

Like a Good Nearest Neighbor: Practical Content Moderation and Text Classification

Anonymous ACL submission

Abstract

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. Inexpensive text classification is important for addressing the problem of domain drift in all classification tasks, and especially in detecting harmful content, which plagues social media platforms. Here, we propose Like a Good Nearest Neighbor (LAGONN), a modification to SetFit that introduces no learnable parameters but alters input text with information from 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 flagging undesirable content and text classification and improves SetFit’s performance. To demonstrate LAGONN’s value, we conduct a thorough study of text classification systems in both the context of content moderation under four label distributions and in a more general classification setting.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) few-shot benchmarks, but most are difficult to use because of their reliance on billion-parameter PLMs, pay-to-use APIs, and/or prompting. Constructing prompts is not trivial and may require domain expertise.

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 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.

SetFit’s performance, speed, and low cost would make it ideal for effective content moderation, however, this type of text classification proves difficult 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,

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

082 appears to be the most difficult dataset; at time of
083 writing, it is the only task where the human base-
084 line has not been surpassed, yet SetFit is among
085 the top ten most performant systems.²

086 Here, we propose a modification to SetFit,
087 called Like a Good Nearest Neighbor (LAGONN).
088 LAGONN introduces no learnable parameters and
089 instead modifies input text by retrieving informa-
090 tion from its nearest neighbors (NN) seen during
091 optimization. Specifically, we append the label,
092 distance, and text of the NNs in the training data
093 to a new instance and encode this modified version
094 with an ST (see Figures 5 and 1 and Table 1). By
095 making input data appear more similar to instances
096 seen during training, we inexpensively exploit the
097 ST’s pretrained or fine-tuned knowledge when con-
098 sidering a novel example. Our method can also be
099 applied to the linear probing of an ST, requiring
100 no expensive fine-tuning of the large embedding
101 model. Finally, we propose a simple alteration to
102 the SetFit training procedure, where we fine-tune
103 the ST on a subset of the training data. This results
104 in a more efficient and performant text classifier
105 that can be used with LAGONN. We summarize
106 our contributions as follows:

- 107 1. We propose LAGONN, an inexpensive mod-
108 ification to Sentence Transformer- or SetFit-
109 based text classification.
- 110 2. We suggest an alternative training procedure
111 to the standard fine-tuning of SetFit, that can
112 be used with or without LAGONN, and results
113 in a cheaper system with similar or improved
114 performance to the more expensive SetFit.
- 115 3. We perform an extensive study of LAGONN,
116 SetFit, and standard transformer fine-tuning
117 in the context of content moderation under dif-
118 ferent label distributions and in a general text
119 classification setting, showing that our method
120 excels at many text classification tasks.

121 2 Related Work

122 There is little work on using sentence embeddings
123 as features for classification despite the pioneering
124 work being five years old (Perone et al., 2018). STs
125 are pretrained with the objective of maximizing
126 the distance between semantically distinct text and
127 minimizing the distance between text that is seman-
128 tically similar in feature space. They are composed

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

129 of a Siamese and triplet architecture that encodes
130 text into dense vectors which can be used as fea-
131 tures for ML. STs were first used to embed text
132 for classification by Piao (2021), however, only
133 pretrained representations were examined.

134 SetFit uses a contrastive learning paradigm
135 (Koch et al., 2015) to optimize the ST embedding
136 model. The ST is fine-tuned with a distance-based
137 loss function, like cosine similarity, such that ex-
138 amples with different labels are separated in fea-
139 ture space. Input text is then encoded with the
140 fine-tuned ST and a classifier, such as logistic re-
141 gression, is trained. This approach creates a strong,
142 few-shot text classification system, transforming
143 the ST from a sentence encoder to a topic encoder.

144 Work done by Xu et al. (2021) showed that re-
145 trieving and concatenating text from training data
146 and external sources, such as ConceptNet (Speer
147 et al., 2017) and the Wiktionary³ definition, can be
148 viewed as a type of external attention that does not
149 alter the architecture of the Transformer in ques-
150 tion answering. Liu et al. (2022b) used PLMs and
151 k -NN lookup to prepend examples that are similar
152 to a GPT-3 query, aiding in prompt engineering
153 for in-context learning. Wang et al. (2022) demon-
154 strated that prepending and appending training data
155 helps PLMs in summarization, language modelling,
156 machine translation, and question answering, us-
157 ing BM25 as their retrieval model (Manning et al.,
158 2008; Robertson and Zaragoza, 2009).

159 We alter the SetFit training procedure by using
160 fewer examples to adapt the embedding model for
161 many-shot learning. LAGONN decorates input text
162 with its NN’s gold label, Euclidean distance, and
163 text from the training data to exploit both the ST’s
164 distance-based pretraining and SetFit’s distance-
165 based fine-tuning objective. Compared to retrieval-
166 based methods, LAGONN uses the same model for
167 both retrieval and encoding, retrieving only infor-
168 mation from the training data for classification.

169 3 Like a Good Nearest Neighbor

170 Xu et al. (2021) formulate a type of external atten-
171 tion, where textual information is retrieved from
172 multiple sources and added to text input to give
173 the model stronger reasoning ability without al-
174 tering the internal architecture. Inspired by this
175 approach, LAGONN exploits pretrained and fine-
176 tuned knowledge through external attention, but the
177 information we retrieve comes only from data used

³<https://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)
DISTANCE	"I love this. [SEP] [0.5]" (0)
LABDIST	"I love this. [SEP] [positive 0.5]" (0)
TEXT	"I love this. [SEP] [positive 0.5] This is great!" (0)
ALL	"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]" (?)
DISTANCE	"So good! [SEP] [1.5]" (?)
LABDIST	"So good! [SEP] [positive 1.5]" (?)
TEXT	"So good! [SEP] [positive 1.5] I love this." (?)
ALL	"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. Text is in quotation marks and the integer label is in parenthesis. In brackets are the gold label or distance from the NN or both. 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 *training instance text*". See Appendix A.9 for examples of LAGONN ALL modified text.

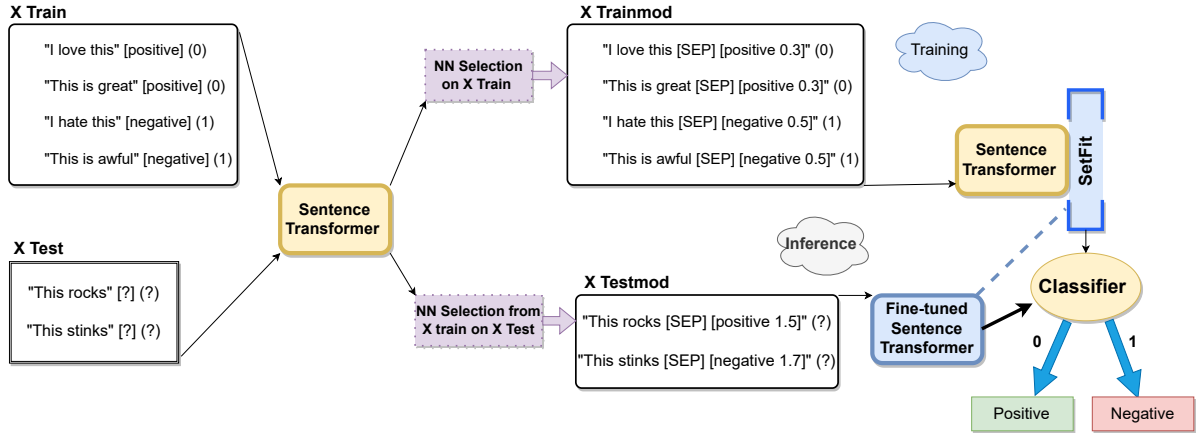


Figure 1: LAGONN LABDIST uses an ST to encode training data, performs NN lookup, appends the NN’s gold label and distance, and optionally SetFit to fine-tune the embedding model. We then embed this new instance and train a classifier. During inference, we use the embedding model to modify the test data with its NN’s gold label and distance from the training data, 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.

during optimization. We consider an embedding function, f , that encodes both training and test data, $f(X_{train})$ and $f(X_{test})$. Considering its success 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.

Encoding and nearest neighbors LAGONN first uses a pretrained Sentence Transformer to em-

bed training text in feature space, $f(X_{train})$, and NN lookup with scikit-learn (Buitinck et al., 2013) on the resulting embeddings.

Nearest neighbor information We extract text from the nearest neighbors 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 of the

194 NN, which we call LABEL. We then consider the
195 Euclidean distance of the NN, which we call DIS-
196 TANCE, giving the model access to a continuous
197 measure of similarity. We then combine these two
198 configurations, appending both the NN’s gold label
199 and Euclidean distance, referring to this as LAB-
200 DIST. Next, we consider the gold label, distance,
201 and the text of the NN, which we refer to as TEXT.
202 Finally, we tried the same format as TEXT but for
203 all possible labels, which we call ALL (see Table
204 1 and Figure 1). Information from the NN is ap-
205 pended to the text following a separator token to
206 indicate this instance is composed of multiple se-
207 quences. See Appendix A.8.1 for a detailed study
208 of and comparison between all LAGONN configu-
209 rations.

210 **Training** LAGONN encodes the modified
211 training data, optionally fine-tunes the embed-
212 ding model via SetFit, and trains a classifier,
213 $CLF(f(X_{trainmod}))$.

214 **Inference** LAGONN uses information from
215 the nearest neighbor in the training data to modify
216 input text. We compute the embeddings of the test
217 data, $f(X_{test})$, and select and extract information
218 from the NN’s training text, decorating the input
219 instance with this information. Finally, we encode
220 the modified data with the embedding model and
221 call the classifier, $CLF(f(X_{testmod}))$.

222 **Intuition** The ST’s pretraining and SetFit’s
223 fine-tuning objective both rely on distance, cre-
224 ating a feature space appropriate for distance-based
225 algorithms, such as our NN-lookup. We hypoth-
226 esize that LAGONN’s modifications make novel
227 data appear semantically similar to their NNs in the
228 training data, that is, more akin to an instance on
229 which the encoder and classifier were optimized.
230 LAGONN’s utilization of distance and clear dis-
231 tinctions between classes inspired our use case of
232 content moderation, where it is realistic to have few
233 labels, harmful or neutral, for example. However,
234 this work demonstrates that LAGONN is useful for
235 general text classification as well.

236 4 Experiments

237 We first study LAGONN’s performance on four
238 binary and one ternary classification dataset related
239 to the task of content moderation. Each dataset is
240 composed of a training, validation, and test split
241 (see Appendix A.1 for details).

242 We study our system by simulating growing
243 training data over ten discrete steps sampled under
244 four different label distributions: extreme, imbal-
245 anced, moderate, and balanced (see Table 4). On
246 each step we add 100 examples (100 on the first,
247 200 on the second, etc.) from the training split sam-
248 pled under one of the four ratios. On each step, we
249 train our method with the sampled data and evalu-
250 ate on the test split. Considering growing training
251 data has two benefits: 1) We can simulate a stream-
252 ing data scenario, where new data are labeled and
253 added for training and 2) We can investigate each
254 method’s sensitivity to the number of training ex-
255 amples. We sampled over five seeds, reporting the
256 mean and standard deviation.

257 4.1 Baselines

258 We compare LAGONN against a number of strong
259 baselines, detailed below. We used default hyper-
260 parameters in all cases unless stated otherwise.

261 **RoBERTa** RoBERTa-base is a pretrained lan-
262 guage model (Liu et al., 2019) that we fine-tuned
263 with the transformers library (Wolf et al., 2020).
264 We select two versions of RoBERTa-base: an ex-
265 pensive version, where we perform standard fine-
266 tuning on each step (RoBERTa_{full}) and a cheaper
267 version, where we freeze the model body after step
268 one and update the classification head on subse-
269 quent steps (RoBERTa_{freeze}). We set the learning
270 rate to $1e^{-5}$, train for a maximum of 70 epochs,
271 and use early stopping, selecting the best model
272 after training. We consider RoBERTa_{full} an upper
273 bound as it has the most trainable parameters and
274 requires the most time to train of all our methods.

275 **Linear probe** We perform linear probing of a
276 pretrained Sentence Transformer by fitting logis-
277 tic regression with default hyperparameters on the
278 training embeddings on each step. We choose this
279 baseline because LAGONN can be applied as a
280 modification in this scenario. We select MPNET
281 (Song et al., 2020) as the ST, for SetFit, and for
282 LAGONN.⁴ We refer to this method as Probe.

283 **SetFit** Here, we perform standard fine-tuning
284 with SetFit on the first step, and then on subsequent
285 steps, freeze the embedding model and retrain only
286 the classification head. We choose this baseline as
287 LAGONN relies on ST/SetFit for its modifications.

⁴<https://huggingface.co/sentence-transformers/paraphrase-mpnet-base-v2>

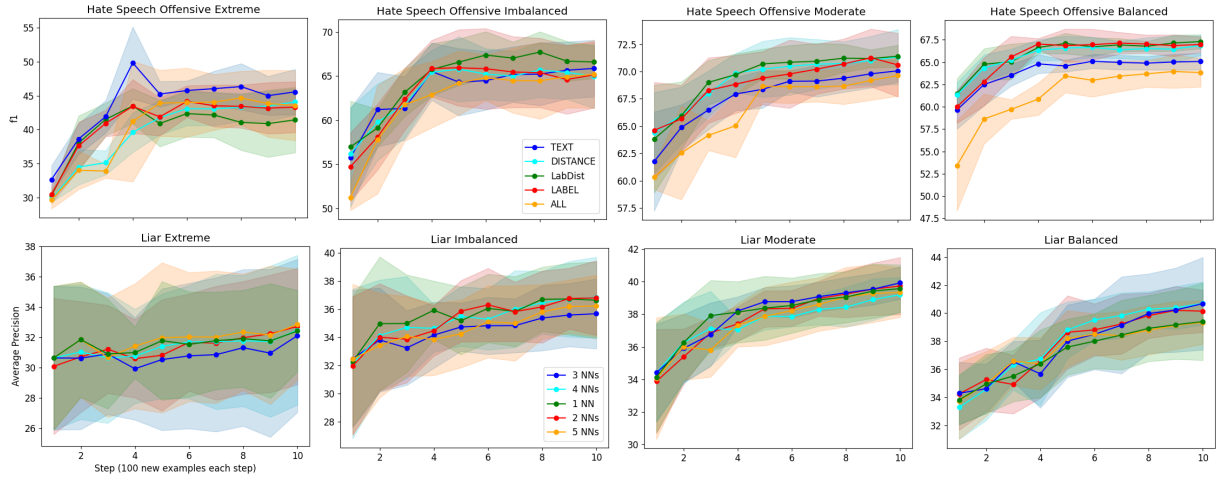


Figure 2: First row: performance for all LAGONN configurations and balance regimes for the Hate Speech Offensive dataset. Second row: LAGONN performance for one to five neighbors for all balance regimes on a collapsed version of the LIAR dataset. We use the LAGONN_{lite} fine-tuning strategy (see Section 5.1).

k -nearest neighbors Similar to the above baseline, we fine-tune the embedding model via SetFit, but swap out the classification head for a k NN 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 k NN.

SetFit expensive For this baseline we perform standard fine-tuning with SetFit on each step. On the first step, this method is equivalent to SetFit. We refer to this as SetFit_{exp}.

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

LAGONN On the first step, we use LAGONN to modify our data and perform standard fine-tuning with SetFit. On subsequent steps, we freeze the embedding model but continue to use it to modify our data. We only fit logistic regression on later steps, referring to this method as LAGONN. It is comparable to SetFit.

LAGONN expensive Here we modify our data and fine-tune the embedding model on each step. We refer to this method as LAGONN_{exp} and it is comparable to SetFit_{exp}. On the first step, this method is equivalent to LAGONN.

4.2 LAGONN configurations

We perform extensive experiments over the different LAGONN configurations. We note that while DISTANCE and LABEL show similar performance, LABDIST in general is the most performant and consistent classifier. TEXT and ALL are arguably the most interesting LAGONN configurations, but are often unstable, low-performing classifiers. In Figure 2, we provide a comparison between the different configurations on the Hate Speech Offensive dataset. As LABDIST is the most performant configuration, it is the version of our method about which we report results hereafter, but detailed ablations can be found in Appendix A.8.1.

4.3 LAGONN k nearest neighbors

To determine how many neighbors we should consider for LAGONN, we perform thorough experiments for one to five neighbors over all datasets, LAGONN configurations, and balance regimes under the LAGONN_{lite} fine-tuning strategy (see Section 5.1). We find that one to three neighbors tends to result in the strongest classifier, but this varies and is a hyperparameter that can be searched over. In Figure 2, we provide a representative example of our NN results for the LABDIST configuration for the LIAR dataset, however, detailed ablations can be found in Appendix A.8.2.

5 Content Moderation Results

Table 2 and Figure 6 show our results. In the cases of the extreme and imbalanced regimes, the performance of SetFit_{exp} steadily increases with

Method	InsincereQs				AmazonCF			
	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 _{exp}	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}
SetFit	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}	24.2 _{9.0}	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}
<i>Balanced</i>								
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 _{exp}	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}
SetFit	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 2: 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.5 for additional results).

the number of training examples. As the label distribution shifts to the balanced regime, however, the performance quickly saturates or even degrades as the number of training examples grows. LAGONN, RoBERTa_{full}, and SetFit, other fine-tuned PLM classifiers, do not exhibit this behavior. LAGONN_{exp}, being based on SetFit_{exp}, exhibits a similar trend, but the performance degradation is mitigated; on the 10th step of Amazon Counterfactual in Table 2 SetFit_{exp}’s performance decreased by 9.7, while LAGONN_{exp} only fell by 3.7. Note that we only consider the first NN here.

LAGONN and LAGONN_{exp} generally outperform SetFit and SetFit_{exp}, 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 on the Insincere Questions dataset for all balance

regimes, despite being the third fastest (see Table 6) and having the second fewest trainable parameters. We attribute this result to the fact that this dataset is composed of questions from Quora⁵ 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.

5.1 SetFit for efficient many-shot learning

Respectively comparing SetFit to SetFit_{exp} and LAGONN to LAGONN_{exp} 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 equiva-

⁵<https://www.quora.com/>

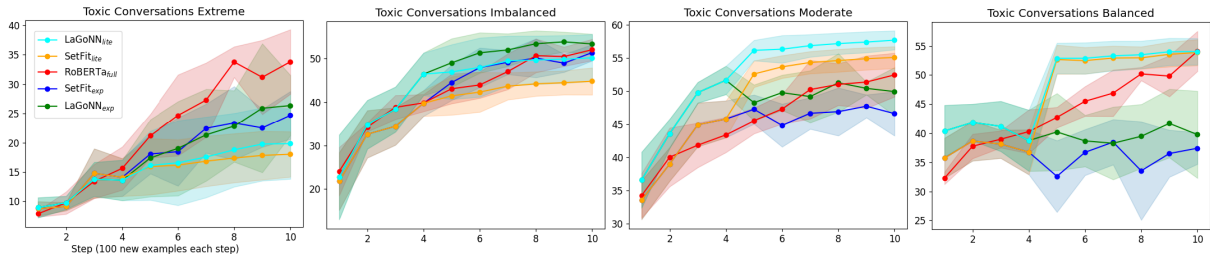


Figure 3: Average performance for all sampling regimes on Toxic Conversations. More expensive models, such as LAGONN_{exp} , SetFit_{exp} , and RoBERTa_{full} perform best when the label distribution is imbalanced. As the distribution becomes more balanced, inexpensive models, such as LAGONN_{lite} , show similar or improved performance. The metric is average precision and we only consider one neighbor for the LAGONN-based methods (see Appendix A.6 for additional results).

lent to SetFit_{exp} and LAGONN_{exp} , except after the fourth step (400 samples), we freeze the encoder and only retrain the classifier on subsequent steps, similar to SetFit and LAGONN .

Figures 3 and 7 show our results with these two new models. As expected, in the cases of extreme and imbalanced distributions, LAGONN_{exp} , SetFit_{exp} , and RoBERTa_{full} , are the strongest performers. We note very different results for both LAGONN_{lite} and SetFit_{lite} compared to LAGONN_{exp} and SetFit_{exp} on Toxic Conversations under the moderate and balanced label distributions. As their expensive counterparts start to plateau or degrade on the fourth step, these two new models dramatically increase, 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.

6 LAGONN as a General Classifier

LAGONN is effective for general text classification. Thus far, we have focused on the important topic of content moderation, but here we turn our attention to general text classification, conducting experiments on six additional datasets (see Appendix A.2 for details). Our experimental setup remains largely the same, but here we restrict ourselves to the balanced sampling regime as it is nontrivial to design sampling strategies for datasets with more than three labels. We respectively compare LAGONN_{lite} against SetFit_{lite} and LAGONN_{exp} against SetFit_{exp} , showing results for one to five

neighbors with LAGONN.

In Figure 4, we demonstrate that LAGONN continues to stabilize and improve SetFit , regardless of the number of neighbors we consider. This is especially clear for IMDB, where in the case of LAGONN_{lite} vs SetFit_{lite} , all versions of our method saturate to an average precision of 98 with 300 fewer training samples. If we consider SetFit_{exp} vs LAGONN_{exp} , consistent with our analysis of other binary datasets, classifier performance begins to degrade if we continue to fine-tune the ST, but LAGONN mitigates this performance drop.

Continuing to fine-tune the embedding model is beneficial when we have many labels. For 20 Newsgroups and Emotion, which have 20 and 28 labels respectively, LAGONN_{exp} is the strongest model and shows no indication of plateauing or degrading, even with 1,000 samples. We attribute this to the relatively high number of labels present in both of these datasets. Our findings related to SST-5 supports this; in intermediate cases when we have five labels, all models saturate quickly and there are minimal performance gains with continued full-model fine-tuning.

7 Discussion

Flagging potentially dangerous text presents a challenge even for state-of-the-art approaches. The content moderation datasets we consider proved more difficult than our general text classification datasets for all models, despite typically having fewer labels. It is imperative that we develop reliable and practical text classifiers for content moderation, such that we can inexpensively re-tune them for novel forms of hate speech, toxicity, and fake news.

Our results suggest that LAGONN_{exp} , a relatively expensive technique, can detect harmful content when dealing with imbalanced label distribu-

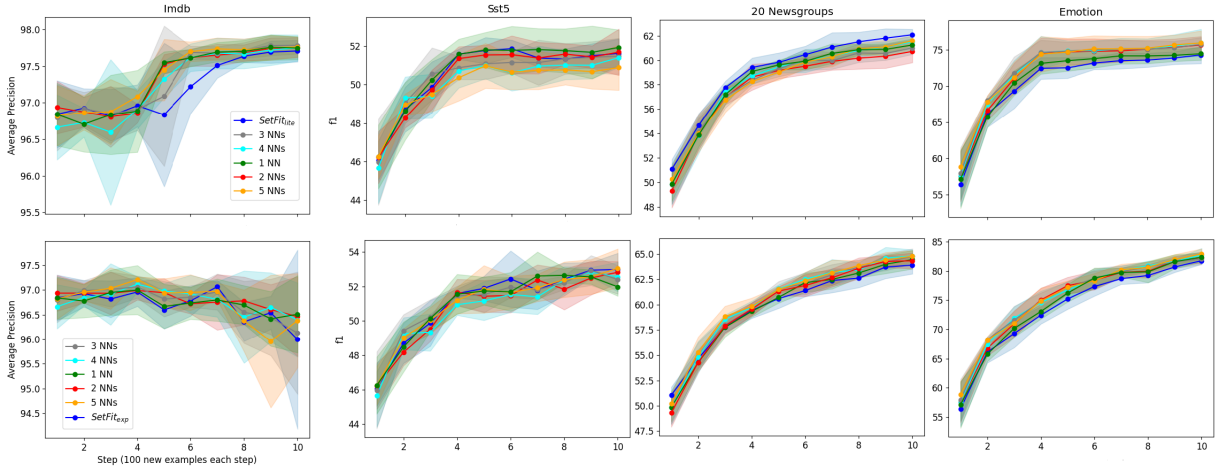


Figure 4: Average performance on four datasets in the balanced sampling regime; the metric is average precision for Imdb, macro-f1 elsewhere. First row: SetFit_{vite} compared to LAGONN_{exp} LABDIST with modifications for one to five neighbors. Second row: SetFit_{exp} compared to LAGONN_{exp}. See Appendix A.7 for additional results.

tions, as is common with realistic datasets. This is intuitive from the perspective that less common instances are more difficult to learn and require more effort. An exception would be our examination of Insincere Questions, where LAGONN_{cheap} excelled in the extreme and balanced settings. 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, directly affecting our own approach. However, we have observed that LAGONN can stabilize SetFit's predictions and reduce its performance drop in many cases. Figures 6, 3, and 4 show that when the label distribution is moderate or balanced (see Table 4), SetFit_{exp} plateaus, yet cheaper systems, such as LAGONN, continue to learn. We believe this is due to SetFit's fine-tuning objective, which optimizes an ST using cosine similarity loss to separate examples belonging to different labels in feature space, assuming independence between labels. This may be too strong an assumption as we fine-tune with more data, which is counter-intuitive for data-hungry transformers; RoBERTa_{full}, optimized with cross-entropy loss, 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, inexpensively boosting or stabilizing SetFit's performance. 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_{exp} (see Figure 3 and Table 6).

8 Conclusion

We have proposed LAGONN, an inexpensive modification to SetFit. LAGONN improves SetFit's performance by modifying text with the nearest neighbors in the training data. To demonstrate the merit of LAGONN, we examined text classification systems for content moderation with different label distributions and for general classification, studying 11 datasets with growing training data. When the training labels are imbalanced, expensive systems, such as LAGONN_{exp} are performant. LAGONN_{exp} also excels on balanced datasets with many labels. However, when the labels are binary or ternary, typical for content moderation, and the distribution is balanced, fine-tuning with SetFit can hurt model performance. We have therefore proposed an alternative but strong training procedure. LAGONN is a practical and simple method for detecting harmful content and text classification.

9 Limitations

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 training 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.

10 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.

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 Warmusley, 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.
- Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. [GoEmotions: A dataset of fine-grained emotions](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4040–4054, Online. Association for Computational Linguistics.

640	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding . In <i>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)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	696
641		697
642		698
643		699
644		700
645		
646		701
647		702
648		703
649	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 . <i>Nature Reviews Psychology</i> , 1(1):13–29.	704
650		705
651		706
652		707
653		708
654		709
655		710
656		711
657		712
658		
659		713
660		714
661		715
662		716
663		
664		717
665		718
666		719
667		720
668		721
669		
670		722
671		723
672		724
673		725
674		726
675		727
676		728
677		729
678		
679		730
680		731
681		
682		732
683		733
684		734
685		
686		735
687		736
688		737
689		738
690		739
691		740
692		741
693		
694		742
695		743
		744
		745
		746
		747
		748
		749
		750
		751
		752

753	Yusuke Shido, Hsien-Chi Liu, and Keisuke Umezawa.	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	810
754	2022. Textual content moderation in C2C market-	Chaumond, Clement Delangue, Anthony Moi, Pier-	811
755	place . In <i>Proceedings of the Fifth Workshop on</i>	ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-	812
756	<i>e-Commerce and NLP (ECNLP 5)</i> , pages 58–62,	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	813
757	Dublin, Ireland. Association for Computational Lin-	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	814
758	guistics.	Teven Le Scao, Sylvain Gugger, Mariama Drame,	815
759	Richard Socher, Alex Perelygin, Jean Wu, Jason	Quentin Lhoest, and Alexander Rush. 2020. Trans-	816
760	Chuang, Christopher D. Manning, Andrew Ng, and	formers: State-of-the-art natural language processing .	817
761	Christopher Potts. 2013. Recursive deep models for	In <i>Proceedings of the 2020 Conference on Empirical</i>	818
762	semantic compositionality over a sentiment treebank .	<i>Methods in Natural Language Processing: System</i>	819
763	In <i>Proceedings of the 2013 Conference on Empiri-</i>	<i>Demonstrations</i> , pages 38–45, Online. Association	820
764	<i>cal Methods in Natural Language Processing</i> , pages	for Computational Linguistics.	821
765	1631–1642, Seattle, Washington, USA. Association		
766	for Computational Linguistics.		
767	Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-	Yichong Xu, Chenguang Zhu, Shuohang Wang, Siqi	822
768	Yan Liu. 2020. Mpnnet: Masked and permuted pre-	Sun, Hao Cheng, Xiaodong Liu, Jianfeng Gao,	823
769	training for language understanding . In <i>Advances in</i>	Pengcheng He, Michael Zeng, and Xuedong Huang.	824
770	<i>Neural Information Processing Systems</i> , volume 33,	2021. Human parity on commonsenseqa: Aug-	825
771	pages 16857–16867. Curran Associates, Inc.	menting self-attention with external attention . <i>arXiv</i>	826
772	Robyn Speer, Joshua Chin, and Catherine Havasi. 2017.	preprint arXiv:2112.03254 , abs/2112.03254.	827
773	Conceptnet 5.5: An open multilingual graph of gen-		
774	eral knowledge . <i>Proceedings of the AAAI Conference</i>	Meng Ye, Karan Sikka, Katherine Atwell, Sabit Hassan,	828
775	<i>on Artificial Intelligence</i> , 31(1).	Ajay Divakaran, and Malihe Alikhani. 2023. Multi-	829
776	Lewis Tunstall, Nils Reimers, Unso Eun Seo Jo, Luke	lingual content moderation: A case study on reddit .	830
777	Bates, Daniel Korat, Moshe Wasserblat, and Oren	<i>arXiv preprint arXiv:2302.09618</i> .	831
778	Pereg. 2022. Efficient few-shot learning without		
779	prompts . <i>arXiv preprint arXiv:2209.11055</i> .	A Appendix	832
780	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	A.1 Content moderation data and balance	833
781	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	regimes	834
782	Kaiser, and Illia Polosukhin. 2017. Attention is all	In this Appendix section, we provide a background	835
783	you need . In <i>Advances in Neural Information Pro-</i>	on the datasets we studied in our experiments and	836
784	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	summarize the label distribution (see Table 3) of	837
785	Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018.	our content moderation datasets and the different	838
786	The spread of true and false news online . <i>Science</i> ,	sampling regimes (see Table 4) we studied in our	839
787	359(6380):1146–1151.	content moderation experiments. LIAR was cre-	840
788	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman-	ated from Politifact ⁶ for fake news detection and is	841
789	preet Singh, Julian Michael, Felix Hill, Omer Levy,	composed of the data fields <i>context</i> , <i>speaker</i> , and	842
790	and Samuel Bowman. 2019. Superglue: A stickier	<i>statement</i> , which are labeled with varying levels of	843
791	benchmark for general-purpose language understand-	truthfulness (Wang, 2017). We used a collapsed	844
792	ing systems . In <i>Advances in Neural Information</i>	version of this dataset where a statement can only	845
793	<i>Processing Systems</i> , volume 32. Curran Associates,	be true or false. We did not use <i>speaker</i> , but did	846
794	Inc.	use <i>context</i> and <i>statement</i> , separated by a separator	847
795	Shuohang Wang, Yichong Xu, Yuwei Fang, Yang Liu,	token. Quora Insincere Questions ⁷ is composed of	848
796	Siqi Sun, Ruochen Xu, Chenguang Zhu, and Michael	neutral and toxic questions, where the author is not	849
797	Zeng. 2022. Training data is more valuable than you	asking in good faith. Hate Speech Offensive ⁸ has	850
798	think: A simple and effective method by retrieving	three labels and is composed of tweets that can con-	851
799	from training data . In <i>Proceedings of the 60th Annual</i>	tain either neutral text, offensive language, or hate	852
800	<i>Meeting of the Association for Computational Lin-</i>	speech (Davidson et al., 2017) ⁹ . Amazon Counter-	853
801	<i>guistics (Volume 1: Long Papers)</i> , pages 3170–3179,	factual ¹⁰ contains sentences from product reviews,	854
802	Dublin, Ireland. Association for Computational Lin-		
803	guistics.		
804	William Yang Wang. 2017. “Liar, liar pants on fire”:	⁶ https://www.politifact.com/	
805	A new benchmark dataset for fake news detection .	⁷ https://www.kaggle.com/c/	
806	In <i>Proceedings of the 55th Annual Meeting of the</i>	quora-insincere-questions-classification	
807	<i>Association for Computational Linguistics (Volume 2:</i>	⁸ https://huggingface.co/datasets/hate_speech_	
808	<i>Short Papers)</i> , pages 422–426, Vancouver, Canada.	offensive	
809	Association for Computational Linguistics.	⁹ For Hate Speech Offensive, 0 and 2 denote undesirable	
		text and 1 denotes neither.	
		¹⁰ https://huggingface.co/datasets/SetFit/	
		amazon_counterfactual_en	

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 Conversations¹¹ is a dataset of comments where the author wrote with unintended bias¹² (see Table 3).

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 3: Summary of content moderation datasets and number of labels. We provide the type of task in parenthesis in unclear cases.

