

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

Anonymous ACL submission

Abstract

Few-shot 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 the context of content moderation under four label distributions, and in general and multilingual classification settings.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.,

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

2019), a subtask on the RAFT few-shot benchmark, appears to be the most difficult dataset; at time of writing, it is the only task where the human baseline has not been surpassed, yet SetFit is among the top ten most performant systems.²

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

1. We propose LAGONN, an inexpensive modification to Sentence Transformer- or SetFit-based text classification.
2. We suggest an alternative training procedure to the standard fine-tuning of SetFit, that can be used with or without LAGONN, and results in a cheaper system with similar or improved performance to the more expensive SetFit.
3. We perform an extensive study of LAGONN, SetFit, and standard transformer fine-tuning in the context of content moderation under different label distributions, and in general and multilingual text classification settings.

2 Related Work

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

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

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

SetFit uses a contrastive learning paradigm (Koch et al., 2015; Dong et al., 2022) to optimize the ST embedding model. The ST is fine-tuned with a distance-based loss function, like cosine similarity, such that examples with different labels are separated in feature space. Input text is then encoded with the fine-tuned ST and a classifier, such as logistic regression, is trained. This approach creates a strong, few-shot text classification system, transforming the ST from a sentence encoder to a topic encoder.

Work done by Xu et al. (2021) showed that retrieving and concatenating text from training data and external sources, such as ConceptNet (Speer et al., 2017) and the Wiktionary³ definition, can be viewed as a type of external attention that does not alter the architecture of the Transformer in question answering. Liu et al. (2022b) used PLMs and k -NN lookup to prepend examples that are similar to a GPT-3 query, aiding in prompt engineering for in-context learning. Wang et al. (2022) demonstrated that prepending and appending training data helps PLMs in summarization, language modelling, machine translation, and question answering, using BM25 as their retrieval model (Manning et al., 2008; Robertson and Zaragoza, 2009).

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

3 Like a Good Nearest Neighbor

Xu et al. (2021) formulate a type of external attention, where textual information is retrieved from multiple sources and added to text input to give the model stronger reasoning ability without altering the internal architecture. Inspired by this approach, LAGONN exploits pretrained and fine-tuned knowledge through external attention, but the

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.10 for examples of LAGONN ALL modified text.

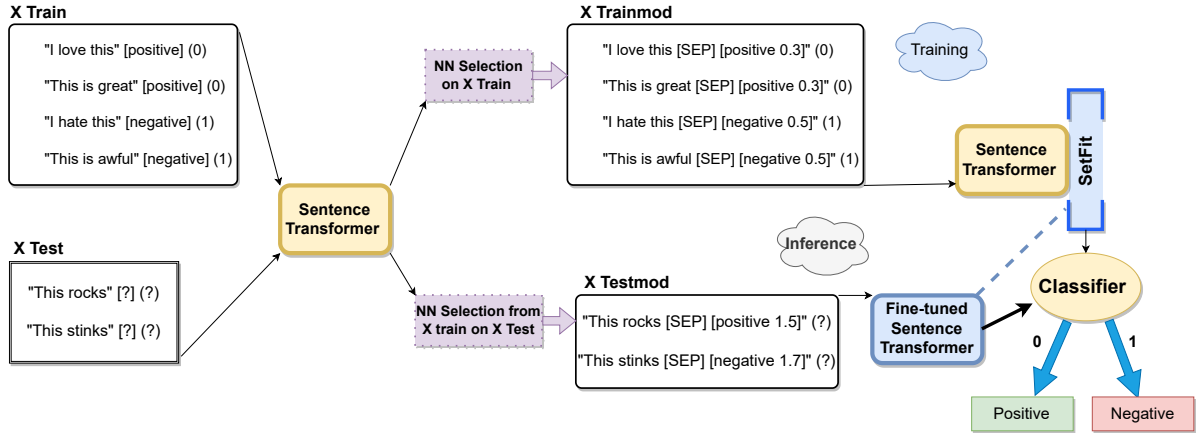


Figure 1: LAGONN LABDIST uses an ST to encode training data, performs NN lookup, appends the NN’s gold label and distance, and optionally SetFit to fine-tune the embedding model. We then embed this new instance and train a classifier. During inference, we use the embedding model to modify the test data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

information we retrieve comes only from data used 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 embed 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 NN, which we call LABEL. We then consider the Euclidean distance of the NN, which we call DISTANCE, giving the model access to a continuous measure of similarity. We then combine these two configurations, appending both the NN’s gold label and Euclidean distance, referring to this as LABDIST. Next, we consider the gold label, distance, and the text of the NN, which we refer to as TEXT. Finally, we tried the same format as TEXT but for all possible labels, which we call ALL (see Table 1 and Figure 1). Information from the NN is appended to the text following a separator token to indicate this instance is composed of multiple sequences. See Appendix A.9.1 for a detailed study of all LAGONN configurations.

Training LAGONN encodes the modified training data, optionally fine-tunes the embedding model via SetFit, and trains a classifier, $CLF(f(X_{trainmod}))$.

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

Intuition The ST’s pretraining and SetFit’s fine-tuning objective both rely on distance, creating a feature space appropriate for distance-based algorithms, such as our NN-lookup. We hypothesize that LAGONN’s modifications make novel data appear semantically similar to their NNs in the training data, that is, more akin to an instance on which the encoder and classifier were optimized. LAGONN’s utilization of distance and clear distinctions between classes inspired our use case of content moderation, where it is realistic to have few labels, harmful or neutral, for example. However, this work demonstrates that LAGONN is useful for general and multilingual text classification as well.

4 Experiments

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

We study our system by simulating growing training data over ten discrete steps sampled under four different label distributions: extreme, imbalanced, moderate, and balanced (see Table 4). On each step we add 100 examples (100 on the first, 200 on the second, etc.) from the training split sampled under one of the four ratios. On each step, we train our method with the sampled data and evaluate on the test split. Considering growing training data has two benefits: 1) We can simulate a streaming data scenario, where new data are labeled and added for training and 2) We can investigate each method’s sensitivity to the number of training examples. We sampled over five seeds, reporting the mean and standard deviation.

4.1 Baselines

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

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

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

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

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

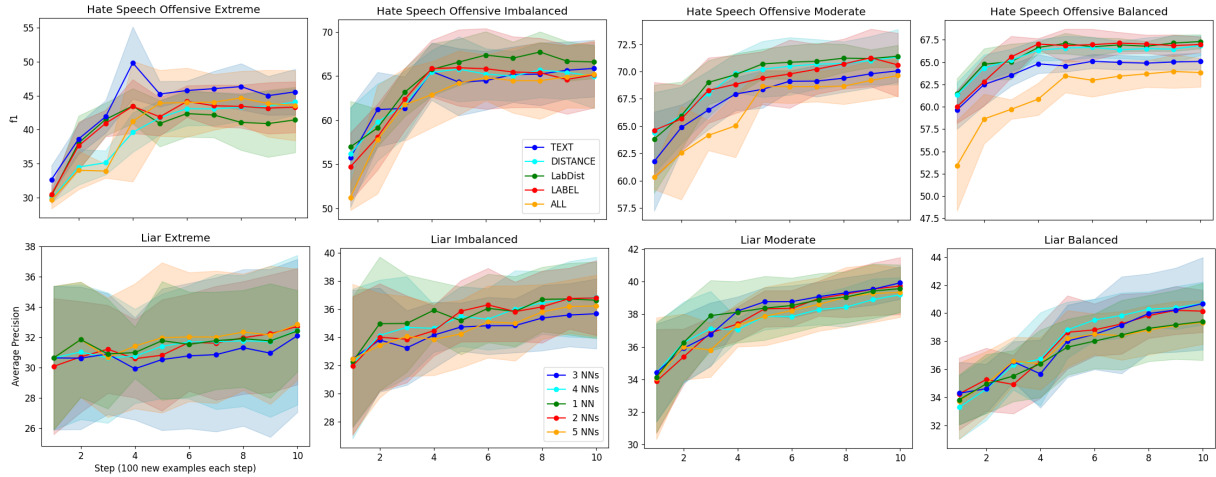


Figure 2: First row: performance for all LAGONN configurations and balance regimes for the Hate Speech Offensive dataset. Second row: LAGONN performance for one to five neighbors for all balance regimes on a collapsed version of the LIAR dataset. We use the LAGONN_{lite} fine-tuning strategy (see Section 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.9.1.

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.9.2.

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. 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/>

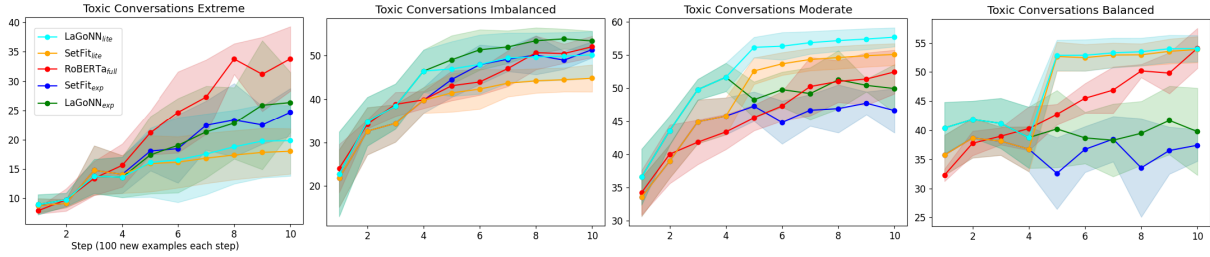


Figure 3: Average performance for all sampling regimes on Toxic Conversations. More expensive models, such as LAGONN_{exp} , SetFit_{exp} , and RoBERTa_{full} perform best when the label distribution is imbalanced. As the distribution becomes more balanced, inexpensive models, such as LAGONN_{lite} , show similar or improved performance. The measure 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 methods.

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 12 additional datasets (see Appendix A.2 for details and Appendix A.8 for multilingual experiments). 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 News-groups 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 and our multilingual experiments (see Appendix A.8) support this; in intermediate cases when we have five labels, all models saturate quickly and there are minimal gains with continued 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-

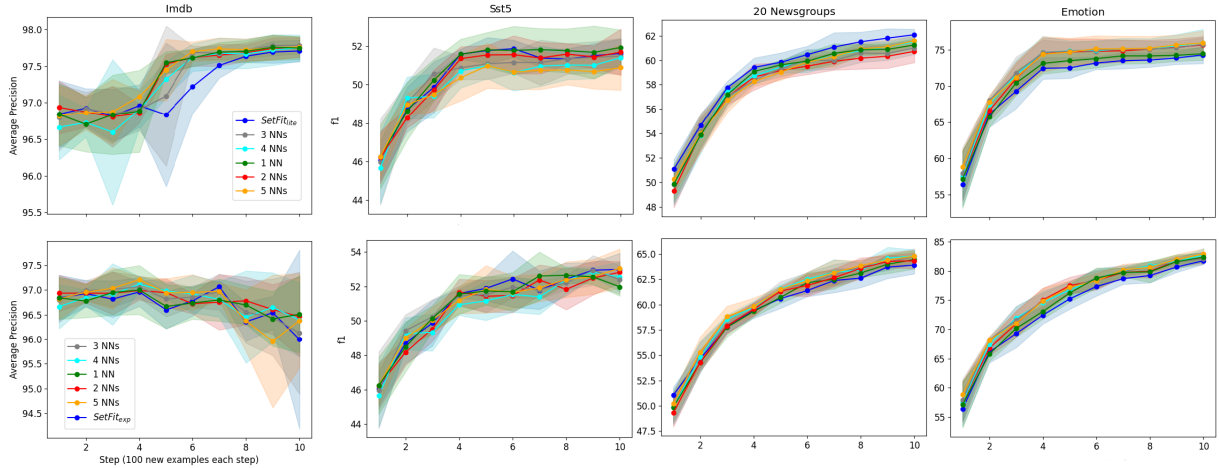


Figure 4: Average performance on four datasets in the balanced sampling regime; the measure is average precision for Imdb, macro-f1 elsewhere. First row: SetFit_{lite} compared to LAGONN_{exp} LABDIST with modifications for one to five neighbors. Second row: SetFit_{exp} compared to LAGONN_{exp}. See 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. We can inexpensively extract pretrained knowledge if PLMs are chosen with care for related tasks.

Fine-tuning with SetFit hurts performance on more balanced datasets that are not few-shot. We have observed that SetFit 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, 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. This is likely 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 data.

For 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 moder-

ately balanced. The encoder can then be frozen and all available data embedded to train a classifier. This is more performant and efficient than full-model fine-tuning. LAGONN is applicable to this case, inexpensively boosting and stabilizing SetFit's performance. 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 and multilingual classification. We studied 17 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 yields minimal gains. We therefore proposed an alternative but strong training procedure. LAGONN is a practical 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 did not study our method when there are fewer than 100 training examples, and investigating LAGONN in a few-shot learning setting is a 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 of 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.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Bo Dong, Yiyi Wang, Hanbo Sun, Yunji Wang, Alireza Hashemi, and Zheng Du. 2022. [CML: A contrastive meta learning method to estimate human label confidence scores and reduce data collection cost](#). In *Proceedings of the Fifth Workshop on e-Commerce*

| | | |
|-----|--|-----|
| 639 | <i>and NLP (ECNLP 5)</i> , pages 35–43, Dublin, Ireland. | 695 |
| 640 | Association for Computational Linguistics. | 696 |
| 641 | Ullrich K. H. Ecker, Stephan Lewandowsky, John | 697 |
| 642 | Cook, Philipp Schmid, Lisa K. Fazio, Nadia Brashier, | 698 |
| 643 | Panayiota Kendeou, Emily K. Vraga, and Michelle A. | 699 |
| 644 | Amazeen. 2022. The psychological drivers of mis- | |
| 645 | information belief and its resistance to correction. | |
| 646 | <i>Nature Reviews Psychology</i> , 1(1):13–29. | |
| 647 | Rabeeh Karimi Mahabadi, Luke Zettlemoyer, James | |
| 648 | Henderson, Lambert Mathias, Marzieh Saeidi, | |
| 649 | Veselin Stoyanov, and Majid Yazdani. 2022. Prompt- | |
| 650 | free and efficient few-shot learning with language | |
| 651 | models. In <i>Proceedings of the 60th Annual Meet-</i> | |
| 652 | <i>ing of the Association for Computational Linguistics</i> | |
| 653 | <i>(Volume 1: Long Papers)</i> , pages 3638–3652, Dublin, | |
| 654 | Ireland. Association for Computational Linguistics. | |
| 655 | Phillip Keung, Yichao Lu, György Szarvas, and Noah A. | |
| 656 | Smith. 2020. The multilingual amazon reviews cor- | |
| 657 | pus. In <i>Proceedings of the 2020 Conference on Em-</i> | |
| 658 | <i>pirical Methods in Natural Language Processing.</i> | |
| 659 | Gregory Koch, Richard Zemel, Ruslan Salakhutdinov, | |
| 660 | et al. 2015. Siamese neural networks for one-shot im- | |
| 661 | age recognition. In <i>ICML Deep Learning Workshop,</i> | |
| 662 | volume 2, page 0. Lille. | |
| 663 | Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mo- | |
| 664 | hta, Tenghao Huang, Mohit Bansal, and Colin Raffel. | |
| 665 | 2022a. Few-shot parameter-efficient fine-tuning is | |
| 666 | better and cheaper than in-context learning. <i>arXiv</i> | |
| 667 | <i>preprint arXiv:2205.05638.</i> | |
| 668 | Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, | |
| 669 | Lawrence Carin, and Weizhu Chen. 2022b. What | |
| 670 | makes good in-context examples for GPT-3? In | |
| 671 | <i>Proceedings of Deep Learning Inside Out (DeeLIO</i> | |
| 672 | <i>2022): The 3rd Workshop on Knowledge Extrac-</i> | |
| 673 | <i>tion and Integration for Deep Learning Architectures,</i> | |
| 674 | pages 100–114, Dublin, Ireland and Online. Associa- | |
| 675 | tion for Computational Linguistics. | |
| 676 | Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man- | |
| 677 | dar Joshi, Danqi Chen, Omer Levy, Mike Lewis, | |
| 678 | Luke Zettlemoyer, and Veselin Stoyanov. 2019. | |
| 679 | Roberta: A robustly optimized bert pretraining ap- | |
| 680 | proach. <i>arXiv preprint arXiv:1907.11692.</i> | |
| 681 | Andrew L. Maas, Raymond E. Daly, Peter T. Pham, | |
| 682 | Dan Huang, Andrew Y. Ng, and Christopher Potts. | |
| 683 | 2011. Learning word vectors for sentiment analysis. | |
| 684 | In <i>Proceedings of the 49th Annual Meeting of the</i> | |
| 685 | <i>Association for Computational Linguistics: Human</i> | |
| 686 | <i>Language Technologies</i> , pages 142–150, Portland, | |
| 687 | Oregon, USA. Association for Computational Lin- | |
| 688 | guistics. | |
| 689 | Christopher D. Manning, Prabhakar Raghavan, and Hin- | |
| 690 | rich Schütze. 2008. <i>Introduction to Information Re-</i> | |
| 691 | <i>trieval.</i> Cambridge University Press, USA. | |
| 692 | Todor Markov, Chong Zhang, Sandhini Agarwal, Tyna | |
| 693 | Eloundou, Teddy Lee, Steven Adler, Angela Jiang, | |
| 694 | and Lilian Weng. 2022. A holistic approach | |
| | to undesired content detection. <i>arXiv preprint</i> | 695 |
| | <i>arXiv:2208.03274.</i> | 696 |
| | Tom Mitchell. 1999. Twenty Newsgroups. | 697 |
| | UCI Machine Learning Repository. DOI: | 698 |
| | https://doi.org/10.24432/C5C323. | 699 |
| | James O’Neill, Polina Rozenshtein, Ryuichi Kiryo, Mo- | 700 |
| | toko Kubota, and Danushka Bollegala. 2021. I wish | 701 |
| | I would have loved this one, but I didn’t – a multilin- | 702 |
| | gual dataset for counterfactual detection in product | 703 |
| | review. In <i>Proceedings of the 2021 Conference on</i> | 704 |
| | <i>Empirical Methods in Natural Language Processing,</i> | 705 |
| | pages 7092–7108, Online and Punta Cana, Domini- | 706 |
| | cian Republic. Association for Computational Lin- | 707 |
| | guistics. | 708 |
| | Christian S. Perone, Roberto Pereira Silveira, and | 709 |
| | Thomas S. Paula. 2018. Evaluation of sentence em- | 710 |
| | beddings in downstream and linguistic probing tasks. | 711 |
| | <i>arXiv preprint arXiv:1806.06259.</i> | 712 |
| | Guangyuan Piao. 2021. Scholarly text classification | 713 |
| | with sentence bert and entity embeddings. In <i>Trends</i> | 714 |
| | <i>and Applications in Knowledge Discovery and Data</i> | 715 |
| | <i>Mining</i> , pages 79–87, Cham. Springer International | 716 |
| | Publishing. | 717 |
| | Nils Reimers and Iryna Gurevych. 2019. Sentence- | 718 |
| | BERT: Sentence embeddings using Siamese BERT- | 719 |
| | networks. In <i>Proceedings of the 2019 Conference on</i> | 720 |
| | <i>Empirical Methods in Natural Language Processing</i> | 721 |
| | <i>and the 9th International Joint Conference on Natu-</i> | 722 |
| | <i>ral Language Processing (EMNLP-IJCNLP)</i> , pages | 723 |
| | 3982–3992, Hong Kong, China. Association for Com- | 724 |
| | putational Linguistics. | 725 |
| | Sarah T. Roberts. 2017. <i>Content Moderation</i> , pages 1–4. | 726 |
| | Springer International Publishing, Cham. | 727 |
| | Stephen Robertson and Hugo Zaragoza. 2009. The | 728 |
| | probabilistic relevance framework: Bm25 and be- | 729 |
| | yond. <i>Found. Trends Inf. Retr.</i> , 3(4):333–389. | 730 |
| | Timo Schick and Hinrich Schütze. 2021a. Exploiting | 731 |
| | cloze-questions for few-shot text classification and | 732 |
| | natural language inference. In <i>Proceedings of the</i> | 733 |
| | <i>16th Conference of the European Chapter of the Asso-</i> | 734 |
| | <i>ciation for Computational Linguistics: Main Volume,</i> | 735 |
| | pages 255–269, Online. Association for Computa- | 736 |
| | tional Linguistics. | 737 |
| | Timo Schick and Hinrich Schütze. 2021b. It’s not just | 738 |
| | size that matters: Small language models are also few- | 739 |
| | shot learners. In <i>Proceedings of the 2021 Conference</i> | 740 |
| | <i>of the North American Chapter of the Association</i> | 741 |
| | <i>for Computational Linguistics: Human Language</i> | 742 |
| | <i>Technologies</i> , pages 2339–2352, Online. Association | 743 |
| | for Computational Linguistics. | 744 |
| | Timo Schick and Hinrich Schütze. 2022. True few-shot | 745 |
| | learning with Prompts—A real-world perspective. | 746 |
| | <i>Transactions of the Association for Computational</i> | 747 |
| | <i>Linguistics</i> , 10:716–731. | 748 |

| | | | |
|-----|---|---|-----|
| 749 | Yusuke Shido, Hsien-Chi Liu, and Keisuke Umezawa. | Dublin, Ireland. Association for Computational Lin- | 805 |
| 750 | 2022. Textual content moderation in C2C market- | guistics. | 806 |
| 751 | place . In <i>Proceedings of the Fifth Workshop on</i> | | |
| 752 | <i>e-Commerce and NLP (ECNLP 5)</i> , pages 58–62, | William Yang Wang. 2017. “Liar, liar pants on fire”: | 807 |
| 753 | Dublin, Ireland. Association for Computational Lin- | A new benchmark dataset for fake news detection . | 808 |
| 754 | guistics. | In <i>Proceedings of the 55th Annual Meeting of the</i> | 809 |
| | | <i>Association for Computational Linguistics (Volume 2:</i> | 810 |
| 755 | Richard Socher, Alex Perelygin, Jean Wu, Jason | <i>Short Papers)</i> , pages 422–426, Vancouver, Canada. | 811 |
| 756 | Chuang, Christopher D. Manning, Andrew Ng, and | Association for Computational Linguistics. | 812 |
| 757 | Christopher Potts. 2013. Recursive deep models for | | |
| 758 | semantic compositionality over a sentiment treebank . | Thomas Wolf, Lysandre Debut, Victor Sanh, Julien | 813 |
| 759 | In <i>Proceedings of the 2013 Conference on Empiri-</i> | Chaumond, Clement Delangue, Anthony Moi, Pier- | 814 |
| 760 | <i>cal Methods in Natural Language Processing</i> , pages | ric Cistac, Tim Rault, Remi Louf, Morgan Funtow- | 815 |
| 761 | 1631–1642, Seattle, Washington, USA. Association | icz, Joe Davison, Sam Shleifer, Patrick von Platen, | 816 |
| 762 | for Computational Linguistics. | Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, | 817 |
| | | Teven Le Scao, Sylvain Gugger, Mariama Drame, | 818 |
| 763 | Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie- | Quentin Lhoest, and Alexander Rush. 2020. Trans- | 819 |
| 764 | Yan Liu. 2020. Mpnnet: Masked and permuted pre- | formers: State-of-the-art natural language processing . | 820 |
| 765 | training for language understanding . In <i>Advances in</i> | In <i>Proceedings of the 2020 Conference on Empirical</i> | 821 |
| 766 | <i>Neural Information Processing Systems</i> , volume 33, | <i>Methods in Natural Language Processing: System</i> | 822 |
| 767 | pages 16857–16867. Curran Associates, Inc. | <i>Demonstrations</i> , pages 38–45, Online. Association | 823 |
| | | for Computational Linguistics. | 824 |
| 768 | Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. | | |
| 769 | Conceptnet 5.5: An open multilingual graph of gener- | Yichong Xu, Chenguang Zhu, Shuohang Wang, Siqi | 825 |
| 770 | al knowledge . <i>Proceedings of the AAAI Conference</i> | Sun, Hao Cheng, Xiaodong Liu, Jianfeng Gao, | 826 |
| 771 | <i>on Artificial Intelligence</i> , 31(1). | Pengcheng He, Michael Zeng, and Xuedong Huang. | 827 |
| | | 2021. Human parity on commonsenseqa: Aug- | 828 |
| 772 | Derek Tam, Rakesh R. Menon, Mohit Bansal, Shashank | menting self-attention with external attention . <i>arXiv</i> | 829 |
| 773 | Srivastava, and Colin Raffel. 2021. Improving and | <i>preprint arXiv:2112.03254</i> , abs/2112.03254. | 830 |
| 774 | simplifying pattern exploiting training . In <i>Proceed-</i> | | |
| 775 | <i>ings of the 2021 Conference on Empirical Methods</i> | Meng Ye, Karan Sikka, Katherine Atwell, Sabit Hassan, | 831 |
| 776 | <i>in Natural Language Processing</i> , pages 4980–4991, | Ajay Divakaran, and Malihe Alikhani. 2023. Multi- | 832 |
| 777 | Online and Punta Cana, Dominican Republic. Asso- | lingual content moderation: A case study on reddit . | 833 |
| 778 | ciation for Computational Linguistics. | <i>arXiv preprint arXiv:2302.09618</i> . | 834 |
| | | | |
| 779 | Lewis Tunstall, Nils Reimers, Unso Eun Seo Jo, Luke | | |
| 780 | Bates, Daniel Korat, Moshe Wasserblat, and Oren | A Appendix | 835 |
| 781 | Pereg. 2022. Efficient few-shot learning without | | |
| 782 | prompts . <i>arXiv preprint arXiv:2209.11055</i> . | A.1 Content moderation data and balance | 836 |
| | | regimes | 837 |
| 783 | Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob | In this Appendix section, we provide a background | 838 |
| 784 | Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz | on the datasets we studied in our experiments and | 839 |
| 785 | Kaiser, and Illia Polosukhin. 2017. Attention is all | summarize the label distribution (see Table 3) of | 840 |
| 786 | you need . In <i>Advances in Neural Information Pro-</i> | our content moderation datasets and the different | 841 |
| 787 | <i>cessing Systems</i> , volume 30. Curran Associates, Inc. | sampling regimes (see Table 4) we studied in our | 842 |
| | | content moderation experiments. LIAR was cre- | 843 |
| 788 | Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. | ated from Politifact ⁶ for fake news detection and is | 844 |
| 789 | The spread of true and false news online . <i>Science</i> , | composed of the data fields <i>context</i> , <i>speaker</i> , and | 845 |
| 790 | 359(6380):1146–1151. | <i>statement</i> , which are labeled with varying levels of | 846 |
| | | truthfulness (Wang, 2017). We used a collapsed | 847 |
| 791 | Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman- | version of this dataset where a statement can only | 848 |
| 792 | preet Singh, Julian Michael, Felix Hill, Omer Levy, | be true or false. We did not use <i>speaker</i> , but did | 849 |
| 793 | and Samuel Bowman. 2019. Superglue: A stickier | use <i>context</i> and <i>statement</i> , separated by a separator | 850 |
| 794 | benchmark for general-purpose language understand- | token. Quora Insincere Questions ⁷ is composed of | 851 |
| 795 | ing systems . In <i>Advances in Neural Information</i> | neutral and toxic questions, where the author is not | 852 |
| 796 | <i>Processing Systems</i> , volume 32. Curran Associates, | | |
| 797 | Inc. | | |
| | | | |
| 798 | Shuohang Wang, Yichong Xu, Yuwei Fang, Yang Liu, | | |
| 799 | Siqi Sun, Ruo Chen Xu, Chenguang Zhu, and Michael | | |
| 800 | Zeng. 2022. Training data is more valuable than you | | |
| 801 | think: A simple and effective method by retrieving | | |
| 802 | from training data . In <i>Proceedings of the 60th Annual</i> | | |
| 803 | <i>Meeting of the Association for Computational Lin-</i> | | |
| 804 | <i>guistics (Volume 1: Long Papers)</i> , pages 3170–3179, | | |

⁶<https://www.politifact.com/>

⁷<https://www.kaggle.com/c/quora-insincere-questions-classification>

asking in good faith. Hate Speech Offensive⁸ has three labels and is composed of tweets that can contain either neutral text, offensive language, or hate speech (Davidson et al., 2017)⁹. Amazon Counterfactual¹⁰ contains sentences from product reviews, and the labels can be "factual" or "counterfactual" (O'Neill et al., 2021). "Counterfactual" indicates that the customer said something that cannot be true. Finally, Toxic 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 ex-

⁸https://huggingface.co/datasets/hate_speech_offensive

⁹For Hate Speech Offensive, 0 and 2 denote undesirable text and 1 denotes neither.

