 254968850738507 > @hjannwilliams whoocoof Chelsey? Fuck her lot. She juss a bitch -s/p> <offensive 2.884238589172859="" <p="" @_ezzylove:="" a="" at="" awwy.="" bad="" bitch,="" get="" it="" language="" let's="" right="" she="" to="">> chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.884238589172859 AT @_Ezzylove: She a bad bitch, let's get to it right awwy. > chether of supposed 54912315 in me and language 2.</offensive>	offensive langu
bless than 3 bad bitters in my bed at a time/p> coffensive language 2.8327350063078717-8 add bitters in the pea make my tops could visy-obless general 3.2539432048797807-1 didn't forsite all other bitches for my wrife to be getting fucked on by another nigga, and you know she married? you gotts die -/p> either 3.470325590055975. Keep those swyrm from fairle by 19 "f bleformen" kund one of my in if for life for someone's fororities and from body to be my disparation.	offensive langu
(Fifted fig. 1) the many int a backstabbing bits that lies and decisions me clys- offenous lengage 2 998188831703873 AT (Bettedfilter life in box may be compared to the class of the class	offensive langu
com, sandroid phiome, Blood Res, xxxx, Exclusivia Squiring pusys and frequent ashable that // 1.co/Revolvivia/ c/p> -center 1.5877334211355459 - xxxx, stock place pain thrst // 1.co/Revolvivia/ c/p> -center 1.587733421135549 - xxxxx stock place pain thrst // 1.co/Revolvivia/ c/p> -center 1.58773421135549 - xxxxx stock place pain thrst // 1.co/Revolvivia/ c/p> -center 1.5877342113549 xxxxxx xxxxxx xxxxxx place pain thrst // 1.co/Revolvivia/ c/p-center xxxxxx xxxxxx xxxxxx xxxxxx xxxxxx	offensive langu
y haby gaines gig is so cate #Addonable - For early reading 15 is no cate #Ad	neither
(例のは、Mon Classifts Philful Nets - Philful	neither
jestatih we could get matching bumer phones and be glated to the for a few months (s/s - shate speech 3.02/#35/27080913; @JAMMI. Byoden buth we can denish reduced; ((flarg dealers)) (s/s - coffensive linguage 3.12525930404683- @JZolly33 @JBlinovich we need to grow mullest together so we can et all the blinkse and (flarmanhsbuds) and an all on our s/s - reafters \$2,70000370900303 #3,70000370003 #3,70000303 #3,70000370003 #3,70000370003 #3,70000370003 #3,70000370003 #3,70000370003 #3,70000370003 #3,70000370003 #3,70000370003 #3,70000370003 #3,70000370003 #3,70000370003 #3,70000370003 #3,7000003 #3,7000003 #3,70	
o u might be stupid if u pay 4.59 for a b…	neither
w White loo Band plays this weedened in Fargo, ND at the Aquarium (21+), Friday (10-29-10) with Charile Part. The next night, Saturdays. < [>> < neither 3.2482401380075849-XT @toddinlife Full @weakenedinachos set (except the last song) from Southern Daviness Fest last month. Who's the age on guitar? typy, 1c.888250, < [>> has been seen 3.3524651527404783> Eagles fuck around & Bamp, lose it'll be kill the crucker at the Sophi crib smft of your continues language 3.51101638885967> My dawag @commanimize told me it's a must the @90112lounge this Saturday ROCKIN that bitch wit Tha 888210, tsp.//t.co/(INV)901410100000000000000000000000000000000	neither
(Edishvett: @FonNews @tjoy? And I don't have any confidence NONWHATSCEVER in you Bursck! You're the sole reason why this country is in this@8230; > -(s) - chate speech 2.695760554449463> RT @veeveeveevee: if I was Obsama id call a priess conference & sit joe bidens neck on live try just 2 show these solers I mean business&88230; -> (s) - chether 2.795964859717> R* @jernafetjerjebubb. 882220; @Burschied and solers are dish of the soler eason why this country is in this@82210; -> (s) - chate speech 2.69576055449463> R* @veeveeveevee: if I was Obsama id call a priess conference & sit joe bidens neck on live try in the call a speech 2.69576055449463> R* @veeveeveevee: if I was Obsama id call a priess conference & sit joe bidens neck on live try in the call a speech 2.69576055449463> R* @veeveeveevee: if I was Obsama id call a priess conference & sit joe bidens neck on live try just 2 show these solers are lived as a speech 2.69576055449463> R* @veeveeveevee: if I was Obsama id call a priess conference & sit joe bidens neck on live try just 2 show these solers are lived as a speech 2.69576055449463> R* @veeveeveevee: if I was Obsama id call a priess conference & sit joe bidens neck on live try just 2 show these solers are lived as a speech 2.69576055449463> R* @veeveeveevee: if I was Obsama id call a priess conference & sit joe bidens neck on live try just and the speech 2.6957605449463> R* @veeveeveevee: if I was Obsama id call a priess conference & sit joe bidens neck on live try just and the speech 2.6957605449463> R* @veeveeveevee: if I was Obsama id call a priess conference & sit joe bidens neck on live try just and the speech as a speech 2.6957605449463> R* @veevee evee it is was Obsama id call a priess conference & sit joe bidens neck as a speech 2.6957605449449	neither

Figure 19: Hate Speech Offensive, step 10.