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

A.2 General text classification data

In this Appendix section, we provide additional information on the datasets we examined in our general text classification experiments. The Internet Movie Database (IMDB) dataset (Maas et al., 2011) is composed of movie reviews that are classified as either positive or negative.¹³ Student Question Categories contains questions from qualifying examinations in India,¹⁴ where the label is the subject the question appeared in and can be from Physics, Chemistry, Biology, or Mathematics.¹⁵ SST5 is an alternative version of the Stanford Sentiment Treebank (Socher et al., 2013) that has five labels, ranging from very positive to very negative.¹⁶ We also include the original version of LIAR, which

¹¹https://huggingface.co/datasets/SetFit/toxic_conversations

¹²<https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification/overview>

¹³<https://huggingface.co/datasets/SetFit/imdb>

¹⁴<https://www.kaggle.com/datasets/mrutyunjaybiswal/iitjee-neet-aims-students-questions-data>

¹⁵<https://huggingface.co/datasets/SetFit/student-question-categories>

¹⁶<https://huggingface.co/datasets/SetFit/sst5>

has six labels of varying levels of truthfulness.¹⁷ We also used 20 Newsgroups¹⁸ (Mitchell, 1999) which contains newspaper articles labeled with the topic they cover.¹⁹ And finally, we ran experiments on GoEmotions (Demszky et al., 2020), a dataset of Reddit comments labeled with 28 classes based on the emotional charge of the post.²⁰

The evaluation metric was average precision in the case of IMDB, macro F1 elsewhere. In cases where the a validation split was not available, we created one by sampling 30% of the test split. Please see Table 5 for a summary regarding the datasets and label information.

Dataset (and Detection Task)	Number of Labels
IMDB (Sentiment Analysis)	2
Student Questions (Question Type)	4
SST5 (Sentiment Analysis)	5
LIAR (Fake News)	6
20 Newsgroups (Topic)	20
GoEmotions (Emotion)	28

Table 5: Summary of datasets and number of labels used in the general text classification experiments. We provide the type of task in parenthesis in unclear cases.

A.3 Observations about LAGONN

Here, at the suggestion of an anonymous reviewer, we include a little background on LAGONN. We originally attempted to use Sentence Transformers/SetFit as a retrieval model that would modify input text and then pass this input to a Transformer-based classifier, such as RoBERTa, instead of back into the ST as in LaGoNN. We experimented with different ST retrieval models and Transformer classifiers, but this system was often beaten by baselines, and performant versions were too expensive to justify their use. The failure of this system is what ultimately inspired LAGONN. We had hoped 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

¹⁷<https://huggingface.co/datasets/LIAR>

¹⁸https://scikit-learn.org/0.19/datasets/twenty_newsgroups.html#the-20-newsgroups-text-dataset

¹⁹https://huggingface.co/datasets/SetFit/20_newsgroups

²⁰https://huggingface.co/datasets/SetFit/go_emotions

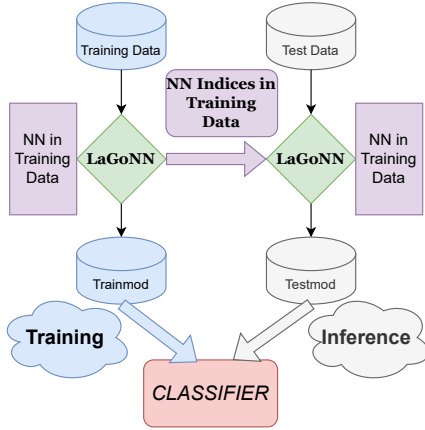


Figure 5: We embed training data, retrieve the text, gold label, and distance for each instance from its nearest neighbor and modify the original text with this information. Then we embed the modified training data and train a classifier. During inference, the NN from the training data is selected, the original text is modified with the text, gold label, and distance from this NN, and the classifier is called.

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.4 LAGoNN’s computational expense

In this Appendix section we discuss and provide results for LAGoNN’s computation time. 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 6, we provide a speed comparison of comparable methods computed on the same hardware.²¹ On average, LAGoNN introduced 24.2 additional seconds of computation compared to its relative counterpart.

Method	Time in seconds
Probe	22.9
LAGoNN _{cheap}	44.2
SetFit	42.9
LAGoNN	63.4
SetFit _{exp}	207.3
LAGoNN _{exp}	238.0
RoBERTa _{full}	446.9

Table 6: Speed comparison between LAGoNN LAB-DIST with one neighbor and comparable methods. Time includes training on 1,000 examples and inference on 51,000 examples.

A.5 Additional results: 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, usually within a point, 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. In the table caption we provide a summary/interpretation of the results for a given setting. 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.

²¹We used a 40 GB NVIDIA A100 Tensor Core GPU.

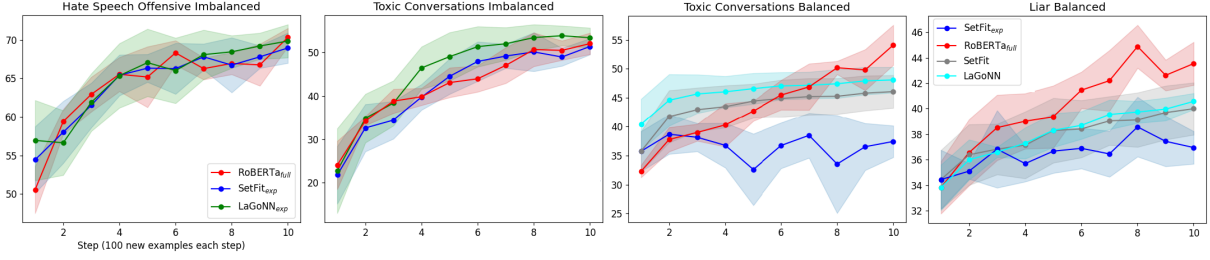


Figure 6: Average performance in the imbalanced and balanced regimes relative to comparable methods. We include RoBERTa_{full} results for reference. The metric is macro-F1 for Hate Speech Offensive, average precision elsewhere.

Method	Insincere-Questions			
	1 st	5 th	10 th	Average
<i>Imbalanced</i>				
RoBERTa _{full}	39.8 _{5.5}	53.1 _{4.6}	55.7 _{1.2}	50.6 _{4.4}
SetFit _{exp}	43.7 _{2.7}	52.2 _{1.9}	53.8 _{0.9}	51.4 _{2.9}
LAGONN _{exp}	44.5 _{4.5}	52.7 _{2.4}	55.4 _{2.0}	51.8 _{3.0}
RoBERTa _{freeze}	39.8 _{5.5}	44.1 _{3.6}	46.3 _{2.4}	44.0 _{2.0}
kNN	23.9 _{2.2}	30.3 _{3.0}	31.6 _{2.4}	30.0 _{2.1}
SetFit	43.7 _{2.7}	47.6 _{1.6}	50.1 _{2.1}	47.6 _{1.8}
LAGONN	44.5 _{4.5}	48.1 _{2.2}	50.3 _{1.7}	48.1 _{1.9}
Probe	40.4 _{4.2}	49.4 _{2.3}	52.3 _{1.7}	49.0 _{3.3}
LAGONN _{cheap}	40.8 _{4.3}	51.1 _{2.4}	54.5 _{1.4}	50.4 _{4.0}

Table 7: LAGONN and LAGONN_{exp} are the strongest performers on the first step, but are overtaken by RoBERTa_{full} on later steps. The average of all steps shows that LAGONN_{exp} is the overall strongest performer, but we note that LAGONN_{cheap} shows comparable performance to RoBERTa_{full} despite being much less expensive.

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

Table 8: LAGONN and LAGONN_{exp} are the strongest performers on the first step, but are overtaken by RoBERTa_{full} on later steps. The average of all steps shows that LAGONN is the overall strongest performer, but we note that LAGONN_{cheap} shows comparable performance to RoBERTa_{full} despite being much less expensive.

Method	Amazon Counterfactual			
	1 st	5 th	10 th	Average
<i>Imbalanced</i>				
RoBERTa _{full}	68.2 _{4.5}	81.0 _{1.7}	82.2 _{1.0}	79.2 _{3.9}
SetFit _{exp}	72.0 _{2.1}	78.4 _{2.8}	78.8 _{1.2}	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}
SetFit	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}

Table 9: LAGONN and LAGONN_{exp} are the strongest performers on the first step, but are overtaken by RoBERTa_{full} on later steps. However, the average of all steps shows that LAGONN_{exp} is the overall strongest performer.

Method	Amazon Counterfactual			
	1 st	5 th	10 th	Average
<i>Moderate</i>				
RoBERTa _{full}	73.9 _{2.5}	80.0 _{1.0}	80.1 _{2.3}	79.1 _{2.1}
SetFit _{exp}	76.5 _{1.6}	77.0 _{2.4}	74.7 _{0.5}	76.5 _{1.0}
LAGONN _{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}
SetFit	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}	63.1 _{4.3}
LAGONN _{cheap}	47.3 _{3.4}	60.7 _{1.5}	65.2 _{1.4}	59.5 _{5.2}

Table 10: LAGONN_{exp} and LAGONN are the strongest performers on the first step, but LAGONN is strongest classifier on subsequent steps and is also the overall strongest performer based on the average over all steps.

Method	Toxic Conversations			
	1 st	5 th	10 th	Average
<i>Extreme</i>				
RoBERTa _{full}	7.9 _{0.5}	21.2 _{3.7}	33.8 _{5.5}	21.9 _{9.3}
SetFit _{exp}	8.8 _{1.2}	18.1 _{3.4}	24.7 _{4.1}	17.6 _{5.5}
LAGONN _{exp}	8.9 _{1.7}	17.4 _{6.6}	26.4 _{5.2}	17.9 _{6.0}
RoBERTa _{freeze}	7.9 _{0.5}	12.8 _{2.4}	19.1 _{3.2}	13.5 _{3.5}
kNN	7.9 _{0.0}	8.7 _{0.4}	8.7 _{0.2}	8.5 _{0.3}
SetFit	8.8 _{1.2}	13.1 _{2.5}	16.3 _{3.0}	13.0 _{2.6}
LAGONN	8.9 _{1.7}	13.8 _{3.9}	17.1 _{4.8}	13.4 _{2.6}
Probe	13.1 _{2.8}	24.6 _{2.6}	30.1 _{2.1}	23.9 _{5.6}
LAGONN _{cheap}	11.3 _{2.2}	21.7 _{2.7}	27.4 _{2.3}	21.3 _{5.3}

Table 11: Probe is strongest performer on every step, except the 10th where it is overtaken by RoBERTa_{full}. If we average over all steps, we see that Probe is the strongest performer. We note, however, that LAGONN and LAGONN_{exp} outperform SetFit and SetFit_{exp} on all steps.

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

Table 12: RoBERTa_{full} and RoBERTa_{freeze} are the strongest performers on the first step, but are overtaken by LAGONN_{exp} for the subsequent steps. The overall strongest performer based on the average over all steps is LAGONN_{exp}.

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

Table 13: LAGONN and LAGONN_{exp} are the strongest performers on the first step and LAGONN_{exp} remains the strongest for subsequent steps, also being the strongest classifier overall based on the average.

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

Table 14: LAGONN and LAGONN_{exp} are the strongest performers on the first step. LAGONN remains the strongest until the 10th, where it is overtaken by RoBERTa_{full}. Overall, LAGONN is the strongest classifier based on the average. Note the performance of SetFit_{exp} and LAGONN_{exp}. While both degrade after the first step, LAGONN_{exp}'s performance drop is dramatically mitigated.

Method	Hate Speech Offensive			
	1 st	5 th	10 th	Average
RoBERTa _{full}	30.2 _{1.4}	43.5 _{2.5}	51.2 _{2.2}	44.3 _{7.4}
SetFit _{exp}	30.3 _{0.8}	44.0 _{1.3}	51.1 _{2.0}	43.8 _{6.5}
LAGONN _{exp}	30.3 _{0.7}	40.7 _{2.9}	49.1 _{4.4}	42.2 _{6.2}
RoBERTa _{freeze}	30.2 _{1.4}	33.5 _{3.1}	34.4 _{3.4}	33.1 _{1.4}
kNN	31.5 _{1.2}	35.9 _{2.7}	37.4 _{2.0}	35.8 _{1.7}
SetFit	30.3 _{0.8}	38.4 _{2.5}	41.1 _{1.5}	37.8 _{3.3}
LAGONN	30.3 _{0.7}	35.7 _{2.6}	39.1 _{2.4}	35.6 _{2.7}
Probe	29.0 _{0.2}	34.7 _{1.5}	40.1 _{2.1}	35.1 _{3.8}
LAGONN _{cheap}	29.0 _{0.1}	36.9 _{1.8}	40.5 _{2.1}	36.2 _{3.7}

Table 15: kNN is the strongest performer on the first step, while SetFit_{exp} is on the 5th, and RoBERTa_{full} is the strongest on the 10th while also being strongest overall performer for all steps. LAGONN-based methods are generally beaten by ST/SetFit-based baselines, with the exception of LAGONN_{cheap} which consistently outperforms Probe.

Method	Hate Speech Offensive			
	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 _{exp}	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}
RoBERTa _{freeze}	50.6 _{3.0}	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}
SetFit	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 16: LAGONN and LAGONN_{exp} are the strongest performers on the first step, with LAGONN_{exp} being the strongest on the 5th and RoBERTa_{full} taking over on the 10th. LAGONN_{exp} is the strongest performer overall based on the average over all steps.

Method	Hate Speech Offensive			
	1 st	5 th	10 th	Average
RoBERTa _{full}	61.9 _{3.4}	70.8 _{1.0}	72.5 _{1.4}	69.9 _{3.2}
SetFit _{exp}	64.3 _{4.2}	70.6 _{2.4}	72.4 _{0.5}	69.8 _{2.8}
LAGONN _{exp}	63.8 _{4.9}	71.0 _{2.1}	72.3 _{1.0}	70.0 _{3.0}
RoBERTa _{freeze}	61.9 _{3.4}	63.2 _{4.1}	64.1 _{4.5}	63.2 _{0.6}
kNN	64.3 _{4.0}	63.3 _{2.9}	63.9 _{2.5}	63.7 _{0.4}
SetFit	64.3 _{4.2}	67.3 _{3.2}	67.6 _{2.3}	66.9 _{1.1}
LAGONN	63.8 _{4.9}	65.0 _{5.3}	66.7 _{5.9}	65.3 _{0.9}
Probe	55.6 _{1.7}	63.8 _{0.8}	66.1 _{0.3}	63.2 _{3.0}
LAGONN _{cheap}	56.0 _{3.6}	62.2 _{1.4}	66.0 _{0.9}	62.3 _{2.9}

Table 17: kNN, SetFit, and SetFit_{exp} start the strongest, but are overtaken by LAGONN_{exp} on the 5th step, which is in turn overtaken by RoBERTa_{full} on the 10th step. Overall LAGONN_{exp} is the strongest performer based on the average.

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

Table 18: LAGONN and LAGONN_{exp} are the strongest performers on the first step, but are overtaken by RoBERTa_{full} on later steps, which also is the strongest overall classifier. We note that LAGONN and LAGONN_{exp} consistently outperform SetFit and SetFit_{exp}, respectively.

Method	LIAR			
	1 st	5 th	10 th	Average
RoBERTa _{full}	32.0 _{2.7}	34.7 _{2.9}	35.1 _{4.3}	33.7 _{1.0}
SetFit _{exp}	31.2 _{3.8}	30.4 _{3.1}	31.8 _{2.9}	31.5 _{0.7}
LAGONN _{exp}	30.6 _{4.7}	30.3 _{2.0}	31.3 _{2.0}	31.1 _{0.6}
RoBERTa _{freeze}	32.0 _{2.7}	32.8 _{4.5}	34.2 _{5.0}	33.2 _{0.7}
kNN	27.0 _{0.5}	27.3 _{0.8}	27.9 _{0.8}	27.4 _{0.3}
SetFit	31.2 _{3.8}	33.7 _{5.1}	35.7 _{5.1}	34.3 _{1.6}
LAGONN	30.6 _{4.7}	32.0 _{4.6}	33.7 _{5.4}	32.6 _{0.9}
Probe	30.7 _{2.0}	30.6 _{3.9}	31.7 _{2.9}	31.1 _{0.4}
LAGONN _{cheap}	30.7 _{2.0}	30.5 _{3.8}	31.4 _{2.6}	31.0 _{0.4}

Table 19: RoBERTa_{freeze} and RoBERTa_{full} start out as the strongest performers but are eventually overtaken by SetFit on the 10th step, and SetFit ends up being the strongest performer over all steps based on the average.

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

Table 20: SetFit, SetFit_{exp}, LAGONN, and LAGONN_{exp} start out as the strongest performers. On the 5th step, SetFit is overtaken by the other systems, but is eventually overtaken by RoBERTa_{full}. Overall SetFit is the strongest system, but we note that LAGONN_{exp} outperforms SetFit_{exp}.

Method	LIAR				
	<i>Moderate</i>	1 st	5 th	10 th	Average
RoBERTa _{full}	33.9 _{3.1}	38.4 _{2.7}	43.9 _{2.2}	39.5 _{3.0}	
SetFit _{exp}	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}	
SetFit	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 21: LAGONN and LAGONN_{exp} start out as the strongest performers and LAGONN_{exp} continues to be strong, until the 10th step where it is overtaken by RoBERTa_{full}, which ends up as the most performant classifier over all steps based on the average.

Method	LIAR				
	<i>Balanced</i>	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 _{exp}	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}	
SetFit	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}	

Table 22: SetFit and SetFit_{exp} are the most performant systems on the first step, but are overtaken by RoBERTa_{full}, the strongest overall classifier. We note that LAGONN outperforms SetFit after the first step and in aggregate.

A.6 Additional results: secondary experiments

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, usually within one point, 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 a given dataset and a given label-balance distribution on the first, fifth, and tenth step followed by the average for all ten steps. In the table caption we provide a summary/interpretation of the results for a given setting. 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	Insincere Questions				
	<i>Extreme</i>	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 _{exp}	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}	24.1 _{6.3}	38.1 _{6.3}	41.1 _{6.5}	35.6 _{5.5}	
LAGONN _{lite}	30.7 _{8.9}	41.8 _{8.3}	43.4 _{8.5}	39.3 _{4.4}	
RoBERTa _{freeze}	19.9 _{8.4}	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}	
SetFit	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}	44.8 _{4.2}	38.3 _{6.2}	
LAGONN _{cheap}	23.6 _{7.8}	40.7 _{5.9}	45.3 _{4.4}	38.6 _{6.6}	

Table 23: LAGONN, LAGONN_{lite}, and LAGONN_{exp} start out as the strongest models, but LAGONN_{lite} remains the most performant by the 10th step. It is also the overall strongest performer based on the average. We note the strength of LAGONN_{cheap} relative to far more expensive methods.

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

Table 24: LAGONN, LAGONN_{lite}, and LAGONN_{exp} start out as the strongest models, but LAGONN_{lite} remains the most performant by the 10th step. It is also the overall strongest performer based on the average. We note the strength of LAGONN_{cheap} relative to far more expensive methods.

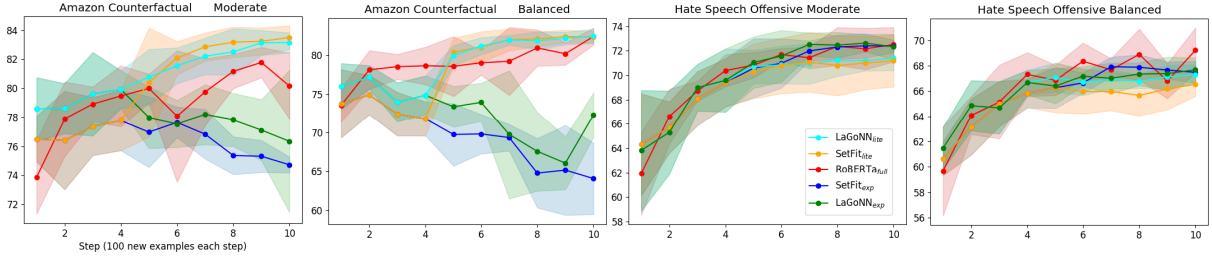


Figure 7: Average performance for all the moderate and balanced sampling regimes on Amazon Counterfactual and Hate Speech Offensive. More expensive models, such as LAGoNN_{exp}, SetFit_{exp}, and RoBERTa_{full} perform best when the label distribution is imbalanced. As the distribution becomes more balanced, inexpensive models, such as LAGoNN_{lite}, show similar or improved performance. The metric is average precision for Amazon Counterfactual and the macro F1 for Hate Speech Offensive. We only consider one neighbor for the LAGoNN-based methods.

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

Table 25: LAGoNN, LAGoNN_{lite}, and LAGoNN_{exp} start out as the strongest models, but SetFit_{lite} overtakes the other methods by the 5th step and is the strongest performer based on the average. We note the strength of LAGoNN_{cheap} relative to far more expensive methods.

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

Table 26: LAGoNN_{cheap} starts out as the strongest model, but SetFit_{lite} overtakes the other methods on the 5th and 10th step. Overall LAGoNN_{cheap} is the strongest model despite being one of the least expensive.

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

Table 27: LAGoNN, LAGoNN_{lite}, and LAGoNN_{exp} are the most performant models on the first step, but only LAGoNN_{exp} remains the most performant on subsequent steps, also being the strongest overall method based on the average over all steps.

Method	Amazon Counterfactual			
	1 st	5 th	10 th	Average
<i>Imbalanced</i>				
RoBERTa _{full}	68.2 _{4.5}	81.0 _{1.7}	82.2 _{1.0}	79.2 _{3.9}
SetFit _{exp}	72.0 _{2.1}	78.4 _{2.8}	78.8 _{1.2}	78.0 _{2.1}
LAGoNN _{exp}	74.3 _{3.8}	80.1 _{1.4}	79.0 _{1.6}	79.5 _{1.9}
SetFit _{lite}	72.0 _{2.1}	79.1 _{1.4}	81.6 _{1.3}	79.1 _{2.7}
LAGoNN _{lite}	74.3 _{3.8}	79.2 _{1.7}	81.9 _{1.1}	80.2 _{2.2}
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}
SetFit	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}

Table 28: On the first step, LAGoNN, LAGoNN_{lite}, and LAGoNN_{exp} start out the strongest but LAGoNN_{lite} performs slightly worse than RoBERTa_{full} on the 5th and 10th step. However, LAGoNN_{lite} is the best overall method based on the average.

Method	Amazon Counterfactual			
	1 st	5 th	10 th	Average
RoBERTa _{full}	73.9 _{2.5}	80.0 _{1.0}	80.1 _{2.3}	79.1 _{2.1}
SetFit _{exp}	76.5 _{1.6}	77.0 _{2.4}	74.7 _{0.5}	76.5 _{1.0}
LAGoNN _{exp}	78.6 _{2.2}	78.0 _{2.1}	76.3 _{4.9}	78.2 _{1.0}
SetFit _{lite}	76.5 _{1.6}	80.4 _{3.8}	83.5 _{0.8}	80.3 _{2.8}
LAGoNN _{lite}	78.6 _{2.2}	80.8 _{1.9}	83.1 _{0.7}	81.0 _{1.7}
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}
SetFit	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}	63.1 _{4.3}
LAGoNN _{cheap}	47.3 _{3.4}	60.7 _{1.5}	65.2 _{1.4}	59.5 _{5.2}

Table 29: On the first step, LAGoNN, LAGoNN_{lite}, and LAGoNN_{exp} start out the strongest. On the 5th step, LAGoNN is the most performant method while on the 10th step it is SetFit_{lite}. However, LAGoNN_{lite} is the best overall method based on the average.

Method	Amazon Counterfactual			
	1 st	5 th	10 th	Average
RoBERTa _{full}	73.6 _{2.1}	78.6 _{3.9}	82.4 _{1.1}	78.9 _{2.2}
SetFit _{exp}	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}	73.8 _{4.4}	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}
SetFit	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	52.4 _{3.4}	64.7 _{2.5}	67.5 _{0.4}	63.4 _{4.4}
LAGoNN _{cheap}	48.1 _{3.4}	62.0 _{2.0}	65.3 _{0.8}	60.5 _{5.0}

Table 30: On the first step, LAGoNN, LAGoNN_{lite}, and LAGoNN_{exp} start out the strongest. On the 5th step, SetFit_{lite} pulls ahead slightly, yet on the 10th step LAGoNN_{lite} is the best performer. Overall, LAGoNN is the best method based on the average.

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

Table 31: Probe is most performant method on all steps and the overall strongest performer. We note, however, that LAGoNN-based methods tend to outperform their SetFit-based counterparts.

Method	Toxic Conversations			
	1 st	5 th	10 th	Average
RoBERTa _{full}	24.1 _{5.6}	43.1 _{3.4}	52.1 _{2.5}	42.4 _{8.2}
SetFit _{exp}	21.8 _{6.6}	44.5 _{4.1}	51.4 _{1.9}	42.1 _{9.3}
LAGoNN _{exp}	22.7 _{9.8}	49.1 _{5.6}	53.4 _{2.3}	45.6 _{9.8}
SetFit _{lite}	21.8 _{6.6}	41.4 _{4.4}	44.8 _{3.1}	39.0 _{7.0}
LAGoNN _{lite}	22.7 _{9.8}	47.0 _{6.3}	50.2 _{5.4}	43.7 _{8.6}
RoBERTa _{freeze}	24.1 _{5.6}	31.2 _{4.4}	34.0 _{4.0}	30.5 _{3.1}
kNN	11.5 _{2.5}	14.7 _{4.0}	15.3 _{3.2}	14.6 _{1.1}
SetFit	21.8 _{6.6}	26.7 _{5.3}	30.2 _{4.0}	26.6 _{2.7}
LAGoNN	22.7 _{9.8}	27.6 _{8.9}	30.3 _{3.7}	27.4 _{2.4}
Probe	23.3 _{2.7}	33.0 _{2.8}	37.1 _{1.8}	32.5 _{4.2}
LAGoNN _{cheap}	20.5 _{3.2}	31.1 _{3.2}	35.6 _{1.8}	30.5 _{4.6}

Table 32: RoBERTa_{full} and RoBERTa_{freeze} start out as the strongest classifiers on the first step, but are overtaken on subsequent steps by LAGoNN_{exp}, which ends up as strongest method overall.

Method	Toxic Conversations			
	1 st	5 th	10 th	Average
RoBERTa _{full}	34.2 _{3.4}	45.5 _{1.9}	52.4 _{3.3}	45.7 _{5.6}
SetFit _{exp}	33.6 _{2.9}	47.2 _{2.2}	46.6 _{3.3}	44.3 _{4.3}
LAGoNN _{exp}	36.6 _{4.2}	48.2 _{2.7}	49.9 _{3.7}	48.0 _{4.4}
SetFit _{lite}	33.6 _{2.9}	52.6 _{2.0}	55.1 _{1.6}	48.8 _{7.3}
LAGoNN _{lite}	36.6 _{4.2}	56.1 _{1.5}	57.7 _{1.4}	52.3 _{6.8}
RoBERTa _{freeze}	34.2 _{3.4}	38.4 _{2.1}	39.5 _{1.8}	38.0 _{1.5}
kNN	19.4 _{1.9}	21.5 _{3.4}	22.4 _{2.9}	21.6 _{0.8}
SetFit	33.6 _{2.9}	39.2 _{2.9}	41.6 _{2.7}	38.6 _{2.4}
LAGoNN	36.6 _{4.2}	42.7 _{3.7}	45.0 _{3.5}	42.0 _{2.5}
Probe	29.0 _{2.7}	36.1 _{1.2}	39.1 _{1.5}	35.5 _{3.3}
LAGoNN _{cheap}	26.1 _{2.7}	34.3 _{1.3}	37.5 _{1.8}	33.6 _{3.6}

Table 33: On the first step, LAGoNN, LAGoNN_{lite}, and LAGoNN_{exp} start out the strongest, but it is LAGoNN_{lite} that remains performant for all other steps. LAGoNN_{lite} is also the strongest overall method based on the average.

Method	Toxic Conversations			
	1 st	5 th	10 th	Average
RoBERTa _{full}	32.3 _{1.1}	42.7 _{1.8}	54.1 _{3.4}	43.8 _{6.3}
SetFit _{exp}	35.7 _{3.4}	32.6 _{6.2}	37.4 _{2.7}	36.5 _{1.9}
LAGoNN _{exp}	40.4 _{4.4}	40.2 _{6.6}	39.8 _{7.5}	40.0 _{1.2}
SetFit _{lite}	35.7 _{3.4}	52.7 _{2.5}	53.9 _{2.2}	46.8 _{7.8}
LAGoNN _{lite}	40.4 _{4.4}	52.9 _{2.6}	54.0 _{2.3}	48.3 _{6.4}
RoBERTa _{freeze}	32.3 _{1.1}	39.2 _{1.5}	41.0 _{0.6}	38.5 _{2.4}
kNN	17.4 _{0.8}	23.7 _{2.6}	24.3 _{2.7}	23.1 _{2.0}
SetFit	35.7 _{3.4}	44.5 _{2.9}	46.1 _{2.8}	43.6 _{2.9}
LAGoNN	40.4 _{4.4}	46.6 _{2.7}	48.1 _{2.2}	46.1 _{2.2}
Probe	29.5 _{2.4}	35.9 _{0.9}	40.2 _{0.9}	36.1 _{3.5}
LAGoNN _{cheap}	26.8 _{2.7}	34.5 _{1.3}	38.5 _{0.8}	34.4 _{3.7}

Table 34: On the first step, LAGoNN, LAGoNN_{lite}, and LAGoNN_{exp} start out the strongest, but it is LAGoNN_{lite} that remains performant for all other steps. LAGoNN_{lite} is also the strongest overall method based on the average.

Method	Hate Speech Offensive			
	1 st	5 th	10 th	Average
<i>Extreme</i>				
RoBERTa _{full}	30.2 _{1.4}	43.5 _{2.5}	51.2 _{2.2}	44.3 _{7.4}
SetFit _{exp}	30.3 _{0.8}	44.0 _{1.3}	51.1 _{2.0}	43.8 _{6.5}
LAGONN _{exp}	30.3 _{0.7}	40.7 _{2.9}	49.1 _{4.4}	42.2 _{6.2}
SetFit _{lite}	30.3 _{0.8}	43.4 _{2.5}	45.5 _{3.4}	41.6 _{4.6}
LAGONN _{lite}	30.3 _{0.7}	40.9 _{3.4}	41.5 _{4.8}	39.1 _{3.6}
RoBERTa _{freeze}	30.2 _{1.4}	33.5 _{3.1}	34.4 _{3.4}	33.1 _{1.4}
kNN	31.5 _{1.2}	35.9 _{2.7}	37.4 _{2.0}	35.8 _{1.7}
SetFit	30.3 _{0.8}	38.4 _{2.5}	41.1 _{1.5}	37.8 _{3.3}
LAGONN	30.3 _{0.7}	35.7 _{2.6}	39.1 _{2.4}	35.6 _{2.7}
Probe	29.0 _{0.2}	34.7 _{1.5}	40.1 _{2.1}	35.1 _{3.8}
LAGONN _{cheap}	29.0 _{0.1}	36.9 _{1.8}	40.5 _{2.1}	36.2 _{3.7}

Table 35: kNN is the strongest method at first, but is overtaken by SetFit_{exp} on the 5th step, which is then overtaken by RoBERTa_{full} on the 10th step. RoBERTa_{full} is overall most performant system based on the average.

Method	Hate Speech Offensive			
	1 st	5 th	10 th	Average
<i>Imbalanced</i>				
RoBERTa _{full}	50.6 _{3.0}	65.2 _{3.9}	70.3 _{1.2}	64.2 _{5.3}
SetFit _{exp}	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}	54.4 _{4.3}	65.5 _{3.0}	65.9 _{3.5}	63.5 _{3.9}
LAGONN _{lite}	57.0 _{5.2}	66.6 _{2.6}	66.6 _{1.9}	64.3 _{4.1}
RoBERTa _{freeze}	50.6 _{3.0}	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}
SetFit	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 36: On the first step, LAGONN, LAGONN_{lite}, and LAGONN_{exp} start out the strongest, and LAGONN_{exp} continues to be performant, but is overtaken on the 10th step by RoBERTa_{full}. LAGONN_{exp} is the strongest overall method based on the average.