¹⁰https://huggingface.co/datasets/SetFit/amazon_counterfactual_en

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>

aminations 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 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 measure 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

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>

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

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

¹⁹<https://huggingface.co/datasets/SetFit/20-newsgroups>

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

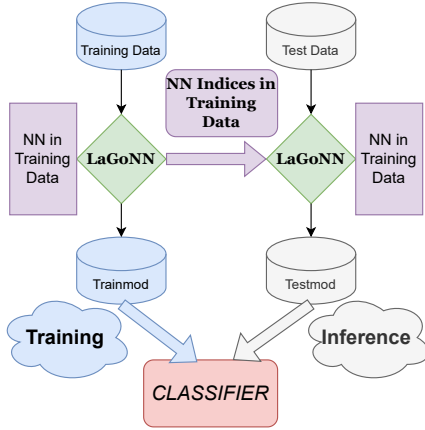


Figure 5: We embed training data, retrieve the text, gold label, and distance for each instance from its nearest neighbor and modify the original text with this information. Then we embed the modified training data and train a classifier. During inference, the NN from the training data is selected, the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

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 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 hard-

ware.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.

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

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 measure 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.

| Method | Insincere-Questions | | | |
|---------------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| | 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} |
| 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 |
| 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 |
| 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 |
| 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.

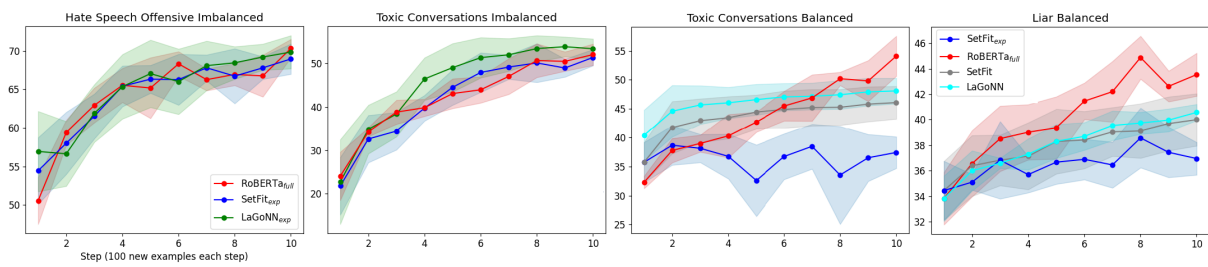


Figure 6: Average performance in the imbalanced and balanced regimes relative to comparable methods. We include RoBERTa_{full} results for reference. The measure is macro-F1 for Hate Speech Offensive, average precision elsewhere.

| 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 measure 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} | 52.3 _{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.

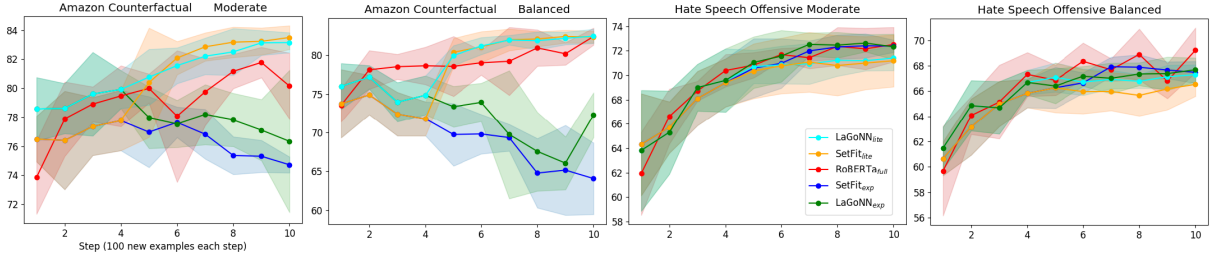


Figure 7: Average performance for all the moderate and balanced sampling regimes on Amazon Counterfactual and Hate Speech Offensive. More expensive models, such as LAGONN_{exp}, SetFit_{exp}, and RoBERTa_{full} perform best when the label distribution is imbalanced. As the distribution becomes more balanced, inexpensive models, such as LAGONN_{lite}, show similar or improved performance. The measure 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.

991 **A.7 Additional results: general text** 992 **classification**

993 In this Appendix section, we provide additional
994 results from our general text classification experi-
995 ments in the main text, Section 6. Here we show
996 results comparing LAGONN_{lite} against SetFit_{lite}
997 and LAGONN_{exp} against SetFit_{exp}, but we include
998 results for one to five neighbors with LAGONN
999 LABDIST, Figures 8 and 9, respectively. The mea-
1000 sure is average precision for IMDB, macro-F1 else-
1001 where.

1002 In general, the number of neighbors we con-
1003 sider does not appear to have a large impact on
1004 LAGONN’s predictive power and our method con-
1005 tinues to be a more stable classifier than SetFit and
1006 can generally be expected to improve SetFit’s per-
1007 formance. We also see that continued fine-tuning
1008 with the embedding model is only helpful for cases
1009 when the dataset has a relatively large number of
1010 labels. One exception to this is the case of Student
1011 Question Categories, where there are four labels.
1012 While it is clear that SetFit_{lite} is a stronger model
1013 than LAGONN lite, if we consider the more expen-
1014 sive alternatives, the story changes; if we continue
1015 to fine-tune, the prediction curves are essentially
1016 the same, and LAGONN_{exp} seems to have a slight
1017 edge on SetFit_{exp} as we add training data.

1018 LIAR, both the collapsed version we consid-
1019 ered in our content moderation experiments and
1020 the original version (Orig Liar) we examine in our
1021 general text classification experiments here, seems
1022 to be a very difficult dataset. Adding examples
1023 or increased fine-tuning does not appear to consis-
1024 tently increase model performance. We observed
1025 this across all experimental settings and balanced
1026 regimes and is a sensible finding, as it should be
1027 very difficult to determine the truth of a specific
1028 statement without additional context.

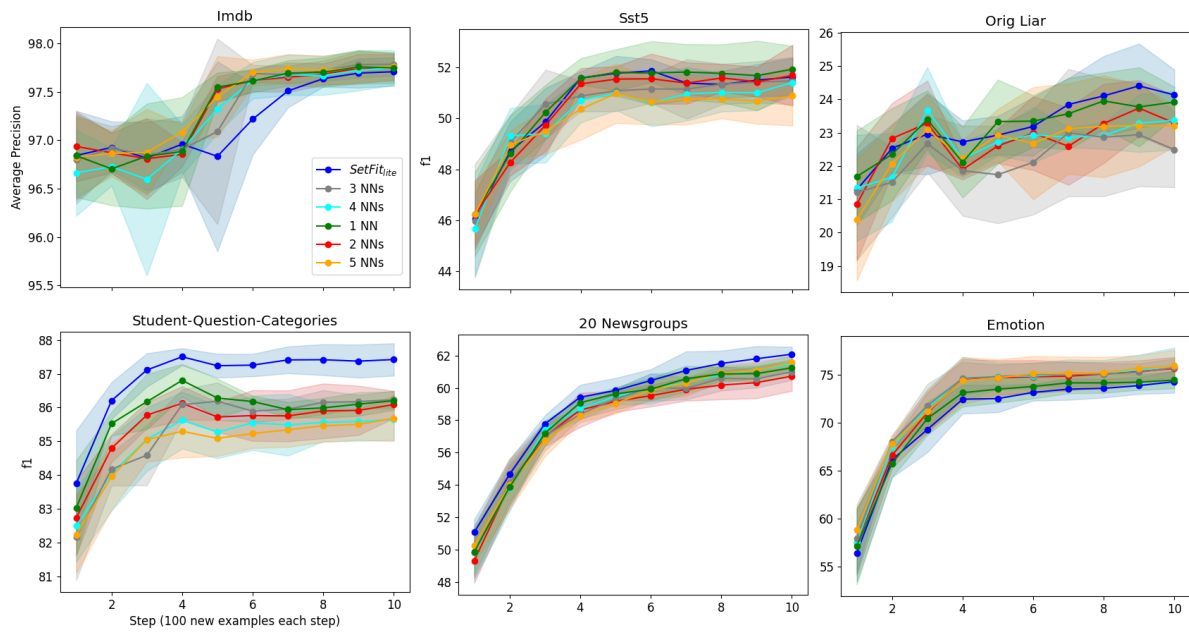


Figure 8: SetFit_{lite} performance compared against one to five neighbors for LAGONN_{lite} LABDIST. The measure is average precision for IMDB, macro-F1 elsewhere.

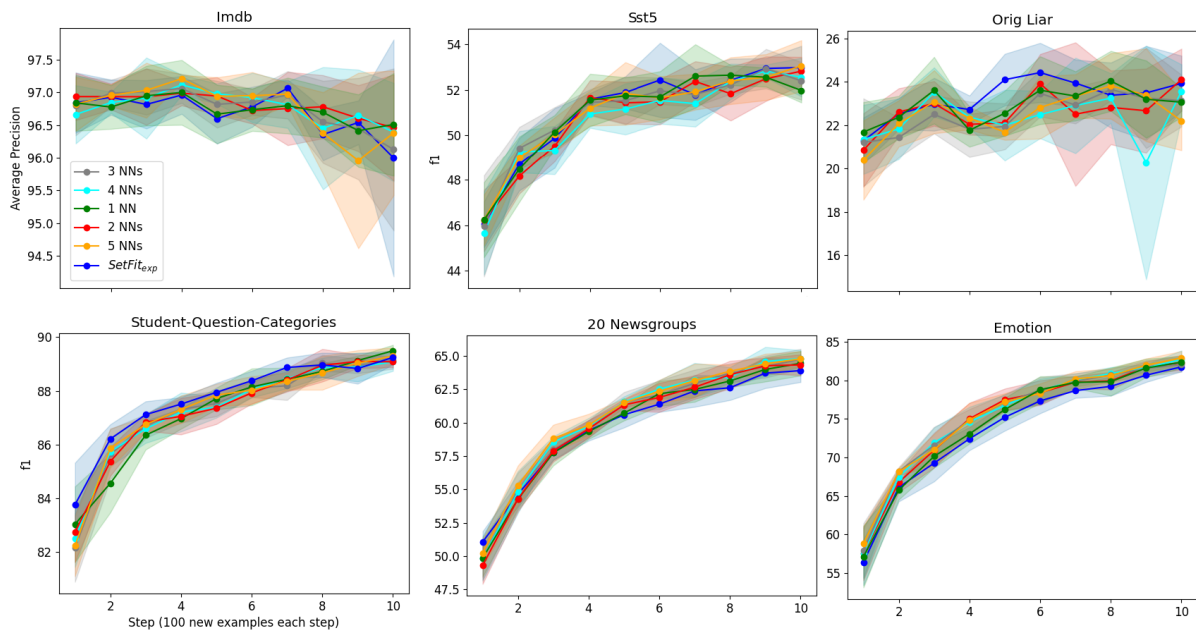


Figure 9: SetFit_{exp} performance compared against one to five neighbors for LAGONN_{exp} LABDIST. The measure is average precision for IMDB, macro-F1 elsewhere.

1029 **A.8 Additional results: multilingual text** 1030 **classification**

1031 In this Appendix section, we provide multilin-
1032 gual text classification results from experiments
1033 where we compare SetFit_{exp} and SetFit_{lite} against
1034 LAGONN_{exp} and LAGONN_{lite} respectively. For
1035 these experiments, we used the Multilingual Ama-
1036 zon Reviews Corpus (Keung et al., 2020), which
1037 has five labels, where each label is a star rating in
1038 Chinese, English, French, German, Japanese, or
1039 Spanish.²² To create the mapping from label to
1040 text, we used code from the ADAPET (Tam et al.,
1041 2021) port in the official SetFit repository.²³ In
1042 these experiments, we used the same multilingual
1043 pretrained Sentence Transformer for all models un-
1044 der the balanced sampling regime.²⁴ In the case of
1045 LAGONN_{exp} and LAGONN_{lite}, we use LABDIST
1046 and search over one to five neighbors, reporting all
1047 results.

1048 Figure 10 shows our results for expensive and
1049 inexpensive models. We note in all cases all mod-
1050 els perform similarly. This supports our assertion
1051 in Section 6 that when the training data is balanced
1052 and we have only a handful of labels or less, it is
1053 sufficient to fine-tune the Sentence Transformer on
1054 only a subset of available training data. A classi-
1055 fier can then be fit on all available data, encoded
1056 with the fine-tuned ST. We observed this for SST-5
1057 and observe it again here, especially clearly on the
1058 Chinese subset of this dataset. SetFit_{exp} plateaus
1059 on the fifth step and stops learning, with different
1060 versions of LAGONN_{exp} outperforming it on later
1061 steps. However, if we move down on row, we see
1062 that all cheaper models continue to learn on all
1063 steps.

²²https://huggingface.co/datasets/amazon_reviews_multi

²³<https://github.com/huggingface/setfit/blob/main/scripts/adapet/ADAPET/utilcode.py>

²⁴<https://huggingface.co/sentence-transformers/paraphrase-multilingual-mpnet-base-v2>

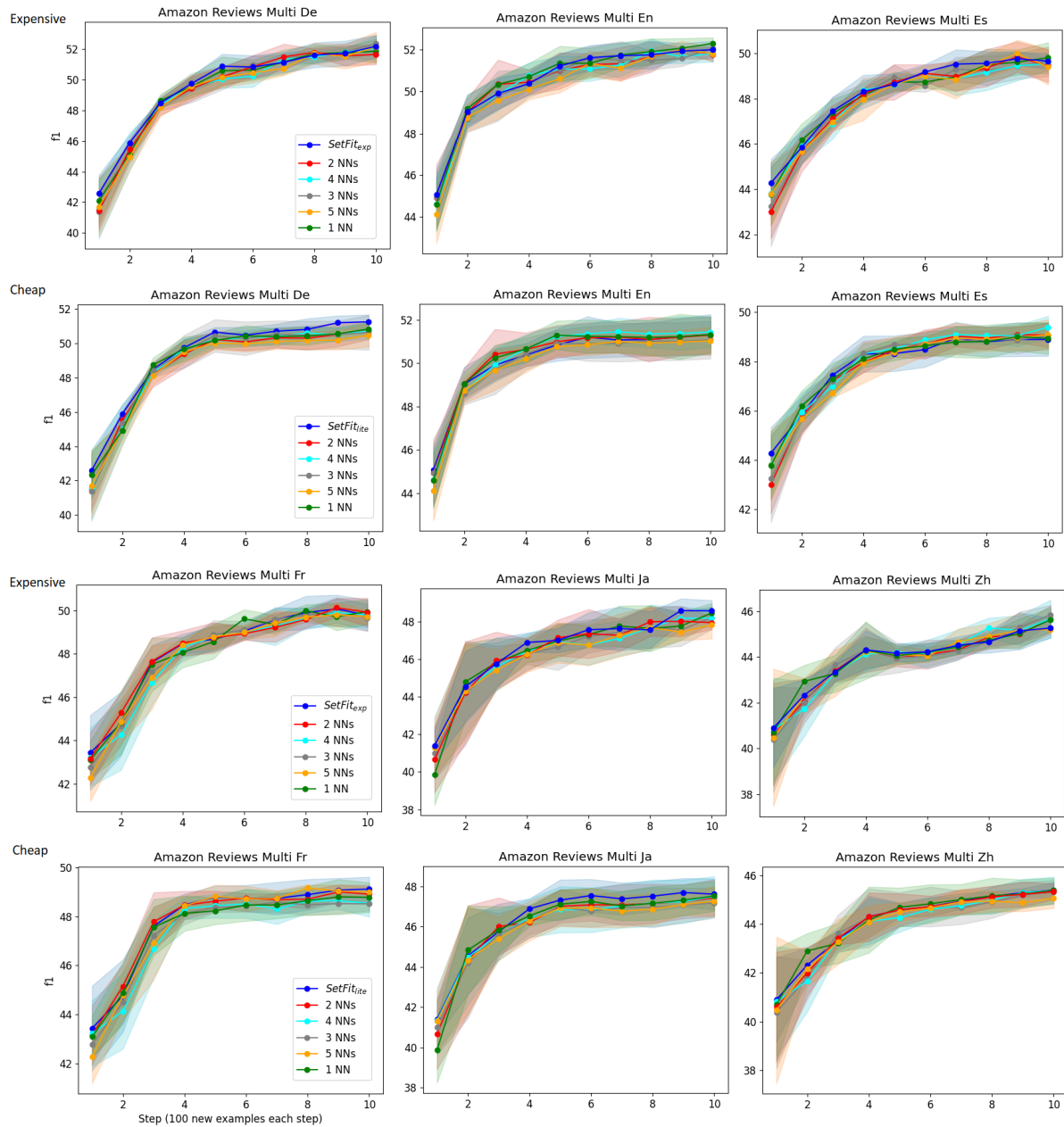


Figure 10: Multilingual classification experiments. In the first row, we display results from expensive models on German, English, Spanish data, with their cheaper counterparts in the following row. In the third and fourth row, we do the same but for French, Japanese, and Chinese. The measure is macro-F1 in all cases.

1064 A.9 Ablations

1065 In this Appendix section, we perform ablation stud-
1066 ies with LAGONN to support our findings in the
1067 main text.

1068 A.9.1 Ablation: LAGONN configurations

1069 Here, we provide an in-depth comparison be-
1070 tween all LAGONN configurations, LABEL, DIS-
1071 TANCE, LABDIST, TEXT, and ALL (see Table
1072 1) for all content moderation datasets, balances, and
1073 levels of expense. The evaluation measure is the
1074 mean average precision ($\times 100$) over five seeds in
1075 all cases except for Hate Speech Offensive where
1076 the measure is the macro-F1.

1077 Below, Figures 11 through 15 are the results
1078 for the LAGONN_{cheap} training strategy, Figures
1079 16 through 20 are the results for LAGONN, Fig-
1080 ures 21 through 25 are the results for LAGONN_{lite},
1081 and Figures 26 through 30 are the results for
1082 LAGONN_{exp}. We place the figures on a new page
1083 for ease of viewing.

1084 In the case of LAGONN_{cheap}, if we do not fine-
1085 tune the embedding model we see little variation in
1086 the standard deviation bands, with the exception of
1087 the LIAR dataset, which seems to be a very difficult
1088 dataset. When we do fine-tune, we see a great deal
1089 of variation, especially in cases of label imbalance,
1090 which is expected as the representations are altered
1091 more. The performance of TEXT and ALL is very
1092 unstable, often being the worst performers, while
1093 sometimes being the best. Interestingly, we note
1094 that DISTANCE, LABEL, and LABDIST often
1095 show very similar performance. In our opinion,
1096 LABDIST seems to be the most consistent and sta-
1097 ble performer, especially in cases when the embed-
1098 ding model is fine-tuned, LAGONN, LAGONN_{lite},
1099 and LAGONN_{exp}.

1100 Overall, we believe that LABDIST is the most
1101 performant/stable configuration of LAGONN, and
1102 it is about this version that we present results in the
1103 main text. We note that we could have presented the
1104 best performer for each evaluation scenario, how-
1105 ever, this is not in the spirit of our work as it adds
1106 yet another hyperparameter to configure, standing
1107 in the way of practical usage and convoluting our
1108 analysis. However, in our codebase, we hope that
1109 we have made it easy for one to change these con-
1110 figurations for their own usage, be it scientific or
1111 otherwise.

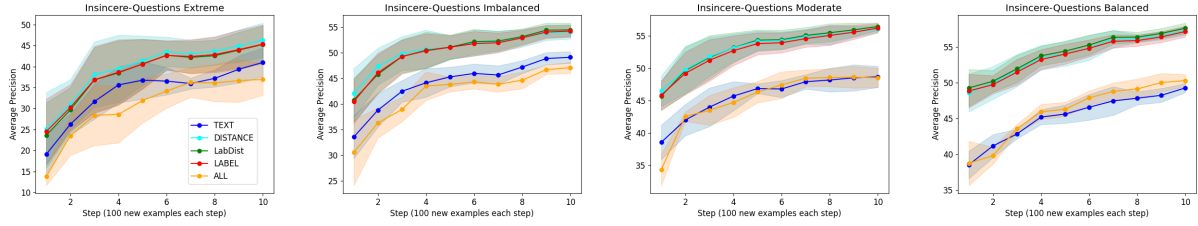


Figure 11: LAGONN_{cheap} performance for all configurations and balance regimes on the Insincere Questions dataset. The relevant balance is in the title of each panel.

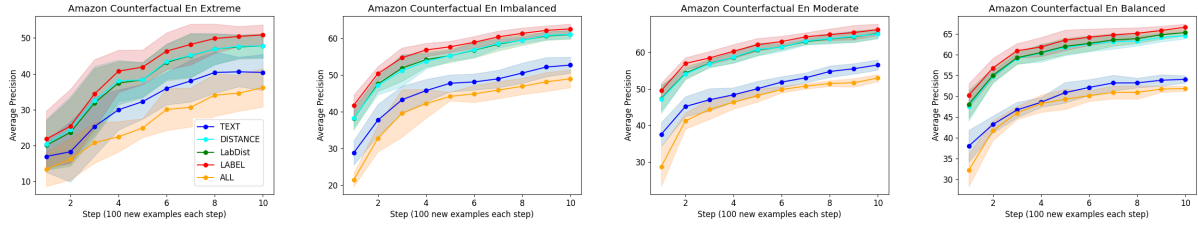


Figure 12: LAGONN_{cheap} performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

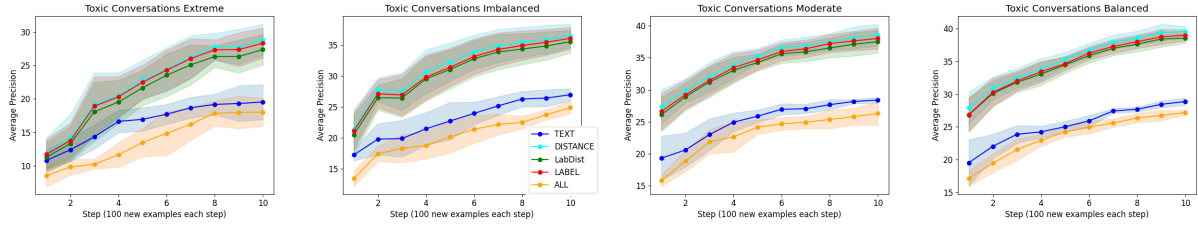


Figure 13: LAGONN_{cheap} performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

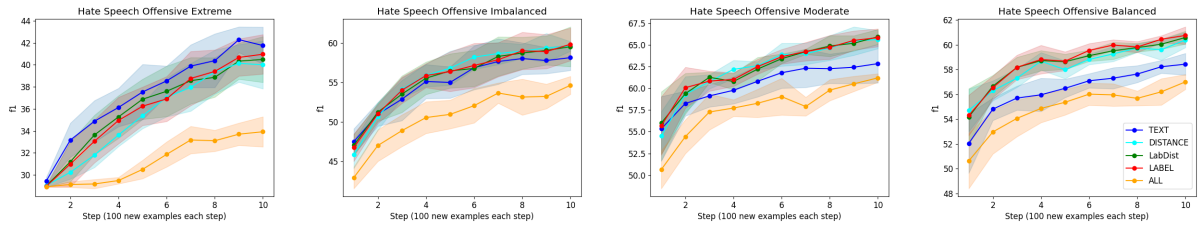


Figure 14: LAGONN_{cheap} performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

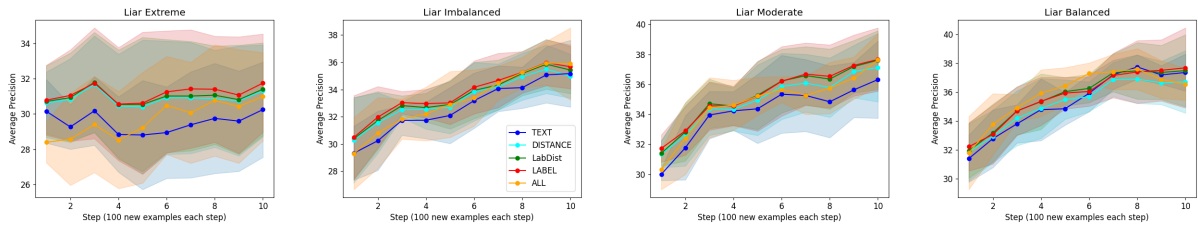


Figure 15: LAGONN_{cheap} performance for all configurations and balance regimes on the LIAR dataset. The relevant balance is in the title of each panel.

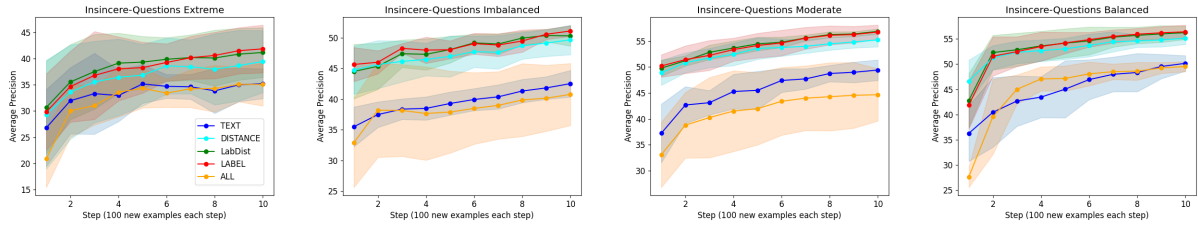


Figure 16: LAGONN performance for all configurations and balance regimes on the InSincere Questions dataset. The relevant balance is in the title of each panel.