Method	Hate Speech Offensive			
	1 st	5 th	10 th	Average
<i>Moderate</i>				
RoBERTa _{full}	61.9 _{3.4}	70.8 _{1.0}	72.5 _{1.4}	69.9 _{3.2}
SetFit _{exp}	64.3 _{4.2}	70.6 _{2.4}	72.4 _{0.5}	69.8 _{2.8}
LAGONN _{exp}	63.8 _{4.9}	71.0 _{2.1}	72.3 _{1.0}	70.0 _{3.0}
SetFit _{lite}	64.3 _{4.2}	70.3 _{2.2}	71.2 _{2.1}	69.3 _{2.3}
LAGONN _{lite}	63.8 _{4.9}	70.7 _{1.4}	71.4 _{1.0}	69.4 _{2.5}
RoBERTa _{freeze}	61.9 _{3.4}	63.2 _{4.1}	64.1 _{4.5}	63.2 _{0.6}
kNN	64.3 _{4.0}	63.3 _{2.9}	63.9 _{2.5}	63.7 _{0.4}
SetFit	64.3 _{4.2}	67.3 _{3.2}	67.6 _{2.3}	66.9 _{1.1}
LAGONN	63.8 _{4.9}	65.0 _{5.3}	66.7 _{5.9}	65.3 _{0.9}
Probe	55.6 _{1.7}	63.8 _{0.8}	66.1 _{0.3}	63.2 _{3.0}
LAGONN _{cheap}	56.0 _{3.6}	62.2 _{1.4}	66.0 _{0.9}	62.3 _{2.9}

Table 37: Similar to the imbalanced setting, on the first step, LAGONN, LAGONN_{lite}, and LAGONN_{exp} start out the strongest, and LAGONN_{exp} continues to be performant, but is overtaken on the 10th step by RoBERTa_{full}. LAGONN_{exp} is the strongest overall method based on the average.

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

Table 38: Similar to the moderate setting, on the first step, LAGONN, LAGONN_{lite}, and LAGONN_{exp} start out the strongest, but RoBERTa_{full} overtakes LAGONN_{lite} by the 10th step. RoBERTa_{full} slightly outperforms LAGONN_{lite} and LAGONN_{exp} as the overall strongest method based on the average.

Method	LIAR			
	1 st	5 th	10 th	Average
<i>Extreme</i>				
RoBERTa _{full}	32.0 _{2.7}	34.7 _{2.9}	35.1 _{4.3}	33.7 _{1.0}
SetFit _{exp}	31.2 _{3.8}	30.4 _{3.1}	31.8 _{2.9}	31.5 _{0.7}
LAGONN _{exp}	30.6 _{4.7}	30.3 _{2.0}	31.3 _{2.0}	31.1 _{0.6}
SetFit _{lite}	31.2 _{3.8}	32.7 _{3.8}	33.5 _{4.2}	32.7 _{0.8}
LAGONN _{lite}	30.6 _{4.7}	31.8 _{3.9}	32.4 _{2.7}	31.6 _{0.6}
RoBERTa _{freeze}	32.0 _{2.7}	32.8 _{4.5}	34.2 _{5.0}	33.2 _{0.7}
kNN	27.0 _{0.5}	27.3 _{0.8}	27.9 _{0.8}	27.4 _{0.3}
SetFit	31.2 _{3.8}	33.7 _{5.1}	35.7 _{5.1}	34.3 _{1.6}
LAGONN	30.6 _{4.7}	32.0 _{4.6}	33.7 _{5.4}	32.6 _{0.9}
Probe	30.7 _{2.0}	30.6 _{3.9}	31.7 _{2.9}	31.1 _{0.4}
LAGONN _{cheap}	30.7 _{2.0}	30.5 _{3.8}	31.4 _{2.6}	31.0 _{0.4}

Table 39: RoBERTa_{freeze} and RoBERTa_{full} start out performant and RoBERTa_{full} continues to be until the 10th step where it is overtaken by SetFit, which ends up being the strongest overall method.

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

Table 40: LAGONN, LAGONN_{lite}, LAGONN_{exp}, SetFit, SetFit_{lite}, and SetFit_{exp} start out as the most performant, but SetFit is the strongest on the 5th step and RoBERTa_{full} on the 10th. Overall, SetFit is strongest method based on the average over all steps.

Method	LIAR			
	1 st	5 th	10 th	Average
<i>Moderate</i>				
RoBERTa _{full}	33.9 _{3.1}	38.4 _{2.7}	43.9 _{2.2}	39.5 _{3.0}
SetFit _{exp}	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}
SetFit _{lite}	33.0 _{2.6}	38.5 _{1.3}	40.4 _{2.0}	38.2 _{2.1}
LAGONN _{lite}	34.1 _{3.4}	38.4 _{2.0}	39.6 _{1.5}	37.9 _{1.6}
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}
SetFit	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 41: LAGONN, LAGONN_{lite}, and LAGONN_{exp} are the most performant classifiers on the first step, while LAGONN_{exp} remains strong until the 10th step where it is overtaken by RoBERTa_{full}. RoBERTa_{full} is the overall strongest method if we aggregate over all steps.

Method	LIAR			
	1 st	5 th	10 th	Average
<i>Balanced</i>				
RoBERTa _{full}	33.8 _{2.1}	39.4 _{2.4}	43.5 _{1.7}	40.2 _{3.2}
SetFit _{exp}	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}
SetFit _{lite}	34.4 _{2.3}	38.7 _{2.3}	40.3 _{2.8}	38.0 _{2.1}
LAGONN _{lite}	33.8 _{1.8}	37.6 _{2.0}	39.4 _{2.8}	37.2 _{1.9}
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}
SetFit	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}

Table 42: SetFit, SetFit_{lite}, and SetFit_{exp} start out the strongest on the first step, but are overtaken by RoBERTa_{full} on the 5th which remains the most performant on the 10th step and if we consider the average over all steps.

988 **A.7 Additional results: general text** 989 **classification**

990 In this Appendix section, we provide additional
991 results from our general text classification experi-
992 ments in the main text, Section 6. Here we show
993 results comparing LAGONN_{lite} against SetFit_{lite}
994 and LAGONN_{exp} against SetFit_{exp}, but we include
995 results for one to five neighbors with LAGONN
996 LABDIST, Figures 8 and 9, respectively. The met-
997 ric is average precision for IMDB, macro-F1 else-
998 where.

999 In general, the number of neighbors we con-
1000 sider does not appear to have a large impact on
1001 LAGONN’s predictive power and our method con-
1002 tinues to be a more stable classifier than SetFit and
1003 can generally be expected to improve SetFit’s per-
1004 formance. We also see that continued fine-tuning
1005 with the embedding model is only helpful for cases
1006 when the dataset has a relatively large number of
1007 labels. One exception to this is the case of Student
1008 Question Categories, where there are four labels.
1009 While it is clear that SetFit_{lite} is a stronger model
1010 than LAGONN lite, if we consider the more expen-
1011 sive alternatives, the story changes; if we continue
1012 to fine-tune, the prediction curves are essentially
1013 the same, and LAGONN_{exp} seems to have a slight
1014 edge on SetFit_{exp} as we add training data.

1015 LIAR, both the collapsed version we consid-
1016 ered in our content moderation experiments and
1017 the original version (Orig Liar) we examine in our
1018 general text classification experiments here, seems
1019 to be a very difficult dataset. Adding examples
1020 or increased fine-tuning does not appear to consis-
1021 tently increase model performance. We observed
1022 this across all experimental settings and balanced
1023 regimes and is a sensible finding, as it should be
1024 very difficult to determine the truth of a specific
1025 statement without additional context.

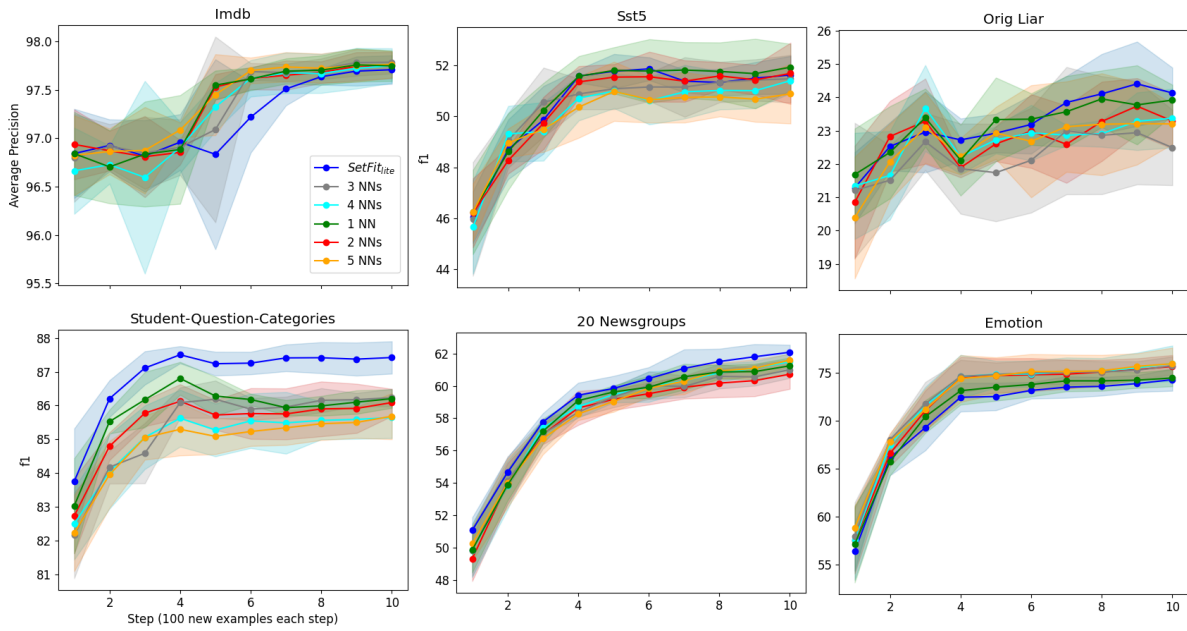


Figure 8: SetFit_{lite} performance compared against one to five neighbors for LAGONN_{lite} LABDIST. The metric is average precision for IMDB, macro-F1 elsewhere.

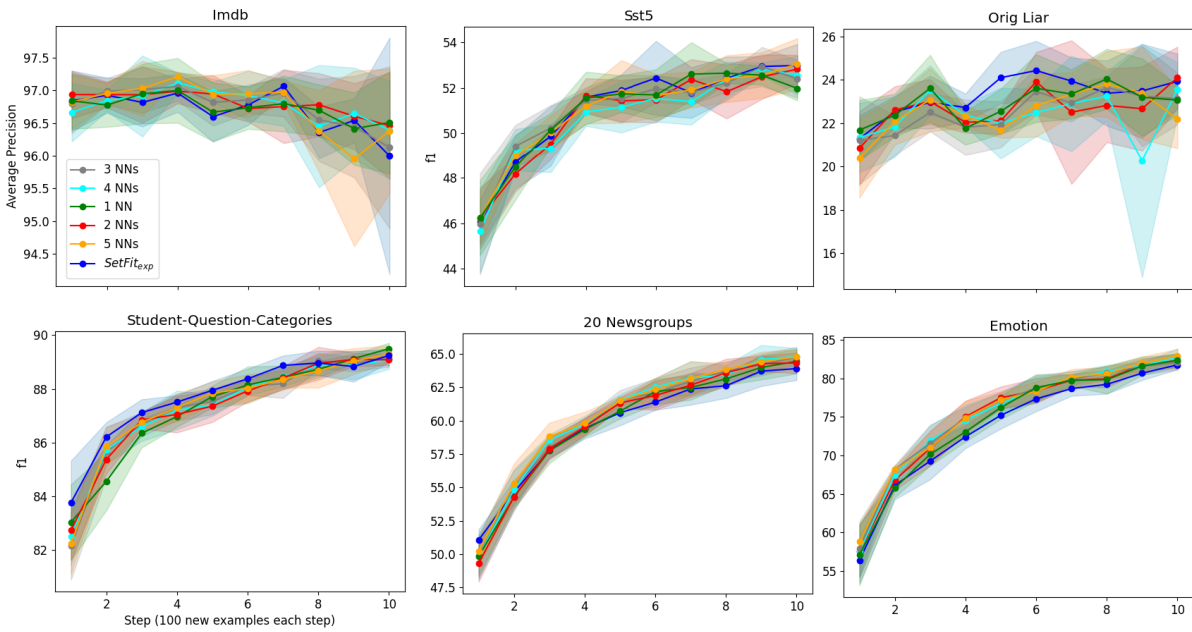


Figure 9: SetFit_{exp} performance compared against one to five neighbors for LAGONN_{exp} LABDIST. The metric is average precision for IMDB, macro-F1 elsewhere.

A.8 Ablations

In this Appendix section, we perform ablation studies with LAGONN to support our findings in the main text.

A.8.1 Ablation: LAGONN configurations

Here, we provide an in-depth comparison between all LAGONN configurations, LABEL, DISTANCE, LABDIST, TEXT, and ALL (see Table 1) for all content moderation datasets, balances, and levels of expense. The evaluation metric is the mean average precision ($\times 100$) over five seeds in all cases except for Hate Speech Offensive where the metric is the macro-F1.

Below, Figures 10 through 14 are the results for the LAGONN_{cheap} training strategy, Figures 15 through 19 are the results for LAGONN, Figures 20 through 24 are the results for LAGONN_{lite}, and Figures 25 through 29 are the results for LAGONN_{exp}. We place the figures on a new page for ease of viewing.

In the case of LAGONN_{cheap}, if we do not fine-tune the embedding model we see little variation in the standard deviation bands, with the exception of the LIAR dataset, which seems to be a very difficult dataset. When we do fine-tune, we see a great deal of variation, especially in cases of label imbalance, which is expected as the representations are altered more. The performance of TEXT and ALL is very unstable, often being the worst performers, while sometimes being the best. Interestingly, we note that DISTANCE, LABEL, and LABDIST often show very similar performance. In our opinion, LABDIST seems to be the most consistent and stable performer, especially in cases when the embedding model is fine-tuned, LAGONN, LAGONN_{lite}, and LAGONN_{exp}.

Overall, we believe that LABDIST is the most performant/stable configuration of LAGONN, and it is about this version that we present results in the main text. We note that we could have presented the best performer for each evaluation scenario, however, this is not in the spirit of our work as it adds yet another hyperparameter to configure, standing in the way of practical usage and convoluting our analysis. However, in our codebase, we hope that we have made it easy for one to change these configurations for their own usage, be it scientific or otherwise.

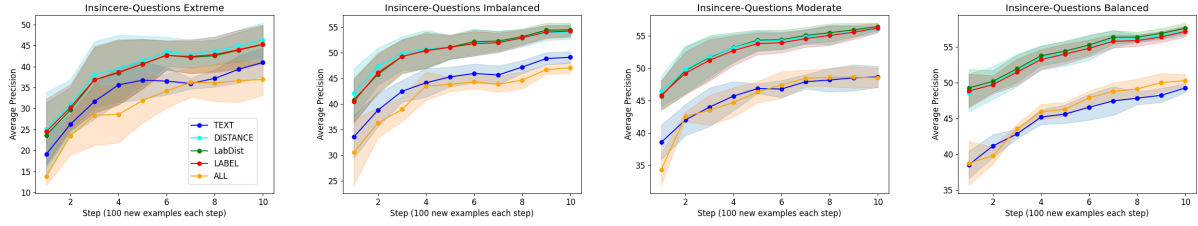


Figure 10: LAGONN_{cheap} performance for all configurations and balance regimes on the Insincere Questions dataset. The relevant balance is in the title of each panel.

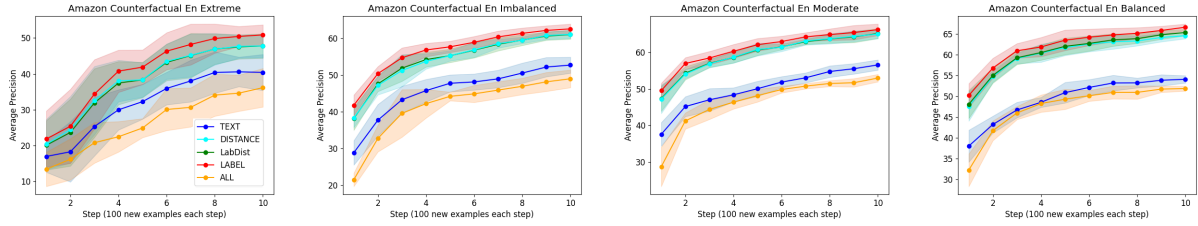


Figure 11: LAGONN_{cheap} performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

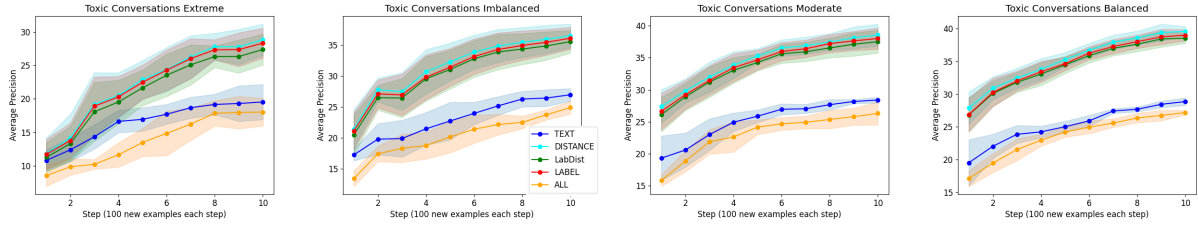


Figure 12: LAGONN_{cheap} performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

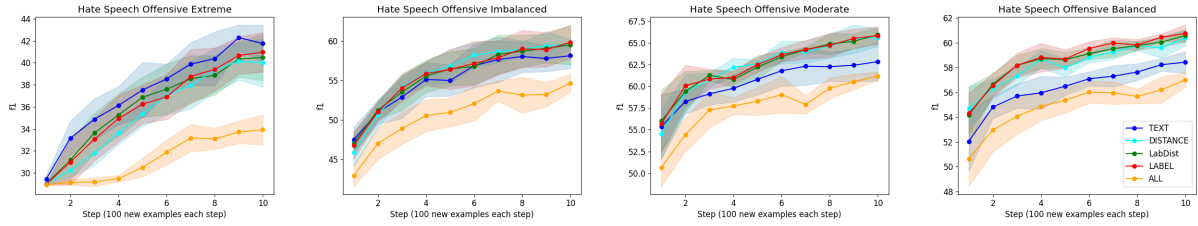


Figure 13: LAGONN_{cheap} performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

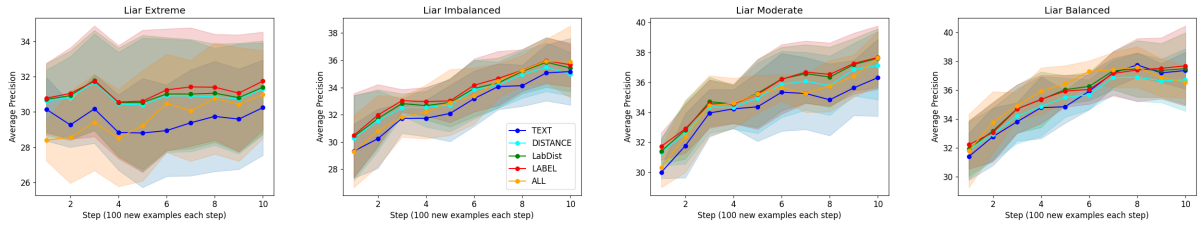


Figure 14: LAGONN_{cheap} performance for all configurations and balance regimes on the LIAR dataset. The relevant balance is in the title of each panel.

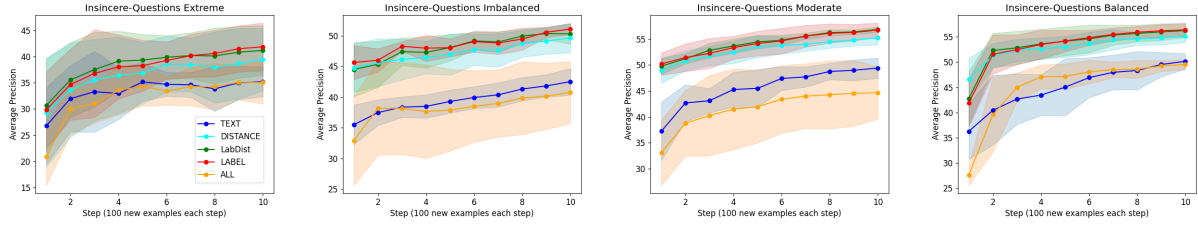


Figure 15: LAGONN performance for all configurations and balance regimes on the Insincere Questions dataset. The relevant balance is in the title of each panel.

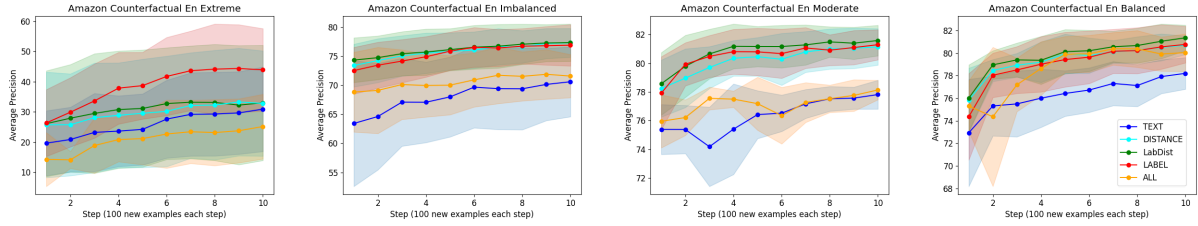


Figure 16: LAGONN performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

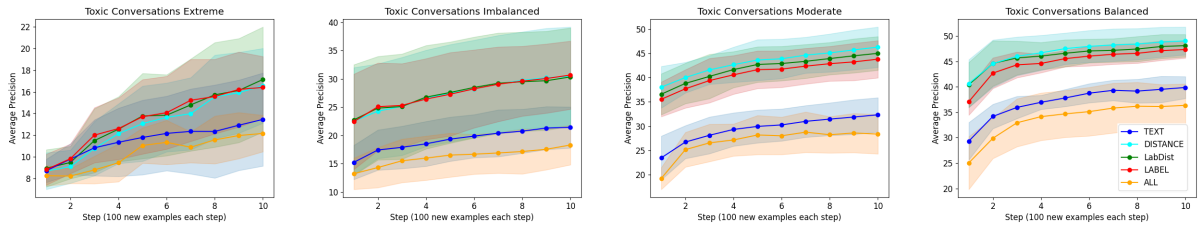


Figure 17: LAGONN performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

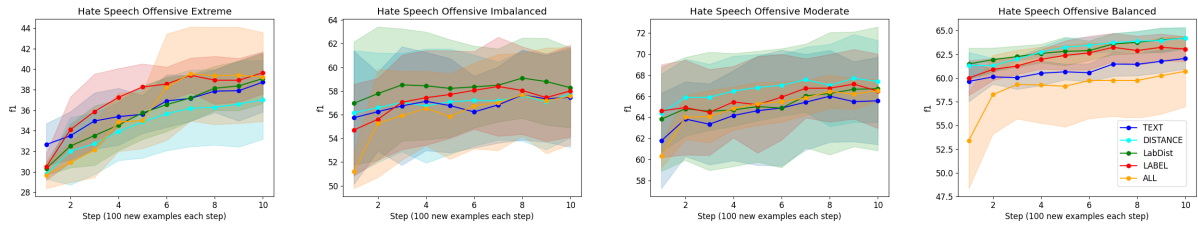


Figure 18: LAGONN performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

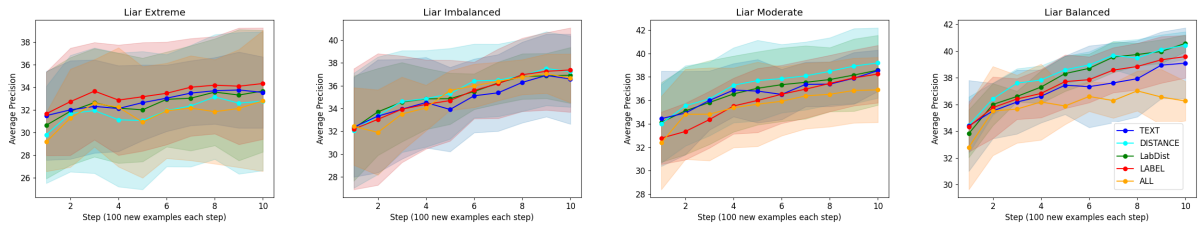


Figure 19: LAGONN performance for all configurations and balance regimes on the LIAR dataset. The relevant balance is in the title of each panel.

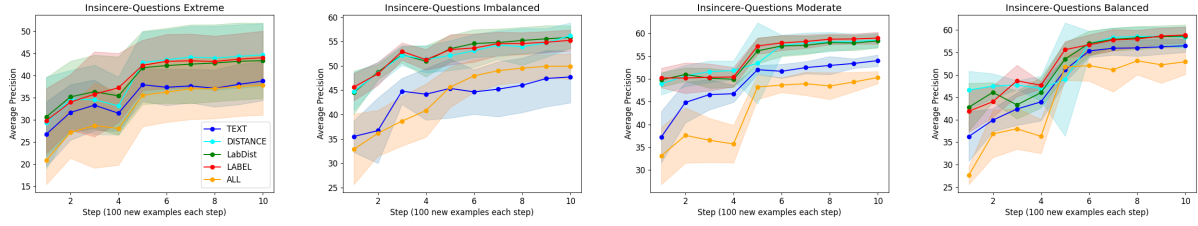


Figure 20: LAGONN_{lite} performance for all configurations and balance regimes on the Insincere Questions dataset. The relevant balance is in the title of each panel.

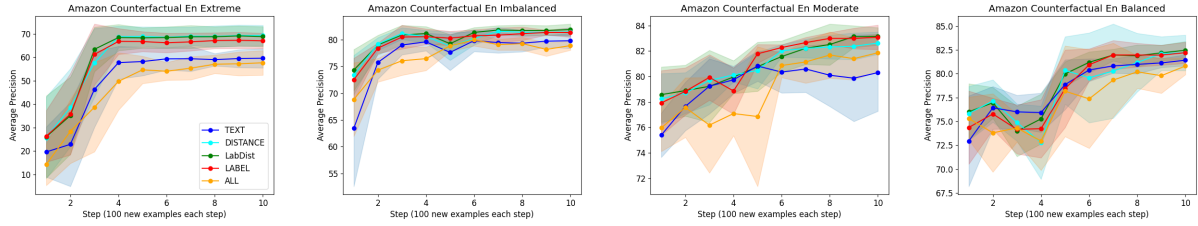


Figure 21: LAGONN_{lite} performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

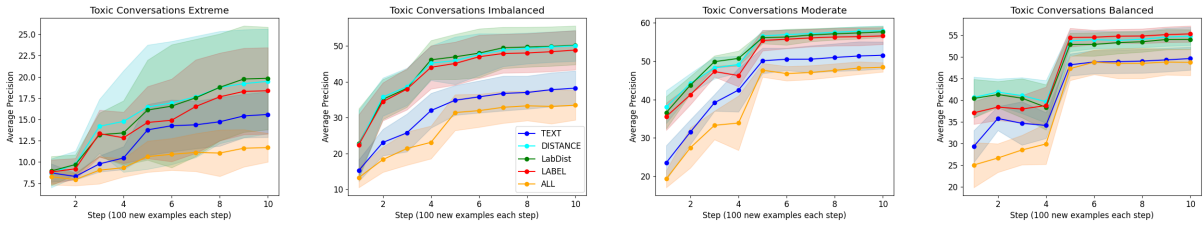


Figure 22: LAGONN_{lite} performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

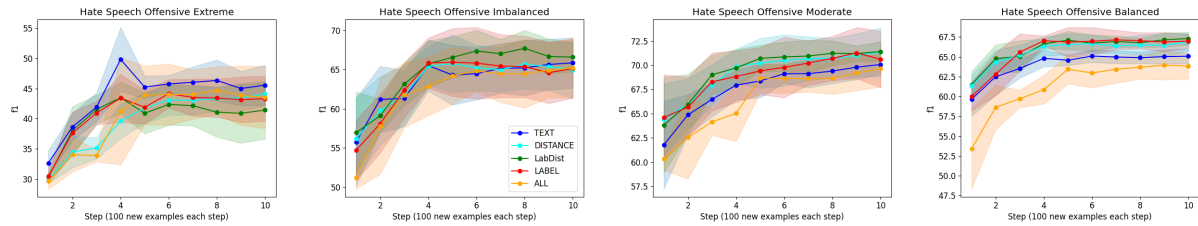


Figure 23: LAGONN_{lite} performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

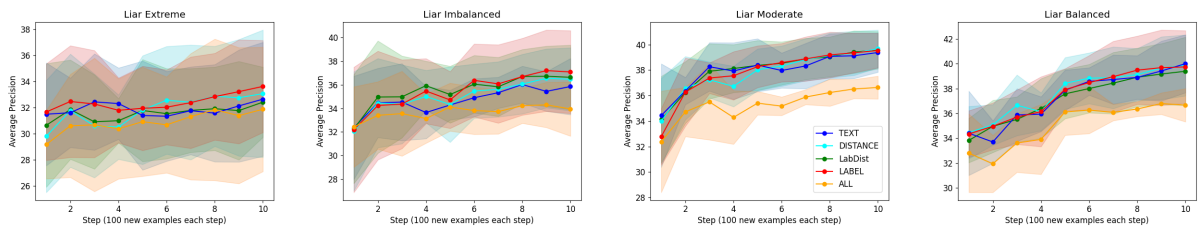


Figure 24: LAGONN_{lite} performance for all configurations and balance regimes on the LIAR dataset. The relevant balance is in the title of each panel.

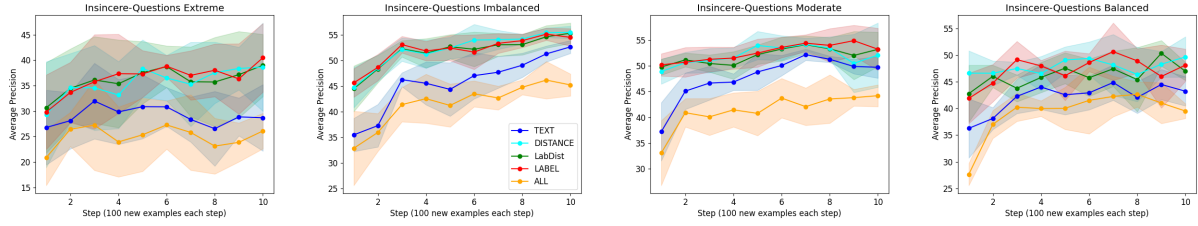


Figure 25: LAGONN_{exp} performance for all configurations and balance regimes on the Insincere Questions dataset. The relevant balance is in the title of each panel.

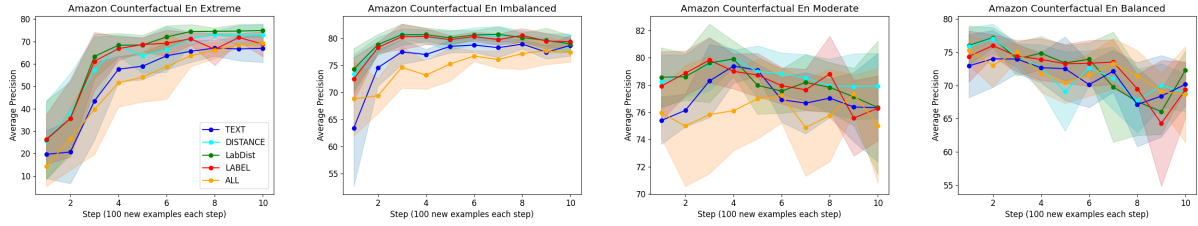


Figure 26: LAGONN_{exp} performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

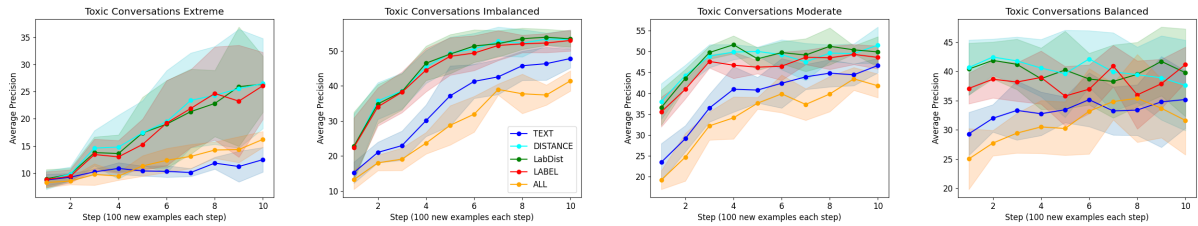


Figure 27: LAGONN_{exp} performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

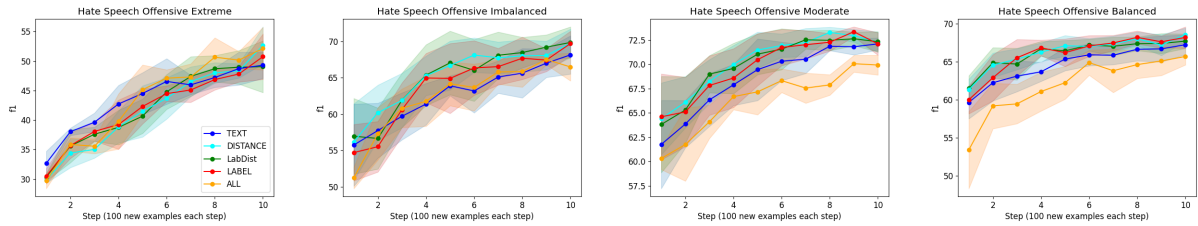


Figure 28: LAGONN_{exp} performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

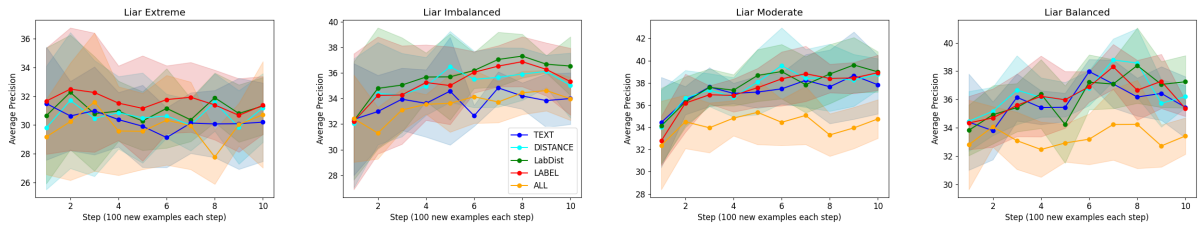


Figure 29: LAGONN_{exp} performance for all configurations and balance regimes on the LIAR dataset. The relevant balance is in the title of each panel.

1074 **A.8.2 Ablation: LAGONN k nearest** 1075 **neighbors**

1076 Here, at the suggestion of an anonymous reviewer,
1077 we present ablation results and analysis of search-
1078 ing over one to five nearest neighbors when modify-
1079 ing input via LAGONN. We present results over all
1080 LAGONN configurations under the LAGONN_{lite}
1081 fine-tuning strategy and with all balance regimes
1082 for the content moderation datasets. For the gen-
1083 eral text classification setting, we present results for
1084 both LAGONN_{lite} and LAGONN_{exp} fine-tuning
1085 under the balanced regime for all datasets with the
1086 LABDIST and TEXT configurations.

1087 If we consider all LAGONN configurations and
1088 balance regimes in the case content moderation,
1089 Figures 30 through 54, the number of neighbors
1090 does not appear to be an important hyperparameter;
1091 the learning curves for a given dataset and balance
1092 regime are very similar. While there is variation,
1093 the trend appears to be that the first NN results in
1094 the stablest, most performant, and most consistent
1095 model.

1096 However, if we only focus on LABDIST (Fig-
1097 ures 30 through 34), the default LAGONN con-
1098 figuration, we see that it can be a very important
1099 hyperparameter to consider in cases of extreme im-
1100 balance or when we have balanced data but few
1101 data points. For example, performance is boosted
1102 by up to five points for Hate Speech Offensive by
1103 the tenth step (1000 examples) with five neighbors
1104 under the extreme balance regime, yet for the bal-
1105 anced regime, the performance curves are roughly
1106 the same. For Toxic Conversations, in the balanced
1107 regime, we see that we can increase performance
1108 by up to seven points on the second step (200 ex-
1109 amples) by considering more neighbors.

1110 Turning our attention now to the general clas-
1111 sification experiments, we see that the number of
1112 neighbors for both the LABDIST and TEXT config-
1113 urations continues to consistently not really make
1114 much of a difference, with all models showing very
1115 similar performance curves for all datasets. We
1116 note however that LABDIST appears to be the most
1117 performant configuration of our method. While
1118 continued fine-tuning on datasets with a large num-
1119 ber of labels does increase performance, we ob-
1120 serve essentially the same boost for all neighbors.
1121 We also observe similar instability and performance
1122 degradation when we fine-tune on a large number
1123 of examples in cases when we have few labels.

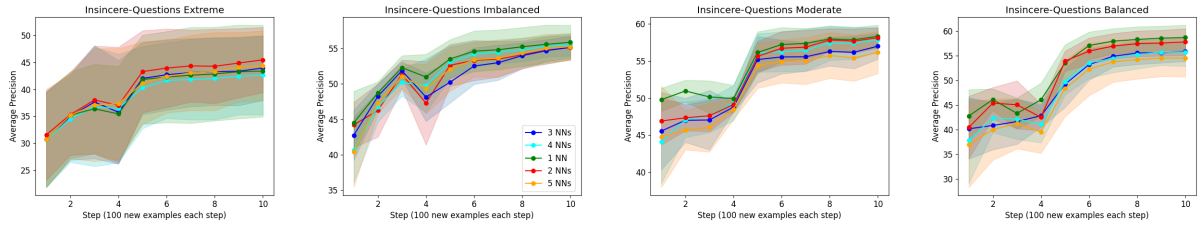


Figure 30: LABDIST results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Inscere Questions dataset. The relevant balance is in the title of each panel.

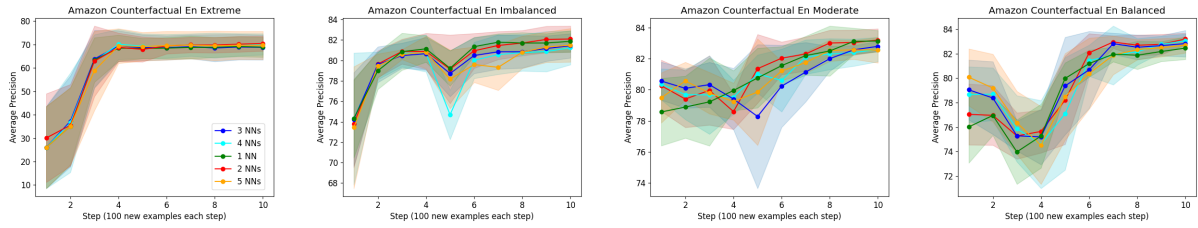


Figure 31: LABDIST results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

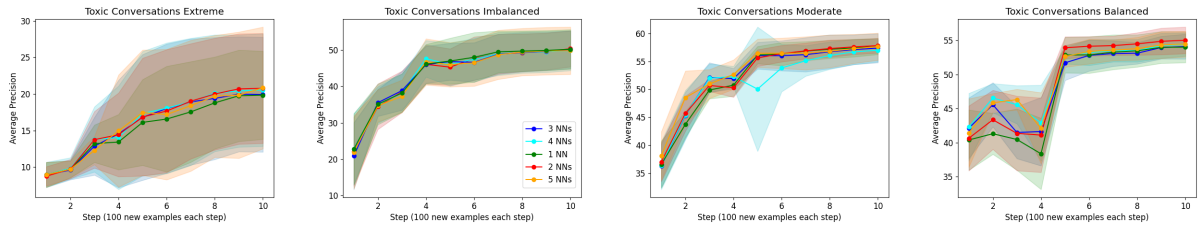


Figure 32: LABDIST results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

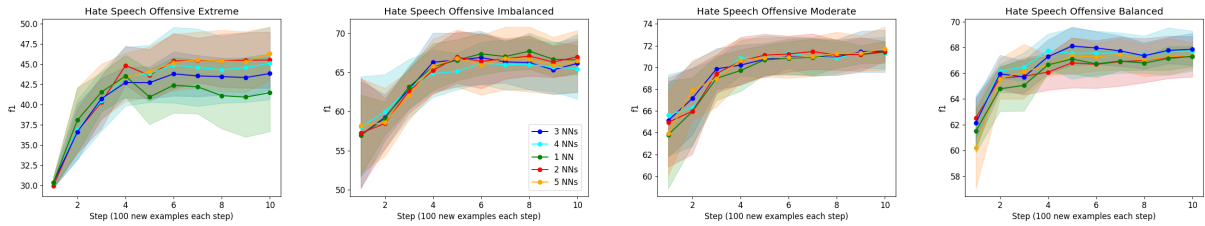


Figure 33: LABDIST results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

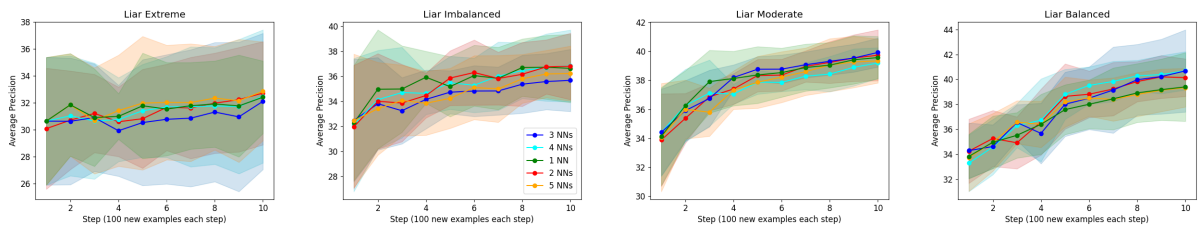


Figure 34: LABDIST results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the LIAR dataset. The relevant balance is in the title of each panel.

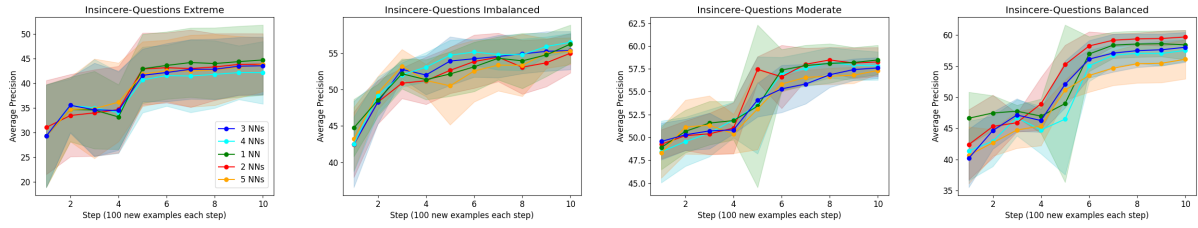


Figure 35: DISTANCE results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the the Insincere Questions dataset. The relevant balance is in the title of each panel.

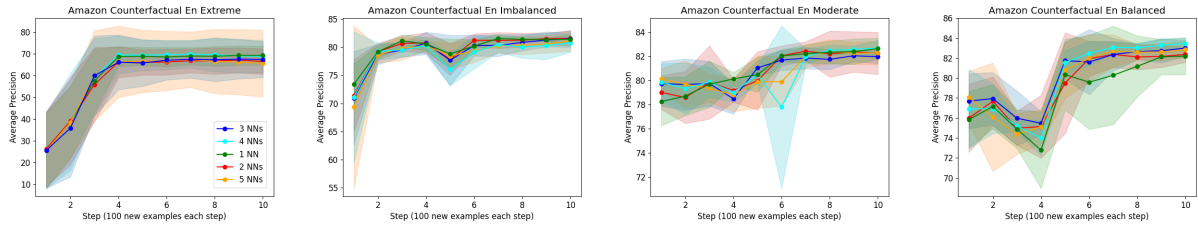


Figure 36: DISTANCE results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

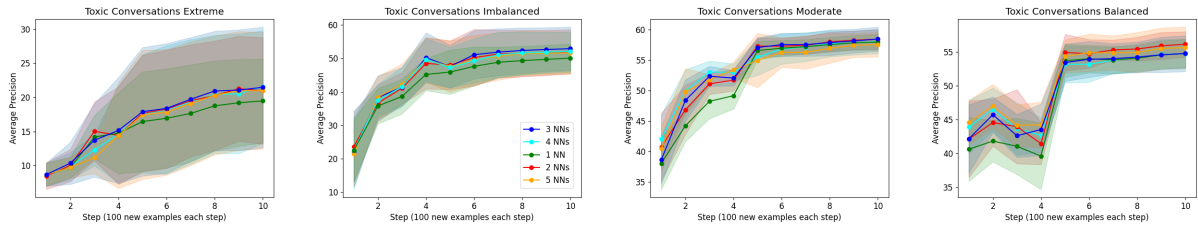


Figure 37: DISTANCE results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

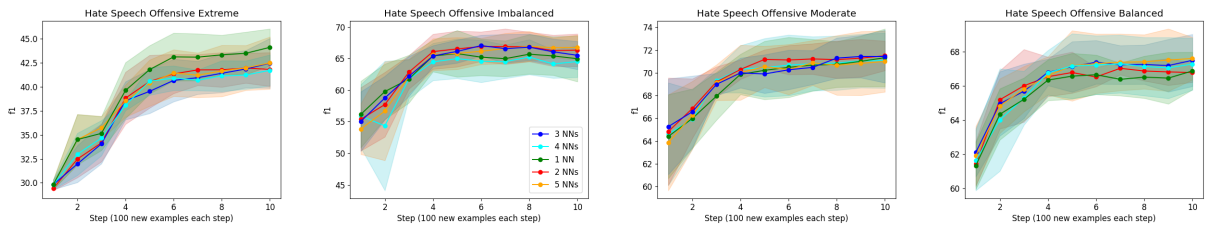


Figure 38: DISTANCE results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

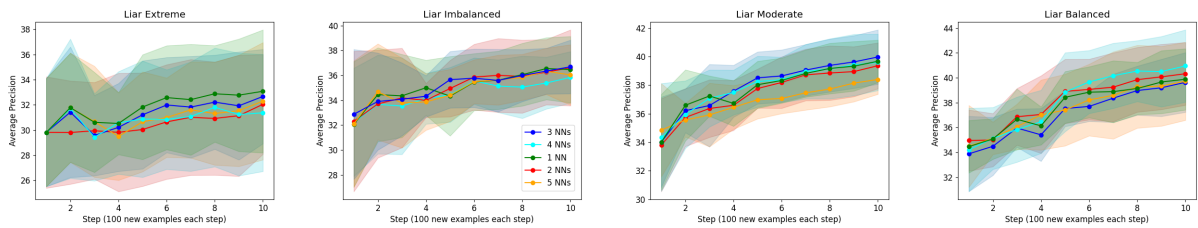


Figure 39: DISTANCE results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the LIAR dataset. The relevant balance is in the title of each panel.

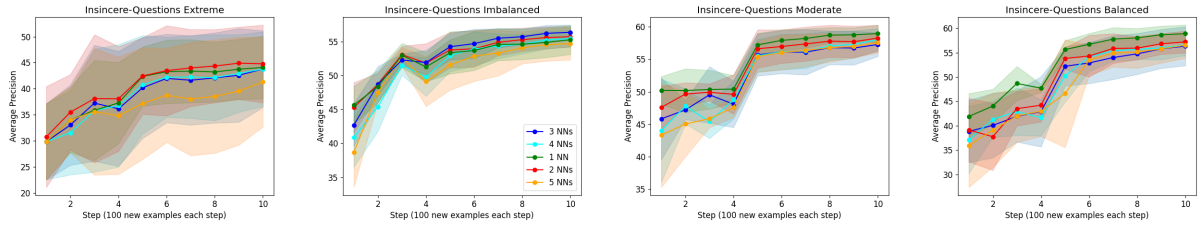


Figure 40: LABEL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Insincere Questions dataset. The relevant balance is in the title of each panel.

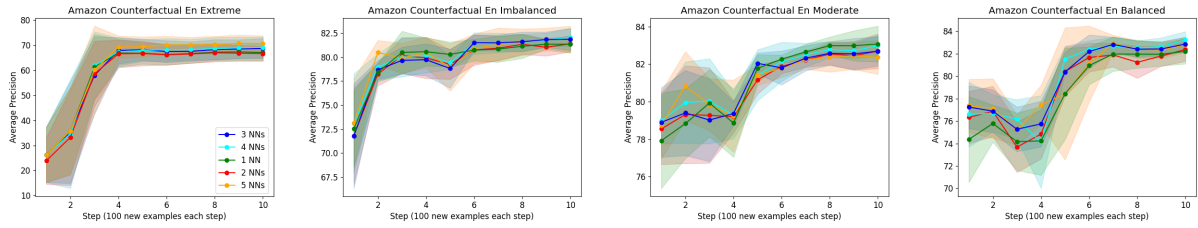


Figure 41: LABEL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

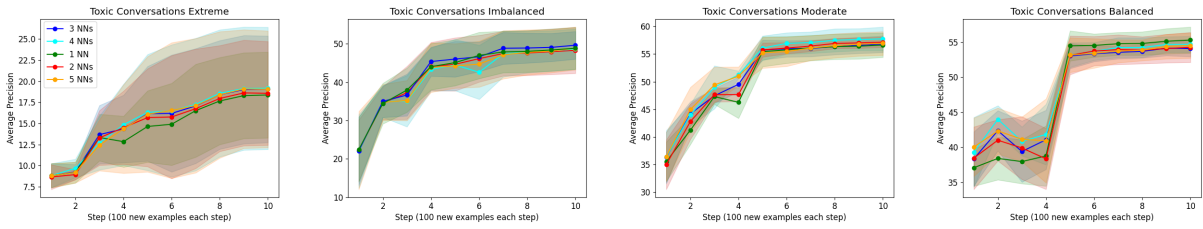


Figure 42: LABEL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

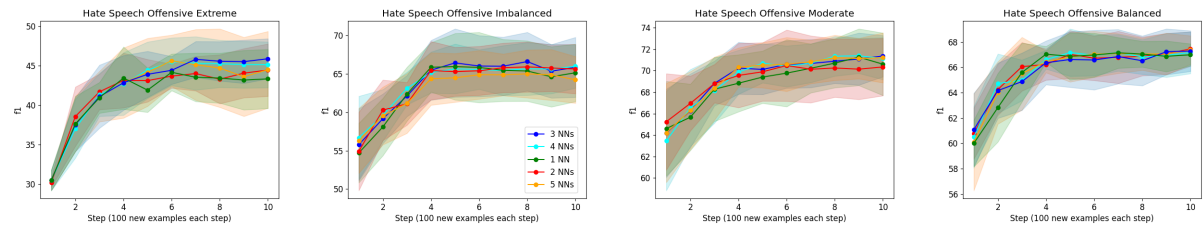


Figure 43: LABEL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

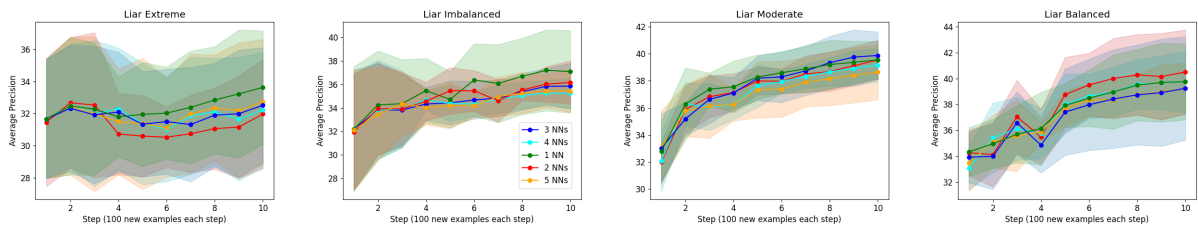


Figure 44: LABEL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the LIAR dataset. The relevant balance is in the title of each panel.

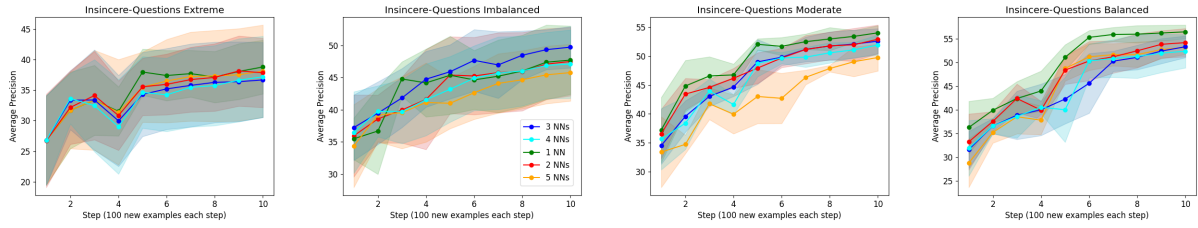


Figure 45: TEXT results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Insincere Questions dataset. The relevant balance is in the title of each panel.

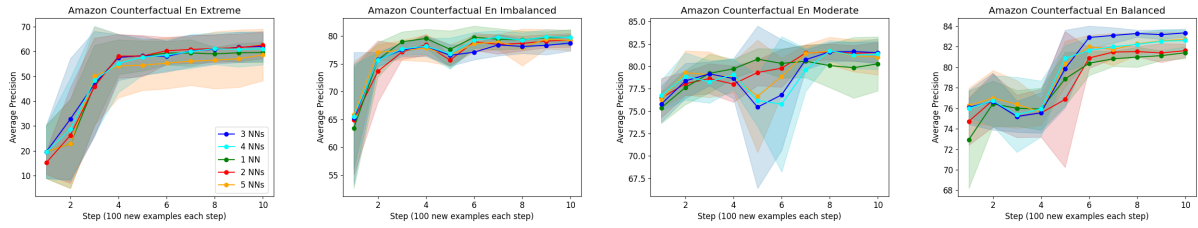


Figure 46: TEXT results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

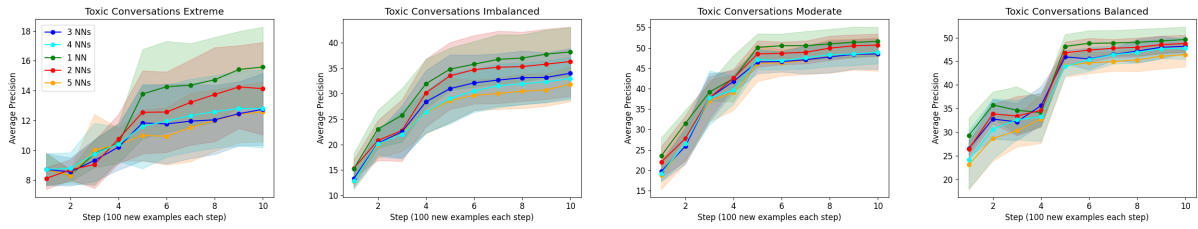


Figure 47: TEXT results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

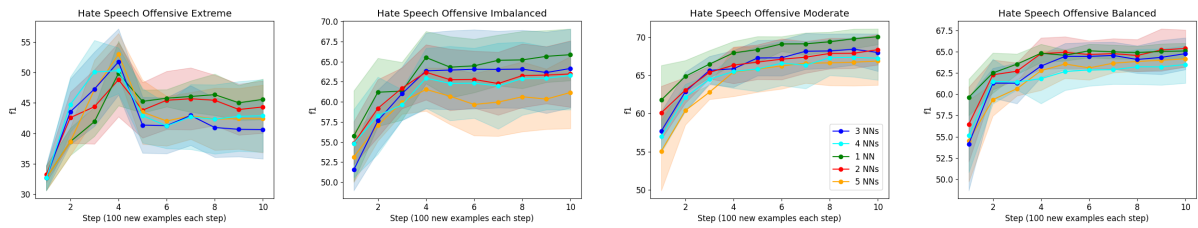


Figure 48: TEXT results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

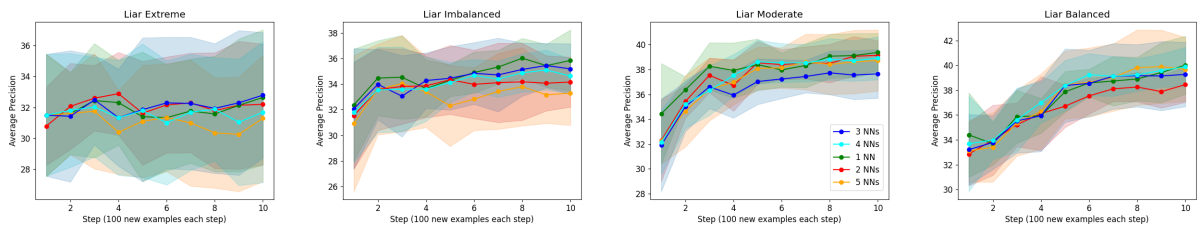


Figure 49: TEXT results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the LIAR dataset. The relevant balance is in the title of each panel.

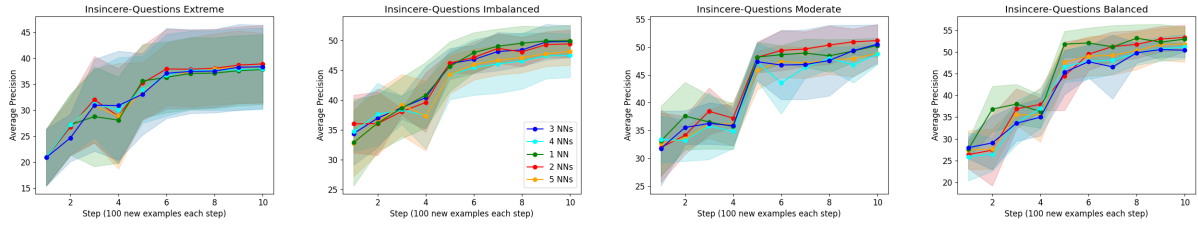


Figure 50: ALL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Insincere Questions dataset. The relevant balance is in the title of each panel.

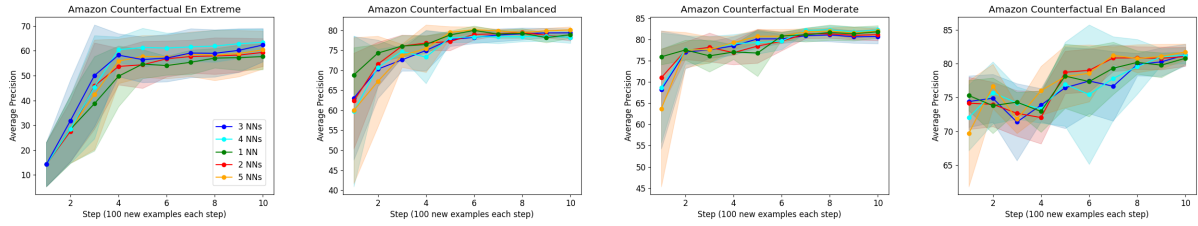


Figure 51: ALL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

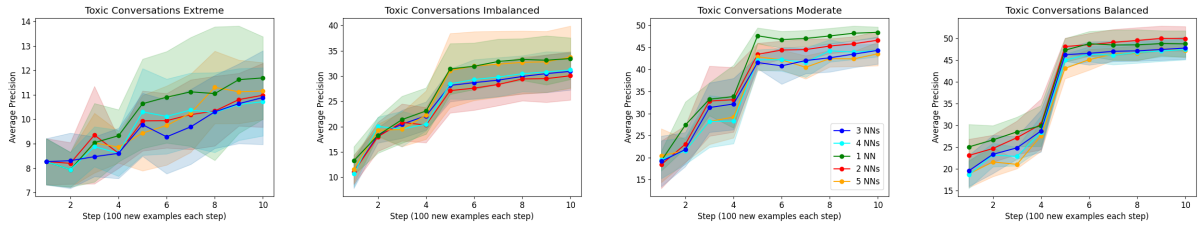


Figure 52: ALL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

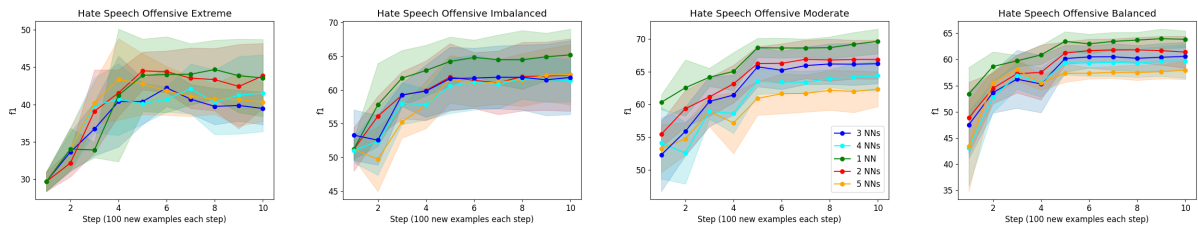


Figure 53: ALL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

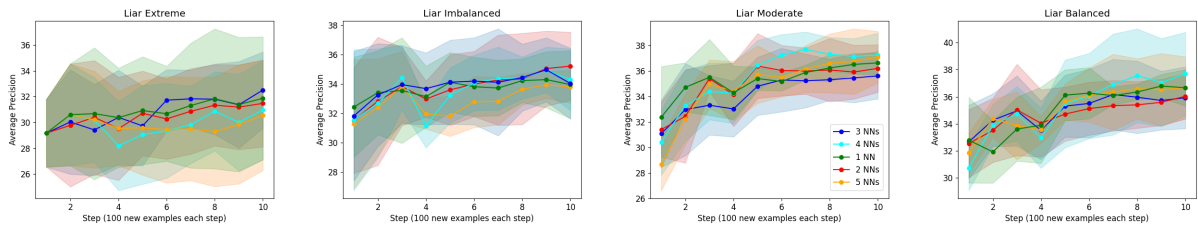


Figure 54: ALL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the LIAR dataset. The relevant balance is in the title of each panel.

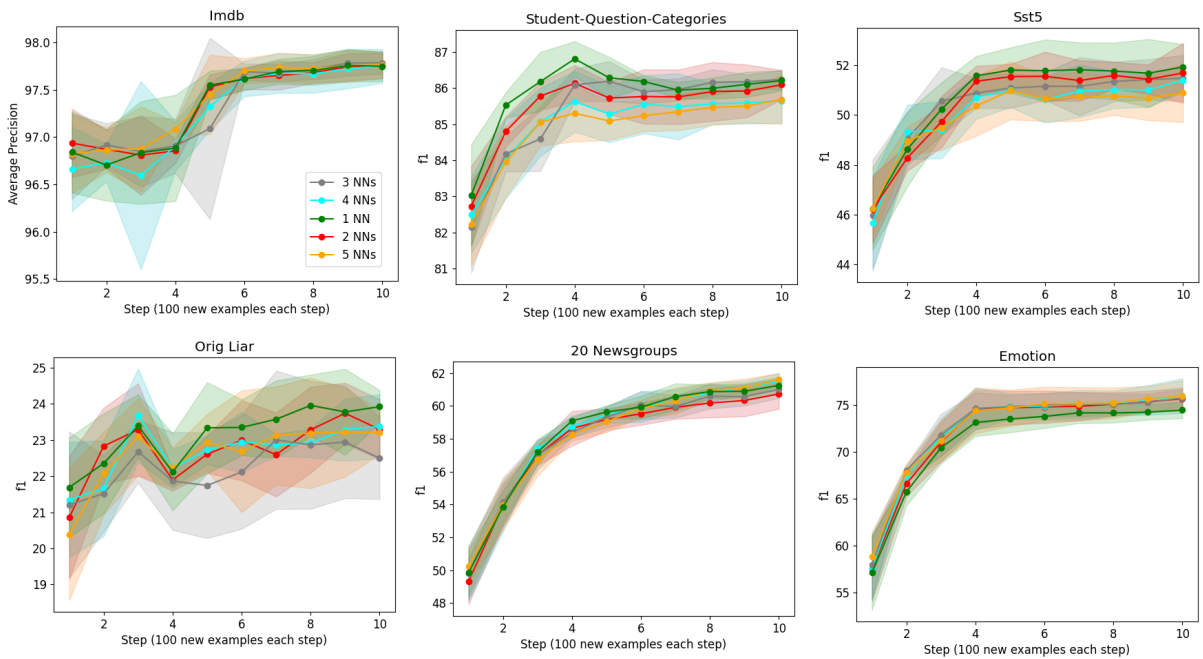


Figure 55: LABDIST results for one to five neighbors under the $LAGONN_{lite}$ fine-tuning strategy over all six general classification datasets. Results are for the balanced sampling regime and the metric is average precision for IMDB, macro-F1 elsewhere.

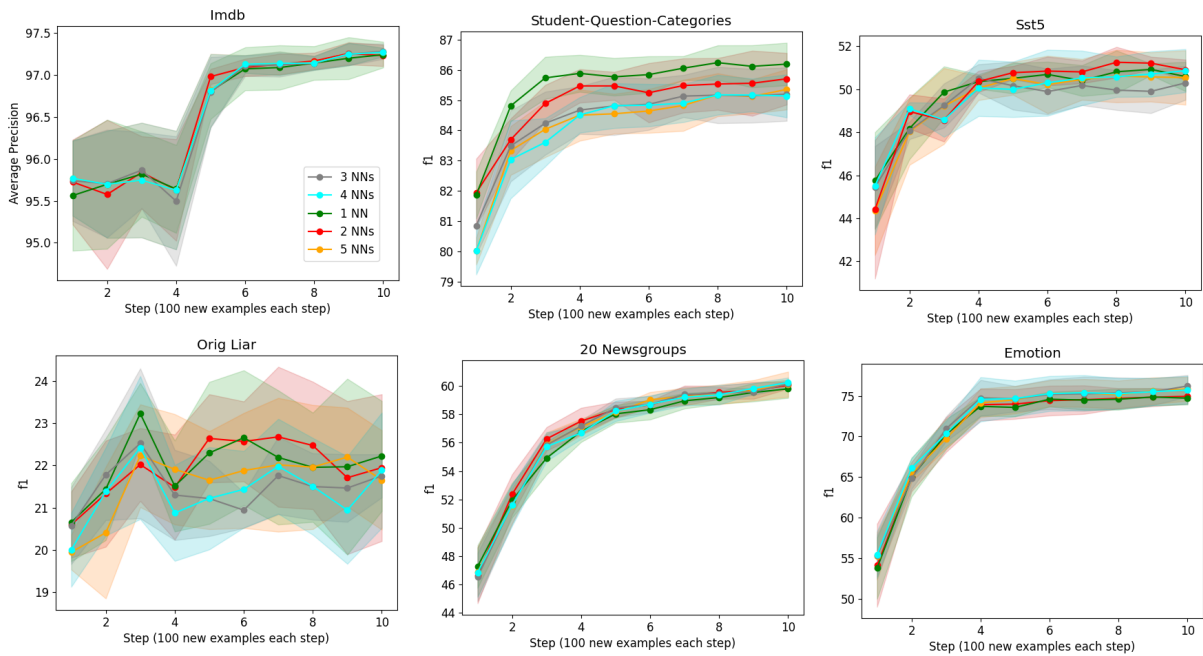


Figure 56: TEXT results for one to five neighbors under the $LAGONN_{lite}$ fine-tuning strategy over all six general classification datasets. Results are for the balanced sampling regime and the metric is average precision for IMDB, macro-F1 elsewhere.