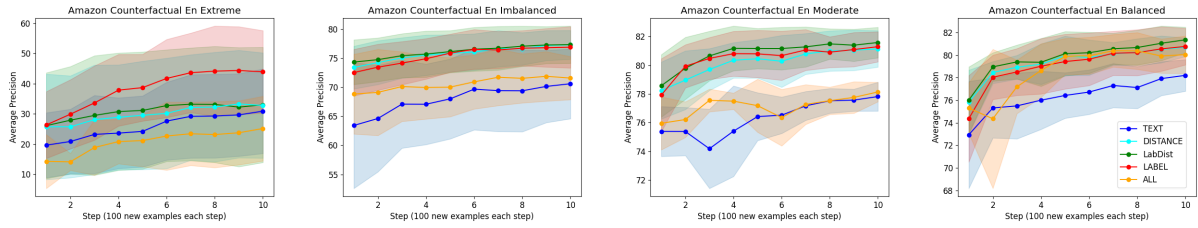


Figure 17: LAGONN performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

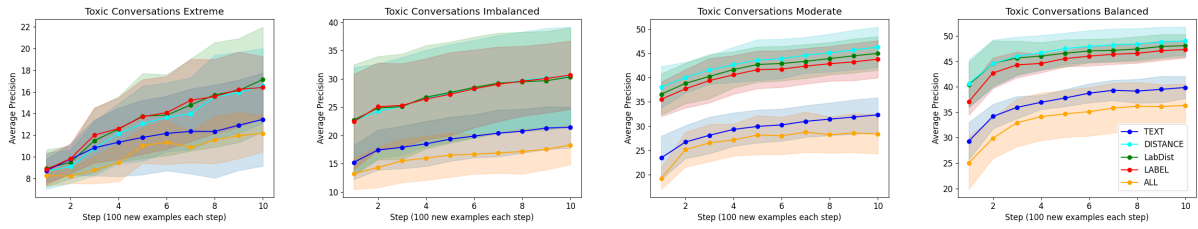


Figure 18: LAGONN performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

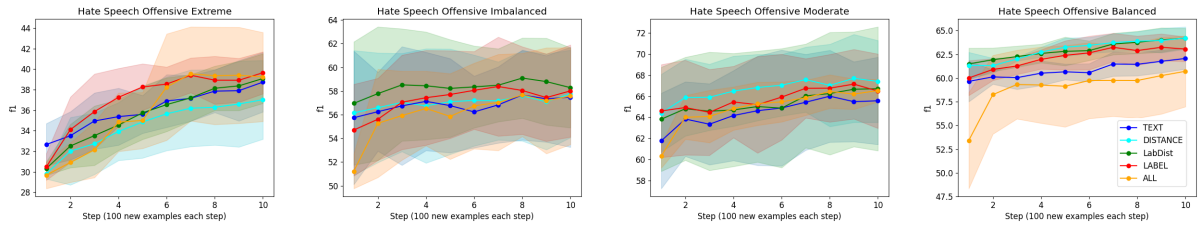


Figure 19: LAGONN performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

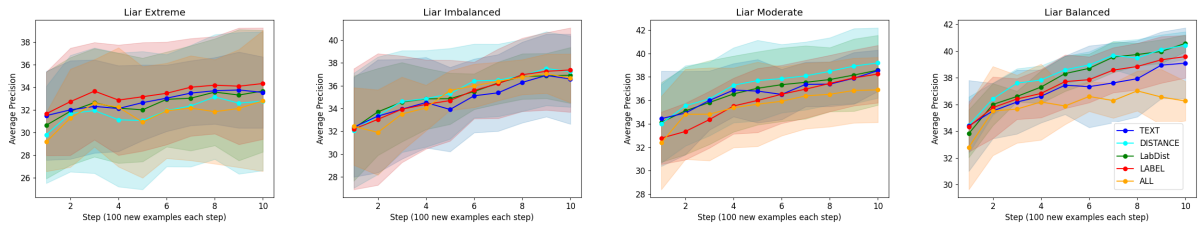


Figure 20: LAGONN performance for all configurations and balance regimes on the LIAR dataset. The relevant balance is in the title of each panel.

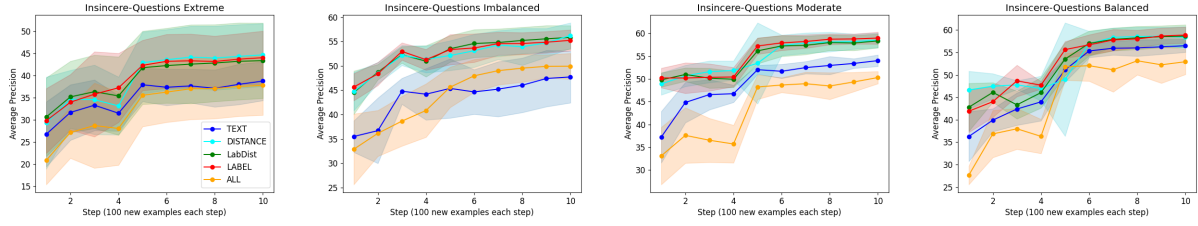


Figure 21: LAGONN_{lite} performance for all configurations and balance regimes on the Insincere Questions dataset. The relevant balance is in the title of each panel.

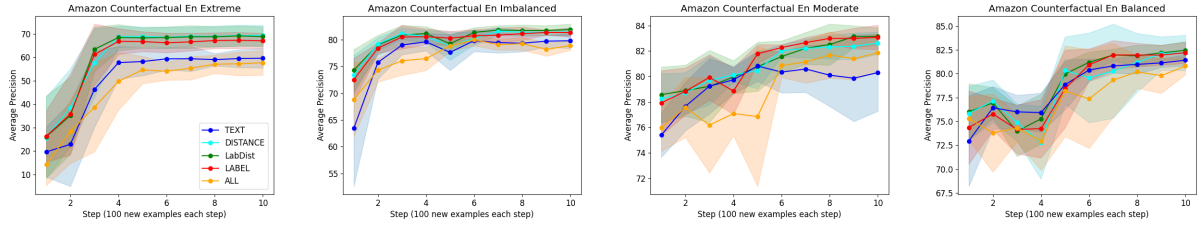


Figure 22: LAGONN_{lite} performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

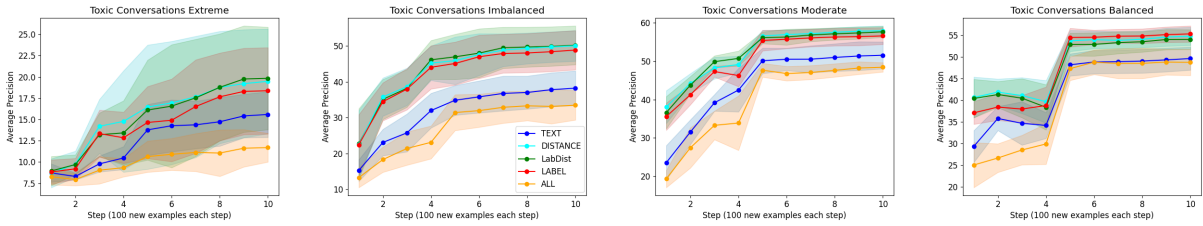


Figure 23: LAGONN_{lite} performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

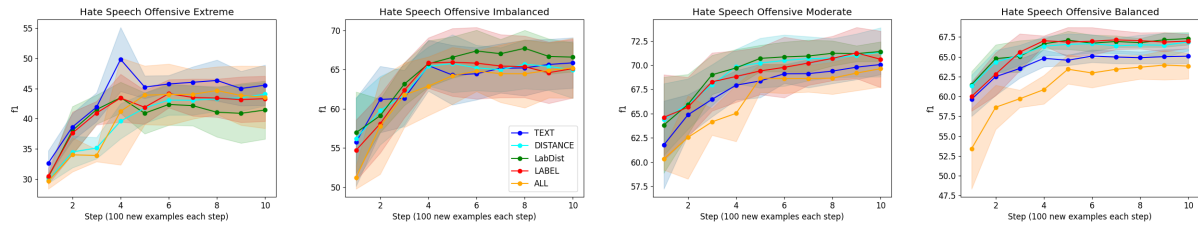


Figure 24: LAGONN_{lite} performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

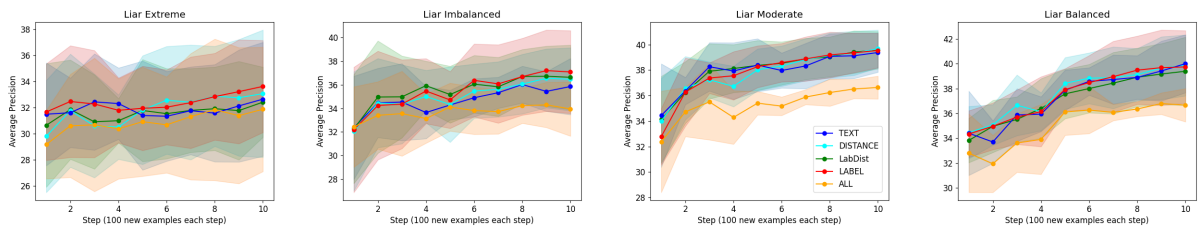


Figure 25: LAGONN_{lite} performance for all configurations and balance regimes on the LIAR dataset. The relevant balance is in the title of each panel.

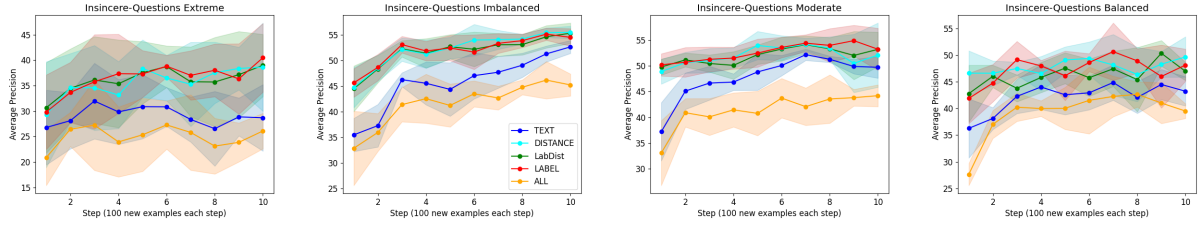


Figure 26: LAGONN_{exp} performance for all configurations and balance regimes on the Insincere Questions dataset. The relevant balance is in the title of each panel.

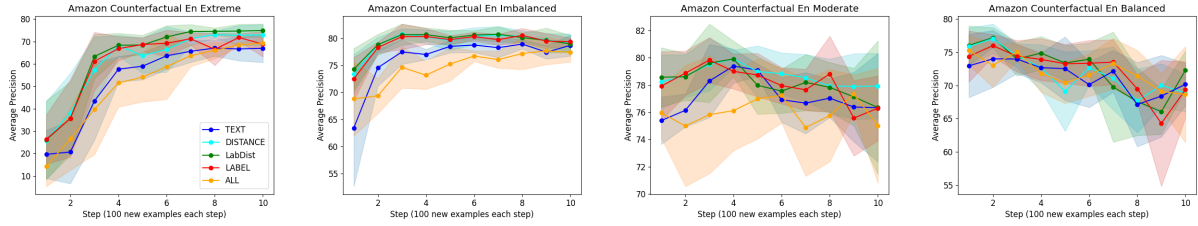


Figure 27: LAGONN_{exp} performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

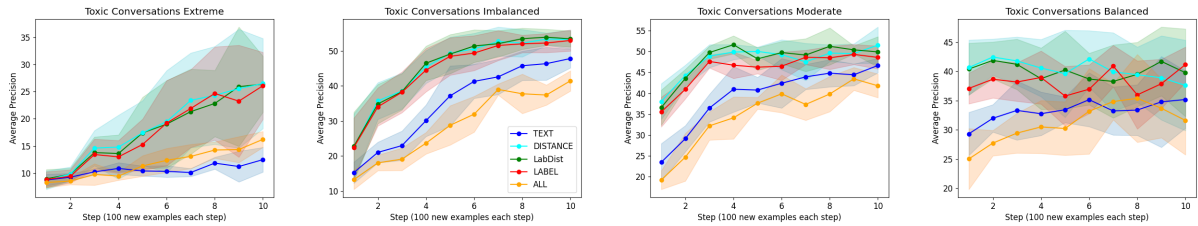


Figure 28: LAGONN_{exp} performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

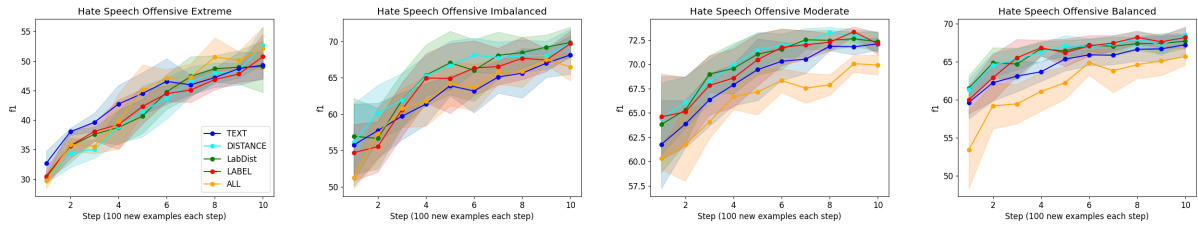


Figure 29: LAGONN_{exp} performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

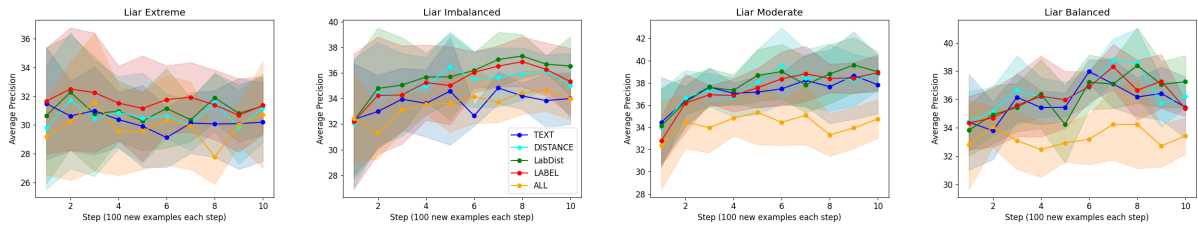


Figure 30: LAGONN_{exp} performance for all configurations and balance regimes on the LIAR dataset. The relevant balance is in the title of each panel.

1112 **A.9.2 Ablation: LAGONN k nearest** 1113 **neighbors**

1114 Here, at the suggestion of an anonymous reviewer,
1115 we present ablation results and analysis of search-
1116 ing over one to five nearest neighbors when modify-
1117 ing input via LAGONN. We present results over all
1118 LAGONN configurations under the LAGONN_{lite}
1119 fine-tuning strategy and with all balance regimes
1120 for the content moderation datasets. For the gen-
1121 eral text classification setting, we present results for
1122 both LAGONN_{lite} and LAGONN_{exp} fine-tuning
1123 under the balanced regime for all datasets with the
1124 LABDIST and TEXT configurations.

1125 If we consider all LAGONN configurations and
1126 balance regimes in the case content moderation,
1127 Figures 31 through 55, the number of neighbors
1128 does not appear to be an important hyperparameter;
1129 the learning curves for a given dataset and balance
1130 regime are very similar. While there is variation,
1131 the trend appears to be that the first NN results in
1132 the stablest, most performant, and most consistent
1133 model.

1134 However, if we only focus on LABDIST (Fig-
1135 ures 31 through 35), the default LAGONN con-
1136 figuration, we see that it can be a very important
1137 hyperparameter to consider in cases of extreme im-
1138 balance or when we have balanced data but few
1139 data points. For example, performance is boosted
1140 by up to five points for Hate Speech Offensive by
1141 the tenth step (1000 examples) with five neighbors
1142 under the extreme balance regime, yet for the bal-
1143 anced regime, the performance curves are roughly
1144 the same. For Toxic Conversations, in the balanced
1145 regime, we see that we can increase performance
1146 by up to seven points on the second step (200 ex-
1147 amples) by considering more neighbors.

1148 Turning our attention now to the general clas-
1149 sification experiments, we see that the number of
1150 neighbors for both the LABDIST and TEXT confi-
1151 gurations continues to consistently not really make
1152 much of a difference, with all models showing very
1153 similar performance curves for all datasets. We
1154 note however that LABDIST appears to be the most
1155 performant configuration of our method. While
1156 continued fine-tuning on datasets with a large num-
1157 ber of labels does increase performance, we ob-
1158 serve essentially the same boost for all neighbors.
1159 We also observe similar instability and performance
1160 degradation when we fine-tune on a large number
1161 of examples in cases when we have few labels.

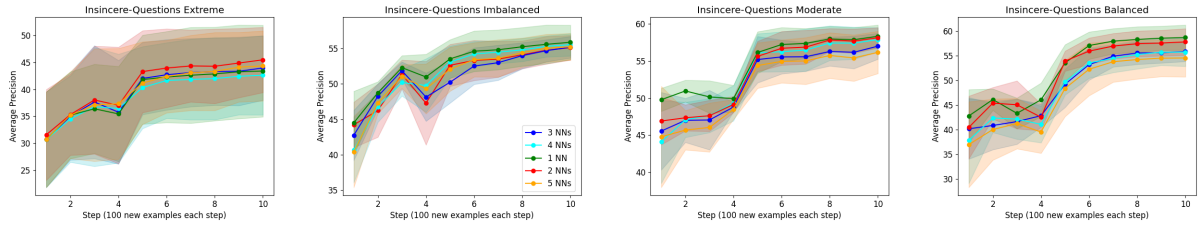


Figure 31: LABDIST results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Inscere Questions dataset. The relevant balance is in the title of each panel.

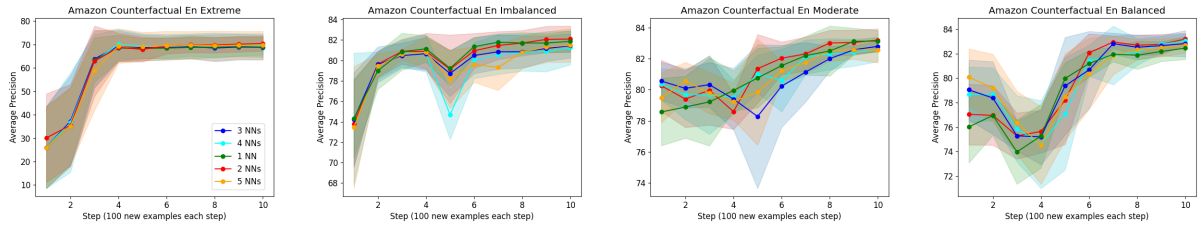


Figure 32: LABDIST results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

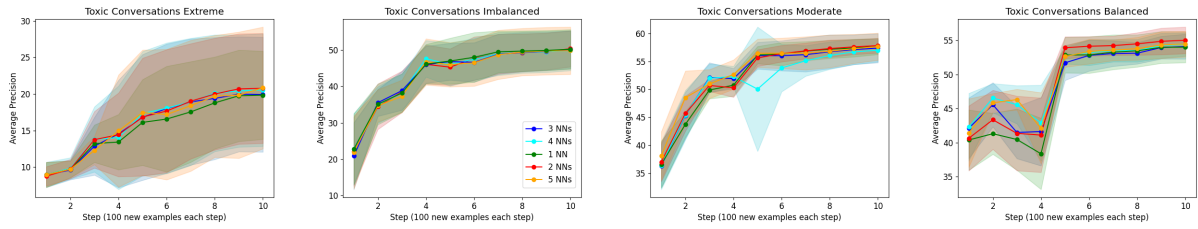


Figure 33: LABDIST results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

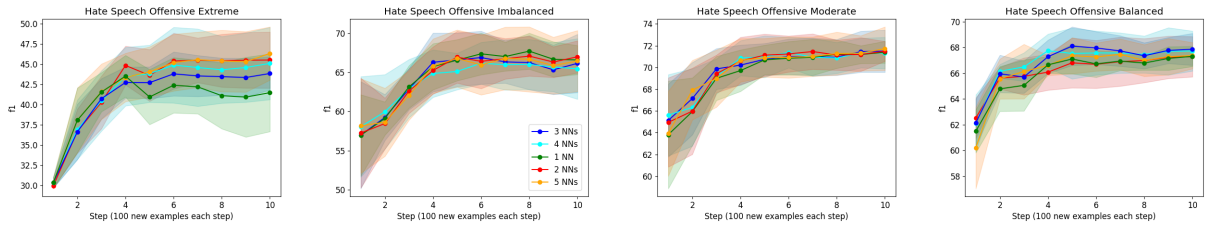


Figure 34: LABDIST results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

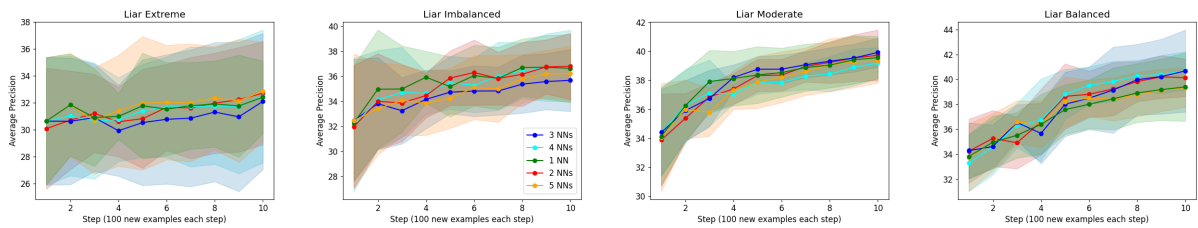


Figure 35: LABDIST results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the LIAR dataset. The relevant balance is in the title of each panel.

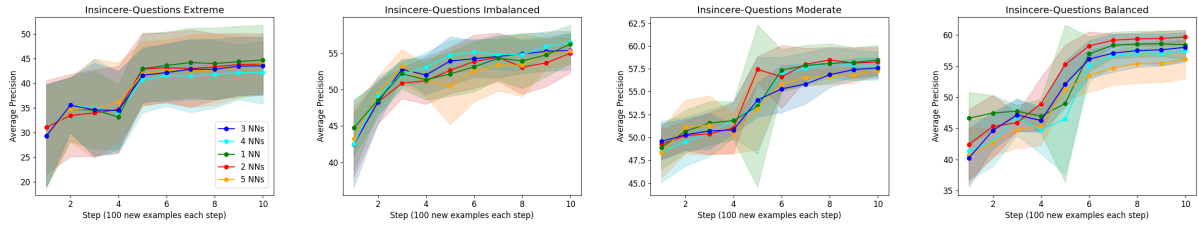


Figure 36: DISTANCE results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the the Insincere Questions dataset. The relevant balance is in the title of each panel.

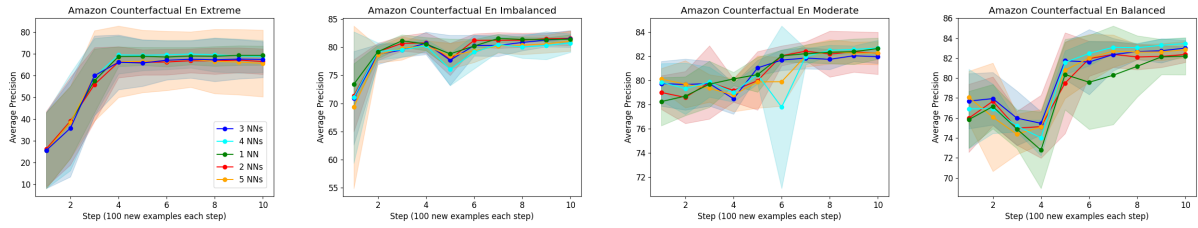


Figure 37: DISTANCE results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

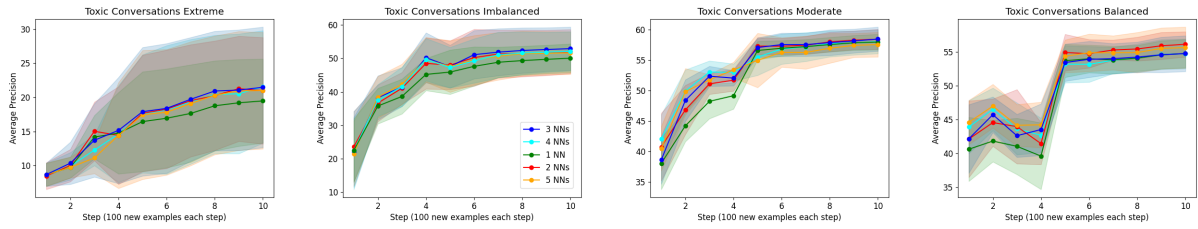


Figure 38: DISTANCE results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

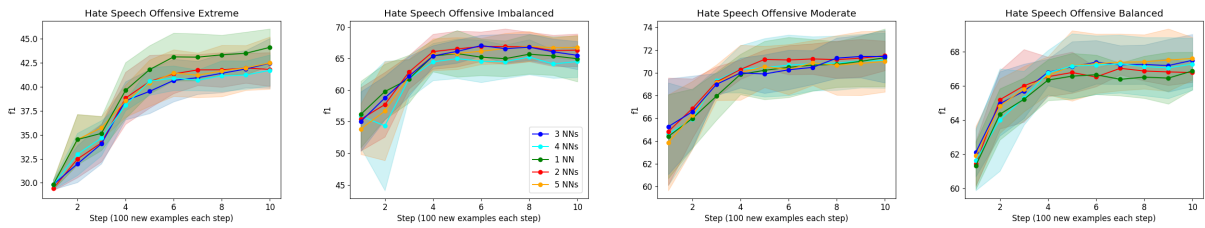


Figure 39: DISTANCE results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

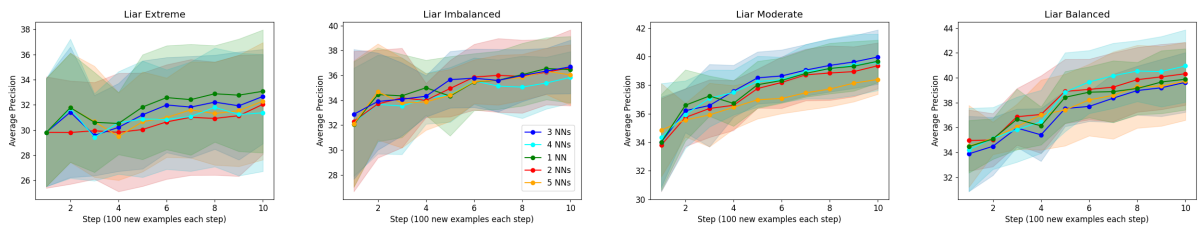


Figure 40: DISTANCE results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the LIAR dataset. The relevant balance is in the title of each panel.