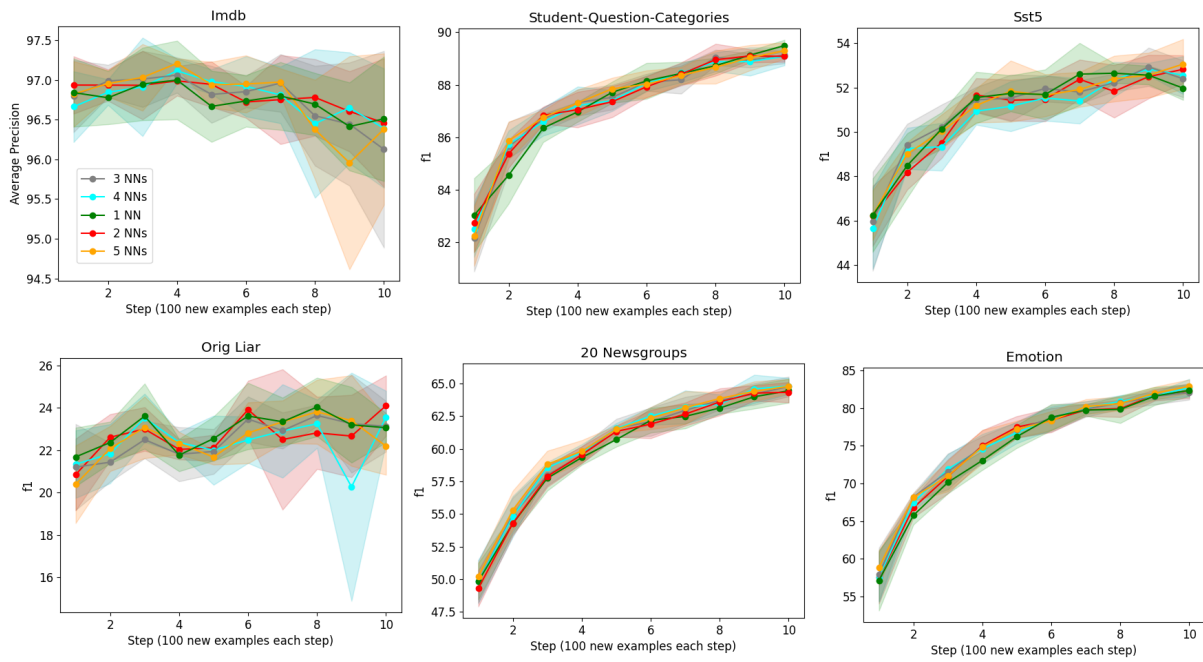


Figure 57: LABDIST results for one to five neighbors under the LAGONN_{exp} fine-tuning strategy over all six general classification datasets. Results are for the balanced sampling regime and the metric is average precision for IMDB, macro-F1 elsewhere.

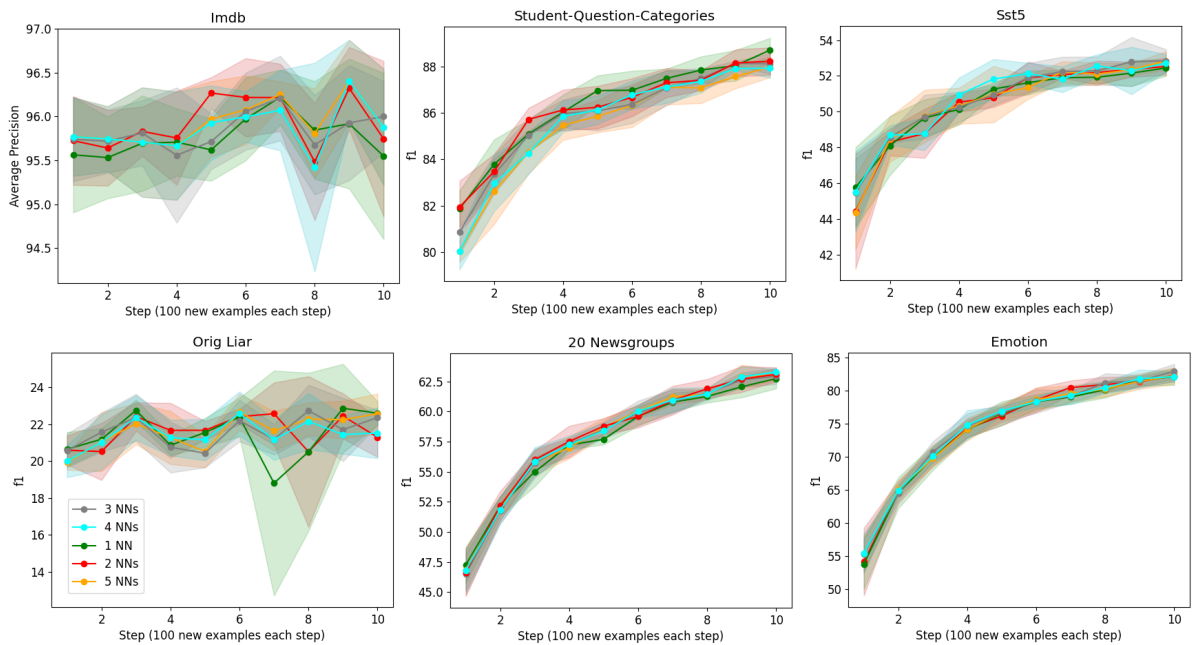


Figure 58: TEXT results for one to five neighbors under the LAGONN_{exp} fine-tuning strategy over all six general classification datasets. Results are for the balanced sampling regime and the metric is average precision for IMDB, macro-F1 elsewhere.

1124 **A.8.3 Ablation: the effect of encoding distance**

1125 Here, at the suggestion of an anonymous reviewer,
1126 we present ablation results and analysis of how en-
1127 coding distance affects LAGONN, because PLMs
1128 often struggle to understand numbers. Note that
1129 during our development stage, we ensured that our
1130 tokenizer was capable of encoding floats with trail-
1131 ing digits. To examine the effect of trailing digits
1132 on LAGONN, we consider the DISTANCE con-
1133 figuration (see Table 1), where we append only
1134 the Euclidean distance to the input text. In this
1135 ablation, however, we round to different levels of
1136 precision. For example, if the distance were a float
1137 of 0.123456789, we round it to the nearest whole
1138 number, 0.0, single digit float, 0.1, three digit float,
1139 0.123, six digit float, 0.123457, and finally keep it
1140 unrounded, that is, the original DISTANCE config-
1141 uration, 0.123456789. The below results are only
1142 for the LAGONN_{lite} training strategy. We chose
1143 LAGONN_{lite} for this ablation because it provides
1144 insight into both how distance affects full-model
1145 fine-tuning and only refitting the classification head.
1146 The results can be seen below in Figures 59 through
1147 63. We place the figures on a new page for ease of
1148 viewing.

1149 Interestingly, we tend to observe very similar per-
1150 formance curves for all rounding precisions. The
1151 exceptions to this would perhaps be Amazon Coun-
1152 terfactual and Hate Speech Offensive in the bal-
1153 anced regime where DISTANCE and rounding
1154 to the third trailing digit respectively exhibit large
1155 instability.

1156 Although not always the case, it appears that
1157 providing the model with the distance rounded to
1158 the nearest whole number tends to result in the
1159 strongest and stablest performer, however, we em-
1160 phasize that in general there does not seem to a
1161 dramatic difference between the rounding preci-
1162 sions we considered. Longer digits slightly worsen
1163 model performance and the model might learn the
1164 most from simpler or abbreviated representations
1165 of distance. This finding motivated us to consider
1166 the ablation in Appendix A.8.4.

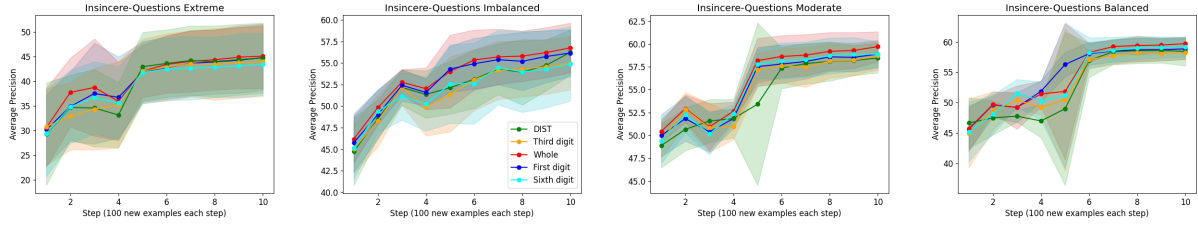


Figure 59: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the InSincere Questions dataset and the relevant balance is in the title of each panel.

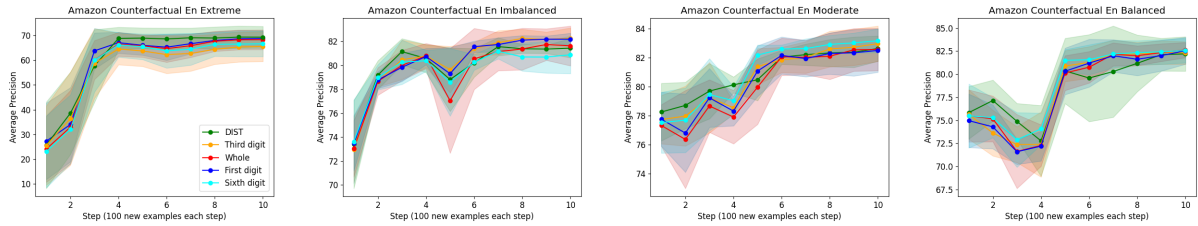


Figure 60: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the Amazon Counterfactual dataset and the relevant balance is in the title of each panel.

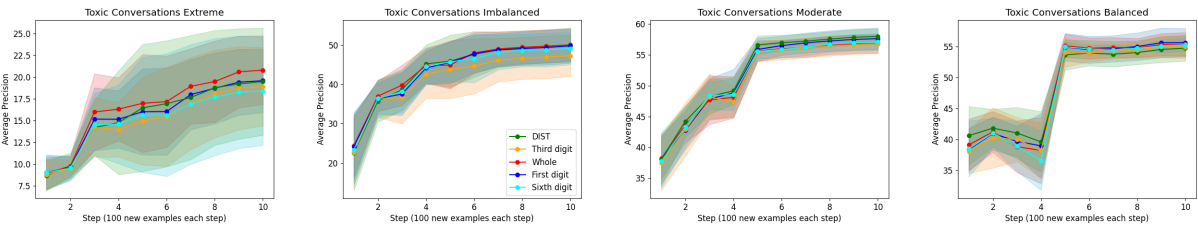


Figure 61: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the Toxic Conversations dataset and the relevant balance is in the title of each panel.

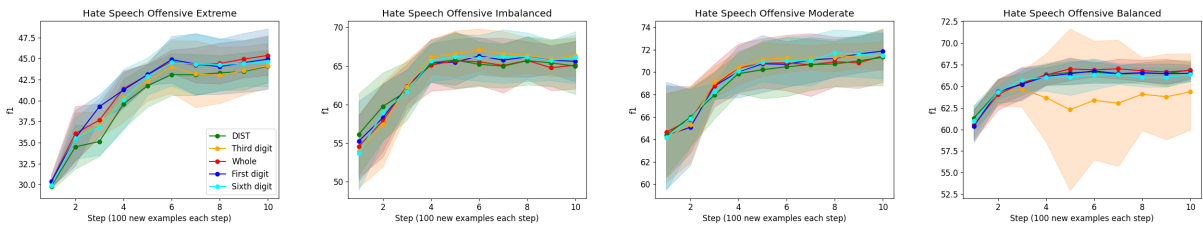


Figure 62: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the Hate Speech Offensive dataset and the relevant balance is in the title of each panel.

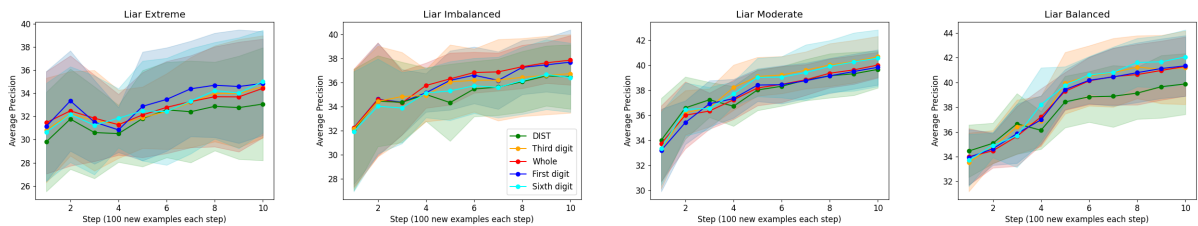


Figure 63: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the LIAR dataset and the relevant balance is in the title of each panel.

1167 **A.8.4 Ablation: support for LABDIST**

1168 The results from the ablation in Appendix A.8.3
1169 suggest that rounding the distance to the nearest
1170 whole number results in a stronger classifier than
1171 appending the unrounded distance. Thus far, we
1172 have asserted that LABDIST, where we append
1173 both the gold label of the NN and unrounded dis-
1174 tance is the most performant version of LAGONN
1175 (see Table 1). To demonstrate that this is reason-
1176 able, in this ablation study, we compare the orig-
1177 inal LABDIST configuration against three mod-
1178 els, namely the LABEL configuration, distance
1179 rounded to near whole number (Whole), and finally
1180 a new configuration similar to LABDIST, but where
1181 we append the gold label and distance rounded to a
1182 whole number, which we refer to as LABROUND.
1183 As in Appendix A.8.3, in this ablation we con-
1184 sider only the LAGONN_{lite} fine-tuning strategy.
1185 We chose LAGONN_{lite} for this ablation because
1186 it provides insight into both how the different con-
1187 figurations affect full-model fine-tuning and only
1188 re-fitting the classification head. The results can be
1189 seen below in Figures 64 through 68. We place the
1190 figures on a new page for ease of viewing.

1191 In general, we note very similar performance
1192 curves for these four models. In the case of Insinc-
1193 ere Questions, appending the distance after round-
1194 ing it to the nearest whole number (Whole, the red
1195 curve), is a strong model, except in the balanced
1196 regime where we note large instability. The results
1197 for Amazon Counterfactual tell a different story,
1198 where rounding the Euclidean distance to the near-
1199 est whole number causes large instability and even
1200 degrades performance on the fifth step.

1201 For the other evaluation scenarios, it is unclear
1202 what is the strongest method as sometimes LAB-
1203 DIST is the best performer and sometimes it is
1204 Whole (the red curve). However, we believe that in
1205 general LABDIST is the most stable model while
1206 also often being the most performant. We therefore
1207 choose it as our default LAGONN configuration as
1208 a compromise between strength and stability. It is
1209 about this configuration which we report results in
1210 the main text. Our interpretation of this is that pass-
1211 ing the model both a discrete prediction (the gold
1212 label of the NN) and a truly continuous measure
1213 of similarity (the unrounded Euclidean distance)
1214 gives it the most consistent and dependable reason-
1215 ing ability.

1216 We note, as we did in Appendix A.8.1, that we
1217 could have presented the best performer for each

1218 evaluation scenario, however, it is not the goal of
1219 our work to create even more hyperparameters that
1220 must be iterated over. However, we hope that our
1221 codebase has made it easy for one to change these
1222 configurations for their own purposes.

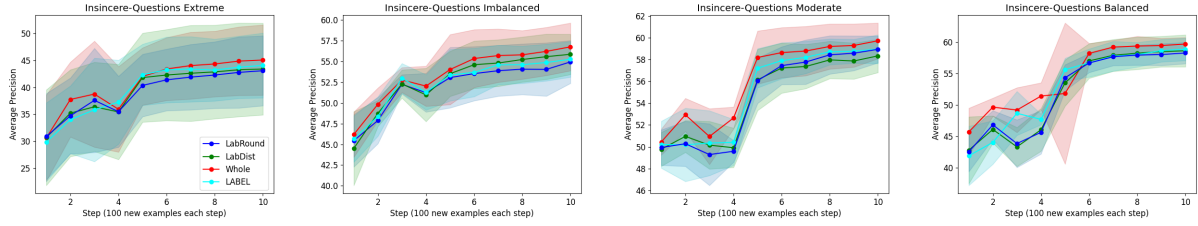


Figure 64: LAGoNN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Insincere Questions dataset and the relevant balance is in the title of each panel.

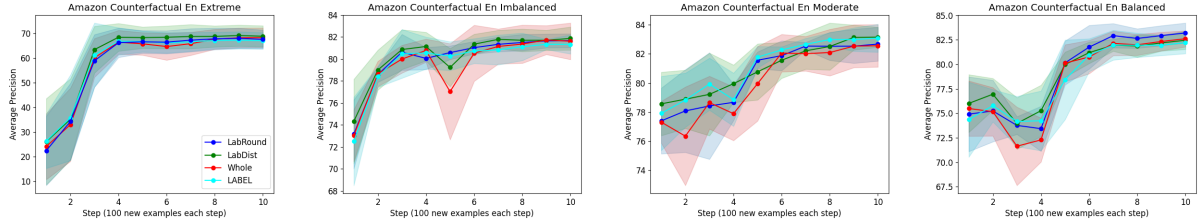


Figure 65: LAGoNN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Amazon Counterfactual dataset and the relevant balance is in the title of each panel.

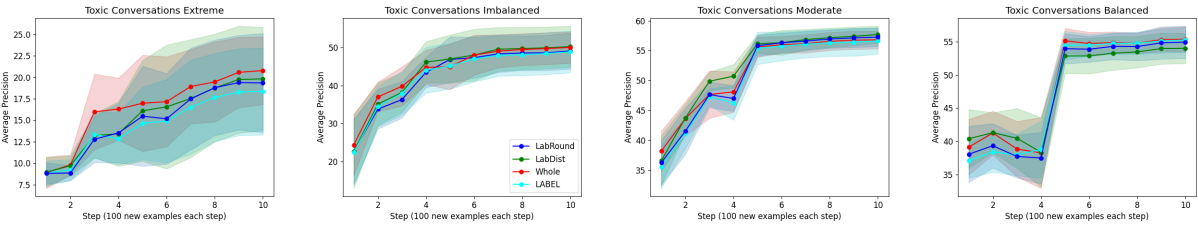


Figure 66: LAGoNN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Toxic Conversations dataset and the relevant balance is in the title of each panel.

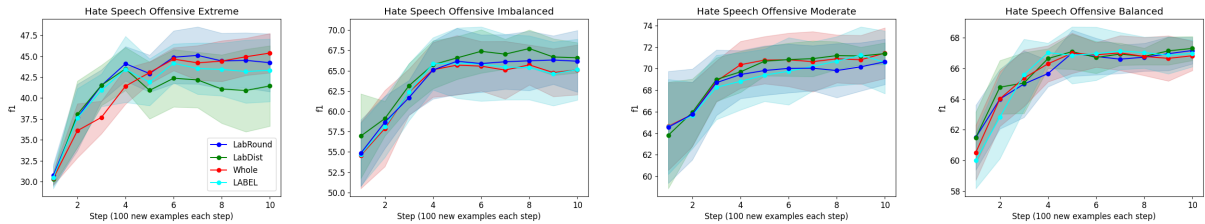


Figure 67: LAGoNN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Hate Speech Offensive dataset and the relevant balance is in the title of each panel.

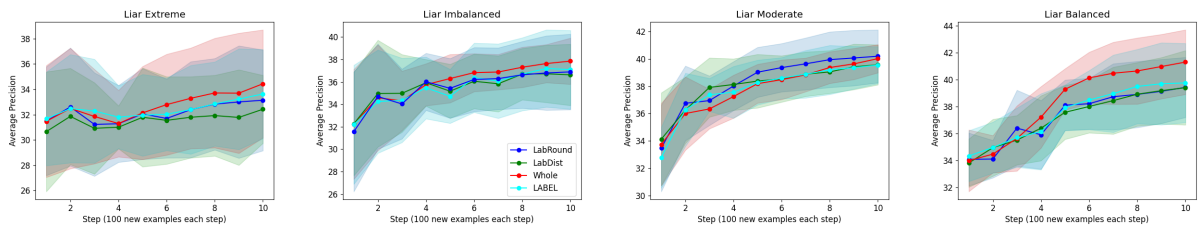


Figure 68: LAGoNN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the LIAR dataset and the relevant balance is in the title of each panel.

1223 **A.9 Examples of LAGONN modified text**

1224 **WARNING:** Some of the examples below are of

1225 an offensive nature. Please view with caution.

1226 In this section, we provide examples of how

1227 LAGONN_{exp} modifies test text from the content

1228 moderation datasets we studied under the ALL con-

1229 figuration. We choose this configuration because

1230 the information it appends from a NN in the train-

1231 ing data to a test instance encapsulates all configu-

1232 rations. LAGONN_{exp} was trained under a balanced

1233 distribution and five examples per label were cho-

1234 sen randomly on the first, fifth, and tenth step to

1235 demonstrate how the same test instance might be

1236 decorated with different training examples as the

1237 training data grow. We have made the .csv files

1238 available with our code and data files. In order to

1239 not break our .pdf generator, we were forced to

1240 remove a handful of symbols from the below text,

1241 but the original modifications remain in-tact in the

1242 .csv files included with our code files. Note that

1243 MPNET’s separator token is `</s>`, not [SEP].

Insincere Questions Step 1 1244

Test Modified What rapper still relevant and 1245
popular today has the best rhyme schemes? `</s>` 1246
`<insincere question 3.859471321105957>` What 1247
would be a good nickname for Trump, Donald 1248
Dumbck, and President Spankovich? `</s>` `<valid` 1249
`question 4.124274253845215>` What are after class 1250
12 courses in commerce stream to choose from? I 1251
have completed my class 12 (expexted 90+) and 1252
aim to do business (not aim to do job). 1253

Label valid question 1254

Test Modified Which books do you sug- 1255
gest to someone who get a free time and will 1256
help him stay motivated? `</s>` `<valid question` 1257
`3.9509353637695312>` What are the best online 1258
courses to learn data science? `</s>` `<insincere ques-` 1259
`tion 4.300448417663574>` What are the more steps 1260
in Career Oriented Education? 1261

Label valid question 1262

Test Modified How will you feel if someone 1263
talks badly about Kunti? `</s>` `<insincere ques-` 1264
`tion 3.5063605308532715>` Why are the UK gov- 1265
ernment and the media (especially the BBC and 1266
the Guardian) demonising ordinary British people, 1267
manipulating buzz words like “alt-right”, “Islama- 1268
phobia”, “racist” to suppress legitimate outrage at 1269
Muslim grooming gangs? `</s>` `<valid question` 1270
`3.6699037551879883>` How do Israelis and Pales- 1271
tinians view Nuseir Yassin? 1272

Label valid question 1273

Test Modified Why is equine HYPP inher- 1274
ited? `</s>` `<valid question 4.066534996032715>` 1275
Can you share some of the pics of hostel 1276
of Indira Gandhi medical college, Shimla, Hi- 1277
machal Pradesh? `</s>` `<insincere question` 1278
`4.231775760650635>` I am an experienced pro- 1279
grammer and in my high school my teacher tried 1280
to make me use python so I said, "No; Trust me, 1281
python is just a language for beginners, thereby 1282
making it not for me." I got sent out. Did I do 1283
anything wrong? 1284

Label valid question 1285

Test Modified How do the Valerie Stevens 1286
leather jackets achieve their quality during the 1287
manufacturing process? `</s>` `<valid question` 1288
`3.9721384048461914>` How are the Lancaster 1289
leather sofas manufactured? `</s>` `<insincere ques-` 1290
`tion 4.3559441566467285>` I am an experienced 1291
programmer and in my high school my teacher tried 1292
to make me use python so I said, "No; Trust me, 1293
python is just a language for beginners, thereby 1294

1295	making it not for me." I got sent out. Did I do	"go fuck yourself," and use the word "pussy"	1346
1296	anything wrong?	to describe women? </s> <insincere question	1347
1297	Label valid question	3.497847080230713> Why are the UK govern-	1348
1298	Test Modified Is Ariana Grande really as mean	ment and the media (especially the BBC and the	1349
1299	and bitchy as she seems? </s> <insincere question	Guardian) demonising ordinary British people, ma-	1350
1300	3.572277545928955> Why is Alia Bhatt so dumb?	nipulating buzz words like "alt-right", "Islama-	1351
1301	</s> <valid question 3.924571990966797> Do you	phobia", "racist" to suppress legitimate outrage	1352
1302	agree with Congressman Steve King's comments	at Muslim grooming gangs? </s> <valid question	1353
1303	on immigrant children in detention centers?	3.845909357070923> Do you agree with Congress-	1354
1304	Label insincere question	man Steve King's comments on immigrant children	1355
1305	Test Modified Do you guys know that aliens	in detention centers?	1356
1306	are real and all those satellites we send up in	Label insincere question	1357
1307	space work as a sort of tracking device for them	Insincere Questions Step 5	1358
1308	so in a few years it will be too late for Earth?	Test Modified What rapper still relevant and	1359
1309	</s> <insincere question 3.6094439029693604>	popular today has the best rhyme schemes? </s>	1360
1310	Have you noticed how conservatives are captur-	<insincere question 3.871907949447632> What	1361
1311	ing the English language and modifying the def-	would be a good nickname for Trump, Donald	1362
1312	initions of political words? </s> <valid question	Dumbck, and President Spankovich? </s> <valid	1363
1313	3.6901655197143555> Do you agree with Con-	question 4.028958797454834> Why does Danc-	1364
1314	gressman Steve King's comments on immigrant	ing with the Stars not include Bachata as one their	1365
1315	children in detention centers?	dance styles?	1366
1316	Label insincere question	Label valid question	1367
1317	Test Modified Is it politically incorrect to say	Test Modified Which books do you sug-	1368
1318	female privilege, but it is a more accurate term to	gest to someone who get a free time and will	1369
1319	say, white female privilege? </s> <insincere ques-	help him stay motivated? </s> <valid question	1370
1320	tion 3.323280096054077> Why are the UK gov-	3.6081225872039795> What is a good degree to	1371
1321	ernment and the media (especially the BBC and	get at community college if you want to explore dif-	1372
1322	the Guardian) demonising ordinary British people,	ferent subjects and figure out your career path? </s>	1373
1323	manipulating buzz words like "alt-right", "Islama-	<insincere question 3.8502604961395264> What	1374
1324	phobia", "racist" to suppress legitimate outrage	are the more steps in Career Oriented Education?	1375
1325	at Muslim grooming gangs? </s> <valid question	Label valid question	1376
1326	3.986680269241333> Do you agree with Congress-	Test Modified How will you feel if someone	1377
1327	man Steve King's comments on immigrant children	talks badly about Kunti? </s> <valid question	1378
1328	in detention centers?	3.5355563163757324> How do I stop feeling bad	1379
1329	Label insincere question	after a girl had a crush on me? </s> <insincere	1380
1330	Test Modified On Mother's Day, is it reasonable	question 3.689171075820923> Why Indian girls	1381
1331	to reflect there is some truth in the unfashionable	go crazy about marrying Shri. Rahul Gandhi ji?	1382
1332	notion than women are more driven by emotion	Label valid question	1383
1333	and men more driven by reason? </s> <insincere	Test Modified Why is equine HYPP inherited?	1384
1334	question 3.499204158782959> Why are the UK	</s> <insincere question 3.6035702228546143>	1385
1335	government and the media (especially the BBC and	Can female animals with male humans sex? </s>	1386
1336	the Guardian) demonising ordinary British people,	<valid question 3.7413032054901123> How long	1387
1337	manipulating buzz words like "alt-right", "Islama-	do guinea pigs live for?	1388
1338	phobia", "racist" to suppress legitimate outrage	Label valid question	1389
1339	at Muslim grooming gangs? </s> <valid question	Test Modified How do the Valerie Stevens	1390
1340	3.771740198135376> Do you agree with Congress-	leather jackets achieve their quality during the	1391
1341	man Steve King's comments on immigrant children	manufacturing process? </s> <valid question	1392
1342	in detention centers?	2.747288227081299> How are the Lancaster	1393
1343	Label insincere question	leather sofas manufactured? </s> <insincere ques-	1394
1344	Test Modified If the U.S. president is a	tion 3.944884777069092> Why don't all Trump	1395
1345	role model, is it acceptable for children to say	supporters buy only made in USA goods, e.g. many	1396

1397	of them have their cars of Asian/European compa-		
1398	nies, shop in places where more than 70 of items		
1399	are not made in USA, eat multi-national cuisine or		
1400	otherwise stop their hypocrisy?		
1401	Label valid question		
1402	Test Modified Is Ariana Grande really as		
1403	mean and bitchy as she seems? </s> <insin-		
1404	cere question 3.3252298831939697> Why is		
1405	Alia Bhatt so dumb? </s> <valid question		
1406	3.7413415908813477> How do I stop feeling bad		
1407	after a girl had a crush on me?		
1408	Label insincere question		
1409	Test Modified Do you guys know that aliens are		
1410	real and all those satellites we send up in space		
1411	work as a sort of tracking device for them so in a		
1412	few years it will be too late for Earth? </s> <insin-		
1413	cere question 3.0673365592956543> Isn't it obvi-		
1414	ous now that walking on the moon by the Amer-		
1415	icans was a hoax, because walking on the bright		
1416	side of the moon, even in a space suit would be		
1417	fatal? </s> <valid question 3.1978228092193604>		
1418	Why do we weunch satellites?		
1419	Label insincere question		
1420	Test Modified Is it politically incorrect to say		
1421	female privilege, but it is a more accurate term to		
1422	say, white female privilege? </s> <insincere ques-		
1423	tion 2.9176812171936035> How does the privi-		
1424	lege of being attractive compare to the privilege		
1425	of being White in the US? </s> <valid question		
1426	3.112481117248535> Is the media wrong for en-		
1427	forcing gender stereotypes?		
1428	Label insincere question		
1429	Test Modified On Mother's Day, is it reasonable		
1430	to reflect there is some truth in the unfashionable		
1431	notion than women are more driven by emotion		
1432	and men more driven by reason? </s> <insincere		
1433	question 3.102353811264038> Do women look		
1434	down on men who are single, even if the man is		
1435	more successful in other aspects of his life? </s>		
1436	<valid question 3.1890125274658203> Why are		
1437	some women uninterested in sex?		
1438	Label insincere question		
1439	Test Modified If the U.S. president is a		
1440	role model, is it acceptable for children to say		
1441	"go fuck yourself," and use the word "pussy"		
1442	to describe women? </s> <insincere question		
1443	3.163693904876709> Is it wrong to take your		
1444	retarded son to a hooker for his 21st birthday?		
1445	</s> <valid question 3.456286907196045> Do you		
1446	agree with Congressman Steve King's comments		
1447	on immigrant children in detention centers?		
	Label insincere question		1448
	Insincere Questions Step 10		1449
	Test Modified What rapper still relevant and		1450
	popular today has the best rhyme schemes? </s>		1451
	<valid question 3.7103171348571777> What is the		1452
	oldest fashion trends running yet? </s> <insincere		1453
	question 3.871907949447632> What would be a		1454
	good nickname for Trump, Donald Dumbck, and		1455
	President Spankovich?		1456
	Label valid question		1457
	Test Modified Which books do you sug-		1458
	gest to someone who get a free time and will		1459
	help him stay motivated? </s> <valid question		1460
	3.1401429176330566> How can I stay motivated		1461
	when learning something new? </s> <insincere		1462
	question 3.7235560417175293> I'm hungry and		1463
	I'm too lazy too get out of bed, should I get a psy-		1464
	chologist or ask you questions?		1465
	Label valid question		1466
	Test Modified How will you feel if some-		1467
	one talks badly about Kunti? </s> <insincere		1468
	question 3.4893462657928467> Does Tamil Isai		1469
	Soundarajan support Vijayendra for disrespect-		1470
	ing the Tamil Anthem? </s> <valid question		1471
	3.5355563163757324> How do I stop feeling bad		1472
	after a girl had a crush on me?		1473
	Label valid question		1474
	Test Modified Why is equine HYPP inher-		1475
	ited? </s> <valid question 3.5067965984344482>		1476
	What disadvantages do animals that don't		1477
	have bones face? </s> <insincere question		1478
	3.6035702228546143> Can female animals with		1479
	male humans sex?		1480
	Label valid question		1481
	Test Modified How do the Valerie Stevens		1482
	leather jackets achieve their quality during the		1483
	manufacturing process? </s> <valid question		1484
	2.747288227081299> How are the Lancaster		1485
	leather sofas manufactured? </s> <insincere		1486
	question 3.9087233543395996> Are Newport		1487
	cigarettes designed to selectively destroy black peo-		1488
	ple's DNA?		1489
	Label valid question		1490
	Test Modified Is Ariana Grande really as mean		1491
	and bitchy as she seems? </s> <valid question		1492
	3.183567762374878> I like this girl who used to		1493
	be quite rude and would run through boyfriends		1494
	very fast. But now that school started again,		1495
	she seems to have gotten a lot nicer through-		1496
	out Summer. Is she faking her politeness, and		1497
	is it worth pursuing her? </s> <insincere ques-		1498