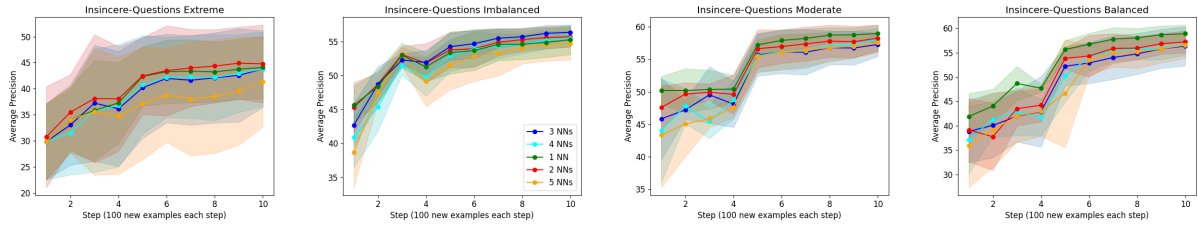


Figure 41: LABEL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Insincere Questions dataset. The relevant balance is in the title of each panel.

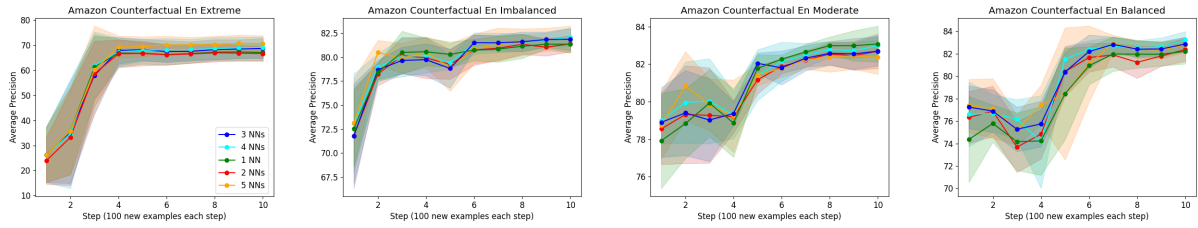


Figure 42: LABEL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

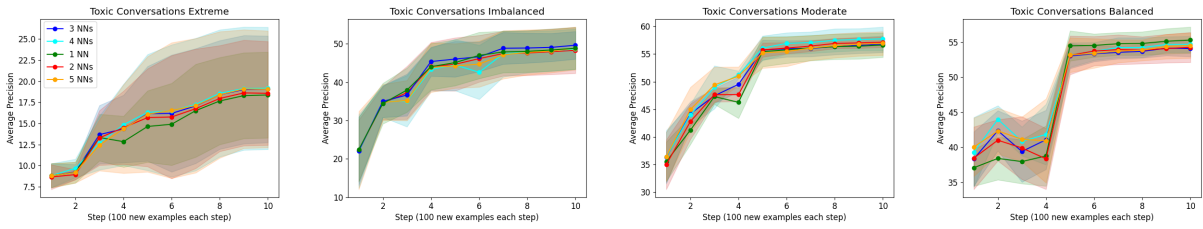


Figure 43: LABEL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

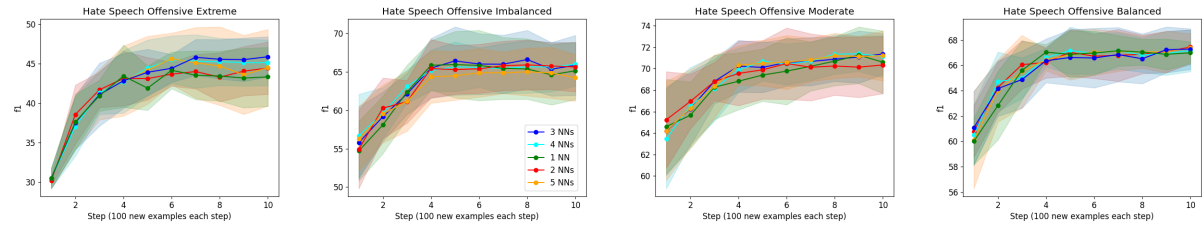


Figure 44: LABEL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

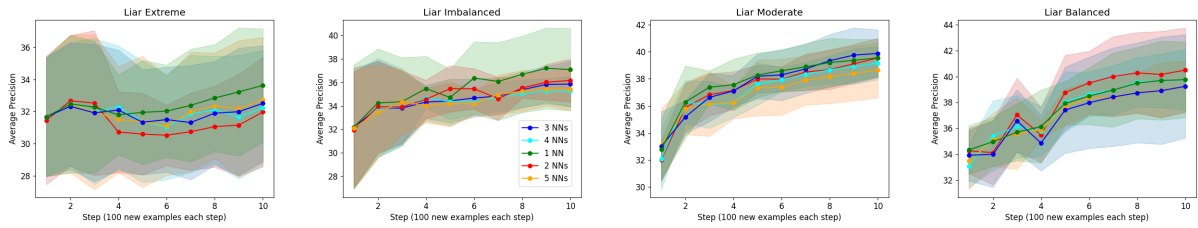


Figure 45: LABEL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the LIAR dataset. The relevant balance is in the title of each panel.

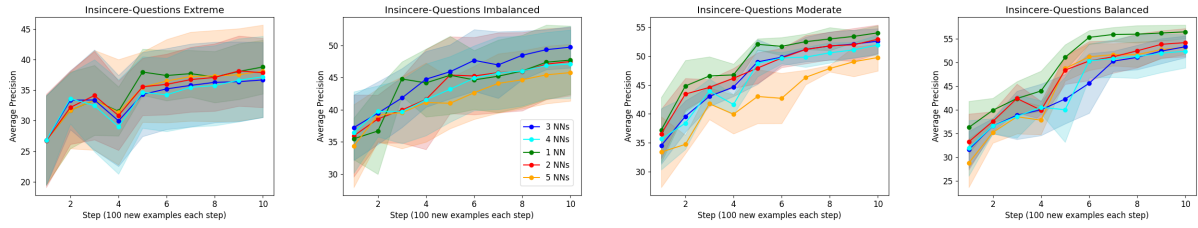


Figure 46: TEXT results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Insincere Questions dataset. The relevant balance is in the title of each panel.

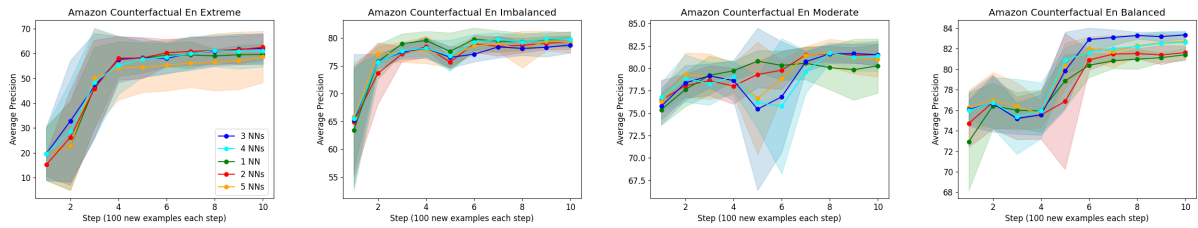


Figure 47: TEXT results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

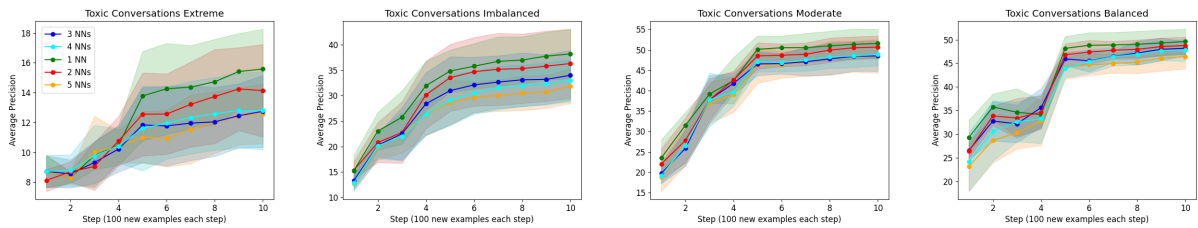


Figure 48: TEXT results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

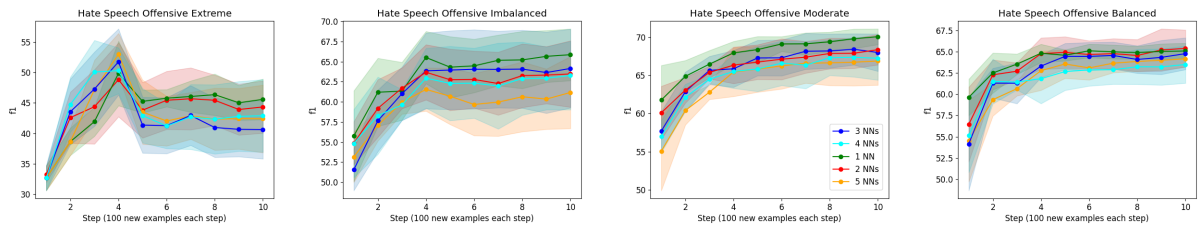


Figure 49: TEXT results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

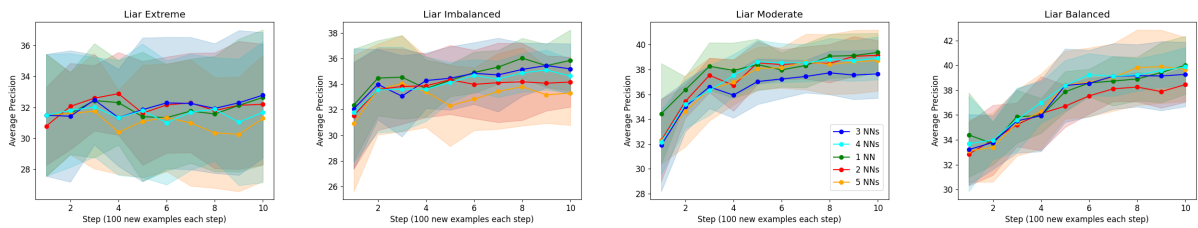


Figure 50: TEXT results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the LIAR dataset. The relevant balance is in the title of each panel.

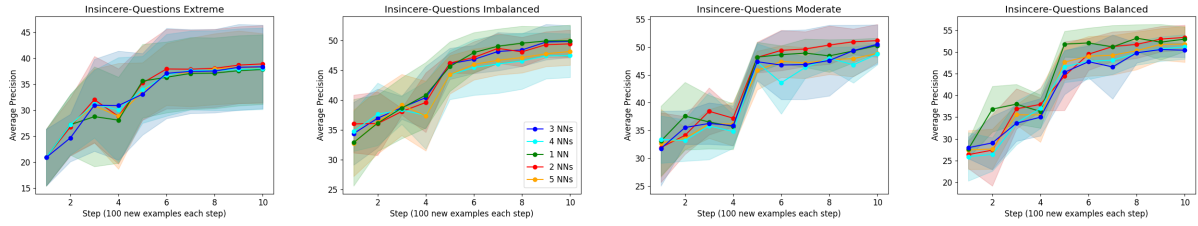


Figure 51: ALL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Insincere Questions dataset. The relevant balance is in the title of each panel.

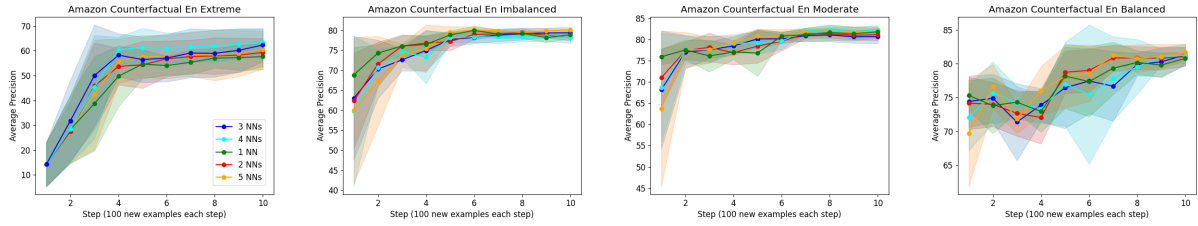


Figure 52: ALL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

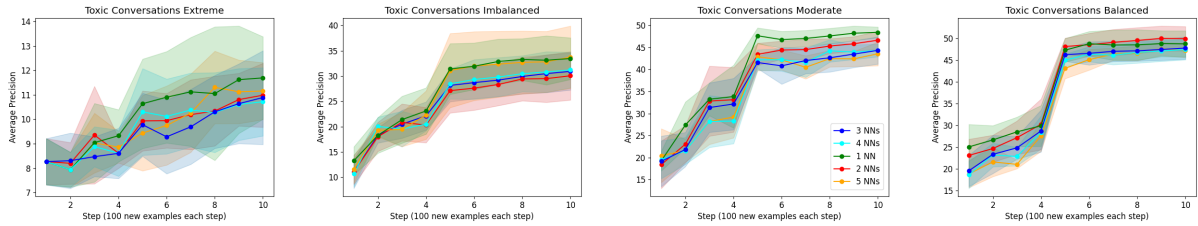


Figure 53: ALL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

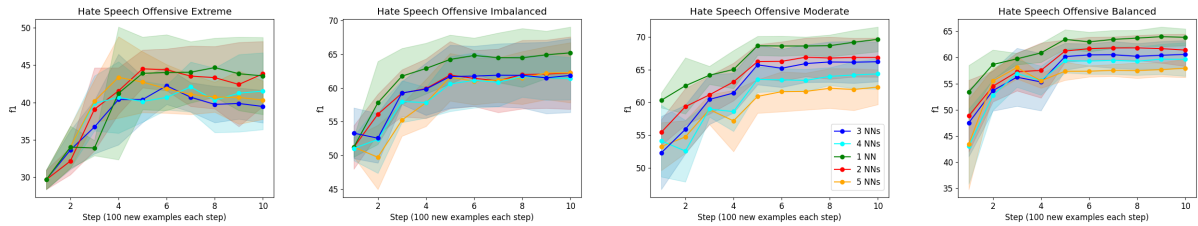


Figure 54: ALL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

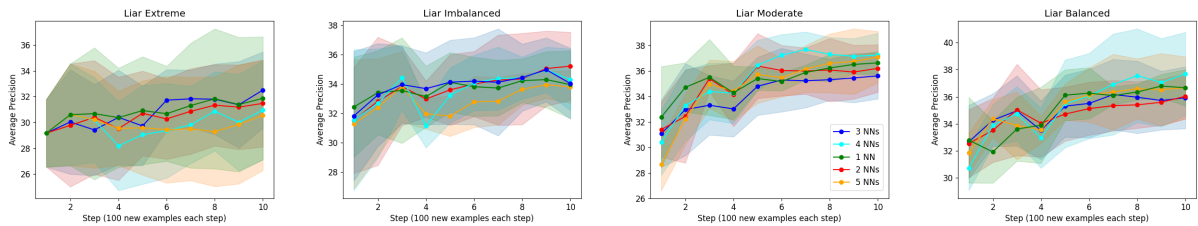


Figure 55: ALL results for one to five neighbors under the LAGONN_{lite} fine-tuning strategy on the LIAR dataset. The relevant balance is in the title of each panel.

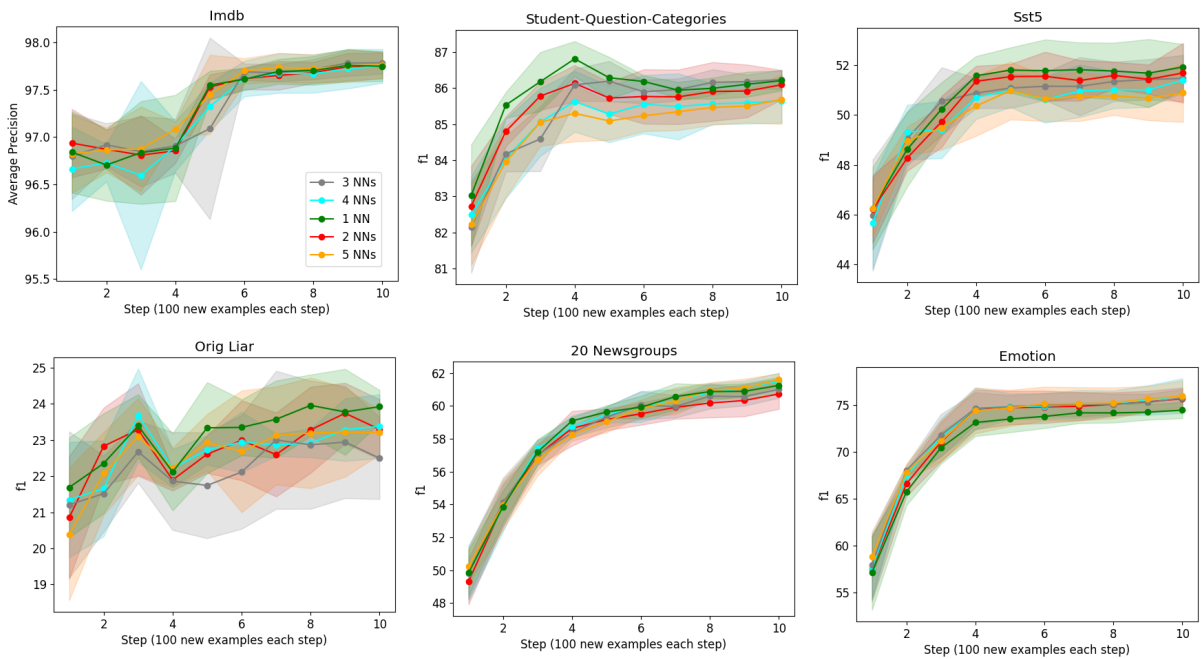


Figure 56: LABDIST results for one to five neighbors under the $LAGONN_{lite}$ fine-tuning strategy over all six general classification datasets. Results are for the balanced sampling regime and the measure is average precision for IMDB, macro-F1 elsewhere.

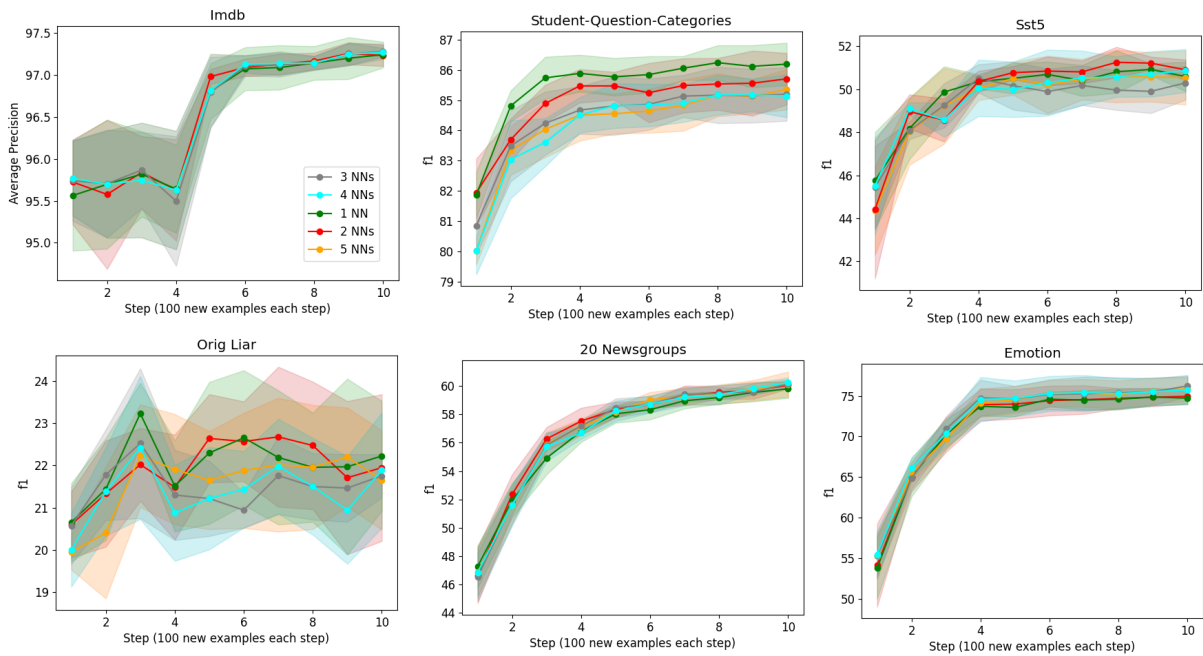


Figure 57: TEXT results for one to five neighbors under the $LAGONN_{lite}$ fine-tuning strategy over all six general classification datasets. Results are for the balanced sampling regime and the measure is average precision for IMDB, macro-F1 elsewhere.

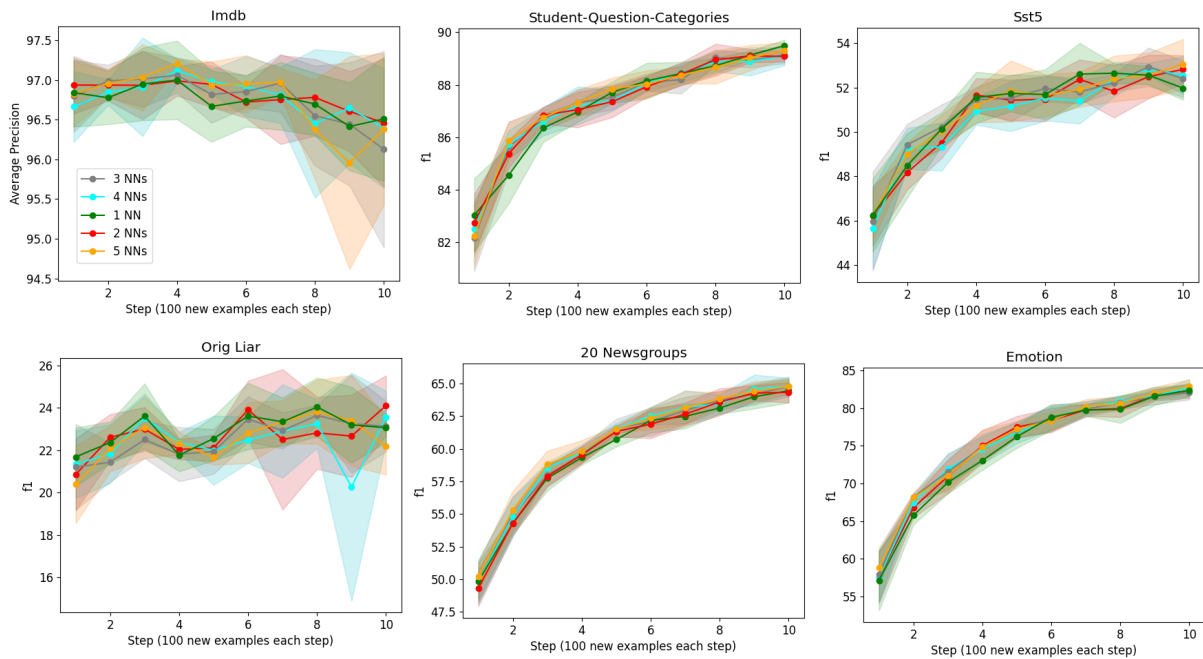


Figure 58: LABDIST results for one to five neighbors under the $LAGONN_{exp}$ fine-tuning strategy over all six general classification datasets. Results are for the balanced sampling regime and the measure is average precision for IMDB, macro-F1 elsewhere.

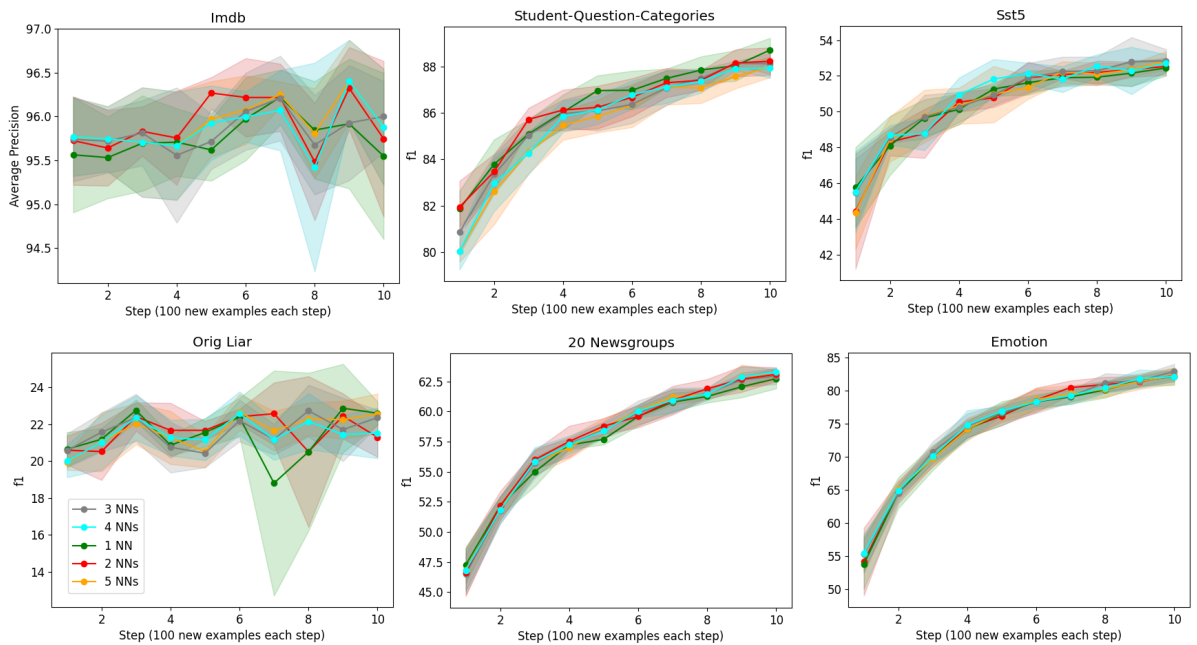


Figure 59: TEXT results for one to five neighbors under the $LAGONN_{exp}$ fine-tuning strategy over all six general classification datasets. Results are for the balanced sampling regime and the measure is average precision for IMDB, macro-F1 elsewhere.

1162 **A.9.3 Ablation: the effect of encoding distance**

1163 Here, at the suggestion of an anonymous reviewer,
1164 we present ablation results and analysis of how en-
1165 coding distance affects LAGONN, because PLMs
1166 often struggle to understand numbers. Note that
1167 during our development stage, we ensured that our
1168 tokenizer was capable of encoding floats with trail-
1169 ing digits. To examine the effect of trailing digits
1170 on LAGONN, we consider the DISTANCE con-
1171 figuration (see Table 1), where we append only
1172 the Euclidean distance to the input text. In this
1173 ablation, however, we round to different levels of
1174 precision. For example, if the distance were a float
1175 of 0.123456789, we round it to the nearest whole
1176 number, 0.0, single digit float, 0.1, three digit float,
1177 0.123, six digit float, 0.123457, and finally keep it
1178 unrounded, that is, the original DISTANCE config-
1179 uration, 0.123456789. The below results are only
1180 for the LAGONN_{lite} training strategy. We chose
1181 LAGONN_{lite} for this ablation because it provides
1182 insight into both how distance affects full-model
1183 fine-tuning and only refitting the classification head.
1184 The results can be seen below in Figures 60 through
1185 64. We place the figures on a new page for ease of
1186 viewing.

1187 Interestingly, we tend to observe very similar per-
1188 formance curves for all rounding precisions. The
1189 exceptions to this would perhaps be Amazon Coun-
1190 terfactual and Hate Speech Offensive in the bal-
1191 anced regime where DISTANCE and rounding
1192 to the third trailing digit respectively exhibit large
1193 instability.

1194 Although not always the case, it appears that
1195 providing the model with the distance rounded to
1196 the nearest whole number tends to result in the
1197 strongest and stablest performer, however, we em-
1198 phasize that in general there does not seem to a
1199 dramatic difference between the rounding preci-
1200 sions we considered. Longer digits slightly worsen
1201 model performance and the model might learn the
1202 most from simpler or abbreviated representations
1203 of distance. This finding motivated us to consider
1204 the ablation in Appendix A.9.4.

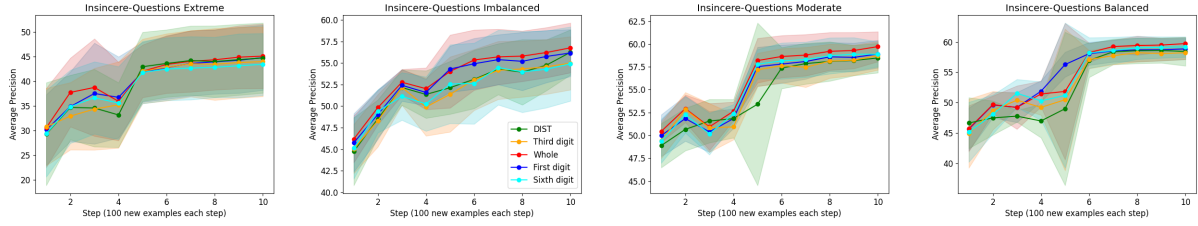


Figure 60: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the InSincere Questions dataset and the relevant balance is in the title of each panel.

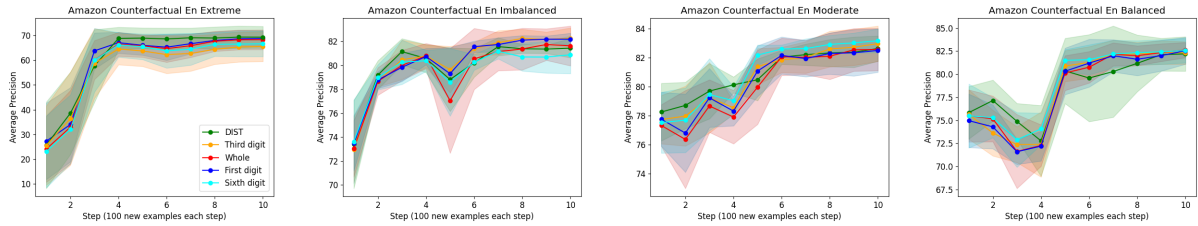


Figure 61: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the Amazon Counterfactual dataset and the relevant balance is in the title of each panel.

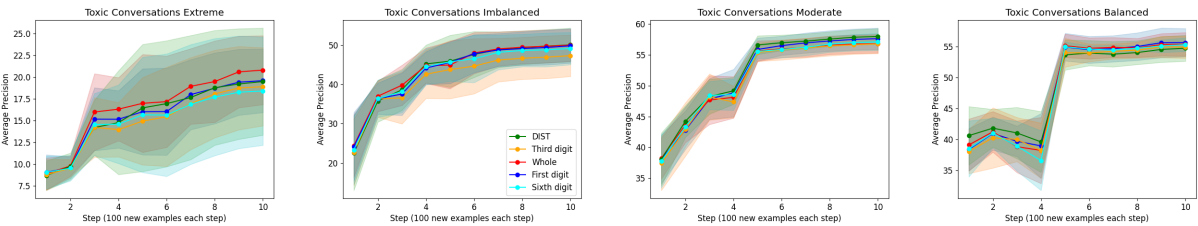


Figure 62: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the Toxic Conversations dataset and the relevant balance is in the title of each panel.

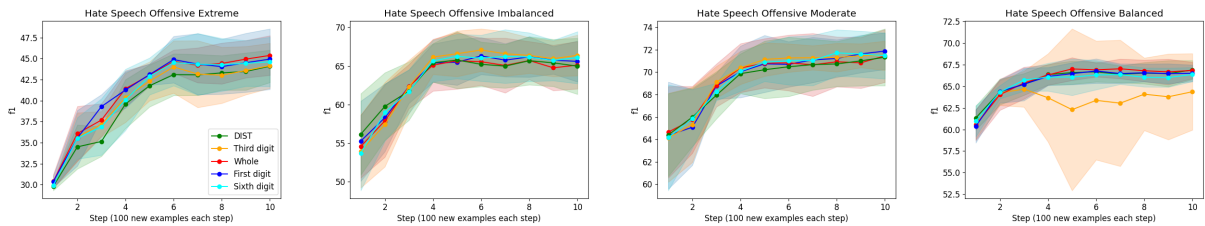


Figure 63: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the Hate Speech Offensive dataset and the relevant balance is in the title of each panel.

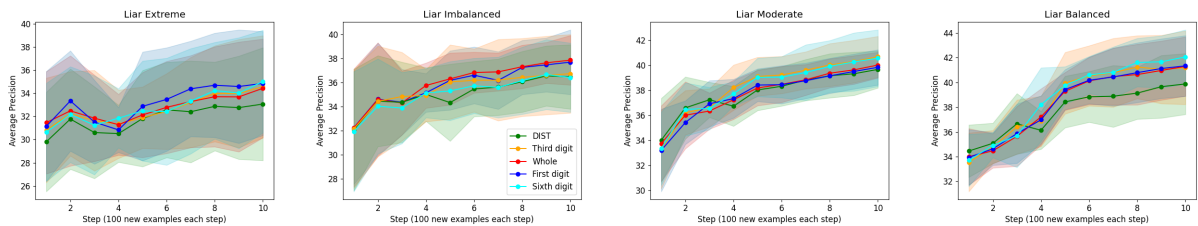


Figure 64: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the LIAR dataset and the relevant balance is in the title of each panel.

1205 **A.9.4 Ablation: support for LABDIST**

1206 The results from the ablation in Appendix A.9.3
1207 suggest that rounding the distance to the nearest
1208 whole number results in a stronger classifier than
1209 appending the unrounded distance. Thus far, we
1210 have asserted that LABDIST, where we append
1211 both the gold label of the NN and unrounded dis-
1212 tance is the most performant version of LAGONN
1213 (see Table 1). To demonstrate that this is reason-
1214 able, in this ablation study, we compare the orig-
1215 inal LABDIST configuration against three mod-
1216 els, namely the LABEL configuration, distance
1217 rounded to near whole number (Whole), and finally
1218 a new configuration similar to LABDIST, but where
1219 we append the gold label and distance rounded to a
1220 whole number, which we refer to as LABROUND.
1221 As in Appendix A.9.3, in this ablation we con-
1222 sider only the LAGONN_{lite} fine-tuning strategy.
1223 We chose LAGONN_{lite} for this ablation because
1224 it provides insight into both how the different con-
1225 figurations affect full-model fine-tuning and only
1226 re-fitting the classification head. The results can be
1227 seen below in Figures 65 through 69. We place the
1228 figures on a new page for ease of viewing.

1229 In general, we note very similar performance
1230 curves for these four models. In the case of Insinc-
1231 ere Questions, appending the distance after round-
1232 ing it to the nearest whole number (Whole, the red
1233 curve), is a strong model, except in the balanced
1234 regime where we note large instability. The results
1235 for Amazon Counterfactual tell a different story,
1236 where rounding the Euclidean distance to the near-
1237 est whole number causes large instability and even
1238 degrades performance on the fifth step.

1239 For the other evaluation scenarios, it is unclear
1240 what is the strongest method as sometimes LAB-
1241 DIST is the best performer and sometimes it is
1242 Whole (the red curve). However, we believe that in
1243 general LABDIST is the most stable model while
1244 also often being the most performant. We therefore
1245 choose it as our default LAGONN configuration as
1246 a compromise between strength and stability. It is
1247 about this configuration which we report results in
1248 the main text. Our interpretation of this is that pass-
1249 ing the model both a discrete prediction (the gold
1250 label of the NN) and a truly continuous measure
1251 of similarity (the unrounded Euclidean distance)
1252 gives it the most consistent and dependable reason-
1253 ing ability.

1254 We note, as we did in Appendix A.9.1, that we
1255 could have presented the best performer for each

1256 evaluation scenario, however, it is not the goal of
1257 our work to create even more hyperparameters that
1258 must be iterated over. However, we hope that our
1259 codebase has made it easy for one to change these
1260 configurations for their own purposes.

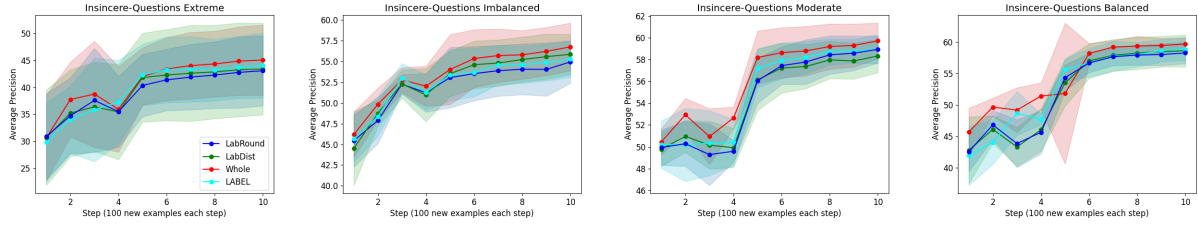


Figure 65: LAGONN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Insincere Questions dataset and the relevant balance is in the title of each panel.

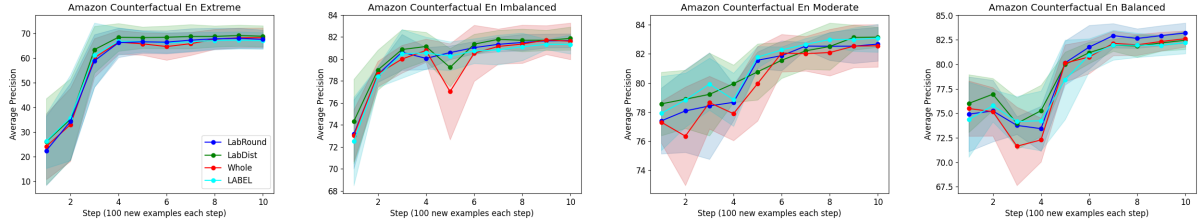


Figure 66: LAGONN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Amazon Counterfactual dataset and the relevant balance is in the title of each panel.

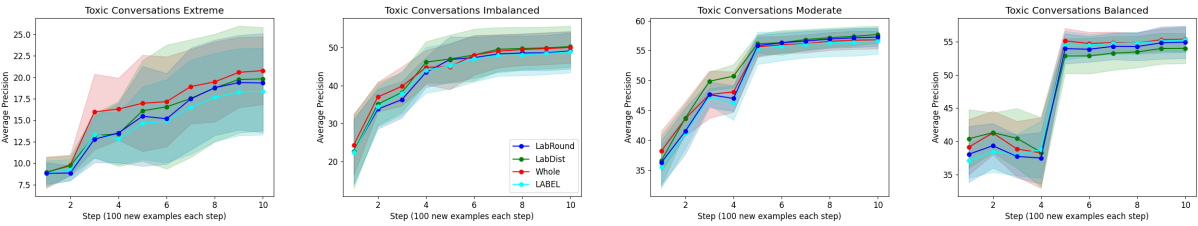


Figure 67: LAGONN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Toxic Conversations dataset and the relevant balance is in the title of each panel.

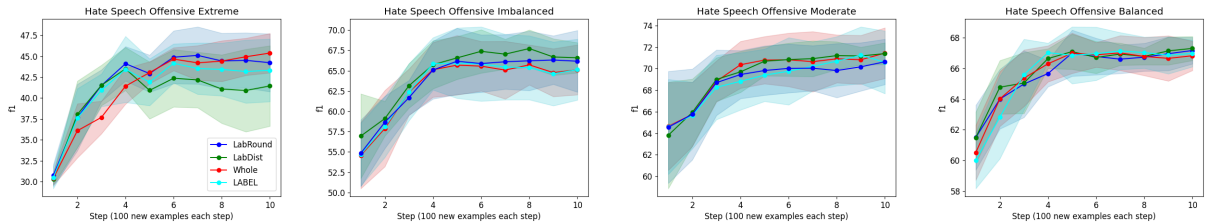


Figure 68: LAGONN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Hate Speech Offensive dataset and the relevant balance is in the title of each panel.

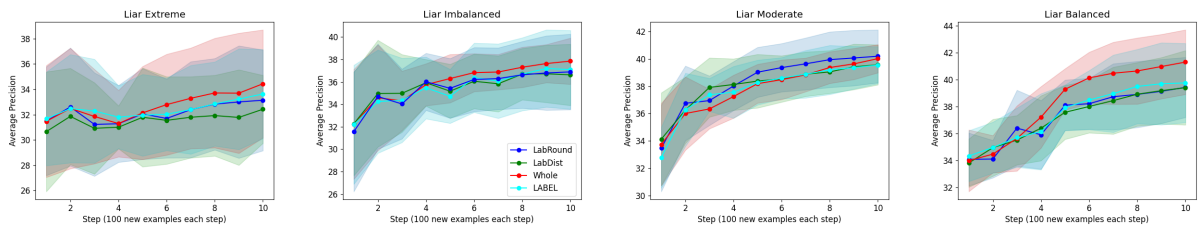


Figure 69: LAGONN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the LIAR dataset and the relevant balance is in the title of each panel.

1261 **A.10 Examples of LAGONN modified text**

1262 **WARNING:** Some of the examples below are of

1263 an offensive nature. Please view with caution.

1264 In this section, we provide examples of how

1265 LAGONN_{exp} modifies test text from the content

1266 moderation datasets we studied under the ALL con-

1267 figuration. We choose this configuration because

1268 the information it appends from a NN in the train-

1269 ing data to a test instance encapsulates all configu-

1270 rations. LAGONN_{exp} was trained under a balanced

1271 distribution and five examples per label were cho-

1272 sen randomly on the first, fifth, and tenth step to

1273 demonstrate how the same test instance might be

1274 decorated with different training examples as the

1275 training data grow. We have made the .csv files

1276 available with our code and data files. In order to

1277 not break our .pdf generator, we were forced to

1278 remove a handful of symbols from the below text,

1279 but the original modifications remain in-tact in the

1280 .csv files included with our code files. Note that

1281 MPNET’s separator token is `</s>`, not [SEP].

Insincere Questions Step 1 1282

Test Modified What rapper still relevant and 1283
popular today has the best rhyme schemes? `</s>` 1284
`<insincere question 3.859471321105957>` What 1285
would be a good nickname for Trump, Donald 1286
Dumbck, and President Spankovich? `</s>` `<valid` 1287
`question 4.124274253845215>` What are after class 1288
12 courses in commerce stream to choose from? I 1289
have completed my class 12 (expexted 90+) and 1290
aim to do business (not aim to do job). 1291

Label valid question 1292

Test Modified Which books do you sug- 1293
gest to someone who get a free time and will 1294
help him stay motivated? `</s>` `<valid question` 1295
`3.9509353637695312>` What are the best online 1296
courses to learn data science? `</s>` `<insincere ques-` 1297
`tion 4.300448417663574>` What are the more steps 1298
in Career Oriented Education? 1299

Label valid question 1300

Test Modified How will you feel if someone 1301
talks badly about Kunti? `</s>` `<insincere ques-` 1302
`tion 3.5063605308532715>` Why are the UK gov- 1303
ernment and the media (especially the BBC and 1304
the Guardian) demonising ordinary British people, 1305
manipulating buzz words like “alt-right”, “Islama- 1306
phobia”, “racist” to suppress legitimate outrage at 1307
Muslim grooming gangs? `</s>` `<valid question` 1308
`3.6699037551879883>` How do Israelis and Pales- 1309
tinians view Nuseir Yassin? 1310

Label valid question 1311

Test Modified Why is equine HYPP inher- 1312
ited? `</s>` `<valid question 4.066534996032715>` 1313
Can you share some of the pics of hostel 1314
of Indira Gandhi medical college, Shimla, Hi- 1315
machal Pradesh? `</s>` `<insincere question` 1316
`4.231775760650635>` I am an experienced pro- 1317
grammer and in my high school my teacher tried 1318
to make me use python so I said, "No; Trust me, 1319
python is just a language for beginners, thereby 1320
making it not for me." I got sent out. Did I do 1321
anything wrong? 1322

Label valid question 1323

Test Modified How do the Valerie Stevens 1324
leather jackets achieve their quality during the 1325
manufacturing process? `</s>` `<valid question` 1326
`3.9721384048461914>` How are the Lancaster 1327
leather sofas manufactured? `</s>` `<insincere ques-` 1328
`tion 4.3559441566467285>` I am an experienced 1329
programmer and in my high school my teacher tried 1330
to make me use python so I said, "No; Trust me, 1331
python is just a language for beginners, thereby 1332

| | | | |
|------|---|---|------|
| 1333 | making it not for me." I got sent out. Did I do | "go fuck yourself," and use the word "pussy" | 1384 |
| 1334 | anything wrong? | to describe women? </s> <insincere question | 1385 |
| 1335 | Label valid question | 3.497847080230713> Why are the UK govern- | 1386 |
| 1336 | Test Modified Is Ariana Grande really as mean | ment and the media (especially the BBC and the | 1387 |
| 1337 | and bitchy as she seems? </s> <insincere question | Guardian) demonising ordinary British people, ma- | 1388 |
| 1338 | 3.572277545928955> Why is Alia Bhatt so dumb? | nipulating buzz words like "alt-right", "Islama- | 1389 |
| 1339 | </s> <valid question 3.924571990966797> Do you | phobia", "racist" to suppress legitimate outrage | 1390 |
| 1340 | agree with Congressman Steve King's comments | at Muslim grooming gangs? </s> <valid question | 1391 |
| 1341 | on immigrant children in detention centers? | 3.845909357070923> Do you agree with Congress- | 1392 |
| 1342 | Label insincere question | man Steve King's comments on immigrant children | 1393 |
| 1343 | Test Modified Do you guys know that aliens | in detention centers? | 1394 |
| 1344 | are real and all those satellites we send up in | Label insincere question | 1395 |
| 1345 | space work as a sort of tracking device for them | Insincere Questions Step 5 | 1396 |
| 1346 | so in a few years it will be too late for Earth? | Test Modified What rapper still relevant and | 1397 |
| 1347 | </s> <insincere question 3.6094439029693604> | popular today has the best rhyme schemes? </s> | 1398 |
| 1348 | Have you noticed how conservatives are captur- | <insincere question 3.871907949447632> What | 1399 |
| 1349 | ing the English language and modifying the def- | would be a good nickname for Trump, Donald | 1400 |
| 1350 | initions of political words? </s> <valid question | Dumbck, and President Spankovich? </s> <valid | 1401 |
| 1351 | 3.6901655197143555> Do you agree with Con- | question 4.028958797454834> Why does Danc- | 1402 |
| 1352 | gressman Steve King's comments on immigrant | ing with the Stars not include Bachata as one their | 1403 |
| 1353 | children in detention centers? | dance styles? | 1404 |
| 1354 | Label insincere question | Label valid question | 1405 |
| 1355 | Test Modified Is it politically incorrect to say | Test Modified Which books do you sug- | 1406 |
| 1356 | female privilege, but it is a more accurate term to | gest to someone who get a free time and will | 1407 |
| 1357 | say, white female privilege? </s> <insincere ques- | help him stay motivated? </s> <valid question | 1408 |
| 1358 | tion 3.323280096054077> Why are the UK gov- | 3.6081225872039795> What is a good degree to | 1409 |
| 1359 | ernment and the media (especially the BBC and | get at community college if you want to explore dif- | 1410 |
| 1360 | the Guardian) demonising ordinary British people, | ferent subjects and figure out your career path? </s> | 1411 |
| 1361 | manipulating buzz words like "alt-right", "Islama- | <insincere question 3.8502604961395264> What | 1412 |
| 1362 | phobia", "racist" to suppress legitimate outrage | are the more steps in Career Oriented Education? | 1413 |
| 1363 | at Muslim grooming gangs? </s> <valid question | Label valid question | 1414 |
| 1364 | 3.986680269241333> Do you agree with Congress- | Test Modified How will you feel if someone | 1415 |
| 1365 | man Steve King's comments on immigrant children | talks badly about Kunti? </s> <valid question | 1416 |
| 1366 | in detention centers? | 3.5355563163757324> How do I stop feeling bad | 1417 |
| 1367 | Label insincere question | after a girl had a crush on me? </s> <insincere | 1418 |
| 1368 | Test Modified On Mother's Day, is it reasonable | question 3.689171075820923> Why Indian girls | 1419 |
| 1369 | to reflect there is some truth in the unfashionable | go crazy about marrying Shri. Rahul Gandhi ji? | 1420 |
| 1370 | notion than women are more driven by emotion | Label valid question | 1421 |
| 1371 | and men more driven by reason? </s> <insincere | Test Modified Why is equine HYPP inherited? | 1422 |
| 1372 | question 3.499204158782959> Why are the UK | </s> <insincere question 3.6035702228546143> | 1423 |
| 1373 | government and the media (especially the BBC and | Can female animals with male humans sex? </s> | 1424 |
| 1374 | the Guardian) demonising ordinary British people, | <valid question 3.7413032054901123> How long | 1425 |
| 1375 | manipulating buzz words like "alt-right", "Islama- | do guinea pigs live for? | 1426 |
| 1376 | phobia", "racist" to suppress legitimate outrage | Label valid question | 1427 |
| 1377 | at Muslim grooming gangs? </s> <valid question | Test Modified How do the Valerie Stevens | 1428 |
| 1378 | 3.771740198135376> Do you agree with Congress- | leather jackets achieve their quality during the | 1429 |
| 1379 | man Steve King's comments on immigrant children | manufacturing process? </s> <valid question | 1430 |
| 1380 | in detention centers? | 2.747288227081299> How are the Lancaster | 1431 |
| 1381 | Label insincere question | leather sofas manufactured? </s> <insincere ques- | 1432 |
| 1382 | Test Modified If the U.S. president is a | tion 3.944884777069092> Why don't all Trump | 1433 |
| 1383 | role model, is it acceptable for children to say | supporters buy only made in USA goods, e.g. many | 1434 |

| | | | |
|------|---|---|------|
| 1435 | of them have their cars of Asian/European compa- | Label insincere question | 1486 |
| 1436 | nies, shop in places where more than 70 of items | Insincere Questions Step 10 | 1487 |
| 1437 | are not made in USA, eat multi-national cuisine or | Test Modified What rapper still relevant and | 1488 |
| 1438 | otherwise stop their hypocrisy? | popular today has the best rhyme schemes? </s> | 1489 |
| 1439 | Label valid question | <valid question 3.7103171348571777> What is the | 1490 |
| 1440 | Test Modified Is Ariana Grande really as | oldest fashion trends running yet? </s> <insincere | 1491 |
| 1441 | mean and bitchy as she seems? </s> <insin- | question 3.871907949447632> What would be a | 1492 |
| 1442 | cere question 3.3252298831939697> Why is | good nickname for Trump, Donald Dumbck, and | 1493 |
| 1443 | Alia Bhatt so dumb? </s> <valid question | President Spankovich? | 1494 |
| 1444 | 3.7413415908813477> How do I stop feeling bad | Label valid question | 1495 |
| 1445 | after a girl had a crush on me? | Test Modified Which books do you sug- | 1496 |
| 1446 | Label insincere question | gest to someone who get a free time and will | 1497 |
| 1447 | Test Modified Do you guys know that aliens are | help him stay motivated? </s> <valid question | 1498 |
| 1448 | real and all those satellites we send up in space | 3.1401429176330566> How can I stay motivated | 1499 |
| 1449 | work as a sort of tracking device for them so in a | when learning something new? </s> <insincere | 1500 |
| 1450 | few years it will be too late for Earth? </s> <insin- | question 3.7235560417175293> I'm hungry and | 1501 |
| 1451 | cere question 3.0673365592956543> Isn't it obvi- | I'm too lazy too get out of bed, should I get a psy- | 1502 |
| 1452 | ous now that walking on the moon by the Ameri- | chologist or ask you questions? | 1503 |
| 1453 | cans was a hoax, because walking on the bright | Label valid question | 1504 |
| 1454 | side of the moon, even in a space suit would be | Test Modified How will you feel if some- | 1505 |
| 1455 | fatal? </s> <valid question 3.1978228092193604> | one talks badly about Kunti? </s> <insincere | 1506 |
| 1456 | Why do we weunch satellites? | question 3.4893462657928467> Does Tamil Isai | 1507 |
| 1457 | Label insincere question | Soundarajan support Vijayendra for disrespect- | 1508 |
| 1458 | Test Modified Is it politically incorrect to say | ing the Tamil Anthem? </s> <valid question | 1509 |
| 1459 | female privilege, but it is a more accurate term to | 3.5355563163757324> How do I stop feeling bad | 1510 |
| 1460 | say, white female privilege? </s> <insincere ques- | after a girl had a crush on me? | 1511 |
| 1461 | tion 2.9176812171936035> How does the privi- | Label valid question | 1512 |
| 1462 | lege of being attractive compare to the privilege | Test Modified Why is equine HYPP inher- | 1513 |
| 1463 | of being White in the US? </s> <valid question | ited? </s> <valid question 3.5067965984344482> | 1514 |
| 1464 | 3.112481117248535> Is the media wrong for en- | What disadvantages do animals that don't | 1515 |
| 1465 | forcing gender stereotypes? | have bones face? </s> <insincere question | 1516 |
| 1466 | Label insincere question | 3.6035702228546143> Can female animals with | 1517 |
| 1467 | Test Modified On Mother's Day, is it reasonable | male humans sex? | 1518 |
| 1468 | to reflect there is some truth in the unfashionable | Label valid question | 1519 |
| 1469 | notion than women are more driven by emotion | Test Modified How do the Valerie Stevens | 1520 |
| 1470 | and men more driven by reason? </s> <insincere | leather jackets achieve their quality during the | 1521 |
| 1471 | question 3.102353811264038> Do women look | manufacturing process? </s> <valid question | 1522 |
| 1472 | down on men who are single, even if the man is | 2.747288227081299> How are the Lancaster | 1523 |
| 1473 | more successful in other aspects of his life? </s> | leather sofas manufactured? </s> <insincere | 1524 |
| 1474 | <valid question 3.1890125274658203> Why are | question 3.9087233543395996> Are Newport | 1525 |
| 1475 | some women uninterested in sex? | cigarettes designed to selectively destroy black peo- | 1526 |
| 1476 | Label insincere question | ple's DNA? | 1527 |
| 1477 | Test Modified If the U.S. president is a | Label valid question | 1528 |
| 1478 | role model, is it acceptable for children to say | Test Modified Is Ariana Grande really as mean | 1529 |
| 1479 | "go fuck yourself," and use the word "pussy" | and bitchy as she seems? </s> <valid question | 1530 |
| 1480 | to describe women? </s> <insincere question | 3.183567762374878> I like this girl who used to | 1531 |
| 1481 | 3.163693904876709> Is it wrong to take your | be quite rude and would run through boyfriends | 1532 |
| 1482 | retarded son to a hooker for his 21st birthday? | very fast. But now that school started again, | 1533 |
| 1483 | </s> <valid question 3.456286907196045> Do you | she seems to have gotten a lot nicer through- | 1534 |
| 1484 | agree with Congressman Steve King's comments | out Summer. Is she faking her politeness, and | 1535 |
| 1485 | on immigrant children in detention centers? | is it worth pursuing her? </s> <insincere ques- | 1536 |