1499	tion 3.3253660202026367> Why is Alia Bhatt so	because It worked."''''	1550
1500	dumb?	Label not-counterfactual	1551
1501	Label insincere question	Test Modified I like these jeans they sit	1552
1502	Test Modified Do you guys know that aliens are	low enough without being inappropriate when	1553
1503	real and all those satellites we send up in space	you sit or bend over. </s> <counterfactual	1554
1504	work as a sort of tracking device for them so in a	3.402600049972534> "But oddly enough, the bot-	1555
1505	few years it will be too late for Earth? </s> <insincere	toms are a little too loose in the waist (37) and could	1556
1506	question 3.0673365592956543> Isn't it obvi-	have used another inch or two in the inseam (I nor-	1557
1507	ously now that walking on the moon by the Ameri-	normally take a 35''' or 36''' in jeans, depending on	1558
1508	cans was a hoax, because walking on the bright	the brand if this helps).'''' </s> <not-counterfactual	1559
1509	side of the moon, even in a space suit would be	3.4201438426971436> These boxer-briefs are very	1560
1510	fatal? </s> <valid question 3.1978228092193604>	soft, very comfortable, and fit like high-end under-	1561
1511	Why do we weunch satellites?	wear the likes of which you might get at, oh, say,	1562
1512	Label insincere question	Calvin Klein for example, but for about half the	1563
1513	Test Modified Is it politically incorrect to say	price.	1564
1514	female privilege, but it is a more accurate term to	Label not-counterfactual	1565
1515	say, white female privilege? </s> <insincere ques-	Test Modified He was very professional and	1566
1516	tion 2.9176158905029297> How does the privi-	wish all transactions I make through Amazon were	1567
1517	lege of being attractive compare to the privilege	this good. </s> <counterfactual 3.4319908618927>	1568
1518	of being White in the US? </s> <valid question	I wish I had had him as an instructor at college.	1569
1519	3.112481117248535> Is the media wrong for en-	</s> <not-counterfactual 4.054030895233154>	1570
1520	forcing gender stereotypes?	I worried that it would be cheap or not fit	1571
1521	Label insincere question	or...whatever...But WOW!	1572
1522	Test Modified On Mother's Day, is it reason-	Label not-counterfactual	1573
1523	able to reflect there is some truth in the unfash-	Test Modified Well written with a twist	1574
1524	ionable notion than women are more driven by	I didn't expect. </s> <not-counterfactual	1575
1525	emotion and men more driven by reason? </s> <in-	3.3257973194122314> "The crossover from the	1576
1526	sincere question 2.9901626110076904> Do you	characters from one novel to others keeps me in-	1577
1527	agree that females think with their brains and	terested; after all, I do hate to miss a Dee-Ann	1578
1528	males with their testicles? </s> <valid question	or Eggie''' appearance.''''' </s> <counterfactual	1579
1529	3.1890125274658203> Why are some women un-	3.6820030212402344> "Had I reviewed this im-	1580
1530	interested in sex?	mediately I would have given this product five stars	1581
1531	Label insincere question	because It worked.'''''	1582
1532	Test Modified If the U.S. president is a	Label not-counterfactual	1583
1533	role model, is it acceptable for children to say	Test Modified Doesn't feel like the quality	1584
1534	"go fuck yourself," and use the word "pussy"	levi's I am used to. </s> <not-counterfactual	1585
1535	to describe women? </s> <insincere question	3.2773308753967285> However, the fabric is not	1586
1536	2.994286298751831> Why do feminists let their	that great, it's cheap scratchy cotton. </s> <counter-	1587
1537	daughters have sex with their boyfriend's at home?	factual 3.746659755706787> The blanket is nice	1588
1538	</s> <valid question 3.456286907196045> Do you	and soft but it is white, so if it doesn't light up it	1589
1539	agree with Congressman Steve King's comments	isn't much use!	1590
1540	on immigrant children in detention centers?	Label not-counterfactual	1591
1541	Label insincere question	Test Modified If we had wall studs, I believe	1592
1542	Amazon Counterfactual Step 1	the enclosed hardware would have been sufficient.	1593
1543	Test Modified Clings to the wall, doesn't flop	</s> <counterfactual 3.4338643550872803> i wish	1594
1544	around when a bag is pulled out, the mess of	the storage compartment was a little bigger and	1595
1545	bags falling out is gone. </s> <not-counterfactual	opened up instead of slidding on and off. </s> <not-	1596
1546	3.6492726802825928> Hopes that it will keep	counterfactual 3.9785308837890625> I worried	1597
1547	it's shape after washing. </s> <counterfactual	that it would be cheap or not fit or...whatever...But	1598
1548	4.012346267700195> "Had I reviewed this im-	WOW!	1599
1549	mediately I would have given this product five stars	Label counterfactual	1600

1601	Test Modified If this ever turns into a film, I	Label not-counterfactual	1652
1602	hope they do it justice! </s> <not-counterfactual	Test Modified I like these jeans they sit	1653
1603	3.5291523933410645> "The crossover from the	low enough without being inappropriate when	1654
1604	characters from one novel to others keeps me inter-	you sit or bend over. </s> <counterfactual	1655
1605	ested; after all, I do hate to miss a Dee-Ann	2.606198310852051> "But oddly enough, the bot-	1656
1606	or Eggie"" appearance."" </s> <counterfactual	toms are a little too loose in the waist (37) and could	1657
1607	3.751143217086792> "Had I reviewed this imme-	have used another inch or two in the inseam (I nor-	1658
1608	diately I would have given this product five stars	normally take a 35"" or 36"" in jeans, depending on	1659
1609	because It worked."" </s> <not-counterfactual	the brand if this helps)."" </s> <not-counterfactual	1660
1610	Label counterfactual	2.6380045413970947> A tad loose but I rather	1661
1611	Test Modified If you don't want a prominent	have it fit this way than too tight.	1662
1612	display this rack is too large for most bed or living	Label not-counterfactual	1663
1613	rooms, it is wider and taller than my tall Broy-	Test Modified He was very professional	1664
1614	hill wardrobe style dresser which was the largest	and wish all transactions I make through Ama-	1665
1615	piece in the room until this shoe rack. </s> <not-	zon were this good. </s> <not-counterfactual	1666
1616	counterfactual 3.865670680999756> "It also vali-	3.3291680812835693> This new speaker was	1667
1617	dates the incorrect"" assumption that we are alone	just what the doctor ordered and I couldn't	1668
1618	in the feelings we suppress when we sense the com-	be more pleased. </s> <counterfactual	1669
1619	plete garbage that is thrown out into society."" </s>	3.4589436054229736> Had the person han-	1670
1620	<counterfactual 4.063361167907715> The blanket	dling the shipping of this item been at all	1671
1621	is nice and soft but it is white, so if it doesn't light	concerned with the use of the product at the end of	1672
1622	up it isn't much use!	the mailing process, the slightest bit of care could	1673
1623	Label counterfactual	have been taken to ensure it's proper delivery.	1674
1624	Test Modified I wish I could have seen all of	Label not-counterfactual	1675
1625	the places he recommends! </s> <counterfactual	Test Modified Well written with a twist	1676
1626	3.5627076625823975> I wish I had had him as	I didn't expect. </s> <not-counterfactual	1677
1627	an instructor at college. </s> <not-counterfactual	2.651658535003662> The book had some interest-	1678
1628	4.141315937042236> I worried that it would be	ing twists that I did see coming and I look forward	1679
1629	cheap or not fit or...whatever...But WOW!	to reading part two of this series. </s> <counterfac-	1680
1630	Label counterfactual	tual 2.8373162746429443> Fun read Could have	1681
1631	Test Modified I wish I could replace just that	been a little longer with more detail.	1682
1632	small stupid piece, since there's nothing wrong	Label not-counterfactual	1683
1633	with the rest of the hose assembly. </s> <count-	Test Modified Doesn't feel like the qual-	1684
1634	erfactual 3.6057372093200684> i wish the stor-	ity levi's I am used to. </s> <counterfactual	1685
1635	age compartment was a little bigger and opened	2.733877182006836> It has the same great com-	1686
1636	up instead of sliding on and off. </s> <not-	fortable flattering features plus the great denim tex-	1687
1637	counterfactual 4.064871311187744> I worried that	ture that Lee has perfected- smoothing and stretchy	1688
1638	it would be cheap or not fit or...whatever...But	without the excessive cling- but I think it must have	1689
1639	WOW!	been designed for people who have a greater sur-	1690
1640	Label counterfactual	plus of belly fat than I. </s> <not-counterfactual	1691
1641	Amazon Counterfactual Step 5	2.856729745864868> Will keep but won't be that	1692
1642	Test Modified Clings to the wall, doesn't flop	casual sexy top you always want to turn to.	1693
1643	around when a bag is pulled out, the mess of	Label not-counterfactual	1694
1644	bags falling out is gone. </s> <not-counterfactual	Test Modified If we had wall studs, I believe the	1695
1645	3.161406993865967> And the dvd cases were	enclosed hardware would have been sufficient. </s>	1696
1646	tightly packed to ensure they didn't move around.	<not-counterfactual 2.6638145446777344> It was	1697
1647	</s> <counterfactual 3.308583974838257> The	a little tricky to find the center of the studs using	1698
1648	case is small, cord seems to always want to stay	my stud finder but once I felt comfortable with	1699
1649	kinked and coiled, plug should be angled and	the lines I had drawn, I drilled the pilot holes and	1700
1650	not straight...which are all items that others have	bolted this thing to the wall. </s> <counterfactual	1701
1651	pointed out.	2.879924774169922> The only thing I would have	1702

1703	like for it to have a hole in the middle so I can put	</s> <counterfactual 3.289605140686035> If I had	1754
1704	the stopper in without removing the mat.	to come up with anything negative, I would say that	1755
1705	Label counterfactual	the attachments don't seem to stay on the vacuum	1756
1706	Test Modified If this ever turns into a film, I	cleaner when not in use - but that could be me not	1757
1707	hope they do it justice! </s> <not-counterfactual	putting them on properly!	1758
1708	2.671574354171753> I read this book because of	Label not-counterfactual	1759
1709	the motion picture that is coming out soon. </s>	Test Modified I like these jeans they sit	1760
1710	<counterfactual 3.1458709239959717> Was a good	low enough without being inappropriate when	1761
1711	story, though there could have been more to it.	you sit or bend over. </s> <not-counterfactual	1762
1712	Label counterfactual	2.447404623031616> These shorts fit really	1763
1713	Test Modified If you don't want a prominent	well and look good too. </s> <counterfactual	1764
1714	display this rack is too large for most bed or	2.550638198852539> The top fits great just wish	1765
1715	living rooms, it is wider and taller than my tall	the bottoms fit too.	1766
1716	Broyhill wardrobe style dresser which was the	Label not-counterfactual	1767
1717	largest piece in the room until this shoe rack. </s>	Test Modified He was very professional	1768
1718	<counterfactual 2.7353768348693848> I bought	and wish all transactions I make through Ama-	1769
1719	this mount because I wanted one that would sit on	zon were this good. </s> <not-counterfactual	1770
1720	three studs instead of two because my TV is quite	3.3291127681732178> This new speaker was	1771
1721	heavy and I would have had a hard time centering	just what the doctor ordered and I couldn't	1772
1722	it on my wall if I didn't have the wide hanging	be more pleased. </s> <counterfactual	1773
1723	rail that this one has. </s> <not-counterfactual	3.3897111415863037> But the author alle-	1774
1724	2.873617172241211> Good for under the bed shoe	viated my concerns quickly with a few well-timed	1775
1725	storage, IF the wife wants to use it.	comments about how it was the man could have	1776
1726	Label counterfactual	known that the arrangement was something Jack	1777
1727	Test Modified I wish I could have seen all of	wanted.	1778
1728	the places he recommends! </s> <counterfactual	Label not-counterfactual	1779
1729	2.799947738647461> I wish I had had him as	Test Modified Well written with a twist	1780
1730	an instructor at college. </s> <not-counterfactual	I didn't expect. </s> <not-counterfactual	1781
1731	3.3013432025909424> And as the ole man isn't	2.557446002960205> "A bit workmanlike, not	1782
1732	any version of slender it was good that he got to try	up to Lord's high standard of A Night to Re-	1783
1733	on some shirts before hand.	member,"" but well-detailed, and a story that	1784
1734	Label counterfactual	not many now know.""" </s> <counterfactual	1785
1735	Test Modified I wish I could replace just that	2.792485475540161> Wow I am really glad I	1786
1736	small stupid piece, since there's nothing wrong	didn't read these reviews BEFORE I read this	1787
1737	with the rest of the hose assembly. </s> <counter-	book because I would have passed on the book	1788
1738	terfactual 2.628289222717285> The only thing I	and missed a really great start to a series that cap-	1789
1739	would have like for it to have a hole in the middle so	tured my attention and made me laugh all the while	1790
1740	I can put the stopper in without removing the mat.	using my imagination and painting a clear picture	1791
1741	</s> <not-counterfactual 2.9200568199157715>	of the author's world she was building for us.	1792
1742	The only downside is my laptop does not have	Label not-counterfactual	1793
1743	the screw holes on it and the screws do not retract	Test Modified Doesn't feel like the qual-	1794
1744	far enough back for me to push the connector all	ity levi's I am used to. </s> <counterfactual	1795
1745	the way in, but a simple smash will rid that issue	2.5612902641296387> i was hoping the pants	1796
1746	(this thing is durable!)	would be thicker but being that it's not too expen-	1797
1747	Label counterfactual	sive it's understandable. </s> <not-counterfactual	1798
1748	Amazon Counterfactual Step 10	2.572395086288452> But it doesn't have a lining	1799
1749	Test Modified Clings to the wall, doesn't flop	like the last couple models I bought.	1800
1750	around when a bag is pulled out, the mess of	Label not-counterfactual	1801
1751	bags falling out is gone. </s> <not-counterfactual	Test Modified If we had wall studs, I believe	1802
1752	3.161406993865967> And the dvd cases were	the enclosed hardware would have been sufficient.	1803
1753	tightly packed to ensure they didn't move around.	</s> <not-counterfactual 2.6638145446777344> It	1804

1805 was a little tricky to find the center of the studs 1856
1806 using my stud finder but once I felt comfortable 1857
1807 with the lines I had drawn, I drilled the pilot holes 1858
1808 and bolted this thing to the wall. </s> <counterfactual 2.771395206451416> Wish it had a little more 1859
1809 padding, otherwise just as advertised. 1860
1810
1811 **Label** counterfactual 1861
1812 **Test Modified** If this ever turns into a film, I 1862
1813 hope they do it justice! </s> <not-counterfactual 2.671574354171753> I read this book because of 1863
1814 the motion picture that is coming out soon. </s> 1864
1815 <counterfactual 3.141676187515259> Wish this 1865
1816 story would have been longer and turned into a 1866
1817 book, with some gut wrenching action, love/hate 1867
1818 lovers quarrels scenes, with a happy ending at the 1868
1819 end... 1869
1820
1821 **Label** counterfactual 1870
1822 **Test Modified** If you don't want a prominent 1871
1823 display this rack is too large for most bed or 1872
1824 living rooms, it is wider and taller than my tall 1873
1825 Broyhill wardrobe style dresser which was the 1874
1826 largest piece in the room until this shoe rack. </s> 1875
1827 <counterfactual 2.7353768348693848> I bought 1876
1828 this mount because I wanted one that would sit on 1877
1829 three studs instead of two because my TV is quite 1878
1830 heavy and I would have had a hard time centering 1879
1831 it on my wall if I didn't have the wide hanging 1880
1832 rail that this one has. </s> <not-counterfactual 2.873617172241211> Good for under the bed shoe 1881
1833 storage, IF the wife wants to use it. 1882
1834
1835 **Label** counterfactual 1883
1836 **Test Modified** I wish I could have seen all of 1884
1837 the places he recommends! </s> <counterfactual 2.7999041080474854> I wish I had had him as 1885
1838 an instructor at college. </s> <not-counterfactual 3.2604622840881348> I wanted to order him a few 1886
1839 affordable hats I wouldn't mind him loosing. 1887
1840
1841 **Label** counterfactual 1888
1842 **Test Modified** I wish I could replace just that 1889
1843 small stupid piece, since there's nothing wrong 1890
1844 with the rest of the hose assembly. </s> <counterfactual 2.474032402038574> I wish I could just 1891
1845 hook up a hose rather than connecting and routing 1892
1846 an tube internally </s> <not-counterfactual 2.9200568199157715> The only downside is my 1893
1847 laptop does not have the screw holes on it and the 1894
1848 screws do not retract far enough back for me to 1895
1849 push the connector all the way in, but a simple 1896
1850 smash will rid that issue (this thing is durable!) 1897
1851
1852 **Label** counterfactual 1898
1853 **Toxic Conversations Step 1** 1899
1854
1855

Test Modified Homos demand that you accept 1856
their fudge packing. But none of us ever will- </s> 1857
<not toxic 3.5619421005249023> Sounds just awful, 1858
working for the state that is. So it begs the ques- 1859
tion, why work for the state if the pay is so bad ver- 1860
sus the private sector? Seems logical to just make 1861
the switch? </s> <toxic 4.056504249572754> I 1862
would imagine it is a miserable existence. Being a 1863
whining sore loser is no way to go through life. 1864
Label not toxic 1865
Test Modified I don't think anyone likes this 1866
health care bill, it stinks for everyone. 50 years 1867
and older are going to get hammered with higher 1868
premiums. People with preexisting conditions will 1869
also see their premiums go through the roof. Even- 1870
tually no one will be able to afford it. They're still 1871
not addressing the reason and the solution. Pre- 1872
scription drugs. Why are Americans paying 2 to 1873
3 times higher for prescription drugs than Cana- 1874
dians are? They should sit down with Democrats 1875
and fix it not play musical chairs with our health 1876
care. Here is the add that AARP put out, and it 1877
has a fact check at the end. It's worth watching. 1878
[http://www.thedenverchannel.com/news/politics](http://www.thedenverchannel.com/news/politics/aarp-ad-says-house-gop-health-care-bill-would-boost-insurance-rates-for-older-americans) 1879
[/aarp-ad-says-house-gop-health-care-bill-](http://www.thedenverchannel.com/news/politics/aarp-ad-says-house-gop-health-care-bill-would-boost-insurance-rates-for-older-americans) 1880
[would-boost-insurance-rates-for-older-americans](http://www.thedenverchannel.com/news/politics/aarp-ad-says-house-gop-health-care-bill-would-boost-insurance-rates-for-older-americans) 1881
</s> <toxic 3.6495330333709717> Calling 1882
everybody that disagrees with you a racist and 1883
anti-Muslim is pretty uncivil. Also intellectually 1884
bankrupt... </s> <not toxic 3.8587634563446045> 1885
Sounds just awful, working for the state that is. 1886
So it begs the question, why work for the state if 1887
the pay is so bad versus the private sector? Seems 1888
logical to just make the switch? 1889
Label not toxic 1890
Test Modified Rather than call you a liar, I'll 1891
just post a link to HART's helicopter footage of 1892
this "most densely populated" ag land. 1893
[https://www.youtube.com/watch?](https://www.youtube.com/watch?v=opJzRWWVDdw) 1894
[v=opJzRWWVDdw](https://www.youtube.com/watch?v=opJzRWWVDdw) 1895
And if you've been following a little closer, 1896
you'd know about the luxury towers and boutique 1897
hotels that developers are planning on fulfilling 1898
their TOD requirement with "cash gifts" to the city 1899
instead of actual affordable TOD. </s> <not toxic 3.8009092807769775> Sounds just awful, working 1900
for the state that is. So it begs the question, why 1901
work for the state if the pay is so bad versus the pri- 1902
vate sector? Seems logical to just make the switch? 1903
</s> <toxic 4.463983535766602> I would imagine 1904
it is a miserable existence. Being a whining sore 1905
loser is no way to go through life. 1906

1907 loser is no way to go through life. 1958
1908 **Label** not toxic 1959
1909 **Test Modified** Get the "real news" from Breit- 1960
1910 bart and Drudge Report? Ahahahahaha! Are you 1961
1911 for real? </s> <not toxic 3.559427499771118> 1962
1912 Sounds just awful, working for the state that is. 1963
1913 So it begs the question, why work for the state 1964
1914 if the pay is so bad versus the private sector? 1965
1915 Seems logical to just make the switch? </s> <toxic 1966
1916 4.239322662353516> I would imagine it is a mis- 1967
1917 erable existence. Being a whining sore loser is no 1968
1918 way to go through life. 1969
1919 **Label** not toxic 1970
1920 **Test Modified** Wud luv to see one of those Amer- 1971
1921 ican Boats with all the planes on it at the bottom 1972
1922 of the sea. </s> <not toxic 3.8268911838531494> 1973
1923 Sounds just awful, working for the state that is. 1974
1924 So it begs the question, why work for the state 1975
1925 if the pay is so bad versus the private sector? 1976
1926 Seems logical to just make the switch? </s> <toxic 1977
1927 4.241443157196045> I would imagine it is a mis- 1978
1928 erable existence. Being a whining sore loser is no 1979
1929 way to go through life. 1980
1930 **Label** not toxic 1981
1931 **Test Modified** Wtf, nowhere in scripture does 1982
1932 it say that Elizabeth was a kid when married to 1983
1933 Zachariah and nowhere does it say that Mary was 1984
1934 a child when married to Joseph. That's a complete 1985
1935 and utter fabrication. In fact, it says in Luke 1:42, 1986
1936 "And she spake out with a loud voice, and said, 1987
1937 Blessed art thou among women, and blessed is the 1988
1938 fruit of thy womb." It says blessed art thou among 1989
1939 WOMEN. </s> <not toxic 3.6674933433532715> 1990
1940 Sounds just awful, working for the state that is. 1991
1941 So it begs the question, why work for the state 1992
1942 if the pay is so bad versus the private sector? 1993
1943 Seems logical to just make the switch? </s> <toxic 1994
1944 4.0054121017456055> I would imagine it is a mis- 1995
1945 erable existence. Being a whining sore loser is no 1996
1946 way to go through life. 1997
1947 **Label** toxic 1998
1948 **Test Modified** Angela Merkel and all other Euro- 1999
1949 pean political leaders who have aided and abetted 2000
1950 the ongoing invasion of Europe by the forces of 2001
1951 the crescent moon death cult should be tried as 2002
1952 accessories to Mr. Urban's murder. </s> <toxic 2003
1953 3.262410879135132> Calling everybody that dis- 2004
1954 agrees with you a racist and anti-Muslim is pretty 2005
1955 uncivil. Also intellectually bankrupt... </s> <not 2006
1956 toxic 3.8916428089141846> It's always important 2007
1957 to remember what can happen when you have so- 2008

ciopaths as leaders and also have compliant fol-
lowers. Some of the younger posters on this site
might want to Google "Jim Jones and Jonestow"...
There were no "checks and balances" in Jonestown;
I fear there are none in North Korea....and I can
only hope those in our country are firmly in place
and functioning. Gary Crum

Label toxic
Test Modified I 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. </s> <toxic
3.0406124591827393> Calling everybody that dis-
agrees with you a racist and anti-Muslim is pretty
uncivil. Also intellectually bankrupt... </s> <not
toxic 4.094666481018066> Christ never said he
would give grace, mercy, and acceptance to those
who determinedly violate Scripture. In fact, he
often spoke of hell.

Label toxic
Test Modified no one cares what a paid liberal
trolling hack like you believes lunatic., </s> <toxic
2.8411786556243896> Calling everybody that dis-
agrees with you a racist and anti-Muslim is pretty
uncivil. Also intellectually bankrupt... </s> <not
toxic 4.034884929656982> Christ never said he
would give grace, mercy, and acceptance to those
who determinedly violate Scripture. In fact, he
often spoke of hell.

Label toxic
Test Modified Ok all you NDP "LEAP" mani-
festo types, where is your hero Naomi Klein? Her
fawning adoration of Chavez and Venezuelan thug-
gery knows no bounds. I'm sure she's awfully
hysterical over the thought that such a pathetic
dictstorship could ever be sanctioned. </s> <toxic
3.3616135120391846> Calling everybody that dis-
agrees with you a racist and anti-Muslim is pretty
uncivil. Also intellectually bankrupt... </s> <not
toxic 3.903903007507324> I have very high re-
spect for teachers that get the job done. Teaching
is an extremely difficult and important job. And
it is quite apparent that we are in desperate need
of teachers that can actually do the job. Pride of
workmanship would have teachers wanting to have
their students periodically evaluated and tested to
show how well they have done their job. We have
some very competent teachers that get the job done
and welcome student testing (in spite of sorry ad-

min and unfair union). But the majority of teachers here instead of doing their jobs they band together wear purple shirts and mob the government for a better contract, and no accountability in the form of testing students. Many teachers are middle class kids that took the path of least resistance in what was expected by their parents (college) and because they lacked drive ended up teachers. That lack of drive shows by what the private sector taxpayers get for their money. Your degrees mean nothing if you don't do your job.

Label toxic

Toxic Conversations Step 5

Test Modified Homos demand that you accept their fudge packing. But none of us ever will- </s><toxic 3.1383495330810547> So you admit you would exterminate inferior humans. </s> <not toxic 3.2954952716827393> Mark MacKinnon and the interests he work for would like us to 'get used to it', because they don't want to do anything practical to stop it.

Label not toxic

Test Modified I don't think anyone likes this health care bill, it stinks for everyone. 50 years and older are going to get hammered with higher premiums. People with preexisting conditions will also see their premiums go through the roof. Eventually no one will be able to afford it. They're still not addressing the reason and the solution. Prescription drugs. Why are Americans paying 2 to 3 times higher for prescription drugs than Canadians are? They should sit down with Democrats and fix it not play musical chairs with our health care. Here is the add that AARP put out, and it has a fact check at the end. It's worth watching. <http://www.thedenverchannel.com/news/politics/aarp-ad-says-house-gop-health-care-bill-would-boost-insurance-rates-for-older-americans> </s><not toxic 2.5086519718170166> so in the mean time tens of thousands of Oregonians go without health insurance which will now be unaffordable to them. And sorry, the republicans have had 8 years to figure out a better system, they aren't going to do it anytime soon. Neither party has any desire to actually find a solution to all this. Hopefully Trump will also soon eliminate the tax penalty for not having insurance so all us folks who buy our own insurance don't get penalized for not being able to afford to buy the insurance we are required to have. Of course he is probably clueless that detail even exists </s> <toxic 2.801957607269287> reducing

number of brackets. Another is lowering corporate tax rates which would be OK if all loopholes, including tax-free political donations for wealthy people only. Another is doubling the earned income tax credit which would help families with children but for people like me, would not make up for loss of the state tax deduction. Essentially the proposed tax "reform bill as it stands is a huge wealth transfer, from working people and the treasury to the offshore accounts of wealthy people and corporations. Median earners like me will be screwed. Poor people with children will be screwed. The rich will get a huge bundle of cash, and the Republicans will lead us further into deep debt.

That's if Ryan and McConnell can pass this huge sack of excrement.

Label not toxic

Test Modified Rather than call you a liar, I'll just post a link to HART's helicopter footage of this "most densely populated" ag land. <https://www.youtube.com/watch?v=opJzR-WWVDdw>

And if you've been following a little closer, you'd know about the luxury towers and boutique hotels that developers are planning on fulfilling their TOD requirement with "cash gifts" to the city instead of actual affordable TOD. </s> <not toxic 3.0351996421813965> So are you saying we should build the road. Lots of doom and gloom but I'm missing your point. </s> <toxic 3.1464040279388428> Isn't that terrible? Caldwell and Inouye do not know the amount a 10-year extension would raise yet they are proposing it. Caldwell asked for two years and he got it. Two years later he is now asking for ten years. Now, which is only weeks later, HART director is claiming 10 years is insufficient. This scenario shows our mayor and HART do not know the scope of this project. Pathetic!

Label not toxic

Test Modified Get the "real news" from Breitbart and Drudge Report? Ahahahahaha! Are you for real? </s> <toxic 2.634126901626587> "If one read the Dispatch one would think Trump is the most evil person on the planet." Not evil, just idiotic. And it would be easy to give his behaviour a pass if he wasn't POTUS. ".....five to eight anti-Trump stories per day. Never any good one's or one's that just stuck to the facts." Well, when there's good a news Trump story to print, I'm sure ADN

2111 will be all over it, problem is, there's been a dearth 2162
2112 of those since his election. Facts? Ok Rich, give 2163
2113 us a list of incorrect facts in the above story. </s> 2164
2114 <not toxic 3.0248677730560303> With dismay I 2165
2115 noticed that ADN had printed yet another column 2166
2116 from Fox commentator Charles Krauthammer but 2167
2117 after reading it I'm glad the editors chose it as the 2168
2118 feature article on the opinion page. Krauthammer 2169
2119 is also a psychiatrist so his analyses of Trump as 2170
2120 a man who has never emotionally, intellectually 2171
2121 developed beyond adolescence holds some weight. 2172
2122 But what does it say about Trumps supporters that 2173
2123 so many millions of them can't see through the 2174
2124 boorish, confrontational attitude of the man? How 2175
2125 can so many Americans have devolved into such 2176
2126 anger, fear and irrationality that they would/could 2177
2127 find redemption in Trump after how he has exposed 2178
2128 his true narcissistic self for all to see. When you've 2179
2129 lost the Jennifer Rubins and Charles Krautham- 2180
2130 mer's of the media world you've lost the battle 2181
2131 yet the Trumpian cult members will soldier on and 2182
2132 then become even angrier and more full of fear af- 2183
2133 ter the election. Something to do with their choice 2184
2134 of "information" sources no doubt. 2185
2135 **Label** not toxic 2186
2136 **Test Modified** Wud luv to see one of those Amer- 2187
2137 ican Boats with all the planes on it at the bottom 2188
2138 of the sea. </s> <toxic 3.3901054859161377> I 2189
2139 bet Regent Seven Seas will never offer Mr Ham- 2190
2140 mond another trip. Wow, what a snarky article. He 2191
2141 makes , I assume, some valid points about food 2192
2142 and atmosphere. However, after discovering the 2193
2143 treats available on his "massive deck" he "blew 2194
2144 off" his remaining restaurant reservations , donned 2195
2145 his comfy bathrobe and ordered-in. He was cer- 2196
2146 tainly not an ideal passenger and, for one floating 2197
2147 on a freebie, he's a total ingrate! </s> <not toxic 2198
2148 3.409156084060669> Now replaced by the sexy 2199
2149 EA-18G Growler! Using a preexisting Military 2200
2150 Operating Area! Get over it!!!!!! 2201
2151 **Label** not toxic 2202
2152 **Test Modified** Wtf, nowhere in scripture does 2203
2153 it say that Elizabeth was a kid when married to 2204
2154 Zachariah and nowhere does it say that Mary 2205
2155 was a child when married to Joseph. That's a 2206
2156 complete and utter fabrication. In fact, it says 2207
2157 in Luke 1:42, "And she spake out with a loud 2208
2158 voice, and said, Blessed art thou among women, 2209
2159 and blessed is the fruit of thy womb." It says 2210
2160 blessed art thou among WOMEN. </s> <not toxic 2211
2161 2.769857406616211> I was informed that my first 2212

grandchild had been conceived the evening of the day when I had inserted a prayer note in the Wailing Wall in Jerusalem that asked God to help my daughter conceive after a year of frustrated attempts. Maybe Elizabeth did the same thing? After all, she was in the same neighborhood. :-) </s> <toxic 3.286393880844116> Christians who support Trump are the most mind-boggling to me. I just don't see how they square the circle between Trump and their moral foundations.

"Beware then of useless grumbling, and keep your tongue from slander; because no secret word is without result, and a lying mouth destroys the soul." (Wisdom 1:11)

If that is the case, then Trump's soul was utterly destroyed decades ago.