| | | | |
|------|--|--|------|
| 1537 | tion 3.3253660202026367> Why is Alia Bhatt so | because It worked.""" | 1588 |
| 1538 | dumb? | | |
| 1539 | Label insincere question | Label not-counterfactual | 1589 |
| 1540 | Test Modified Do you guys know that aliens are | Test Modified I like these jeans they sit | 1590 |
| 1541 | real and all those satellites we send up in space | low enough without being inappropriate when | 1591 |
| 1542 | work as a sort of tracking device for them so in a | you sit or bend over. </s> <counterfactual | 1592 |
| 1543 | few years it will be too late for Earth? </s> <insincere | 3.402600049972534> "But oddly enough, the bot- | 1593 |
| 1544 | question 3.0673365592956543> Isn't it obvi- | oms are a little too loose in the waist (37) and could | 1594 |
| 1545 | ous now that walking on the moon by the Amer- | have used another inch or two in the inseam (I nor- | 1595 |
| 1546 | icans was a hoax, because walking on the bright | mally take a 35"" or 36"" in jeans, depending on | 1596 |
| 1547 | side of the moon, even in a space suit would be | the brand if this helps)."" </s> <not-counterfactual | 1597 |
| 1548 | fatal? </s> <valid question 3.1978228092193604> | 3.4201438426971436> These boxer-briefs are very | 1598 |
| 1549 | Why do we weunch satellites? | soft, very comfortable, and fit like high-end under- | 1599 |
| 1550 | Label insincere question | wear the likes of which you might get at, oh, say, | 1600 |
| 1551 | Test Modified Is it politically incorrect to say | Calvin Klein for example, but for about half the | 1601 |
| 1552 | female privilege, but it is a more accurate term to | price. | 1602 |
| 1553 | say, white female privilege? </s> <insincere ques- | Label not-counterfactual | 1603 |
| 1554 | tion 2.9176158905029297> How does the privi- | Test Modified He was very professional and | 1604 |
| 1555 | lege of being attractive compare to the privilege | wish all transactions I make through Amazon were | 1605 |
| 1556 | of being White in the US? </s> <valid question | this good. </s> <counterfactual 3.4319908618927> | 1606 |
| 1557 | 3.112481117248535> Is the media wrong for en- | I wish I had had him as an instructor at college. | 1607 |
| 1558 | forcing gender stereotypes? | </s> <not-counterfactual 4.054030895233154> | 1608 |
| 1559 | Label insincere question | I worried that it would be cheap or not fit | 1609 |
| 1560 | Test Modified On Mother's Day, is it reason- | or...whatever...But WOW! | 1610 |
| 1561 | able to reflect there is some truth in the unfash- | Label not-counterfactual | 1611 |
| 1562 | ionable notion than women are more driven by | Test Modified Well written with a twist | 1612 |
| 1563 | emotion and men more driven by reason? </s> <in- | I didn't expect. </s> <not-counterfactual | 1613 |
| 1564 | sincere question 2.9901626110076904> Do you | 3.3257973194122314> "The crossover from the | 1614 |
| 1565 | agree that females think with their brains and | characters from one novel to others keeps me in- | 1615 |
| 1566 | males with their testicles? </s> <valid question | terested; after all, I do hate to miss a Dee-Ann | 1616 |
| 1567 | 3.1890125274658203> Why are some women un- | or Eggie"" appearance."" </s> <counterfactual | 1617 |
| 1568 | interested in sex? | 3.6820030212402344> "Had I reviewed this im- | 1618 |
| 1569 | Label insincere question | mediately I would have given this product five stars | 1619 |
| 1570 | Test Modified If the U.S. president is a | because It worked.""" | 1620 |
| 1571 | role model, is it acceptable for children to say | Label not-counterfactual | 1621 |
| 1572 | "go fuck yourself," and use the word "pussy" | Test Modified Doesn't feel like the quality | 1622 |
| 1573 | to describe women? </s> <insincere question | levi's I am used to. </s> <not-counterfactual | 1623 |
| 1574 | 2.994286298751831> Why do feminists let their | 3.2773308753967285> However, the fabric is not | 1624 |
| 1575 | daughters have sex with their boyfriend's at home? | that great, it's cheap scratchy cotton. </s> <counter- | 1625 |
| 1576 | </s> <valid question 3.456286907196045> Do you | factual 3.746659755706787> The blanket is nice | 1626 |
| 1577 | agree with Congressman Steve King's comments | and soft but it is white, so if it doesn't light up it | 1627 |
| 1578 | on immigrant children in detention centers? | isn't much use! | 1628 |
| 1579 | Label insincere question | Label not-counterfactual | 1629 |
| 1580 | Amazon Counterfactual Step 1 | Test Modified If we had wall studs, I believe | 1630 |
| 1581 | Test Modified Clings to the wall, doesn't flop | the enclosed hardware would have been sufficient. | 1631 |
| 1582 | around when a bag is pulled out, the mess of | </s> <counterfactual 3.4338643550872803> i wish | 1632 |
| 1583 | bags falling out is gone. </s> <not-counterfactual | the storage compartment was a little bigger and | 1633 |
| 1584 | 3.6492726802825928> Hopes that it will keep | opened up instead of slidding on and off. </s> <not- | 1634 |
| 1585 | it's shape after washing. </s> <counterfactual | counterfactual 3.9785308837890625> I worried | 1635 |
| 1586 | 4.012346267700195> "Had I reviewed this im- | that it would be cheap or not fit or...whatever...But | 1636 |
| 1587 | mediately I would have given this product five stars | WOW! | 1637 |
| | | Label counterfactual | 1638 |

| | | | |
|------|--|--|------|
| 1639 | Test Modified If this ever turns into a film, I | Label not-counterfactual | 1690 |
| 1640 | hope they do it justice! </s> <not-counterfactual | Test Modified I like these jeans they sit | 1691 |
| 1641 | 3.5291523933410645> "The crossover from the | low enough without being inappropriate when | 1692 |
| 1642 | characters from one novel to others keeps me inter- | you sit or bend over. </s> <counterfactual | 1693 |
| 1643 | ested; after all, I do hate to miss a Dee-Ann | 2.606198310852051> "But oddly enough, the bot- | 1694 |
| 1644 | or Eggie"" appearance."" </s> <counterfactual | toms are a little too loose in the waist (37) and could | 1695 |
| 1645 | 3.751143217086792> "Had I reviewed this imme- | have used another inch or two in the inseam (I nor- | 1696 |
| 1646 | diately I would have given this product five stars | normally take a 35"" or 36"" in jeans, depending on | 1697 |
| 1647 | because It worked."" </s> <not-counterfactual | the brand if this helps)."" </s> <not-counterfactual | 1698 |
| 1648 | Label counterfactual | 2.6380045413970947> A tad loose but I rather | 1699 |
| 1649 | Test Modified If you don't want a prominent | have it fit this way than too tight. | 1700 |
| 1650 | display this rack is too large for most bed or living | Label not-counterfactual | 1701 |
| 1651 | rooms, it is wider and taller than my tall Broy- | Test Modified He was very professional | 1702 |
| 1652 | hill wardrobe style dresser which was the largest | and wish all transactions I make through Ama- | 1703 |
| 1653 | piece in the room until this shoe rack. </s> <not- | zon were this good. </s> <not-counterfactual | 1704 |
| 1654 | counterfactual 3.865670680999756> "It also vali- | 3.3291680812835693> This new speaker was | 1705 |
| 1655 | dates the incorrect"" assumption that we are alone | just what the doctor ordered and I couldn't | 1706 |
| 1656 | in the feelings we suppress when we sense the com- | be more pleased. </s> <counterfactual | 1707 |
| 1657 | plete garbage that is thrown out into society."" </s> | 3.4589436054229736> Had the person han- | 1708 |
| 1658 | <counterfactual 4.063361167907715> The blanket | dling the shipping of this item been at all | 1709 |
| 1659 | 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 | 1710 |
| 1660 | up it isn't much use! | the mailing process, the slightest bit of care could | 1711 |
| 1661 | Label counterfactual | have been taken to ensure it's proper delivery. | 1712 |
| 1662 | Test Modified I wish I could have seen all of | Label not-counterfactual | 1713 |
| 1663 | the places he recommends! </s> <counterfactual | Test Modified Well written with a twist | 1714 |
| 1664 | 3.5627076625823975> I wish I had had him as | I didn't expect. </s> <not-counterfactual | 1715 |
| 1665 | an instructor at college. </s> <not-counterfactual | 2.651658535003662> The book had some interest- | 1716 |
| 1666 | 4.141315937042236> I worried that it would be | ing twists that I did see coming and I look forward | 1717 |
| 1667 | cheap or not fit or...whatever...But WOW! | to reading part two of this series. </s> <counterfac- | 1718 |
| 1668 | Label counterfactual | tual 2.8373162746429443> Fun read Could have | 1719 |
| 1669 | Test Modified I wish I could replace just that | been a little longer with more detail. | 1720 |
| 1670 | small stupid piece, since there's nothing wrong | Label not-counterfactual | 1721 |
| 1671 | with the rest of the hose assembly. </s> <count- | Test Modified Doesn't feel like the qual- | 1722 |
| 1672 | erfactual 3.6057372093200684> i wish the stor- | ity levi's I am used to. </s> <counterfactual | 1723 |
| 1673 | age compartment was a little bigger and opened | 2.733877182006836> It has the same great com- | 1724 |
| 1674 | up instead of sliding on and off. </s> <not- | fortable flattering features plus the great denim tex- | 1725 |
| 1675 | counterfactual 4.064871311187744> I worried that | ture that Lee has perfected- smoothing and stretchy | 1726 |
| 1676 | it would be cheap or not fit or...whatever...But | without the excessive cling- but I think it must have | 1727 |
| 1677 | WOW! | been designed for people who have a greater sur- | 1728 |
| 1678 | Label counterfactual | plus of belly fat than I. </s> <not-counterfactual | 1729 |
| 1679 | Amazon Counterfactual Step 5 | 2.856729745864868> Will keep but won't be that | 1730 |
| 1680 | Test Modified Clings to the wall, doesn't flop | casual sexy top you always want to turn to. | 1731 |
| 1681 | around when a bag is pulled out, the mess of | Label not-counterfactual | 1732 |
| 1682 | bags falling out is gone. </s> <not-counterfactual | Test Modified If we had wall studs, I believe the | 1733 |
| 1683 | 3.161406993865967> And the dvd cases were | enclosed hardware would have been sufficient. </s> | 1734 |
| 1684 | tightly packed to ensure they didn't move around. | <not-counterfactual 2.6638145446777344> It was | 1735 |
| 1685 | </s> <counterfactual 3.308583974838257> The | a little tricky to find the center of the studs using | 1736 |
| 1686 | case is small, cord seems to always want to stay | my stud finder but once I felt comfortable with | 1737 |
| 1687 | kinked and coiled, plug should be angled and | the lines I had drawn, I drilled the pilot holes and | 1738 |
| 1688 | not straight...which are all items that others have | bolted this thing to the wall. </s> <counterfactual | 1739 |
| 1689 | pointed out. | 2.879924774169922> The only thing I would have | 1740 |

| | | | |
|------|--|---|------|
| 1741 | like for it to have a hole in the middle so I can put | </s> <counterfactual 3.289605140686035> If I had | 1792 |
| 1742 | the stopper in without removing the mat. | to come up with anything negative, I would say that | 1793 |
| 1743 | Label counterfactual | the attachments don't seem to stay on the vacuum | 1794 |
| 1744 | Test Modified If this ever turns into a film, I | cleaner when not in use - but that could be me not | 1795 |
| 1745 | hope they do it justice! </s> <not-counterfactual | putting them on properly! | 1796 |
| 1746 | 2.671574354171753> I read this book because of | Label not-counterfactual | 1797 |
| 1747 | the motion picture that is coming out soon. </s> | Test Modified I like these jeans they sit | 1798 |
| 1748 | <counterfactual 3.1458709239959717> Was a good | low enough without being inappropriate when | 1799 |
| 1749 | story, though there could have been more to it. | you sit or bend over. </s> <not-counterfactual | 1800 |
| 1750 | Label counterfactual | 2.447404623031616> These shorts fit really | 1801 |
| 1751 | Test Modified If you don't want a prominent | well and look good too. </s> <counterfactual | 1802 |
| 1752 | display this rack is too large for most bed or | 2.550638198852539> The top fits great just wish | 1803 |
| 1753 | living rooms, it is wider and taller than my tall | the bottoms fit too. | 1804 |
| 1754 | Broyhill wardrobe style dresser which was the | Label not-counterfactual | 1805 |
| 1755 | largest piece in the room until this shoe rack. </s> | Test Modified He was very professional | 1806 |
| 1756 | <counterfactual 2.7353768348693848> I bought | and wish all transactions I make through Ama- | 1807 |
| 1757 | this mount because I wanted one that would sit on | zon were this good. </s> <not-counterfactual | 1808 |
| 1758 | three studs instead of two because my TV is quite | 3.3291127681732178> This new speaker was | 1809 |
| 1759 | heavy and I would have had a hard time centering | just what the doctor ordered and I couldn't | 1810 |
| 1760 | it on my wall if I didn't have the wide hanging | be more pleased. </s> <counterfactual | 1811 |
| 1761 | rail that this one has. </s> <not-counterfactual | 3.3897111415863037> But the author alle- | 1812 |
| 1762 | 2.873617172241211> Good for under the bed shoe | viated my concerns quickly with a few well-timed | 1813 |
| 1763 | storage, IF the wife wants to use it. | comments about how it was the man could have | 1814 |
| 1764 | Label counterfactual | known that the arrangement was something Jack | 1815 |
| 1765 | Test Modified I wish I could have seen all of | wanted. | 1816 |
| 1766 | the places he recommends! </s> <counterfactual | Label not-counterfactual | 1817 |
| 1767 | 2.799947738647461> I wish I had had him as | Test Modified Well written with a twist | 1818 |
| 1768 | an instructor at college. </s> <not-counterfactual | I didn't expect. </s> <not-counterfactual | 1819 |
| 1769 | 3.3013432025909424> And as the ole man isn't | 2.557446002960205> "A bit workmanlike, not | 1820 |
| 1770 | any version of slender it was good that he got to try | up to Lord's high standard of A Night to Re- | 1821 |
| 1771 | on some shirts before hand. | member,"" but well-detailed, and a story that | 1822 |
| 1772 | Label counterfactual | not many now know.""" </s> <counterfactual | 1823 |
| 1773 | Test Modified I wish I could replace just that | 2.792485475540161> Wow I am really glad I | 1824 |
| 1774 | small stupid piece, since there's nothing wrong | didn't read these reviews BEFORE I read this | 1825 |
| 1775 | with the rest of the hose assembly. </s> <counter- | book because I would have passed on the book | 1826 |
| 1776 | terfactual 2.628289222717285> The only thing I | and missed a really great start to a series that cap- | 1827 |
| 1777 | would have like for it to have a hole in the middle so | tured my attention and made me laugh all the while | 1828 |
| 1778 | I can put the stopper in without removing the mat. | using my imagination and painting a clear picture | 1829 |
| 1779 | </s> <not-counterfactual 2.9200568199157715> | of the author's world she was building for us. | 1830 |
| 1780 | The only downside is my laptop does not have | Label not-counterfactual | 1831 |
| 1781 | the screw holes on it and the screws do not retract | Test Modified Doesn't feel like the qual- | 1832 |
| 1782 | far enough back for me to push the connector all | ity levi's I am used to. </s> <counterfactual | 1833 |
| 1783 | the way in, but a simple smash will rid that issue | 2.5612902641296387> i was hoping the pants | 1834 |
| 1784 | (this thing is durable!) | would be thicker but being that it's not too expen- | 1835 |
| 1785 | Label counterfactual | sive it's understandable. </s> <not-counterfactual | 1836 |
| 1786 | Amazon Counterfactual Step 10 | 2.572395086288452> But it doesn't have a lining | 1837 |
| 1787 | Test Modified Clings to the wall, doesn't flop | like the last couple models I bought. | 1838 |
| 1788 | around when a bag is pulled out, the mess of | Label not-counterfactual | 1839 |
| 1789 | bags falling out is gone. </s> <not-counterfactual | Test Modified If we had wall studs, I believe | 1840 |
| 1790 | 3.161406993865967> And the dvd cases were | the enclosed hardware would have been sufficient. | 1841 |
| 1791 | tightly packed to ensure they didn't move around. | </s> <not-counterfactual 2.6638145446777344> It | 1842 |

1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893

was a little tricky to find the center of the studs using my stud finder but once I felt comfortable with the lines I had drawn, I drilled the pilot holes and bolted this thing to the wall. </s> <counterfactual 2.771395206451416> Wish it had a little more padding, otherwise just as advertised.

Label counterfactual

Test Modified If this ever turns into a film, I hope they do it justice! </s> <not-counterfactual 2.671574354171753> I read this book because of the motion picture that is coming out soon. </s> <counterfactual 3.141676187515259> Wish this story would have been longer and turned into a book, with some gut wrenching action, love/hate lovers quarrels scenes, with a happy ending at the end...

Label counterfactual

Test Modified If you don't want a prominent display this rack is too large for most bed or living rooms, it is wider and taller than my tall Broyhill wardrobe style dresser which was the largest piece in the room until this shoe rack. </s> <counterfactual 2.7353768348693848> I bought this mount because I wanted one that would sit on three studs instead of two because my TV is quite heavy and I would have had a hard time centering it on my wall if I didn't have the wide hanging rail that this one has. </s> <not-counterfactual 2.873617172241211> Good for under the bed shoe storage, IF the wife wants to use it.

Label counterfactual

Test Modified I wish I could have seen all of the places he recommends! </s> <counterfactual 2.7999041080474854> I wish I had had him as an instructor at college. </s> <not-counterfactual 3.2604622840881348> I wanted to order him a few affordable hats I wouldn't mind him loosing.

Label counterfactual

Test Modified I wish I could replace just that small stupid piece, since there's nothing wrong with the rest of the hose assembly. </s> <counterfactual 2.474032402038574> I wish I could just hook up a hose rather than connecting and routing an tube internally </s> <not-counterfactual 2.9200568199157715> The only downside is my laptop does not have the screw holes on it and the screws do not retract far enough back for me to push the connector all the way in, but a simple smash will rid that issue (this thing is durable!)

Label counterfactual

Toxic Conversations Step 1

Test Modified Homos demand that you accept their fudge packing. But none of us ever will- </s> <not toxic 3.5619421005249023> Sounds just awful, working for the state that is. So it begs the question, why work for the state if the pay is so bad versus the private sector? Seems logical to just make the switch? </s> <toxic 4.056504249572754> I would imagine it is a miserable existence. Being a whining sore loser is no way to go through life.

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> <toxic 3.6495330333709717> Calling everybody that disagrees with you a racist and anti-Muslim is pretty uncivil. Also intellectually bankrupt... </s> <not toxic 3.8587634563446045> Sounds just awful, working for the state that is. So it begs the question, why work for the state if the pay is so bad versus the private sector? Seems logical to just make the switch?

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=opJzRWWVDdw>

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 3.8009092807769775> Sounds just awful, working for the state that is. So it begs the question, why work for the state if the pay is so bad versus the private sector? Seems logical to just make the switch? </s> <toxic 4.463983535766602> I would imagine it is a miserable existence. Being a whining sore

1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944

1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995

loser is no way to go through life.

Label not toxic

Test Modified Get the "real news" from Breitbart and Drudge Report? Ahahahahaha! Are you for real? </s> <not toxic 3.559427499771118> Sounds just awful, working for the state that is. So it begs the question, why work for the state if the pay is so bad versus the private sector? Seems logical to just make the switch? </s> <toxic 4.239322662353516> I would imagine it is a miserable existence. Being a whining sore loser is no way to go through life.

Label not toxic

Test Modified Wud luv to see one of those American Boats with all the planes on it at the bottom of the sea. </s> <not toxic 3.8268911838531494> Sounds just awful, working for the state that is. So it begs the question, why work for the state if the pay is so bad versus the private sector? Seems logical to just make the switch? </s> <toxic 4.241443157196045> I would imagine it is a miserable existence. Being a whining sore loser is no way to go through life.

Label not toxic

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 3.6674933433532715> Sounds just awful, working for the state that is. So it begs the question, why work for the state if the pay is so bad versus the private sector? Seems logical to just make the switch? </s> <toxic 4.0054121017456055> I would imagine it is a miserable existence. Being a whining sore loser is no way to go through life.

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 3.262410879135132> Calling everybody that disagrees with you a racist and anti-Muslim is pretty uncivil. Also intellectually bankrupt... </s> <not toxic 3.8916428089141846> It's always important to remember what can happen when you have so-

ciopaths as leaders and also have compliant followers. 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 disagrees 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 disagrees 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" 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 3.3616135120391846> Calling everybody that disagrees with you a racist and anti-Muslim is pretty uncivil. Also intellectually bankrupt... </s> <not toxic 3.903903007507324> I have very high respect 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-

1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046

2047 min and unfair union). But the majority of teachers 2098
2048 here instead of doing their jobs they band together 2099
2049 wear purple shirts and mob the government for a 2100
2050 better contract, and no accountability in the form 2101
2051 of testing students. Many teachers are middle class 2102
2052 kids that took the path of least resistance in what 2103
2053 was expected by their parents (college) and because 2104
2054 they lacked drive ended up teachers. That lack of 2105
2055 drive shows by what the private sector taxpayers 2106
2056 get for their money. Your degrees mean nothing if 2107
2057 you don't do your job. 2108

2058 **Label** toxic 2109

2059 **Toxic Conversations Step 5** 2110

2060 **Test Modified** Homos demand that you accept 2111
2061 their fudge packing. But none of us ever will- </s> 2112
2062 <toxic 3.1383495330810547> So you admit you 2113
2063 would exterminate inferior humans. </s> <not 2114
2064 toxic 3.2954952716827393> Mark MacKinnon 2115
2065 and the interests he work for would like us to 'get 2116
2066 used to it', because they don't want to do anything 2117
2067 practical to stop it. 2118

2068 **Label** not toxic 2119

2069 **Test Modified** I don't think anyone likes this 2120
2070 health care bill, it stinks for everyone. 50 years 2121
2071 and older are going to get hammered with higher 2122
2072 premiums. People with preexisting conditions will 2123
2073 also see their premiums go through the roof. Even- 2124
2074 tually no one will be able to afford it. They're still 2125
2075 not addressing the reason and the solution. Pre- 2126
2076 scription drugs. Why are Americans paying 2 to 2127
2077 3 times higher for prescription drugs than Cana- 2128
2078 dians are? They should sit down with Democrats 2129
2079 and fix it not play musical chairs with our health 2130
2080 care. Here is the add that AARP put out, and it 2131
2081 has a fact check at the end. It's worth watching. 2132
2082 [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) 2133
2083 [aarp-ad-says-house-gop-health-care-bill-would-](http://www.thedenverchannel.com/news/politics/aarp-ad-says-house-gop-health-care-bill-would-boost-insurance-rates-for-older-americans) 2134
2084 [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) </s> 2135
2085 <not toxic 2.5086519718170166> so in the mean 2136
2086 time tens of thousands of Oregonians go without 2137
2087 health insurance which will now be unaffordable to 2138
2088 them. And sorry, the republicans have had 8 years 2139
2089 to figure out a better system, they aren't going to 2140
2090 do it anytime soon. Neither party has any desire to 2141
2091 actually find a solution to all this. Hopefully Trump 2142
2092 will also soon eliminate the tax penalty for not 2143
2093 having insurance so all us folks who buy our own 2144
2094 insurance don't get penalized for not being able to 2145
2095 afford to buy the insurance we are required to have. 2146
2096 Of course he is probably clueless that detail even 2147
2097 exists </s> <toxic 2.801957607269287> reducing 2148

number of brackets. Another is lowering corporate 2098
tax rates which would be OK if all loopholes, 2099
including tax-free political donations for wealthy 2100
people only. Another is doubling the earned 2101
income tax credit which would help families with 2102
children but for people like me, would not make 2103
up for loss of the state tax deduction. Essentially 2104
the proposed tax "reform bill as it stands is a 2105
huge wealth transfer, from working people and 2106
the treasury to the offshore accounts of wealthy 2107
people and corporations. Median earners like me 2108
will be screwed. Poor people with children will be 2109
screwed. The rich will get a huge bundle of cash, 2110
and the Republicans will lead us further into deep 2111
debt. 2112

That's if Ryan and McConnell can pass this huge 2113
sack of excrement. 2114

Label not toxic 2115

Test Modified Rather than call you a liar, I'll 2116
just post a link to HART's helicopter footage of 2117
this "most densely populated" ag land. 2118
[https://www.youtube.com/watch?v=opJzR-](https://www.youtube.com/watch?v=opJzR-WWVDdw) 2119
[WWVDdw](https://www.youtube.com/watch?v=opJzR-WWVDdw) 2120

And if you've been following a little closer, 2121
you'd know about the luxury towers and boutique 2122
hotels that developers are planning on fulfilling 2123
their TOD requirement with "cash gifts" to the 2124
city instead of actual affordable TOD. </s> <not 2125
toxic 3.0351996421813965> So are you saying 2126
we should build the road. Lots of doom and 2127
gloom but I'm missing your point. </s> <toxic 2128
3.1464040279388428> Isn't that terrible? Cald- 2129
well and Inouye do not know the amount a 10-year 2130
extension would raise yet they are proposing it. 2131
Caldwell asked for two years and he got it. Two 2132
years later he is now asking for ten years. Now, 2133
which is only weeks later, HART director is claim- 2134
ing 10 years is insufficient. This scenario shows 2135
our mayor and HART do not know the scope of 2136
this project. Pathetic! 2137