Label toxic

Test Modified Angela Merkel and all other European political leaders who 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. </s> <toxic 3.2037758827209473> that's what happens when you betray the people of your country for foreign bs. let's go Le Pen, Geert Wilders. If the media refuses to mention the muslim crisis the total incompatibility of primitive, uneducated muslim males swarming countries and turning them into misogynistic fundamentalist religious areas then we need these people to save us from YOU! </s> <not toxic 3.2326502799987793> your first mistake is believing what a politician says because generally it has nothing to do with what they do.

The 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 Butts will ram through Ranked Ballot

Label toxic

Test Modified I 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. </s> <toxic 2.8730525970458984> Ms. Van Brocklin: You state that there is Payne's conduct is inexcusable, then proceed to use the rest of your space to justify him and his actions. You have denigrated the proud courage of countless people who took up causes

2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263

via civil disobedience. I marched in the non-violent peaceful Civil Rights protest. So the dogs and the firehoses used by a certain southern sheriff were justified, by your logic. So were the citizens beaten by Chicago police during the Democratic Convention Police riots. Resolved in 20 minutes? Nonsense. If Miss Wubbels hadn't protested as she did, she'd likely have ended up in a back room, somewhere, and who know when she would have been granted the presence of a lawyer and what she would have gone through prior to that. You are completely ignorant of the shortage of nurses in this country - in some cases, critical shortages. And why would anyone want to be a nurse when they are disrespected by a former state and federal prosecutor such as you. </s> <not toxic 2.92986798286438> Acquit her, then commit her. This womens cheese has slid so far off the cracker she's a danger to herself and to others.

Animal rights activism is a just cause, but her and her group have gone off the deep end into radical extremist territory.

Label toxic

Test Modified no one cares what a paid liberal trolling hack like you believes lunatic,, </s> <toxic 2.6076736450195312> aa another hate filled left winger again! save the stupid nonsense sheep, trump is not causing anything, our weak leadership is. </s> <not toxic 2.810284376144409> Ouch... didn't see that one coming. A liberal stealing my own line... just like they take everything else they like.

Label toxic

Test Modified Ok all you NDP "LEAP" manifesto types, where is your hero Naomi Klein? Her fawning adoration of Chavez and Venezuelan thugery knows no bounds. I'm sure she's awfully hysterical over the thought that such a pathetic dictstorship could ever be sanctioned. </s> <toxic 2.941824436187744> Ms. Van Brocklin: You state that there is Payne's conduct is inexcusable, then proceed to use the rest of your space to justify him and his actions. You have denigrated the proud courage of countless people who took up causes via civil disobedience. I marched in the non-violent peaceful Civil Rights protest. So the dogs and the firehoses used by a certain southern sheriff were justified, by your logic. So were the citizens beaten by Chicago police during the Democratic Convention Police riots. Resolved in 20 minutes? Nonsense. If Miss Wubbels hadn't protested as she did, she'd

likely have ended up in a back room, somewhere, and who know when she would have been granted the presence of a lawyer and what she would have gone through prior to that. You are completely ignorant of the shortage of nurses in this country - in some cases, critical shortages. And why would anyone want to be a nurse when they are disrespected by a former state and federal prosecutor such as you. </s> <not toxic 3.0424022674560547> So anyone who doesn't share your opinion is a "raving loony?"

And as for my post being "speculation?" - which part - that the Liberals 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 clique, as you apparently are.

Label toxic

Toxic Conversations Step 10

Test Modified Homos demand that you accept their fudge packing. But none of us ever will- </s> <toxic 3.1383132934570312> So you admit you would exterminate inferior humans. </s> <not toxic 3.295428514480591> Mark MacKinnon and the interests he work for would like us to 'get used to it', because they don't want to do anything practical to stop it.

Label not toxic

Test Modified I don't think anyone likes this health care bill, it stinks for everyone. 50 years and older are going to get hammered with higher premiums. People with preexisting conditions will also see their premiums go through the roof. Eventually no one will be able to afford it. They're still not addressing the reason and the solution. Prescription drugs. Why are Americans paying 2 to 3 times higher for prescription drugs than Canadians are? They should sit down with Democrats and fix it not play musical chairs with our health care. Here is the add that AARP put out, and it has a fact check at the end. It's worth watching. <http://www.thedenverchannel.com/news/politics/aarp-ad-says-house-gop-health-care-bill-would-boost-insurance-rates-for-older-americans> </s> <not toxic 2.5086519718170166> so in the mean time tens of thousands of Oregonians go without health insurance which will now be unaffordable to them. And sorry, the republicans have had 8 years to figure out a better system, they aren't going to do it anytime soon. Neither party has any desire to actually find a solution to all this. Hopefully Trump

2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314

2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365

will also soon eliminate the tax penalty for not having insurance so all us folks who buy our own insurance don't get penalized for not being able to afford to buy the insurance we are required to have. Of course he is probably clueless that detail even exists </s> <toxic 2.5147175788879395> hate to bust the bubble but over 60 of people including all those trump voters never liked or wanted obamacare and dont want it now. , trump has NOT told the gop to back off you lying sack of bs.. he wants it gone and replaced period and the gop are doing just that together WITH trump. , you will continue to turn reality into stupidity

Label not toxic

Test Modified Rather than call you a liar, I'll just post a link to HART's helicopter footage of this "most densely populated" ag land.

<https://www.youtube.com/watch?v=opJzR-WWVDdw>

And if you've been following a little closer, you'd know about the luxury towers and boutique hotels that developers are planning on fulfilling their TOD requirement with "cash gifts" to the city instead of actual affordable TOD. </s> <not toxic 2.9202661514282227> I suppose you just support urban sprawl then with that logic. </s> <toxic 2.9730660915374756> Why don't you go and live in one of their buildings and see what they're like? "Deadbeats" - you're an idiot. They're my neighbours.

Label not toxic

Test Modified Get the "real news" from Breitbart and Drudge Report? Ahahahahaha! Are you for real? </s> <toxic 2.634126901626587> "If one read the Dispatch one would think Trump is the most evil person on the planet." Not evil, just idiotic. And it would be easy to give his behaviour a pass if he wasn't POTUS. ".....five to eight anti-Trump stories per day. Never any good one's or one's that just stuck to the facts." Well, when there's good a news Trump story to print, I'm sure ADN will be all over it, problem is, there's been a dearth of those since his election. Facts? Ok Rich, give us a list of incorrect facts in the above story. </s> <not toxic 2.9902079105377197> "a gift from the political gods when the struggling effort to pass a health bill dominates the headlines."

It was a gift from media that isn'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. 2366

Every stupid mistake they make gives an even larger advantage to Trump and strengthens his supporters that already believe the MSM is biased against him and makes fence sitters begin to question what's news and what's crap. Like I've been saying, the MSM is slitting it's own throat. 2367
2368
2369
2370
2371
2372

Label not toxic 2373

Test Modified Wud luv to see one of those American Boats with all the planes on it at the bottom of the sea. </s> <toxic 3.390166759490967> I bet Regent Seven Seas will never offer Mr Hammond another trip. Wow, what a snarky article. He makes , I assume, some valid points about food and atmosphere. However, after discovering the treats available on his "massive deck" he "blew off" his remaining restaurant reservations , donned his comfy bathrobe and ordered-in. He was certainly not an ideal passenger and, for one floating on a freebie, he's a total ingrate! </s> <not toxic 3.409156084060669> Now replaced by the sexy EA-18G Growler! Using a preexisting Military Operating Area! Get over it!!!!!! 2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388

Label not toxic 2389

Test Modified Wtf, nowhere in scripture does it say that Elizabeth was a kid when married to Zachariah and nowhere does it say that Mary was a child when married to Joseph. That's a complete and utter fabrication. In fact, it says in Luke 1:42, "And she spake out with a loud voice, and said, Blessed art thou among women, and blessed is the fruit of thy womb." It says blessed art thou among WOMEN. </s> <not toxic 2.769857406616211> I was informed that my first grandchild had been conceived the evening of the day when I had inserted a prayer note in the Wailing Wall in Jerusalem that asked God to help my daughter conceive after a year of frustrated attempts. Maybe Elizabeth did the same thing? After all, she was in the same neighborhood. :-) </s> <toxic 2.9830095767974854> I don't know if the bish-ops ever told priests that if they rape a child and celebrate Mass they are committing sacrilege. But even if they didn't, the priest should know that...it is pretty basic Catholicism that receiving Communion with a mortal sin on your soul is sacrilege. 2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410

By the logic of your second paragraph, there can never be a mortally sinful sexual act, since sex acts by definition occur in a state of passion. Which begs the question, why in this case would the Scriptures go through the trouble of condemn- 2411
2412
2413
2414
2415
2416

2417 ing sexual immorality? This sounds like something
2418 your example of a rapist priest could say to himself
2419 to justify himself taking Communion in a state of
2420 mortal sin.

2421 **Label toxic**

2422 **Test Modified** Angela Merkel and all other Euro-
2423 pean political leaders who have aided and abetted
2424 the ongoing invasion of Europe by the forces of
2425 the crescent moon death cult should be tried as
2426 accessories to Mr. Urban's murder. </s> <toxic
2427 3.024935007095337> About 415 million Euro-
2428 peans cannot compete with Google, Amazon, Face-
2429 book, Oracle, Intel, Apple, etc. and the socialist Eu-
2430 ropean welfare states need more revenue because
2431 they are running out of other peoples' money.

2432 Thus the Euro-socialist-bureaucrats pick the low-
2433 hanging fruit with litigious persecution of Ameri-
2434 can firms which dominate because unlike their pa-
2435 thetic Euro-competitors, the U.S. firms are clever,
2436 hard-working, and well-capitalized.

2437 If the the Europeans wish to engage in this trans-
2438 parent financial inquisition, then the US should
2439 respond with counter litigation for trillions against
2440 corrupt scofflaws like VW (think diesel fiddle!)
2441 as well as UBS/Credit Suisse/HSBC/Credit Lyon-
2442 naise (think tax cheats!)and sue/litigate them out
2443 of existence.

2444 If the lazy, corrupt, incompetent Euros want to
2445 play with fire, then let them be financially inciner-
2446 ated! </s> <not toxic 3.2237966060638428> Rome
2447 should never have made such inane pronounce-
2448 ments at Trent in their attempt to define the sub-
2449 stance of holy Eucharist. Most reasonable people
2450 understand that perfectly well. That Rome also
2451 made their pronouncements (faith and morals) "in-
2452 fallible" is equally tragic, for the simple reason that
2453 so-called infallible statements cannot be retracted
2454 without calling into question other so-called infalli-
2455 ble statements.

2456 Sincere question for you: If Jesus and his follow-
2457 ers celebrated Eucharist as a communal meal seated
2458 around a table, what gives Rome the right to alter
2459 this simple act of worship (perhaps "fellowship" is
2460 a better word—more suited toward love of God and
2461 neighbor), given to us by the Lord himself?

2462 **Label toxic**

2463 **Test Modified** I hope you don't have kids if you
2464 see this woman's actions as acceptable. And I ap-
2465 plaud the den for kicking the kid out. She brought
2466 unwanted negative attention upon them. However,
2467 she will, and is already likely, pay the the price for

her stupid stunt. </s> <toxic 2.873025417327881> 2468
Ms. Van Brocklin: You state that there is Payne's 2469
conduct is inexcusable, then proceed to use the rest 2470
of your space to justify him and his actions. You 2471
have denigrated the proud courage of countless peo- 2472
ple who took up causes via civil disobedience. I 2473
marched in the non-violent peaceful Civil Rights 2474
protest. So the dogs and the firehoses used by a cer- 2475
tain southern sheriff were justified, by your logic. 2476
So were the citizens beaten by Chicago police dur- 2477
ing the Democratic Convention Police riots. Res- 2478
olved in 20 minutes? Nonsense. If Miss Wubbels 2479
hadn't protested as she did, she'd likely have ended 2480
up in a back room, somewhere, and who know 2481
when she would have been granted the presence of 2482
a lawyer and what she would have gone through 2483
prior to that. You are completely ignorant of the 2484
shortage of nurses in this country - in some cases, 2485
critical shortages. And why would anyone want 2486
to be a nurse when they are disrespected by a for- 2487
mer state and federal prosecutor such as you. </s> 2488
<not toxic 2.8961446285247803> Well, I can't very 2489
well respect or fear an imaginary sky-being. As for 2490
my concept of character, it was good enough for 2491
the Alaska Judicial Council and Governor Knowles. 2492
But that was long ago. I've gotten older and, crikey, 2493
maybe I am going downhill. You're right about 2494
the inappropriateness of my comment. First Lady 2495
Walker's piece is very laudable and I shouldn't have 2496
taken it as an occasion to rant. (But look on the 2497
bright side: my misplaced comment gave you yet 2498
another occasion to rant about how your Fosterism 2499
is saving civilization.) 2500

2501 **Label toxic**

2502 **Test Modified** no one cares what a paid liberal
2503 trolling hack like you believes lunatic., </s> <toxic
2504 2.5876824855804443> It always amuses me when
2505 a troll gets on, they like their own comments and
2506 simply assert everyone else is wrong. Never any
2507 evidence to rebut it just blind assertions. </s> <not
2508 toxic 2.810284376144409> Ouch... didn't see that
2509 one coming. A liberal stealing my own line... just
2510 like they take everything else they like.

2511 **Label toxic**

2512 **Test Modified** Ok all you NDP "LEAP" mani-
2513 festo types, where is your hero Naomi Klein? Her
2514 fawning adoration of Chavez and Venezuelan thug-
2515 gery knows no bounds. I'm sure she's awfully
2516 hysterical over the thought that such a pathetic dict-
2517 storship could ever be sanctioned. </s> <toxic
2518 2.883418321609497> Um, no. The major left-

2519 wing Labour party was decimated; Mr. Rutte lost 2570
2520 8 seats; and Mr. Wilders Freedom party GAINED 2571
2521 4 seats. Now Mr. Rutte will have to garner favour 2572
2522 among RIGHT-WING parties to cobble together 2573
2523 his coalition. And as to your glib little comment 2574
2524 about people "embracing left-leaning parties", you 2575
2525 need only look to Canada to see the fiasco that 2576
2526 results... The corrupt and divisive Trudeau gov- 2577
2527 ernment was elected on a slate of fuzzy, alt-left 2578
2528 promises; only to renege on most of them ! </s> 2579
2529 <not toxic 3.0424022674560547> So anyone who 2580
2530 doesn't share your opinion is a "raving loony?" 2581
2531 And as for my post being "speculation? - which 2582
2532 part - that the Liberals are the party in power, or 2583
2533 that this involves money? 2584
2534 As for me not knowing what is going on, you 2585
2535 are correct, I am not a member of the Liberal party 2586
2536 insider clique, as you apparently are. 2587
2537 **Label toxic** 2588
2538 **Hate Speech Offensive Step 1** 2589
2539 **Test Modified** ' If this ugly ass nig- 2590
2540 gah can be with her , I can too . Tf 2591
2541 she with Chuu Baka for anyways this nig- 2592
2542 gah look like a... <http://t.co/IVNDRDPpQT> 2593
2543 </s> <hate speech 3.680654525756836> AtBlack- 2594
2544 face MrMooncricket UncleUnco BobbyBotsods 2595
2545 FuckTheMallcops RacistNegro69 ape? Fuck 2596
2546 u broke ass racist </s> <offensive language 2597
2547 3.9933786392211914> Tamela X But a wannabe 2598
2548 ghetto white girl will deal with an unaccomplished 2599
2549 black man for more than sex. </s> <neither 2600
2550 4.20903205871582> VNNForum These Marxist 2601
2551 kikes promoted Anti-Semitism by creating this film, 2602
2552 what was their incentive? They just promote their 2603
2553 enemies. 2604
2554 **Label hate speech** 2605
2555 **Test Modified** 8220; iamkrause: No need to 2606
2556 thank me, killing niggers is it's own reward. 8221; 2607
2557 Dis my pal from our graduation when we was 26 2608
2558 </s> <neither 3.5232555866241455> LanceTHES- 2609
2559 POKEN Statement of the Day: Ppl who RT 2610
2560 ChiefsatWar are even dumber than he is. The 2611
2561 ole birds of a feather type of deal. </s> <hate 2612
2562 speech 3.5796432495117188> Brenddan they said 2613
2563 some bullshit. Damn niggers </s> <offensive lan- 2614
2564 guage 3.7068209648132324> BRUH lmaoo 8220; 2615
2565 NoBeetch: Them hoes was nasty but I kept eating 2616
2566 them " roblokk: What the fuck happened to these? 2617
2567 <http://t.co/G9w10SmQdr>" 8221; 2618
2568 **Label hate speech** 2619
2569 **Test Modified** RT IsaidNick: niggers are 2620

so ignorant <http://t.co/P9jDdVsRVb> </s> <hate 2570
speech 3.2627384662628174> AtBlackface Mr- 2571
Mooncricket UncleUnco BobbyBotsods FuckThe- 2572
Mallcops RacistNegro69 ape? Fuck u broke 2573
ass racist </s> <neither 3.7804527282714844> 2574
VNNForum These Marxist kikes promoted Anti- 2575
Semitism by creating this film, what was their in- 2576
centive? They just promote their enemies. </s> <of- 2577
fensive language 3.9285027980804443> Tamela X 2578
But a wannabe ghetto white girl will deal with an 2579
unaccomplished black man for more than sex. 2580
Label hate speech 2581
Test Modified RT RosieZaya1: Ur 2582
fucking white trash </s> <hate speech 2583
2.951470136642456> AtBlackface MrMoon- 2584
cricket UncleUnco BobbyBotsods FuckTheMall- 2585
cops RacistNegro69 ape? Fuck u broke ass racist 2586
</s> <offensive language 3.6144936084747314> 2587
Tamela X But a wannabe ghetto white girl will 2588
deal with an unaccomplished black man for more 2589
than sex. </s> <neither 3.7668633460998535> 2590
VNNForum These Marxist kikes promoted 2591
Anti-Semitism by creating this film, what was their 2592
incentive? They just promote their enemies. 2593
Label hate speech 2594
Test Modified mike ray7 congratulations, you 2595
are officially fucking retarded. </s> <nei- 2596
ther 3.4077796936035156> RT JakeG Based- 2597
God: "Never go full retard" </s> <hate speech 2598
3.479813575744629> Brenddan they said some 2599
bullshit. Damn niggers </s> <offensive lan- 2600
guage 3.601623773574829> BRUH lmaoo 8220; 2601
NoBeetch: Them hoes was nasty but I kept eating 2602
them " roblokk: What the fuck happened to these? 2603
<http://t.co/G9w10SmQdr>" 8221; 2604
Label hate speech 2605
Test Modified gonna have them pussies mix up 2606
some concrete today. teach them to pose like me. 2607
I am a badass motherfucker. and I will let you 2608
be too (: </s> <hate speech 3.227602243423462> 2609
AtBlackface MrMooncricket UncleUnco Bobby- 2610
Botsods FuckTheMallcops RacistNegro69 ape? 2611
Fuck u broke ass racist </s> <offensive lan- 2612
guage 3.5520639419555664> BRUH lmaoo 8220; 2613
NoBeetch: Them hoes was nasty but I kept 2614
eating them " roblokk: What the fuck hap- 2615
pened to these? <http://t.co/G9w10SmQdr>" 8221; 2616
</s> <neither 3.809976816177368> RT Venus 2617
Lynn: 8220; dylxnl: look ghetto but it work 2618
<http://t.co/chrwW9dPca> 8221; 128557; 128557; 2619
128557; 128557; 128557; 128557; 128557; 2620

2621	128557; 128557; 128557;		
2622	Label offensive language		
2623	Test Modified maddievaans u call ur bestie a	Test Modified My baby guinea pig is so cute	2672
2624	bitch I'm guessing she's a dog that barks too much	Adorable </s> <hate speech 3.82438588142395>	2673
2625	</s> <offensive language 3.2484242916107178>	whiteangelss84 y dont u end us monkeys then? U	2674
2626	RT Im Amy Bitches: I'm not always a	hate us because were better than crackers amp; I	2675
2627	bitch, sometimes I sleep </s> <hate speech	I love it. Deep down u know weve done nothin 2 u	2676
2628	3.3403964042663574> Women who are feminist	</s> <offensive language 3.8302881717681885> I	2677
2629	are the ugly bitches who cant find a man for them-	always wanted a bull dog them hoes clean fuck a	2678
2630	selves </s> <neither 3.8210768699645996> Gary's	pit </s> <neither 3.8650126457214355> Break-	2679
2631	girl was a tranny. RT Em 99car: I wonder what	fast fried chicken jerk chicken Tater tots white	2680
2632	would have become of rickspringfield if he'd writ-	rice nd press yellow rice nd beans Mac nd cheese	2681
2633	ten 'Gary's Girl'.	http://t.co/Usz8gJnZl0	2682
2634	Label offensive language	Label neither	2683
2635	Test Modified No less than 3 bad bitches in	Test Modified RT Kick Man: Giants- Pitiful	2684
2636	my bed at a time... </s> <offensive language	.. Jets-Pitiful .. Mets- Pitiful .. Yankees-Pitiful	2685
2637	3.2010550498962402> Then again... i shoot	.. Nets- Pitiful .. Knicks-Pitiful ... Ny sports-	2686
2638	bitches in the ass for doing some stupid shit like	Pitiful </s> <neither 3.8152754306793213> You	2687
2639	key my car or bust my windows. </s> <hate	know I'm not big on the NFL, but I'm so sick of	2688
2640	speech 3.666656732559204> bitch kill yoself, go	hearing all of this "Black and yellow" shit. LOL	2689
2641	on to the bathroom and eat the pills bitch, all	bandwagon fans and hell, GO PACKERS! </s>	2690
2642	of em </s> <neither 4.249817848205566> 9733;	<offensive language 4.071953773498535> BRUH	2691
2643	BEST ASIAN MASSAGE ON THE park slope	lmaoo 8220; NoBeeetch: Them hoes was nasty but	2692
2644	-TOP RATED SPA 9733; 718-622-0221 - 24	I kept eating them " roblokk: What the fuck hap-	2693
2645	http://t.co/ZsAAzFL0p5	pened to these? http://t.co/G9w10SmQdr" 8221;	2694
2646	Label offensive language	</s> <hate speech 4.095180511474609> whitean-	2695
2647	Test Modified RT TheDrugTribe: mary isn't a	angelss84 y dont u end us monkeys then? U hate us	2696
2648	backstabbing bitch that lies and deceives me </s>	because were better than crackers amp; I love it.	2697
2649	<offensive language 3.4536943435668945> RT Im	Deep down u know weve done nothin 2 u	2698
2650	Amy Bitches: I'm not always a bitch, sometimes	Label neither	2699
2651	I sleep </s> <hate speech 3.6065785884857178>	Test Modified jesstoth we could get matching	2700
2652	vinny2vicious faggot I knew you weren't really	burner phones and be ghetto fab for a few months	2701
2653	my friend. </s> <neither 3.638406753540039>	</s> <hate speech 3.4667954444885254> whitean-	2702
2654	Gary's girl was a tranny. RT Em 99car: I wonder	angelss84 y dont u end us monkeys then? U hate	2703
2655	what would have become of rickspringfield if he'd	us because were better than crackers amp; I love	2704
2656	written 'Gary's Girl'.	it. Deep down u know weve done nothin 2 u	2705
2657	Label offensive language	</s> <offensive language 3.595543622970581> RT	2706
2658	Test Modified porn, android, iphone, ipad, sex,	NickBratton3: I wish my parents bought me a	2707
2659	xxx, CloseUp Squirtng pussy and fingered ass-	car man.. People bitch about not getting what	2708
2660	hole http://t.co/bKYeoUwWv2 </s> <offensive lan-	car they want when they want it, and its free	2709
2661	guage 3.5574071407318115> BRUH lmaoo 8220;	8230; </s> <neither 3.6075007915496826> RT	2710
2662	NoBeeetch: Them hoes was nasty but I kept eat-	Venus Lynn: 8220; dylxnl: look ghetto but it work	2711
2663	ing them " roblokk: What the fuck happened	http://t.co/chrvW9dPca 8221; 128557; 128557;	2712
2664	to these? http://t.co/G9w10SmQdr" 8221; </s>	128557; 128557; 128557; 128557; 128557;	2713
2665	<neither 3.6928675174713135> DegenerateArtist	128557; 128557; 128557;	2714
2666	Sniffs whiffy balls involuntary, cuz a FAIRY walks	Label neither	2715
2667	DOWNTOWN HAIRY, climatechange all our scary	Test Modified Thw White Iron Band plays	2716
2668	sense of hillbilly sneeze gt; </s> <hate speech	this weekend in Fargo,ND at the Aquar-	2717
2669	3.7489092350006104> Wake up you stoopid dyke	ium(21+) ,Friday(10-29-10) with Charlie	2718
2670	RandyPawlik	Parr. The next night,Saturday... </s> <hate	2719
2671	Label offensive language	speech 3.8393898010253906> 128514; 128514;	2720
		128514;RT kwagiheath: Them 1st 48 Charlotte	2721
		goon niccas gon Fkkk around and burn Bank Of	2722

2723 America stadium down 128293; </s> <offensive
2724 language 3.9061098098754883> I be telling
2725 Mcgirt music ain't enough.You gotta have a non
2726 music related agenda.They crackers want to
2727 sell something with your face. </s> <neither
2728 3.980140447616577> You know I'm not big on
2729 the NFL, but I'm so sick of hearing all of this
2730 "Black and yellow" shit. LOL bandwagon fans and
2731 hell, GO PACKERS!

2732 **Label** neither 2779

2733 **Test Modified** RT dsrtvet: FoxNews tjoy7 2780
2734 And I don't have any confidence NONWHAT- 2781
2735 SOEVER in you Barack! You're the sole 2782
2736 reason why this country is in thi 8220; 2783
2737 </s> <neither 3.78818678855896> RT lach- 2784
2738 lan: So who wants to tell the Washing- 2785
2739 ton Post that Charlie Crist is a Democrat? 2786
2740 http://t.co/LGzzYusEKZ http://t.co/2jCVv7qxqf 2787
2741 </s> <hate speech 3.942119598388672> whitean- 2788
2742 gels84 y dont u end us monkeys then? U hate 2789
2743 us because were better than crackers amp; I love 2790
2744 it. Deep down u know weve done nothin 2 u </s> 2791
2745 <offensive language 3.984135150909424> BRUH 2792
2746 lmaoo 8220; NoBeetch: Them hoes was nasty but 2793
2747 I kept eating them " roblockk: What the fuck hap- 2794
2748 pened to these? http://t.co/G9w10SmQdr" 8221; 2795
2749 **Label** neither 2796

2750 **Hate Speech Offensive Step 5** 2797

2751 **Test Modified** ' If this ugly ass niggah can 2798
2752 be with her , I can too . Tf she with 2799
2753 Chuu Baka for anyways this niggah look like 2800
2754 a... http://t.co/lVNDRDPPpQT </s> <hate speech 2801
2755 2.7132835388183594> RT WaavyLee: His balls 2802
2756 ashy RT Yattabing: Trelaire1st: Real women do 2803
2757 this http://t.co/VG5DBqH8aT 8221; and real fag- 2804
2758 gots let em do that.. Smh </s> <offensive language 2805
2759 2.7854363918304443> RT JayyMfCarter: If you 2806
2760 gotta nigga or bitch.... PLEASE STAY OUT MY 2807
2761 DM's I'm tired of y'all taken girls </s> <neither 2808
2762 3.0996105670928955> NeonTreezz PoCBeauty 2809
2763so.....white people selling African art? NOT 2810
2764 COOL? Okay. blackish redskins doublestandard 2811
2765 **Label** hate speech 2812

2766 **Test Modified** 8220; iamkrause: No need to 2813
2767 thank me, killing niggers is it's own reward. 8221; 2814
2768 Dis my pal from our graduation when we 2815
2769 was 26 </s> <hate speech 2.6058554649353027> 2816
2770 samzbikowski some negro amigo pulled a gun 2817
2771 on Nate and I a few weeks ago. I was 2818
2772 STOKED! 128299; 128299; 128299; </s> <neither 2819
2773 2.670192003250122> RT jennaferjenbub: 8220; 2820

BarryClerjuste: "Anything below a A+ and we 2774
disown you ling ling" http://t.co/m1QiWK4xZg 2775
8221; AustinBedsaul </s> <offensive language 2776
2.740609645843506> 8220; Alondra Lu: Ain't 2777
that a bitch 8221; 2778

Label hate speech 2779

Test Modified RT IsaidNick: niggers are 2780
so ignorant http://t.co/P9jDdVsRVb </s> <hate 2781
speech 2.057800769805908> RT WhitesOnly 1: 2782
niggers! http://t.co/Hb3uJaLky2 </s> <neither 2783
2.7749483585357666> amp; thots are wearing 2784
Uggs RT BigBootyJudy814: ItsFallBecause ne- 2785
gros are pulling out their Timbs" </s> <offensive 2786
language 2.926729440689087> RT Jayy Gee96: 2787
Dumb bitches 2788

Label hate speech 2789

Test Modified RT RosieZaya1: Ur 2790
fucking white trash </s> <hate speech 2791
2.422173500061035> FrankieJGrande fugly 2792
queer white trash </s> <offensive language 2793
2.6756434440612793> RT Jayy Gee96: Dumb 2794
bitches </s> <neither 2.783188819885254> RT 2795
BeardedNixon: Poont gotta be trash 2796

Label hate speech 2797

Test Modified mike ray7 congratulations, you 2798
are officially fucking retarded. </s> <hate 2799
speech 2.4854748249053955> darthdanaa Yes 2800
you do retard. </s> <offensive language 2801
2.851564645767212> Lol!! 8220; ItzSweetz 2802
Bitch: Ooop! QT TIFFANY PORSCHE: You little 2803
twats. 8221; </s> <neither 2.8971688747406006> 2804
RT jennaferjenbub: 8220; BarryClerjuste: "Any- 2805
thing below a A+ and we disown you ling ling" 2806
http://t.co/m1QiWK4xZg 8221; AustinBedsaul 2807

Label hate speech 2808

Test Modified gonna have them pussies mix up 2809
some concrete today. teach them to pose like me. I 2810
am a badass motherfucker. and I will let you be too 2811
(: </s> <offensive language 2.7589027881622314> 2812
40oz VAN IYCM I. I can't get any work done if 2813
you keep showin off your bitches. </s> <hate 2814
speech 2.8690829277038574> SlightlyAdjusted 2815
RT CapoToHeaven Alls niggers wanna do is fuck, 2816
tweet, and drink pineapple soda all day </s> <nei- 2817
ther 3.0193798542022705> cakedjake We're lay- 2818
ing rock around our lake. You're welcome to join a 2819
redneck workout. muscles 128170; 128513; 2820

Label offensive language 2821

Test Modified maddieevaans u call ur bestie a 2822
bitch I'm guessing she's a dog that barks too much 2823
</s> <offensive language 3.0170230865478516> 2824

2825 Lol!! 8220; ItzSweetz Bitch: Ooop! QT TiFFANY 2876
2826 P0RSCH3: You little twats. 8221; </s> <hate 2877
2827 speech 3.0850884914398193> Princesslexii16 2878
2828 Fucking coon </s> <neither 3.21132493019104> 2879
2829 lmaoooo RT ComedyTruth: Girls, don't let a guy 2880
2830 treat you like a yellow starburst. You are a pink 2881
2831 starburst. 2882
2832 **Label** offensive language 2883
2833 **Test Modified** No less than 3 bad bitches in 2884
2834 my bed at a time... </s> <offensive language 2885
2835 3.023017406463623> Don't lose sleep bout these 2886
2836 bitches bc they come and go 128076; </s> <hate 2887
2837 speech 3.2786214351654053> bitch kill yoself, go 2888
2838 on to the bathroom and eat the pills bitch, all of em 2889
2839 </s> <neither 3.4171059131622314> Keep those 2890
2840 away frm Charlie Day RT JhonenV: Just once in 2891
2841 my life I'd like for someone's favorite part of my 2892
2842 body to be my disgusting knees. 2893
2843 **Label** offensive language 2894
2844 **Test Modified** RT TheDrugTribe: mary isn't 2895
2845 a backstabbing bitch that lies and deceives me 2896
2846 </s> <offensive language 2.991457939147949> 2897
2847 RT StevStiffler: If her bio says "Only God can 2898
2848 judge me" she's a hoe. </s> <hate speech 2899
2849 3.098494529724121> RT sorryimalex: I got 2900
2850 called a faggot for buying girl toms so now I'm 2901
2851 gonna fuck that person in the ass </s> <neither 2902
2852 3.3266849517822266> lmaoooo RT ComedyTruth: 2903
2853 Girls, don't let a guy treat you like a yellow star- 2904
2854 burst. You are a pink starburst. 2905
2855 **Label** offensive language 2906
2856 **Test Modified** porn, android, iphone, ipad, sex, 2907
2857 xxx, | CloseUp | Squirting pussy and fingered 2908
2858 asshole http://t.co/bKYeoUwWv2 </s> <neither 2909
2859 1.5677733421325684> porn, android, iphone, ipad, 2910
2860 sex, xxx, | Desi | paki http://t.co/XxcdQvzi9t </s> 2911
2861 <hate speech 2.8566393852233887> RT mitchman- 2912
2862 cuso: BrantPrintup:What straight guys take a pic- 2913
2863 ture of themselves naked in a hot tun.What fags 2914
2864 jakesiw Ryan Murphy3 Randy 8230; </s> <offen- 2915
2865 sive language 2.932191848754883> Lol!! 8220; 2916
2866 ItzSweetz Bitch: Ooop! QT TiFFANY P0RSCH3: 2917
2867 You little twats. 8221; 2918
2868 **Label** offensive language 2919
2869 **Test Modified** My baby guinea pig is so cute 2920
2870 Adorable </s> <neither 3.1643435955047607> 2921
2871 Our female guinea pig is pregnant 127882; 127881; 2922
2872 127873; 128525; 128525; 128525; </s> <offen- 2923
2873 sive language 3.4907007217407227> I impress da 2924
2874 young white girl next doe by taking out my gi- 2925
2875 ant negro thang and usin it to flip da hamburgers 2926

2927 **Test Modified** RT dsrtvet: FoxNews tjoy7 2978
2928 And I don't have any confidence NONWHAT- 2979
2929 SOEVER in you Barack! You're the sole rea- 2980
2930 son why this country is in thi 8230; </s> <hate 2981
2931 speech 2.696760654449463> RT veeveeveeveevee: 2982
2932 If I was Obama Id call a press conference amp; 2983
2933 slit joe bidens neck on live tv just 2 show these 2984
2934 crackers I mean business 8230; </s> <neither 2985
2935 2.762817144393921> RT jennaferjenbub: 8220; 2986
2936 BarryClerjuste: "Anything below a A+ and we 2987
2937 disown you ling ling" http://t.co/m1QiWK4xZg 2988
2938 8221; AustinBedsaul </s> <offensive language 2989
2939 2.9660327434539795> RT CoffyBrownChi: If he 2990
2940 don't believe you, no refunds hoe. 2991
2941 **Label** neither 2992
2942 **Hate Speech Offensive Step 10** 2993
2943 **Test Modified** ' If this ugly ass niggah can 2994
2944 be with her , I can too . Tf she with 2995
2945 Chuu Baka for anyways this niggah look like 2996
2946 a... http://t.co/IVNDRDPpQT </s> <offensive 2997
2947 language 2.6535706520080566> RT CurrenSy 2998
2948 Spitta: And if a bitch can't respect a nigga 2999
2949 wit some paper and a fresh pair of bball shorts 3000
2950 then she was raised terribly.. </s> <hate speech 3001
2951 2.7132835388183594> RT WaavyLee: His balls 3002
2952 ashy RT Yattabing: Trelaire1st: Real women 3003
2953 do this http://t.co/VG5DBqH8aT 8221; and real 3004
2954 faggots let em do that.. Smh </s> <neither 3005
2955 3.0996105670928955> NeonTreezz PoCBeauty 3006
2956so.....white people selling African art? NOT 3007
2957 COOL? Okay. blackish redskins doublestandard 3008
2958 **Label** hate speech 3009
2959 **Test Modified** 8220; iamkrause: No need to 3010
2960 thank me, killing niggers is it's own reward. 8221; 3011
2961 Dis my pal from our graduation when we was 3012
2962 26 </s> <hate speech 2.545886278152466> RT 3013
2963 Tae Rhodes: 8220; kim92493: 8220; Tae Rhodes: 3014
2964 kim92493 patpatbush uhhh you've been judged 3015
2965 8221; it happens. whitepower...I'll hang you nig- 3016
2966 ger 8221; wo 8230; </s> <offensive language 3017
2967 2.6044790744781494> 8220; NoRapist: on my 3018
2968 way to fuck ur bitch http://t.co/SgVBBrwOg2 8221; 3019
2969 mckinley719 </s> <neither 2.670259714126587> 3020
2970 RT jennaferjenbub: 8220; BarryClerjuste: "Any- 3021
2971 thing below a A+ and we disown you ling ling" 3022
2972 http://t.co/m1QiWK4xZg 8221; AustinBedsaul 3023
2973 **Label** hate speech 3024
2974 **Test Modified** RT IsaidNick: niggers are 3025
2975 so ignorant http://t.co/P9jDdVsRVb </s> <hate 3026
2976 speech 2.057800769805908> RT WhitesOnly 1: 3027
2977 niggers! http://t.co/Hb3uJaLky2 </s> <neither 3028

2.7749483585357666> amp; thots are wearing 2978
Uggs RT BigBootyJudy814: ItsFallBecause ne- 2979
gros are pulling out their Timbs" </s> <offensive 2980
language 2.8298287391662598> This Uncle Tom 2981
mother fucking wants to invoke 3/5 a man in his 2982
speech? Dude you ain't white no matter how 2983
much... http://t.co/3yrcyC9ezc 2984
Label hate speech 2985
Test Modified RT RosieZaya1: Ur 2986
fucking white trash </s> <hate speech 2987
2.422173500061035> FrankieJGrande fugly 2988
queer white trash </s> <offensive language 2989
2.6756434440612793> RT Jayy Gee96: Dumb 2990
bitches </s> <neither 2.783188819885254> RT 2991
BeardedNixon: Poont gotta be trash 2992
Label hate speech 2993
Test Modified mike ray7 congratulations, you 2994
are officially fucking retarded. </s> <hate 2995
speech 2.4854748249053955> darthdanaa Yes 2996
you do retard. </s> <offensive language 2997
2.8516175746917725> Lol!! 8220; ItzSweetz 2998
Bitch: Ooop! QT TiFFANY PORSCHE: You little 2999
twats. 8221; </s> <neither 2.8972203731536865> 3000
RT jennaferjenbub: 8220; BarryClerjuste: "Any- 3001
thing below a A+ and we disown you ling ling" 3002
http://t.co/m1QiWK4xZg 8221; AustinBedsaul 3003
Label hate speech 3004
Test Modified gonna have them pussies mix up 3005
some concrete today. teach them to pose like me. I 3006
am a badass motherfucker. and I will let you be too 3007
(: </s> <offensive language 2.758687734603882> 3008
40oz VAN IYCM I. I can't get any work done if 3009
you keep showin off your bitches. </s> <hate 3010
speech 2.8036365509033203> Just to get u mad 3011
go on your search bar on here and search up 3012
"stupid niggers" amp; hop on somebodys head 3013
then mention me lol stonethegreat23 </s> <nei- 3014
ther 3.012741804122925> charloosss keepitplur 3015
nicoleariel I'll chug my tall can . but homegirl 3016
won't approve lol 3017
Label offensive language 3018
Test Modified maddievaans u call ur bestie 3019
a bitch I'm guessing she's a dog that barks too 3020
much </s> <hate speech 2.8469488620758057> 3021
RylannWilliams whooooo? Chelsey? Fuck her 3022
lol. She juss a bitch </s> <offensive language 3023
2.8842358589172363> RT Ezzzylove: She a bad 3024
bitch, let's get to it right away . </s> <neither 3025
3.0819990634918213> charliesheen Charlie, im 3026
an old lady. don't EVER SAY UGLY THINGS 3027
ABOUT UR CHILDRENS MOM.. I GET IT!!!, 3028

3029 **JUS DONT! BIG HUG**

3030 **Label** offensive language

3031 **Test Modified** No less than 3 bad bitches

3032 in my bed at a time... </s> <offensive lan-

3033 guage 2.8522520065307617> Bad bitches in

3034 the pen make my toes curl </s> <hate speech

3035 3.2539432048797607> I didn't forsake all other

3036 bitches for my wife to be getting fucked on by

3037 another nigga. and you know she married? you

3038 gotta die. </s> <neither 3.4170782566070557>

3039 Keep those away frm Charlie Day RT JhonenV:

3040 Just once in my life I'd like for someone's favorite

3041 part of my body to be my disgusting knees.

3042 **Label** offensive language

3043 **Test Modified** RT TheDrugTribe: mary isn't

3044 a backstabbing bitch that lies and deceives me

3045 </s> <offensive language 2.9916186332702637>

3046 RT StevStiffler: If her bio says "Only God can

3047 judge me" she's a hoe. </s> <hate speech

3048 3.02489972114563> triple6em96 Hunglikerobby

3049 bitch you watch your fucking mouth you dirty

3050 whore. I swear to god that's a thin line </s>

3051 <neither 3.1058743000030518> RT shakiraevanss:

3052 Criticize Amanda for saying the n word, sure, but

3053 don't make jokes about her sexual assault, don't be

3054 trash.

3055 **Label** offensive language

3056 **Test Modified** porn, android, iphone, ipad, sex,

3057 xxx, | CloseUp | Squirting pussy and fingered

3058 asshole http://t.co/bKYeoUwWv2 </s> <neither

3059 1.5677733421325684> porn, android, iphone, ipad,

3060 sex, xxx, | Desi | paki http://t.co/XxcdQvzI9t

3061 </s> <offensive language 2.8408925533294678>

3062 RT FunnyPicsDepot: bitches be like "I'm a vir-

3063 gin" http://t.co/mFDwXmg8ic </s> <hate speech

3064 2.8566393852233887> RT mitchmancuso: Brant-

3065 Printup:What straight guys take a picture of them-

3066 selves naked in a hot tun.What fags jakesiw Ryan

3067 Murphy3 Randy 8230;

3068 **Label** offensive language

3069 **Test Modified** My baby guinea pig is so cute

3070 Adorable </s> <neither 3.1643435955047607>

3071 Our female guinea pig is pregnant 127882; 127881;

3072 127873; 128525; 128525; 128525; </s> <offen-

3073 sive language 3.4907007217407227> I impress da

3074 young white girl next doe by taking out my gi-

3075 ant negro thang and usin it to flip da hamburgers

3076 for da KoolQueefTribute 160; </s> <hate speech

3077 3.5861661434173584> What a wetback looks like

3078 when he gets caught crossing the border. Ilovebamf

3079 http://t.co/j3Uf1TYubO

Label neither

Test Modified RT Kick Man: Giants- Piti-

ful .. Jets-Pitiful .. Mets- Pitiful .. Yankees-

Pitiful .. Nets- Pitiful .. Knicks-Pitiful ... Ny

sports- Pitiful </s> <neither 3.0366640090942383>

Buster ESPN Huh.....last 10 games..Tampa 8-2/Balt

7-3/Yanks 6-4...and they lost their best pitcher.

Please explain your logic. </s> <hate speech

3.3251242637634277> RT J R: Smh nigga is

mildly retarded RT Thotcho: LMFAO RT JustDoJ:

If Griff wasn 8217;t injuries we 8217;d legit be 6-1

</s> <offensive language 3.3433356285095215>

Them shits ugly hoe. RT SirRocObama: RT

BurgerKing: All these nuggets amp; u still actin

chicken. http://t.co/tRy8Lvyo9O

Label neither

Test Modified jesstoth we could get match-

ing burner phones and be ghetto fab for a few

months </s> <hate speech 3.034785270690918>

SAMMI boyden bruh we can finally roll like red-

necks (: ((drug dealers)) </s> <offensive lan-

guage 3.122525930404663> JZolly23 JBilovich

we need to grow mullets together so we can get all

the bitches and HannahKubiak can hate on us </s>

<neither 3.271000623703003> RT sassytbh: a girl

tweeted "you might be ghetto if u bring food from

outside into the movies"

no u might be stupid if u pay 4.99 for a b 8230;

Label neither

Test Modified Thw White Iron Band plays

this weekend in Fargo,ND at the Aquar-

ium(21+) ,Friday(10-29-10) with Charlie Parr.

The next night,Saturday... </s> <neither

3.2462401390075684> RT toddknife: Full weak-

enednachos set (except the last song) from South-

ern Darkness Fest last month. Who's the ape

on guitar? https://t.c 8230; </s> <hate speech

3.3524651527404785> Eagles fuck around amp;

lose it'll be kill the cracker at the Sophi crib smfh

</s> <offensive language 3.511016368865967>

My dawg ceomiamimike told me it's a must I be

901k2lounge this Saturday ROCKIN that bitch wit

Tha 8230; http://t.co/0NV9cHtwOs

Label neither

Test Modified RT dsrtvet: FoxNews tjoy7

And I don't have any confidence NONWHAT-

SOEVER in you Barack! You're the sole rea-

son why this country is in thi 8230; </s> <hate

speech 2.696760654449463> RT veeveeveeveevee:

If I was Obama Id call a press conference amp;

slit joe bidens neck on live tv just 2 show these

3131 crackers I mean business 8230; </s> <neither
3132 2.762908458709717> RT jennaferjenbub: 8220;
3133 BarryClerjuste: "Anything below a A+ and we
3134 disown you ling ling" http://t.co/m1QiWK4xZg
3135 8221; AustinBedsaul </s> <offensive language
3136 2.881894588470459> 8220; LongMoneyTonny:
3137 vintage monroe DONT Say Shit Else ! Just Stfu
3138 ! 8221;bitch we can do a lot more off this Twitter
3139 shit you can come see me
3140 **Label** neither
3141 **LIAR (collapsed) Step 1**
3142 **Test Modified** Afscome says In labor negotia-
3143 tions with city employees, Milwaukee Mayor Tom
3144 Barrett demanded concessions that went beyond
3145 those mandated by Gov. Scott Walkers collective
3146 bargaining law </s> a letter to members </s> <true
3147 statement 3.833270311355591> Donald Trump
3148 says Libya Ambassador (Christopher) Stevens sent
3149 600 requests for help in Benghazi. </s> the sec-
3150 ond 2016 presidential debate </s> <false state-
3151 ment 4.013778209686279> Donald Trump says
3152 The federal government is sending refugees to
3153 states with governors who are Republicans, not
3154 to the Democrats. </s> an interview on Laura In-
3155 graham's radio show
3156 **Label** true statement
3157 **Test Modified** Rick Scott says All Aboard
3158 Florida is a 100 percent private venture. There
3159 is no state money involved. </s> a TV interview
3160 </s> <false statement 3.664231777191162> Don-
3161 ald Trump says The federal government is sending
3162 refugees to states with governors who are Republi-
3163 cans, not to the Democrats. </s> an interview on
3164 Laura Ingraham's radio show </s> <true statement
3165 3.831820011138916> Patrick Murphy says Marco
3166 Rubio opposes immigration reform. Worse, Rubio
3167 supports Donald Trump. His plan would deport
3168 800,000 children, destroying families. </s> a TV
3169 ad
3170 **Label** true statement
3171 **Test Modified** Julie Pace says The Obama ad-
3172 ministration is using as its legal justification for
3173 these airstrikes (on the Islamic State), an autho-
3174 rization for military force that the president him-
3175 self has called for repeal of. </s> a question to
3176 White House Press Secretary Josh Earnest </s>
3177 <false statement 3.5803754329681396> Donald
3178 Trump says Hillary Clinton invented ISIS with
3179 her stupid policies. She is responsible for ISIS.
3180 </s> an interview on 60 Minutes </s> <true state-
3181 ment 3.869307518005371> Donald Trump says

Libya Ambassador (Christopher) Stevens sent 600
requests for help in Benghazi. </s> the second
2016 presidential debate
Label true statement
Test Modified John Kasich says We are now
eighth in the nation in job creation . . . we are
No. 1 in the Midwest. </s> a news conference </s>
<true statement 3.851958990097046> Jorge Elorza
says In the last six years of Ciancis administration
violent crime was down in the United States. It
was down in the region. It was down in Rhode
Island. But it was up in Providence. </s> a debate
</s> <false statement 4.010262966156006> Don-
ald Trump says The federal government is sending
refugees to states with governors who are Republi-
cans, not to the Democrats. </s> an interview on
Laura Ingraham's radio show
Label true statement
Test Modified Mike Pence says It was Hillary
Clinton who left Americans in harms way in Beng-
hazi and after four Americans fell said, What dif-
ference at this point does it make? </s> the Re-
publican national convention </s> <true statement
3.7440342903137207> Jorge Elorza says In the last
six years of Ciancis administration violent crime
was down in the United States. It was down in the
region. It was down in Rhode Island. But it was up
in Providence. </s> a debate </s> <false statement
3.746598958969116> Donald Trump says You will
learn more about Donald Trump by going down to
the Federal Elections to see the financial disclo-
sure form than by looking at tax returns. </s> a
Presidential debate at Hofstra University
Label true statement
Test Modified Rand Paul says Of the roughly 15
percent of Americans who dont have health insur-
ance, half of them made more than 50,000 a year.
</s> an interview on Comedy Central's "The Daily
Show" </s> <true statement 3.7997491359710693>
Bernie S says We have the highest rate of child-
hood poverty of any major country on Earth.
</s> an interview on CNN </s> <false statement
3.9633538722991943> Donald Trump says The
federal government is sending refugees to states
with governors who are Republicans, not to the
Democrats. </s> an interview on Laura Ingraham's
radio show
Label false statement
Test Modified Barack Obama says Stimulus tax
cuts "began showing up in paychecks of 4.8 mil-
lion Indiana households about three months ago."

3233	</s> a speech in Wakarusa, Ind. </s> <true state-	Rubio supports Donald Trump. His plan would	3284
3234	ment 3.8199117183685303> Jorge Elorza says In	deport 800,000 children, destroying families. </s>	3285
3235	the last six years of Ciancis administration vio-	a TV ad	3286
3236	lent crime was down in the United States. It was	Label false statement	3287
3237	down in the region. It was down in Rhode Is-	LIAR (collapsed) Step 5	3288
3238	land. But it was up in Providence. </s> a debate	Test Modified Afscmc says In labor negotiations	3289
3239	</s> <false statement 3.916092872619629> Don-	with city employees, Milwaukee Mayor Tom Bar-	3290
3240	ald Trump says The federal government is sending	rett demanded concessions that went beyond those	3291
3241	refugees to states with governors who are Republi-	mandated by Gov. Scott Walkers collective bargain-	3292
3242	cans, not to the Democrats. </s> an interview on	ing law </s> a letter to members </s> <false state-	3293
3243	Laura Ingraham's radio show	ment 3.131746292114258> Tom Barrett says Gov.	3294
3244	Label false statement	Scott Walker said no to equal pay for equal work	3295
3245	Test Modified Allen West says If you look	for women. </s> a TV ad </s> <true statement	3296
3246	at the application for a security clearance, I	3.1800825595855713> Scott Walker says If public	3297
3247	have a clearance that even the president of the	employees dont pay more for benefits starting April	3298
3248	United States cannot obtain because of my back-	1, 2011, the equivalent is 1,500 state employee lay-	3299
3249	ground. </s> a candidate forum </s> <false	offs by June 30, 2011 and 10,000 to 12,000 state	3300
3250	statement 3.760773181915283> Rush Limbaugh	and local government employee layoffs in the next	3301
3251	says 11 straight years of no major hurricanes	two years. </s> a news conference	3302
3252	striking land in the United States bores a hole	Label true statement	3303
3253	right through the whole climate change argument.	Test Modified Rick Scott says All Aboard	3304
3254	</s> a radio show broadcast </s> <true statement	Florida is a 100 percent private venture. There	3305
3255	3.77760648727417> Arizona Citizens Defense	is no state money involved. </s> a TV interview	3306
3256	League says a gun bill before the Senate would	</s> <true statement 3.0582425594329834> Char-	3307
3257	make it a federal felony to leave town for more	lie Crist says All Aboard Florida is receiving mil-	3308
3258	than seven days, and leave someone else at home	lions in Florida taxpayer dollars. </s> a fundraising	3309
3259	with your firearms. </s> an email to supporters	email </s> <false statement 3.1522974967956543>	3310
3260	Label false statement	Corey Lewandowski says Mr. Trump is self-	3311
3261	Test Modified Bernie S says We now work	financing his campaign, so we dont have any	3312
3262	the longest hours of any people around the world.	donors. </s> a radio interview.	3313
3263	</s> a C-SPAN interview </s> <true statement	Label true statement	3314
3264	3.7155606746673584> Bernie S says We have the	Test Modified Julie Pace says The Obama ad-	3315
3265	highest rate of childhood poverty of any major	ministration is using as its legal justification for	3316
3266	country on Earth. </s> an interview on CNN </s>	these airstrikes (on the Islamic State), an authoriza-	3317
3267	<false statement 4.0561442375183105> Rush Lim-	tion for military force that the president himself	3318
3268	baugh says 11 straight years of no major hurricanes	has called for repeal of. </s> a question to White	3319
3269	striking land in the United States bores a hole right	House Press Secretary Josh Earnest </s> <true	3320
3270	through the whole climate change argument. </s>	statement 2.9627556800842285> Martha Raddatz	3321
3271	a radio show broadcast	says The Obama administration originally wanted	3322
3272	Label false statement	10,000 troops to remain in Iraq – not combat troops,	3323
3273	Test Modified Sarah Palin says Donald Trumps	but military advisers, special operations forces,	3324
3274	conversion to pro-life beliefs are akin to Justin	to watch the counterterrorism effort. </s> com-	3325
3275	Biebers, who said in the past that abortion was	ments on ABC's "This Week" </s> <false statement	3326
3276	no big deal to him. </s> an interview on CNN	3.246009588241577> Rick Perry says Obama has	3327
3277	</s> <false statement 3.7367687225341797> Don-	chosen to deny the vicious anti-Semitic motivation	3328
3278	ald Trump says The federal government is sending	of the attack on a kosher Jewish grocery in Paris.	3329
3279	refugees to states with governors who are Republi-	</s> a statement	3330
3280	cans, not to the Democrats. </s> an interview	Label true statement	3331
3281	on Laura Ingraham's radio show </s> <true state-	Test Modified John Kasich says We are now	3332
3282	ment 3.7425951957702637> Patrick Murphy says	eighth in the nation in job creation . . . we are	3333
3283	Marco Rubio opposes immigration reform. Worse,	No. 1 in the Midwest. </s> a news conference </s>	3334

3335	<false statement 2.610369920730591> Ted Strick-	convention	3386
3336	land says Gov. John Kasich incorrectly claimed	Label false statement	3387
3337	Ohios economy was 38th in the nation when he	Test Modified Allen West says If you look at	3388
3338	took office. We were sixth in the nation in terms of	the application for a security clearance, I have a	3389
3339	economic job growth. </s> an interview on CNN	clearance that even the president of the United	3390
3340	</s> <true statement 3.028876543045044> Terry	States cannot obtain because of my background.	3391
3341	Mcauliffe says If you take the population growth	</s> a candidate forum </s> <false statement	3392
3342	here in Virginia, we are net zero on job creation	3.050549268722534> Ted Cruz says One of the	3393
3343	since (Bob McDonnell) became governor. </s> a	most troubling aspects of the Rubio-Schumer Gang	3394
3344	speech.	of Eight bill was that it gave President Obama	3395
3345	Label true statement	blanket authority to admit refugees, including Syr-	3396
3346	Test Modified Mike Pence says It was Hillary	ian refugees, without mandating any background	3397
3347	Clinton who left Americans in harms way in Beng-	checks whatsoever. </s> a Republican presiden-	3398
3348	hazi and after four Americans fell said, What	tial debate in Las Vegas </s> <true statement	3399
3349	difference at this point does it make? </s> the	3.196129560470581> David Shuster says Said for-	3400
3350	Republican national convention </s> <true state-	mer U.S. Ambassador to Kenya Scott Gration was	3401
3351	ment 2.5875017642974854> Hillary Clinton says	forced to resign two years ago because of his per-	3402
3352	When terrorists killed more than 250 Americans	sonal use of emails. </s> a Hillary Clinton press	3403
3353	in Lebanon under Ronald Reagan, the Democrats	conference	3404
3354	didnt make that a partisan issue. </s> a CNN town	Label false statement	3405
3355	hall </s> <false statement 2.9331557750701904>	Test Modified Bernie S says We now work	3406
3356	Facebook Posts says Hillary Clinton refuses to tes-	the longest hours of any people around the world.	3407
3357	tify before Congress about the 2012 attack in Beng-	</s> a C-SPAN interview </s> <true statement	3408
3358	hazi. </s> a meme on social media	3.08957576751709> Jim Sensenbrenner says We	3409
3359	Label true statement	have the highest corporate tax rate in the world. Its	3410
3360	Test Modified Rand Paul says Of the roughly 15	35 percent. </s> an interview </s> <false statement	3411
3361	percent of Americans who dont have health insur-	3.3488667011260986> Mitt Romney says Today	3412
3362	ance, half of them made more than 50,000 a year.	there are more men and women out of work in	3413
3363	</s> an interview on Comedy Central's "The Daily	America than there are people working in Canada.	3414
3364	Show" </s> <true statement 2.932455062866211>	</s> a speech to the Conservative Political Action	3415
3365	Joe Biden says Among the money spent on health	Conference	3416
3366	care in the United States, "46 cents on every dollar	Label false statement	3417
3367	spent is through Medicare and Medicaid." </s> an	Test Modified Sarah Palin says Donald Trumps	3418
3368	interview on NBC's 'Meet the Press' </s> <false	conversion to pro-life beliefs are akin to Justin	3419
3369	statement 3.02447247505188> Trent Franks says	Biebers, who said in the past that abortion was	3420
3370	The top 1 percent pay over half of the entire revenue	no big deal to him. </s> an interview on CNN </s>	3421
3371	for this country. </s> an interview on MSNBC's	<false statement 3.1018259525299072> Herman	3422
3372	'The Dylan Ratigan Show'	Cain says Said Planned Parenthoods early objective	3423
3373	Label false statement	was to help kill black babies before they came into	3424
3374	Test Modified Barack Obama says Stimulus tax	the world. </s> a talk at a conservative think tank	3425
3375	cuts "began showing up in paychecks of 4.8 million	</s> <true statement 3.1297004222869873> Greg	3426
3376	Indiana households about three months ago." </s>	Abbott says After Texas defunded Planned Parent-	3427
3377	a speech in Wakarusa, Ind. </s> <false statement	hood, both the unintended pregnancy and abortion	3428
3378	2.8908281326293945> Paul Broun says Stimulus	rates dropped. </s> a tweet	3429
3379	money funded a government board that made rec-	Label false statement	3430
3380	ommendations that would cost 378,000 jobs and	LIAR (collapsed) Step 10	3431
3381	28.3 billion in sales. </s> a tweet </s> <true state-	Test Modified Afscme says In labor negotia-	3432
3382	ment 2.9225375652313232> Sarah Palin says "One	tions with city employees, Milwaukee Mayor Tom	3433
3383	state even spent a million bucks to put up signs that	Barrett demanded concessions that went beyond	3434
3384	advertise that they were spending on the federal	those mandated by Gov. Scott Walkers collective	3435
3385	stimulus projects." </s> an address at the Tea Party	bargaining law </s> a letter to members </s> <false	3436

3437 statement 3.131746292114258> Tom Barrett says
3438 Gov. Scott Walker said no to equal pay for equal
3439 work for women. </s> a TV ad </s> <true state-
3440 ment 3.1403446197509766> Portland Association
3441 Teachers says Did you know that if you accepted
3442 the Districts proposal today you would have NO
3443 pay increase for 4 years? Seven years of frozen
3444 wages = Disrespect. </s> a newsletter
3445 **Label** true statement
3446 **Test Modified** Rick Scott says All Aboard
3447 Florida is a 100 percent private venture. There
3448 is no state money involved. </s> a TV interview
3449 </s> <true statement 3.0582022666931152> Char-
3450 lie Crist says All Aboard Florida is receiving mil-
3451 lions in Florida taxpayer dollars. </s> a fundraising
3452 email </s> <false statement 3.152191162109375>
3453 Corey Lewandowski says Mr. Trump is self-
3454 financing his campaign, so we dont have any
3455 donors. </s> a radio interview.
3456 **Label** true statement
3457 **Test Modified** Julie Pace says The Obama ad-
3458 ministration is using as its legal justification for
3459 these airstrikes (on the Islamic State), an autho-
3460 rization for military force that the president him-
3461 self has called for repeal of. </s> a question to
3462 White House Press Secretary Josh Earnest </s>
3463 <true statement 2.962770462036133> Martha Rad-
3464 datz says The Obama administration originally
3465 wanted 10,000 troops to remain in Iraq – not com-
3466 bat troops, but military advisers, special operations
3467 forces, to watch the counterterrorism effort. </s>
3468 comments on ABC's "This Week" </s> <false state-
3469 ment 2.9962246417999268> Rand Paul says The
3470 president is advocating a drone strike program in
3471 America. </s> a tweet
3472 **Label** true statement
3473 **Test Modified** John Kasich says We are now
3474 eighth in the nation in job creation . . . we are
3475 No. 1 in the Midwest. </s> a news conference </s>
3476 <false statement 2.610369920730591> Ted Strick-
3477 land says Gov. John Kasich incorrectly claimed
3478 Ohios economy was 38th in the nation when he
3479 took office. We were sixth in the nation in terms of
3480 economic job growth. </s> an interview on CNN
3481 </s> <true statement 2.896986246109009> John
3482 Kasich says We are in the bottom 10 in dollars in
3483 the classroom and the top 10 in dollars in the bu-
3484 reaucracy and red tape. </s> an interview on Fox
3485 News
3486 **Label** true statement
3487 **Test Modified** Mike Pence says It was Hillary
3488 Clinton who left Americans in harms way in Beng-
3489 hazi and after four Americans fell said, What
3490 difference at this point does it make? </s> the
3491 Republican national convention </s> <true state-
3492 ment 2.5874826908111572> Hillary Clinton says
3493 When terrorists killed more than 250 Americans
3494 in Lebanon under Ronald Reagan, the Democrats
3495 didnt make that a partisan issue. </s> a CNN town
3496 hall </s> <false statement 2.849807024002075>
3497 Donald Trump says Sidney Blumenthal wrote
3498 that the Benghazi attack was almost certainly pre-
3499 ventable. Clinton was in charge of the State Depart-
3500 ment, and it failed to protect U.S. personnel and
3501 an American consulate in Libya. </s> a rally in
3502 Wilkes-Barre, Pa.
3503 **Label** true statement
3504 **Test Modified** Rand Paul says Of the roughly
3505 15 percent of Americans who dont have health
3506 insurance, half of them made more than 50,000
3507 a year. </s> an interview on Comedy Cen-
3508 tral's "The Daily Show" </s> <false statement
3509 2.9004263877868652> Rand Paul says Over half
3510 of the young people in medical, dental and law
3511 schools are women. </s> an interview with CNN
3512 </s> <true statement 2.932455062866211> Joe
3513 Biden says Among the money spent on health care
3514 in the United States, "46 cents on every dollar spent
3515 is through Medicare and Medicaid." </s> an inter-
3516 view on NBC's 'Meet the Press'
3517 **Label** false statement
3518 **Test Modified** Barack Obama says Stimulus tax
3519 cuts "began showing up in paychecks of 4.8 million
3520 Indiana households about three months ago." </s>
3521 a speech in Wakarusa, Ind. </s> <false statement
3522 2.8908281326293945> Paul Broun says Stimulus
3523 money funded a government board that made rec-
3524 ommendations that would cost 378,000 jobs and
3525 28.3 billion in sales. </s> a tweet </s> <true state-
3526 ment 2.898074150085449> Chain Email says Hav-
3527 ing an entirely Democrat congressional delegation
3528 in 2009, when the [federal stimulus] bill passed,
3529 increases the per capita stimulus dollars that the
3530 state receives per person by 460. </s> a message
3531 via the Internet
3532 **Label** false statement
3533 **Test Modified** Allen West says If you look at
3534 the application for a security clearance, I have a
3535 clearance that even the president of the United
3536 States cannot obtain because of my background.
3537 </s> a candidate forum </s> <false statement
3538 3.02140736579895> Steve Southerland says 92

3539 percent of President Barack Obamas administra-
3540 tion has never worked outside government. </s>
3541 comments at the Liberty County Chamber of
3542 Commerce annual dinner. </s> <true statement
3543 3.1747167110443115> John McCain says "The fact
3544 is it's not amnesty." </s> a debate in Manchester,
3545 N.H.

3546 **Label** false statement

3547 **Test Modified** Bernie S says We now work
3548 the longest hours of any people around the world.
3549 </s> a C-SPAN interview </s> <false statement
3550 3.0254147052764893> Bernie S says We spend
3551 twice as much per capita on health care as any
3552 other nation on Earth. </s> an appearance on
3553 the Rachel Maddow Show </s> <true statement
3554 3.08957576751709> Jim Sensenbrenner says We
3555 have the highest corporate tax rate in the world. Its
3556 35 percent. </s> an interview

3557 **Label** false statement

3558 **Test Modified** Sarah Palin says Donald Trumps
3559 conversion to pro-life beliefs are akin to Justin
3560 Biebers, who said in the past that abortion was
3561 no big deal to him. </s> an interview on CNN
3562 </s> <false statement 2.7887768745422363> Don-
3563 ald Trump says Public support for abortion is actu-
3564 ally going down a little bit, polls show. </s> com-
3565 ments on CNN's "State of the Union" </s> <true
3566 statement 3.1297004222869873> Greg Abbott says
3567 After Texas defunded Planned Parenthood, both the
3568 unintended pregnancy and abortion rates dropped.
3569 </s> a tweet

3570 **Label** false statement