Label not toxic 2138

Test Modified Get the "real news" from Breit- 2139
bart and Drudge Report? Ahahahahaha! Are you 2140
for real? </s> <toxic 2.634126901626587> "If one 2141
read the Dispatch one would think Trump is the 2142
most evil person on the planet." Not evil, just id- 2143
iotic. And it would be easy to give his behaviour 2144
a pass if he wasn't POTUS. ".....five to eight anti- 2145
Trump stories per day. Never any good one's or 2146
one's that just stuck to the facts." Well, when there's 2147
good a news Trump story to print, I'm sure ADN 2148

2149 will be all over it, problem is, there's been a dearth 2200
2150 of those since his election. Facts? Ok Rich, give 2201
2151 us a list of incorrect facts in the above story. </s> 2202
2152 <not toxic 3.0248677730560303> With dismay I 2203
2153 noticed that ADN had printed yet another column 2204
2154 from Fox commentator Charles Krauthammer but 2205
2155 after reading it I'm glad the editors chose it as the 2206
2156 feature article on the opinion page. Krauthammer 2207
2157 is also a psychiatrist so his analyses of Trump as 2208
2158 a man who has never emotionally, intellectually 2209
2159 developed beyond adolescence holds some weight. 2210
2160 But what does it say about Trumps supporters that 2211
2161 so many millions of them can't see through the 2212
2162 boorish, confrontational attitude of the man? How 2213
2163 can so many Americans have devolved into such 2214
2164 anger, fear and irrationality that they would/could 2215
2165 find redemption in Trump after how he has exposed 2216
2166 his true narcissistic self for all to see. When you've 2217
2167 lost the Jennifer Rubins and Charles Krauthammer's 2218
2168 of the media world you've lost the battle 2219
2169 yet the Trumpian cult members will soldier on and 2220
2170 then become even angrier and more full of fear after 2221
2171 the election. Something to do with their choice 2222
2172 of "information" sources no doubt. 2223
2173 **Label** not toxic 2224
2174 **Test Modified** Wud luv to see one of those American 2225
2175 Boats with all the planes on it at the bottom 2226
2176 of the sea. </s> <toxic 3.3901054859161377> I 2227
2177 bet Regent Seven Seas will never offer Mr Hammond 2228
2178 another trip. Wow, what a snarky article. He 2229
2179 makes , I assume, some valid points about food 2230
2180 and atmosphere. However, after discovering the 2231
2181 treats available on his "massive deck" he "blew 2232
2182 off" his remaining restaurant reservations , donned 2233
2183 his comfy bathrobe and ordered-in. He was certainly 2234
2184 not an ideal passenger and, for one floating 2235
2185 on a freebie, he's a total ingrate! </s> <not toxic 2236
2186 3.409156084060669> Now replaced by the sexy 2237
2187 EA-18G Growler! Using a preexisting Military 2238
2188 Operating Area! Get over it!!!!!! 2239
2189 **Label** not toxic 2240
2190 **Test Modified** Wtf, nowhere in scripture does 2241
2191 it say that Elizabeth was a kid when married to 2242
2192 Zachariah and nowhere does it say that Mary 2243
2193 was a child when married to Joseph. That's a 2244
2194 complete and utter fabrication. In fact, it says 2245
2195 in Luke 1:42, "And she spake out with a loud 2246
2196 voice, and said, Blessed art thou among women, 2247
2197 and blessed is the fruit of thy womb." It says 2248
2198 blessed art thou among WOMEN. </s> <not toxic 2249
2199 2.769857406616211> I was informed that my first 2250

grandchild had been conceived the evening of the 2200
day when I had inserted a prayer note in the Wailing 2201
Wall in Jerusalem that asked God to help my 2202
daughter conceive after a year of frustrated attempts. 2203
Maybe Elizabeth did the same thing? After all, she was 2204
in the same neighborhood. :-) </s> 2205
<toxic 3.286393880844116> Christians who support 2206
Trump are the most mind-boggling to me. I just don't see 2207
how they square the circle between Trump and their moral 2208
foundations. 2209
"Beware then of useless grumbling, and keep 2210
your tongue from slander; because no secret word 2211
is without result, and a lying mouth destroys the 2212
soul." (Wisdom 1:11) 2213
If that is the case, then Trump's soul was utterly 2214
destroyed decades ago. 2215
Label toxic 2216
Test Modified Angela Merkel and all other European 2217
political leaders who have aided and abetted 2218
the ongoing invasion of Europe by the forces of 2219
the crescent moon death cult should be tried as 2220
accessories to Mr. Urban's murder. </s> <toxic 2221
3.2037758827209473> that's what happens when 2222
you betray the people of your country for foreign 2223
bs. let's go Le Pen, Geert Wilders. If the media 2224
refuses to mention the muslim crisis the total 2225
incompatibility of primitive, uneducated muslim 2226
males swarming countries and turning them into 2227
misogynistic fundamentalist religious areas then 2228
we need these people to save us from YOU! </s> 2229
<not toxic 3.2326502799987793> your first mistake 2230
is believing what a politician says because 2231
generally it has nothing to do with what they do. 2232
The Libs will be happy to let this die because 2233
Monsef is now a very poor salesman given her own 2234
immigration dishonesty. That said if the election 2235
prospects sour significantly for the Libs I have no 2236
doubts that PM Butts will ram through Ranked 2237
Ballot 2238
Label toxic 2239
Test Modified I hope you don't have kids if 2240
you see this woman's actions as acceptable. And 2241
I applaud the den for kicking the kid out. She 2242
brought unwanted negative attention upon them. 2243
However, she will, and is already likely, pay 2244
the the price for her stupid stunt. </s> <toxic 2245
2.8730525970458984> Ms. Van Brocklin: You 2246
state that there is Payne's conduct is inexcusable, 2247
then proceed to use the rest of your space to justify 2248
him and his actions. You have denigrated the proud 2249
courage of countless people who took up causes 2250

2251 via civil disobedience. I marched in the non-violent
2252 peaceful Civil Rights protest. So the dogs and the
2253 firehoses used by a certain southern sheriff were just-
2254 ified, by your logic. So were the citizens beaten by
2255 Chicago police during the Democratic Convention
2256 Police riots. Resolved in 20 minutes? Nonsense.
2257 If Miss Wubbels hadn't protested as she did, she'd
2258 likely have ended up in a back room, somewhere,
2259 and who know when she would have been granted
2260 the presence of a lawyer and what she would have
2261 gone through prior to that. You are completely ig-
2262 norant of the shortage of nurses in this country - in
2263 some cases, critical shortages. And why would any-
2264 one want to be a nurse when they are disrespected
2265 by a former state and federal prosecutor such as
2266 you. </s> <not toxic 2.92986798286438> Acquit
2267 her, then commit her. This womens cheese has slid
2268 so far off the cracker she's a danger to herself and
2269 to others.

2270 Animal rights activism is a just cause, but her
2271 and her group have gone off the deep end into
2272 radical extremist territory.

2273 **Label toxic**

2274 **Test Modified** no one cares what a paid lib-
2275 eral trolling hack like you believes lunatic,, </s>
2276 <toxic 2.6076736450195312> aa another hate filled
2277 left winger again! save the stupid nonsense sheep,
2278 trump is not causing anything, our weak leadership
2279 is. </s> <not toxic 2.810284376144409> Ouch...
2280 didn't see that one coming. A liberal stealing my
2281 own line... just like they take everything else they
2282 like.

2283 **Label toxic**

2284 **Test Modified** Ok all you NDP "LEAP" mani-
2285 festo types, where is your hero Naomi Klein? Her
2286 fawning adoration of Chavez and Venezuelan thug-
2287 gery knows no bounds. I'm sure she's awfully
2288 hysterical over the thought that such a pathetic
2289 dictstorship could ever be sanctioned. </s> <toxic
2290 2.941824436187744> Ms. Van Brocklin: You state
2291 that there is Payne's conduct is inexcusable, then
2292 proceed to use the rest of your space to justify him
2293 and his actions. You have denigrated the proud
2294 courage of countless people who took up causes
2295 via civil disobedience. I marched in the non-violent
2296 peaceful Civil Rights protest. So the dogs and the
2297 firehoses used by a certain southern sheriff were just-
2298 ified, by your logic. So were the citizens beaten by
2299 Chicago police during the Democratic Convention
2300 Police riots. Resolved in 20 minutes? Nonsense.
2301 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 ig-
norant of the shortage of nurses in this country - in
some cases, critical shortages. And why would any-
one 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 prac-
tical 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. Even-
tually no one will be able to afford it. They're still
not addressing the reason and the solution. Pre-
scription drugs. Why are Americans paying 2 to
3 times higher for prescription drugs than Cana-
dians 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](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

2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403

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. 2404

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. 2405
2406
2407
2408
2409
2410

Label not toxic 2411

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!!!!!! 2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426

Label not toxic 2427

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. 2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449

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- 2450
2451
2452
2453
2454

2455 ing sexual immorality? This sounds like something
2456 your example of a rapist priest could say to himself
2457 to justify himself taking Communion in a state of
2458 mortal sin.

2459 **Label toxic**

2460 **Test Modified** Angela Merkel and all other Euro-
2461 pean political leaders who have aided and abetted
2462 the ongoing invasion of Europe by the forces of
2463 the crescent moon death cult should be tried as
2464 accessories to Mr. Urban's murder. </s> <toxic
2465 3.024935007095337> About 415 million Euro-
2466 peans cannot compete with Google, Amazon, Face-
2467 book, Oracle, Intel, Apple, etc. and the socialist Eu-
2468 ropean welfare states need more revenue because
2469 they are running out of other peoples' money.

2470 Thus the Euro-socialist-bureaucrats pick the low-
2471 hanging fruit with litigious persecution of Ameri-
2472 can firms which dominate because unlike their pa-
2473 thetic Euro-competitors, the U.S. firms are clever,
2474 hard-working, and well-capitalized.

2475 If the the Europeans wish to engage in this trans-
2476 parent financial inquisition, then the US should
2477 respond with counter litigation for trillions against
2478 corrupt scofflaws like VW (think diesel fiddle!)
2479 as well as UBS/Credit Suisse/HSBC/Credit Lyon-
2480 naise (think tax cheats!)and sue/litigate them out
2481 of existence.

2482 If the lazy, corrupt, incompetent Euros want to
2483 play with fire, then let them be financially inciner-
2484 ated! </s> <not toxic 3.2237966060638428> Rome
2485 should never have made such inane pronounce-
2486 ments at Trent in their attempt to define the sub-
2487 stance of holy Eucharist. Most reasonable people
2488 understand that perfectly well. That Rome also
2489 made their pronouncements (faith and morals) "in-
2490 fallible" is equally tragic, for the simple reason that
2491 so-called infallible statements cannot be retracted
2492 without calling into question other so-called infalli-
2493 ble statements.

2494 Sincere question for you: If Jesus and his follow-
2495 ers celebrated Eucharist as a communal meal seated
2496 around a table, what gives Rome the right to alter
2497 this simple act of worship (perhaps "fellowship" is
2498 a better word—more suited toward love of God and
2499 neighbor), given to us by the Lord himself?

2500 **Label toxic**

2501 **Test Modified** I hope you don't have kids if you
2502 see this woman's actions as acceptable. And I ap-
2503 plaud the den for kicking the kid out. She brought
2504 unwanted negative attention upon them. However,
2505 she will, and is already likely, pay the the price for

her stupid stunt. </s> <toxic 2.873025417327881> 2506
Ms. Van Brocklin: You state that there is Payne's 2507
conduct is inexcusable, then proceed to use the rest 2508
of your space to justify him and his actions. You 2509
have denigrated the proud courage of countless peo- 2510
ple who took up causes via civil disobedience. I 2511
marched in the non-violent peaceful Civil Rights 2512
protest. So the dogs and the firehoses used by a cer- 2513
tain southern sheriff were justified, by your logic. 2514
So were the citizens beaten by Chicago police dur- 2515
ing the Democratic Convention Police riots. Res- 2516
olved in 20 minutes? Nonsense. If Miss Wubbels 2517
hadn't protested as she did, she'd likely have ended 2518
up in a back room, somewhere, and who know 2519
when she would have been granted the presence of 2520
a lawyer and what she would have gone through 2521
prior to that. You are completely ignorant of the 2522
shortage of nurses in this country - in some cases, 2523
critical shortages. And why would anyone want 2524
to be a nurse when they are disrespected by a for- 2525
mer state and federal prosecutor such as you. </s> 2526
<not toxic 2.8961446285247803> Well, I can't very 2527
well respect or fear an imaginary sky-being. As for 2528
my concept of character, it was good enough for 2529
the Alaska Judicial Council and Governor Knowles. 2530
But that was long ago. I've gotten older and, crikey, 2531
maybe I am going downhill. You're right about 2532
the inappropriateness of my comment. First Lady 2533
Walker's piece is very laudable and I shouldn't have 2534
taken it as an occasion to rant. (But look on the 2535
bright side: my misplaced comment gave you yet 2536
another occasion to rant about how your Fosterism 2537
is saving civilization.) 2538

2539 **Label toxic**

2540 **Test Modified** no one cares what a paid liberal
2541 trolling hack like you believes lunatic., </s> <toxic
2542 2.5876824855804443> It always amuses me when
2543 a troll gets on, they like their own comments and
2544 simply assert everyone else is wrong. Never any
2545 evidence to rebut it just blind assertions. </s> <not
2546 toxic 2.810284376144409> Ouch... didn't see that
2547 one coming. A liberal stealing my own line... just
2548 like they take everything else they like.

2549 **Label toxic**

2550 **Test Modified** Ok all you NDP "LEAP" mani-
2551 festo types, where is your hero Naomi Klein? Her
2552 fawning adoration of Chavez and Venezuelan thug-
2553 gery knows no bounds. I'm sure she's awfully
2554 hysterical over the thought that such a pathetic dict-
2555 storship could ever be sanctioned. </s> <toxic
2556 2.883418321609497> Um, no. The major left-

2557 wing Labour party was decimated; Mr. Rutte lost 2608
2558 8 seats; and Mr. Wilders Freedom party GAINED 2609
2559 4 seats. Now Mr. Rutte will have to garner favour 2610
2560 among RIGHT-WING parties to cobble together 2611
2561 his coalition. And as to your glib little comment 2612
2562 about people "embracing left-leaning parties", you 2613
2563 need only look to Canada to see the fiasco that 2614
2564 results... The corrupt and divisive Trudeau gov- 2615
2565 ernment was elected on a slate of fuzzy, alt-left 2616
2566 promises; only to renege on most of them ! </s> 2617
2567 <not toxic 3.0424022674560547> So anyone who 2618
2568 doesn't share your opinion is a "raving loony?" 2619
2569 And as for my post being "speculation? - which 2620
2570 part - that the Liberals are the party in power, or 2621
2571 that this involves money? 2622
2572 As for me not knowing what is going on, you 2623
2573 are correct, I am not a member of the Liberal party 2624
2574 insider clique, as you apparently are. 2625
2575 **Label toxic** 2626
2576 **Hate Speech Offensive Step 1** 2627
2577 **Test Modified** ' If this ugly ass nig- 2628
2578 gah can be with her , I can too . Tf 2629
2579 she with Chuu Baka for anyways this nig- 2630
2580 gah look like a... <http://t.co/IVNDRDPpQT> 2631
2581 </s> <hate speech 3.680654525756836> AtBlack- 2632
2582 face MrMooncricket UncleUnco BobbyBotsods 2633
2583 FuckTheMallcops RacistNegro69 ape? Fuck 2634
2584 u broke ass racist </s> <offensive language 2635
2585 3.9933786392211914> Tamela X But a wannabe 2636
2586 ghetto white girl will deal with an unaccomplished 2637
2587 black man for more than sex. </s> <neither 2638
2588 4.20903205871582> VNNForum These Marxist 2639
2589 kikes promoted Anti-Semitism by creating this film, 2640
2590 what was their incentive? They just promote their 2641
2591 enemies. 2642
2592 **Label hate speech** 2643
2593 **Test Modified** 8220; iamkrause: No need to 2644
2594 thank me, killing niggers is it's own reward. 8221; 2645
2595 Dis my pal from our graduation when we was 26 2646
2596 </s> <neither 3.5232555866241455> LanceTHES- 2647
2597 POKEN Statement of the Day: Ppl who RT 2648
2598 ChiefsatWar are even dumber than he is. The 2649
2599 ole birds of a feather type of deal. </s> <hate 2650
2600 speech 3.5796432495117188> Brenddan they said 2651
2601 some bullshit. Damn niggers </s> <offensive lan- 2652
2602 guage 3.7068209648132324> BRUH lmaoo 8220; 2653
2603 NoBeetch: Them hoes was nasty but I kept eating 2654
2604 them " roblokk: What the fuck happened to these? 2655
2605 <http://t.co/G9w10SmQdr>" 8221; 2656
2606 **Label hate speech** 2657
2607 **Test Modified** RT IsaidNick: niggers are 2658

so ignorant <http://t.co/P9jDdVsRVb> </s> <hate 2608
speech 3.2627384662628174> AtBlackface Mr- 2609
Mooncricket UncleUnco BobbyBotsods FuckThe- 2610
Mallcops RacistNegro69 ape? Fuck u broke 2611
ass racist </s> <neither 3.7804527282714844> 2612
VNNForum These Marxist kikes promoted Anti- 2613
Semitism by creating this film, what was their in- 2614
centive? They just promote their enemies. </s> <of- 2615
fensive language 3.9285027980804443> Tamela X 2616
But a wannabe ghetto white girl will deal with an 2617
unaccomplished black man for more than sex. 2618

Label hate speech 2619

Test Modified RT RosieZaya1: Ur 2620
fucking white trash </s> <hate speech 2621
2.951470136642456> AtBlackface MrMoon- 2622
cricket UncleUnco BobbyBotsods FuckTheMall- 2623
cops RacistNegro69 ape? Fuck u broke ass racist 2624
</s> <offensive language 3.6144936084747314> 2625
Tamela X But a wannabe ghetto white girl will 2626
deal with an unaccomplished black man for more 2627
than sex. </s> <neither 3.7668633460998535> 2628
VNNForum These Marxist kikes promoted 2629
Anti-Semitism by creating this film, what was their 2630
incentive? They just promote their enemies. 2631

Label hate speech 2632

Test Modified mike ray7 congratulations, you 2633
are officially fucking retarded. </s> <nei- 2634
ther 3.4077796936035156> RT JakeG Based- 2635
God: "Never go full retard" </s> <hate speech 2636
3.479813575744629> Brenddan they said some 2637
bullshit. Damn niggers </s> <offensive lan- 2638
guage 3.601623773574829> BRUH lmaoo 8220; 2639
NoBeetch: Them hoes was nasty but I kept eating 2640
them " roblokk: What the fuck happened to these? 2641
<http://t.co/G9w10SmQdr>" 8221; 2642

Label hate speech 2643

Test Modified gonna have them pussies mix up 2644
some concrete today. teach them to pose like me. 2645
I am a badass motherfucker. and I will let you 2646
be too (: </s> <hate speech 3.227602243423462> 2647
AtBlackface MrMooncricket UncleUnco Bobby- 2648
Botsods FuckTheMallcops RacistNegro69 ape? 2649
Fuck u broke ass racist </s> <offensive lan- 2650
guage 3.5520639419555664> BRUH lmaoo 8220; 2651
NoBeetch: Them hoes was nasty but I kept 2652
eating them " roblokk: What the fuck hap- 2653
pened to these? <http://t.co/G9w10SmQdr>" 8221; 2654
</s> <neither 3.809976816177368> RT Venus 2655
Lynn: 8220; dylxnl: look ghetto but it work 2656
<http://t.co/chrwW9dPca> 8221; 128557; 128557; 2657
128557; 128557; 128557; 128557; 128557; 2658

| | | | |
|------|--|---|------|
| 2659 | 128557; 128557; 128557; | | |
| 2660 | Label offensive language | | |
| 2661 | Test Modified maddievaans u call ur bestie a | Test Modified My baby guinea pig is so cute | 2710 |
| 2662 | bitch I'm guessing she's a dog that barks too much | Adorable </s> <hate speech 3.82438588142395> | 2711 |
| 2663 | </s> <offensive language 3.2484242916107178> | whiteangelss84 y dont u end us monkeys then? U | 2712 |
| 2664 | RT Im Amy Bitches: I'm not always a | hate us because were better than crackers amp; I | 2713 |
| 2665 | bitch, sometimes I sleep </s> <hate speech | I love it. Deep down u know weve done nothin 2 u | 2714 |
| 2666 | 3.3403964042663574> Women who are feminist | </s> <offensive language 3.8302881717681885> I | 2715 |
| 2667 | are the ugly bitches who cant find a man for them- | always wanted a bull dog them hoes clean fuck a | 2716 |
| 2668 | selves </s> <neither 3.8210768699645996> Gary's | pit </s> <neither 3.8650126457214355> Break- | 2717 |
| 2669 | girl was a tranny. RT Em 99car: I wonder what | fast fried chicken jerk chicken Tater tots white | 2718 |
| 2670 | would have become of rickspringfield if he'd writ- | rice nd press yellow rice nd beans Mac nd cheese | 2719 |
| 2671 | ten 'Gary's Girl'. | http://t.co/Usz8gJnZl0 | 2720 |
| 2672 | Label offensive language | Label neither | 2721 |
| 2673 | Test Modified No less than 3 bad bitches in | Test Modified RT Kick Man: Giants- Pitiful | 2722 |
| 2674 | my bed at a time... </s> <offensive language | .. Jets-Pitiful .. Mets- Pitiful .. Yankees-Pitiful | 2723 |
| 2675 | 3.2010550498962402> Then again... i shoot | .. Nets- Pitiful .. Knicks-Pitiful ... Ny sports- | 2724 |
| 2676 | bitches in the ass for doing some stupid shit like | Pitiful </s> <neither 3.8152754306793213> You | 2725 |
| 2677 | key my car or bust my windows. </s> <hate | know I'm not big on the NFL, but I'm so sick of | 2726 |
| 2678 | speech 3.666656732559204> bitch kill yoself, go | hearing all of this "Black and yellow" shit. LOL | 2727 |
| 2679 | on to the bathroom and eat the pills bitch, all | bandwagon fans and hell, GO PACKERS! </s> | 2728 |
| 2680 | of em </s> <neither 4.249817848205566> 9733; | <offensive language 4.071953773498535> BRUH | 2729 |
| 2681 | BEST ASIAN MASSAGE ON THE park slope | lmaoo 8220; NoBeeetch: Them hoes was nasty but | 2730 |
| 2682 | -TOP RATED SPA 9733; 718-622-0221 - 24 | I kept eating them " robloccc: What the fuck hap- | 2731 |
| 2683 | http://t.co/ZsAAzFL0p5 | pened to these? http://t.co/G9w10SmQdr" 8221; | 2732 |
| 2684 | Label offensive language | </s> <hate speech 4.095180511474609> whitean- | 2733 |
| 2685 | Test Modified RT TheDrugTribe: mary isn't a | angelss84 y dont u end us monkeys then? U hate us | 2734 |
| 2686 | backstabbing bitch that lies and deceives me </s> | because were better than crackers amp; I love it. | 2735 |
| 2687 | <offensive language 3.4536943435668945> RT Im | Deep down u know weve done nothin 2 u | 2736 |
| 2688 | Amy Bitches: I'm not always a bitch, sometimes | Label neither | 2737 |
| 2689 | I sleep </s> <hate speech 3.6065785884857178> | Test Modified jesstoth we could get matching | 2738 |
| 2690 | vinny2vicious faggot I knew you weren't really | burner phones and be ghetto fab for a few months | 2739 |
| 2691 | my friend. </s> <neither 3.638406753540039> | </s> <hate speech 3.4667954444885254> whitean- | 2740 |
| 2692 | Gary's girl was a tranny. RT Em 99car: I wonder | angelss84 y dont u end us monkeys then? U hate | 2741 |
| 2693 | what would have become of rickspringfield if he'd | us because were better than crackers amp; I love | 2742 |
| 2694 | written 'Gary's Girl'. | it. Deep down u know weve done nothin 2 u | 2743 |
| 2695 | Label offensive language | </s> <offensive language 3.595543622970581> RT | 2744 |
| 2696 | Test Modified porn, android, iphone, ipad, sex, | NickBratton3: I wish my parents bought me a | 2745 |
| 2697 | xxx, CloseUp Squirtng pussy and fingered ass- | car man.. People bitch about not getting what | 2746 |
| 2698 | hole http://t.co/bKYeoUwWv2 </s> <offensive lan- | car they want when they want it, and its free | 2747 |
| 2699 | guage 3.5574071407318115> BRUH lmaoo 8220; | 8230; </s> <neither 3.6075007915496826> RT | 2748 |
| 2700 | NoBeeetch: Them hoes was nasty but I kept eat- | Venus Lynn: 8220; dylxnl: look ghetto but it work | 2749 |
| 2701 | ing them " robloccc: What the fuck happened | http://t.co/chrvW9dPca 8221; 128557; 128557; | 2750 |
| 2702 | to these? http://t.co/G9w10SmQdr" 8221; </s> | 128557; 128557; 128557; 128557; 128557; | 2751 |
| 2703 | <neither 3.6928675174713135> DegenerateArtist | 128557; 128557; 128557; | 2752 |
| 2704 | Sniffs whiffy balls involuntary, cuz a FAIRY walks | Label neither | 2753 |
| 2705 | DOWNTOWN HAIRY, climatechange all our scary | Test Modified Thw White Iron Band plays | 2754 |
| 2706 | sense of hillbilly sneeze gt; </s> <hate speech | this weekend in Fargo,ND at the Aquar- | 2755 |
| 2707 | 3.7489092350006104> Wake up you stoopid dyke | ium(21+) ,Friday(10-29-10) with Charlie | 2756 |
| 2708 | RandyPawlik | Parr. The next night,Saturday... </s> <hate | 2757 |
| 2709 | Label offensive language | speech 3.8393898010253906> 128514; 128514; | 2758 |
| | | 128514;RT kwagiheath: Them 1st 48 Charlotte | 2759 |
| | | goon niccas gon Fkkk around and burn Bank Of | 2760 |

2761 America stadium down 128293; </s> <offensive
2762 language 3.9061098098754883> I be telling
2763 Mcgirt music ain't enough.You gotta have a non
2764 music related agenda.They crackers want to
2765 sell something with your face. </s> <neither
2766 3.980140447616577> You know I'm not big on
2767 the NFL, but I'm so sick of hearing all of this
2768 "Black and yellow" shit. LOL bandwagon fans and
2769 hell, GO PACKERS!

2770 **Label** neither 2817

2771 **Test Modified** RT dsrtvet: FoxNews tjoy7 2818
2772 And I don't have any confidence NONWHAT- 2819
2773 SOEVER in you Barack! You're the sole 2820
2774 reason why this country is in thi 8220; 2821
2775 </s> <neither 3.78818678855896> RT lach- 2822
2776 lan: So who wants to tell the Washing- 2823
2777 ton Post that Charlie Crist is a Democrat? 2824
2778 http://t.co/LGzzYusEKZ http://t.co/2jCVv7qxqf 2825
2779 </s> <hate speech 3.942119598388672> whitean- 2826
2780 gelss84 y dont u end us monkeys then? U hate 2827
2781 us because were better than crackers amp; I love 2828
2782 it. Deep down u know weve done nothin 2 u </s> 2829
2783 <offensive language 3.984135150909424> BRUH 2830
2784 lmaoo 8220; NoBeetch: Them hoes was nasty but 2831
2785 I kept eating them " roblokk: What the fuck hap- 2832
2786 pened to these? http://t.co/G9w10SmQdr" 8221; 2833

2787 **Label** neither 2834

2788 **Hate Speech Offensive Step 5** 2835

2789 **Test Modified** ' If this ugly ass niggah can 2836
2790 be with her , I can too . Tf she with 2837
2791 Chuu Baka for anyways this niggah look like 2838
2792 a... http://t.co/lVNDRDPPpQT </s> <hate speech 2839
2793 2.7132835388183594> RT WaavyLee: His balls 2840
2794 ashy RT Yattabing: Trelaire1st: Real women do 2841
2795 this http://t.co/VG5DBqH8aT 8221; and real fag- 2842
2796 gots let em do that.. Smh </s> <offensive language 2843
2797 2.7854363918304443> RT JayyMfCarter: If you 2844
2798 gotta nigga or bitch.... PLEASE STAY OUT MY 2845
2799 DM's I'm tired of y'all taken girls </s> <neither 2846
2800 3.0996105670928955> NeonTreezz PoCBeauty 2847
2801so.....white people selling African art? NOT 2848
2802 COOL? Okay. blackish redskins doublestandard 2849
2803 **Label** hate speech 2850

2804 **Test Modified** 8220; iamkrause: No need to 2851
2805 thank me, killing niggers is it's own reward. 8221; 2852
2806 Dis my pal from our graduation when we 2853
2807 was 26 </s> <hate speech 2.6058554649353027> 2854
2808 samzbikowski some negro amigo pulled a gun 2855
2809 on Nate and I a few weeks ago. I was 2856
2810 STOKED! 128299; 128299; 128299; </s> <neither 2857
2811 2.670192003250122> RT jennaferjenbub: 8220; 2858

BarryClerjuste: "Anything below a A+ and we 2812
disown you ling ling" http://t.co/m1QiWK4xZg 2813
8221; AustinBedsaul </s> <offensive language 2814
2.740609645843506> 8220; Alondra Lu: Ain't 2815
that a bitch 8221; 2816

Label hate speech 2817

Test Modified RT IsaidNick: niggers are 2818
so ignorant http://t.co/P9jDdVsRVb </s> <hate 2819
speech 2.057800769805908> RT WhitesOnly 1: 2820
niggers! http://t.co/Hb3uJaLky2 </s> <neither 2821
2.7749483585357666> amp; thots are wearing 2822
Uggs RT BigBootyJudy814: ItsFallBecause ne- 2823
gros are pulling out their Timbs" </s> <offensive 2824
language 2.926729440689087> RT Jayy Gee96: 2825
Dumb bitches 2826

Label hate speech 2827

Test Modified RT RosieZaya1: Ur 2828
fucking white trash </s> <hate speech 2829
2.422173500061035> FrankieJGrande fugly 2830
queer white trash </s> <offensive language 2831
2.6756434440612793> RT Jayy Gee96: Dumb 2832
bitches </s> <neither 2.783188819885254> RT 2833
BeardedNixon: Poont gotta be trash 2834

Label hate speech 2835

Test Modified mike ray7 congratulations, you 2836
are officially fucking retarded. </s> <hate 2837
speech 2.4854748249053955> darthdanaa Yes 2838
you do retard. </s> <offensive language 2839
2.851564645767212> Lol!! 8220; ItzSweetz 2840
Bitch: Ooop! QT TIFFANY PORSCHE: You little 2841
twats. 8221; </s> <neither 2.8971688747406006> 2842
RT jennaferjenbub: 8220; BarryClerjuste: "Any- 2843
thing below a A+ and we disown you ling ling" 2844
http://t.co/m1QiWK4xZg 8221; AustinBedsaul 2845

Label hate speech 2846

Test Modified gonna have them pussies mix up 2847
some concrete today. teach them to pose like me. I 2848
am a badass motherfucker. and I will let you be too 2849
(: </s> <offensive language 2.7589027881622314> 2850
40oz VAN IYCM I. I can't get any work done if 2851
you keep showin off your bitches. </s> <hate 2852
speech 2.8690829277038574> SlightlyAdjusted 2853
RT CapoToHeaven Alls niggers wanna do is fuck, 2854
tweet, and drink pineapple soda all day </s> <nei- 2855
ther 3.0193798542022705> cakedjake We're lay- 2856
ing rock around our lake. You're welcome to join a 2857
redneck workout. muscles 128170; 128513; 2858

Label offensive language 2859

Test Modified maddieevaans u call ur bestie a 2860
bitch I'm guessing she's a dog that barks too much 2861
</s> <offensive language 3.0170230865478516> 2862

| | | | |
|------|--|--|------|
| 2863 | Lol!! 8220; ItzSweetz Bitch: Ooop! QT TiFFANY | for da KoolQueefTribute 160; </s> <hate speech | 2914 |
| 2864 | P0RSCHHE: You little twats. 8221; </s> <hate | 3.5861661434173584> What a wetback looks like | 2915 |
| 2865 | speech 3.0850884914398193> Princesslexii16 | when he gets caught crossing the border. Ilovebamf | 2916 |
| 2866 | Fucking coon </s> <neither 3.21132493019104> | http://t.co/j3Uf1TYubO | 2917 |
| 2867 | lmaoooo RT ComedyTruth: Girls, don't let a guy | Label neither | 2918 |
| 2868 | treat you like a yellow starburst. You are a pink | Test Modified RT Kick Man: Giants- Pitiful | 2919 |
| 2869 | starburst. | .. Jets-Pitiful .. Mets- Pitiful .. Yankees-Pitiful .. | 2920 |
| 2870 | Label offensive language | Nets- Pitiful .. Knicks-Pitiful ... Ny sports- Pitiful | 2921 |
| 2871 | Test Modified No less than 3 bad bitches in | </s> <hate speech 3.3251242637634277> RT J R: | 2922 |
| 2872 | my bed at a time... </s> <offensive language | Smh nigga is mildly retarded RT Thotcho: LMFAO | 2923 |
| 2873 | 3.023017406463623> Don't lose sleep bout these | RT JustDoJ: If Griff was 8217;t injuries we 8217;d | 2924 |
| 2874 | bitches bc they come and go 128076; </s> <hate | legit be 6-1 </s> <neither 3.328580856323242> | 2925 |
| 2875 | speech 3.2786214351654053> bitch kill yoself, go | Don't follow the astros they said. They're trash they | 2926 |
| 2876 | on to the bathroom and eat the pills bitch, all of em | said. Well now look at them astros </s> <offen- | 2927 |
| 2877 | </s> <neither 3.4171059131622314> Keep those | sive language 3.3433356285095215> Them shits | 2928 |
| 2878 | away frm Charlie Day RT JhonenV: Just once in | ugly hoe. RT SirRocObama: RT BurgerK- | 2929 |
| 2879 | my life I'd like for someone's favorite part of my | ing: All these nuggets amp; u still actin chicken. | 2930 |
| 2880 | body to be my disgusting knees. | http://t.co/tRy8Lvyo9O | 2931 |
| 2881 | Label offensive language | Label neither | 2932 |
| 2882 | Test Modified RT TheDrugTribe: mary isn't | Test Modified jesstoth we could get matching | 2933 |
| 2883 | a backstabbing bitch that lies and deceives me | burner phones and be ghetto fab for a few months | 2934 |
| 2884 | </s> <offensive language 2.991457939147949> | </s> <offensive language 3.122525930404663> | 2935 |
| 2885 | RT StevStiffler: If her bio says "Only God can | JZolly23 JBilinovich we need to grow mullets | 2936 |
| 2886 | judge me" she's a hoe. </s> <hate speech | together so we can get all the bitches and Han- | 2937 |
| 2887 | 3.098494529724121> RT sorryimalex: I got | nahKubiak can hate on us </s> <hate speech | 2938 |
| 2888 | called a faggot for buying girl toms so now I'm | 3.291858434677124> RT NoWomanIsRight: You | 2939 |
| 2889 | gonna fuck that person in the ass </s> <neither | can be a good girl all you want and those hoes still | 2940 |
| 2890 | 3.3266849517822266> lmaoooo RT ComedyTruth: | gonna get us niggas attention from time to time | 2941 |
| 2891 | Girls, don't let a guy treat you like a yellow star- | </s> <neither 3.3428289890289307> RT Venus | 2942 |
| 2892 | burst. You are a pink starburst. | Lynn: 8220; dylxnl: look ghetto but it work | 2943 |
| 2893 | Label offensive language | http://t.co/chrvW9dPca 8221; 128557; 128557; | 2944 |
| 2894 | Test Modified porn, android, iphone, ipad, sex, | 128557; 128557; 128557; 128557; 128557; | 2945 |
| 2895 | xxx, CloseUp Squirting pussy and fingered | 128557; 128557; 128557; | 2946 |
| 2896 | asshole http://t.co/bKYeoUwWv2 </s> <neither | Label neither | 2947 |
| 2897 | 1.5677733421325684> porn, android, iphone, ipad, | Test Modified Thw White Iron Band plays | 2948 |
| 2898 | sex, xxx, Desi paki http://t.co/XxcdQvzI9t </s> | this weekend in Fargo,ND at the Aquar- | 2949 |
| 2899 | <hate speech 2.8566393852233887> RT mitchman- | ium(21+) ,Friday(10-29-10) with Charlie Parr. | 2950 |
| 2900 | cuso: BrantPrintup:What straight guys take a pic- | The next night,Saturday... </s> <neither | 2951 |
| 2901 | ture of themselves naked in a hot tun.What fags | 3.4018728733062744> Lmaooo naw man RT | 2952 |
| 2902 | jakesiw Ryan Murphy3 Randy 8230; </s> <offen- | DipOnline Yo want in RT HumbltonBanks: | 2953 |
| 2903 | sive language 2.932191848754883> Lol!! 8220; | U serious bro?? lol RT CheezMoeJenk- | 2954 |
| 2904 | ItzSweetz Bitch: Ooop! QT TiFFANY P0RSCHHE: | inz 2-3:10am early bird special </s> <hate | 2955 |
| 2905 | You little twats. 8221; | speech 3.544551372528076> 128514; 128514; | 2956 |
| 2906 | Label offensive language | 128514;RT kwagiheath: Them 1st 48 Charlotte | 2957 |
| 2907 | Test Modified My baby guinea pig is so cute | goon nicas gon Fkkk around and burn Bank Of | 2958 |
| 2908 | Adorable </s> <neither 3.1643435955047607> | America stadium down 128293; </s> <offensive | 2959 |
| 2909 | Our female guinea pig is pregnant 127882; 127881; | language 3.711003065109253> I be telling Mcgirt | 2960 |
| 2910 | 127873; 128525; 128525; 128525; </s> <offen- | music ain't enough.You gotta have a non music re- | 2961 |
| 2911 | sive language 3.4907007217407227> I impress da | lated agenda.They crackers want to sell something | 2962 |
| 2912 | young white girl next doe by taking out my gi- | with your face. | 2963 |
| 2913 | ant negro thang and usin it to flip da hamburgers | Label neither | 2964 |

2965 **Test Modified** RT dsrtvet: FoxNews tjoy7 2.7749483585357666> amp; thots are wearing 3016

2966 And I don't have any confidence NONWHAT- 3017

2967 SOEVER in you Barack! You're the sole rea- 3018

2968 son why this country is in thi 8230; </s> <hate 3019

2969 speech 2.696760654449463> RT veeveeveeveeve: 3020

2970 If I was Obama Id call a press conference amp; 3021

2971 slit joe bidens neck on live tv just 2 show these 3022

2972 crackers I mean business 8230; </s> <neither 3023

2973 2.762817144393921> RT jennaferjenbub: 8220; 3024

2974 BarryClerjuste: "Anything below a A+ and we 3025

2975 disown you ling ling" <http://t.co/m1QiWK4xZg> 3026

2976 8221; AustinBedsaul </s> <offensive language 3027

2977 2.9660327434539795> RT CoffyBrownChi: If he 3028

2978 don't believe you, no refunds hoe. 3029

2979 **Label** neither 3030

2980 **Hate Speech Offensive Step 10** 3031

2981 **Test Modified** ' If this ugly ass niggah can 3032

2982 be with her , I can too . Tf she with 3033

2983 Chuu Baka for anyways this niggah look like 3034

2984 a... <http://t.co/IVNDRDPpQT> </s> <offensive 3035

2985 language 2.6535706520080566> RT CurrenSy 3036

2986 Spitta: And if a bitch can't respect a nigga 3037

2987 wit some paper and a fresh pair of bball shorts 3038

2988 then she was raised terribly.. </s> <hate speech 3039

2989 2.7132835388183594> RT WaavyLee: His balls 3040

2990 ashy RT Yattabing: Trelaire1st: Real women 3041

2991 do this <http://t.co/VG5DBqH8aT> 8221; and real 3042

2992 faggots let em do that.. Smh </s> <neither 3043

2993 3.0996105670928955> NeonTreezz PoCBeauty 3044

2994so.....white people selling African art? NOT 3045

2995 COOL? Okay. blackish redskins doublestandard 3046

2996 **Label** hate speech 3047

2997 **Test Modified** 8220; iamkrause: No need to 3048

2998 thank me, killing niggers is it's own reward. 8221; 3049

2999 Dis my pal from our graduation when we was 3050

3000 26 </s> <hate speech 2.545886278152466> RT 3051

3001 Tae Rhodes: 8220; kim92493: 8220; Tae Rhodes: 3052

3002 kim92493 patpatbush uhhh you've been judged 3053

3003 8221; it happens. whitepower...I'll hang you nig- 3054

3004 ger 8221; wo 8230; </s> <offensive language 3055

3005 2.6044790744781494> 8220; NoRapist: on my 3056

3006 way to fuck ur bitch <http://t.co/SgVBBrwOg2> 8221; 3057

3007 mckinley719 </s> <neither 2.670259714126587> 3058

3008 RT jennaferjenbub: 8220; BarryClerjuste: "Any- 3059

3009 thing below a A+ and we disown you ling ling" 3060

3010 <http://t.co/m1QiWK4xZg> 8221; AustinBedsaul 3061

3011 **Label** hate speech 3062

3012 **Test Modified** RT IsaidNick: niggers are 3063

3013 so ignorant <http://t.co/P9jDdVsRVb> </s> <hate 3064

3014 speech 2.057800769805908> RT WhitesOnly 1: 3065

3015 niggers! <http://t.co/Hb3uJaLky2> </s> <neither 3066

2.7749483585357666> amp; thots are wearing 3016

Uggs RT BigBootyJudy814: ItsFallBecause ne- 3017

gros are pulling out their Timbs" </s> <offensive 3018

language 2.8298287391662598> This Uncle Tom 3019

mother fucking wants to invoke 3/5 a man in his 3020

speech? Dude you ain't white no matter how 3021

much... <http://t.co/3yrcyC9ezc> 3022

Label hate speech 3023

Test Modified RT RosieZaya1: Ur 3024

fucking white trash </s> <hate speech 3025

2.422173500061035> FrankieJGrande fugly 3026

queer white trash </s> <offensive language 3027

2.6756434440612793> RT Jayy Gee96: Dumb 3028

bitches </s> <neither 2.783188819885254> RT 3029

BeardedNixon: Poont gotta be trash 3030

Label hate speech 3031

Test Modified mike ray7 congratulations, you 3032

are officially fucking retarded. </s> <hate 3033

speech 2.4854748249053955> darthdanaa Yes 3034

you do retard. </s> <offensive language 3035

2.8516175746917725> Lol!! 8220; ItzSweetz 3036

Bitch: Ooop! QT TiFFANY PORSCHE: You little 3037

twats. 8221; </s> <neither 2.8972203731536865> 3038

RT jennaferjenbub: 8220; BarryClerjuste: "Any- 3039

thing below a A+ and we disown you ling ling" 3040

<http://t.co/m1QiWK4xZg> 8221; AustinBedsaul 3041

Label hate speech 3042

Test Modified gonna have them pussies mix up 3043

some concrete today. teach them to pose like me. I 3044

am a badass motherfucker. and I will let you be too 3045

(: </s> <offensive language 2.758687734603882> 3046

40oz VAN IYCM I. I can't get any work done if 3047

you keep showin off your bitches. </s> <hate 3048

speech 2.8036365509033203> Just to get u mad 3049

go on your search bar on here and search up 3050

"stupid niggers" amp; hop on somebodys head 3051

then mention me lol stonethegreat23 </s> <nei- 3052

ther 3.012741804122925> charloosss keepitplur 3053

nicoleariel I'll chug my tall can . but homegirl 3054

won't approve lol 3055

Label offensive language 3056

Test Modified maddievaans u call ur bestie 3057

a bitch I'm guessing she's a dog that barks too 3058

much </s> <hate speech 2.8469488620758057> 3059

RylannWilliams whooooo? Chelsey? Fuck her 3060

lol. She juss a bitch </s> <offensive language 3061

2.8842358589172363> RT Ezzzylove: She a bad 3062

bitch, let's get to it right away . </s> <neither 3063

3.0819990634918213> charliesheen Charlie, im 3064

an old lady. don't EVER SAY UGLY THINGS 3065

ABOUT UR CHILDRENS MOM.. I GET IT!!!, 3066

| | | | |
|------|--|---|------|
| 3067 | JUS DONT! BIG HUG | | |
| 3068 | Label offensive language | Label neither | 3118 |
| 3069 | Test Modified No less than 3 bad bitches | Test Modified RT Kick Man: Giants- Piti- | 3119 |
| 3070 | in my bed at a time... </s> <offensive lan- | ful .. Jets-Pitiful .. Mets- Pitiful .. Yankees- | 3120 |
| 3071 | guage 2.8522520065307617> Bad bitches in | Pitiful .. Nets- Pitiful .. Knicks-Pitiful ... Ny | 3121 |
| 3072 | the pen make my toes curl </s> <hate speech | sports- Pitiful </s> <neither 3.0366640090942383> | 3122 |
| 3073 | 3.2539432048797607> I didn't forsake all other | Buster ESPN Huh.....last 10 games..Tampa 8-2/Balt | 3123 |
| 3074 | bitches for my wife to be getting fucked on by | 7-3/Yanks 6-4...and they lost their best pitcher. | 3124 |
| 3075 | another nigga. and you know she married? you | Please explain your logic. </s> <hate speech | 3125 |
| 3076 | gotta die. </s> <neither 3.4170782566070557> | 3.3251242637634277> RT J R: Smh nigga is | 3126 |
| 3077 | Keep those away frm Charlie Day RT JhonenV: | mildly retarded RT Thotcho: LMFAO RT JustDoJ: | 3127 |
| 3078 | Just once in my life I'd like for someone's favorite | If Griff wasn 8217;t injuries we 8217;d legit be 6-1 | 3128 |
| 3079 | part of my body to be my disgusting knees. | </s> <offensive language 3.3433356285095215> | 3129 |
| 3080 | Label offensive language | Them shits ugly hoe. RT SirRocObama: | 3130 |
| 3081 | Test Modified RT TheDrugTribe: mary isn't | RT BurgerKing: All these nuggets amp; u still actin | 3131 |
| 3082 | a backstabbing bitch that lies and deceives me | chicken. http://t.co/tRy8Lvyo9O | 3132 |
| 3083 | </s> <offensive language 2.9916186332702637> | Label neither | 3133 |
| 3084 | RT StevStiffler: If her bio says "Only God can | Test Modified jesstoth we could get match- | 3134 |
| 3085 | judge me" she's a hoe. </s> <hate speech | ing burner phones and be ghetto fab for a few | 3135 |
| 3086 | 3.02489972114563> triple6em96 Hunglikerobby | months </s> <hate speech 3.034785270690918> | 3136 |
| 3087 | bitch you watch your fucking mouth you dirty | SAMMI boyden bruh we can finally roll like red- | 3137 |
| 3088 | whore. I swear to god that's a thin line </s> | necks (: ((drug dealers)) </s> <offensive lan- | 3138 |
| 3089 | <neither 3.1058743000030518> RT shakiraevanss: | guage 3.122525930404663> JZolly23 JBilovich | 3139 |
| 3090 | Criticize Amanda for saying the n word, sure, but | we need to grow mullets together so we can get all | 3140 |
| 3091 | don't make jokes about her sexual assault, don't be | the bitches and HannahKubiak can hate on us </s> | 3141 |
| 3092 | trash. | <neither 3.271000623703003> RT sassytbh: a girl | 3142 |
| 3093 | Label offensive language | tweeted "you might be ghetto if u bring food from | 3143 |
| 3094 | Test Modified porn, android, iphone, ipad, sex, | outside into the movies" | 3144 |
| 3095 | xxx, CloseUp Squirting pussy and fingered | no u might be stupid if u pay 4.99 for a b 8230; | 3145 |
| 3096 | asshole http://t.co/bKYeoUwWv2 </s> <neither | Label neither | 3146 |
| 3097 | 1.5677733421325684> porn, android, iphone, ipad, | Test Modified Thw White Iron Band plays | 3147 |
| 3098 | sex, xxx, Desi paki http://t.co/XxcdQvzI9t | this weekend in Fargo,ND at the Aquar- | 3148 |
| 3099 | </s> <offensive language 2.8408925533294678> | ium(21+) ,Friday(10-29-10) with Charlie Parr. | 3149 |
| 3100 | RT FunnyPicsDepot: bitches be like "I'm a vir- | The next night,Saturday... </s> <neither | 3150 |
| 3101 | gin" http://t.co/mFDwXmg8ic </s> <hate speech | 3.2462401390075684> RT toddknife: Full weak- | 3151 |
| 3102 | 2.8566393852233887> RT mitchmancuso: Brant- | enednachos set (except the last song) from South- | 3152 |
| 3103 | Printup:What straight guys take a picture of them- | ern Darkness Fest last month. Who's the ape | 3153 |
| 3104 | selves naked in a hot tun.What fags jakesiw Ryan | on guitar? https://t.c 8230; </s> <hate speech | 3154 |
| 3105 | Murphy3 Randy 8230; | 3.3524651527404785> Eagles fuck around amp; | 3155 |
| 3106 | Label offensive language | lose it'll be kill the cracker at the Sophi crib smfh | 3156 |
| 3107 | Test Modified My baby guinea pig is so cute | </s> <offensive language 3.511016368865967> | 3157 |
| 3108 | Adorable </s> <neither 3.1643435955047607> | My dawg ceomiamimike told me it's a must I be | 3158 |
| 3109 | Our female guinea pig is pregnant 127882; 127881; | 901k2lounge this Saturday ROCKIN that bitch wit | 3159 |
| 3110 | 127873; 128525; 128525; 128525; </s> <offen- | Tha 8230; http://t.co/0NV9cHtwOs | 3160 |
| 3111 | sive language 3.4907007217407227> I impress da | Label neither | 3161 |
| 3112 | young white girl next doe by taking out my gi- | Test Modified RT dsrtvet: FoxNews tjoy7 | 3162 |
| 3113 | ant negro thang and usin it to flip da hamburgers | And I don't have any confidence NONWHAT- | 3163 |
| 3114 | for da KoolQueefTribute 160; </s> <hate speech | SOEVER in you Barack! You're the sole rea- | 3164 |
| 3115 | 3.5861661434173584> What a wetback looks like | son why this country is in thi 8230; </s> <hate | 3165 |
| 3116 | when he gets caught crossing the border. Ilovebamf | speech 2.696760654449463> RT veeveeveeveevee: | 3166 |
| 3117 | http://t.co/j3Uf1TYubO | If I was Obama Id call a press conference amp; | 3167 |
| | | slit joe bidens neck on live tv just 2 show these | 3168 |

| | | | |
|------|--|---|------|
| 3169 | crackers I mean business 8230; </s> <neither | Libya Ambassador (Christopher) Stevens sent 600 | 3220 |
| 3170 | 2.762908458709717> RT jennaferjenbub: 8220; | requests for help in Benghazi. </s> the second | 3221 |
| 3171 | BarryClerjuste: "Anything below a A+ and we | 2016 presidential debate | 3222 |
| 3172 | disown you ling ling" http://t.co/m1QiWK4xZg | Label true statement | 3223 |
| 3173 | 8221; AustinBedsaul </s> <offensive language | Test Modified John Kasich says We are now | 3224 |
| 3174 | 2.881894588470459> 8220; LongMoneyTonny: | eighth in the nation in job creation . . . we are | 3225 |
| 3175 | vintage monroe DONT Say Shit Else ! Just Stfu | No. 1 in the Midwest. </s> a news conference </s> | 3226 |
| 3176 | ! 8221;bitch we can do a lot more off this Twitter | <true statement 3.851958990097046> Jorge Elorza | 3227 |
| 3177 | shit you can come see me | says In the last six years of Ciancis administration | 3228 |
| 3178 | Label neither | violent crime was down in the United States. It | 3229 |
| 3179 | LIAR (collapsed) Step 1 | was down in the region. It was down in Rhode | 3230 |
| 3180 | Test Modified Afscome says In labor negotia- | Island. But it was up in Providence. </s> a debate | 3231 |
| 3181 | tions with city employees, Milwaukee Mayor Tom | </s> <false statement 4.010262966156006> Don- | 3232 |
| 3182 | Barrett demanded concessions that went beyond | ald Trump says The federal government is sending | 3233 |
| 3183 | those mandated by Gov. Scott Walkers collective | refugees to states with governors who are Republi- | 3234 |
| 3184 | bargaining law </s> a letter to members </s> <true | cans, not to the Democrats. </s> an interview on | 3235 |
| 3185 | statement 3.833270311355591> Donald Trump | Laura Ingraham's radio show | 3236 |
| 3186 | says Libya Ambassador (Christopher) Stevens sent | Label true statement | 3237 |
| 3187 | 600 requests for help in Benghazi. </s> the sec- | Test Modified Mike Pence says It was Hillary | 3238 |
| 3188 | ond 2016 presidential debate </s> <false state- | Clinton who left Americans in harms way in Beng- | 3239 |
| 3189 | ment 4.013778209686279> Donald Trump says | hazi and after four Americans fell said, What dif- | 3240 |
| 3190 | The federal government is sending refugees to | ference at this point does it make? </s> the Re- | 3241 |
| 3191 | states with governors who are Republicans, not | publican national convention </s> <true statement | 3242 |
| 3192 | to the Democrats. </s> an interview on Laura In- | 3.7440342903137207> Jorge Elorza says In the last | 3243 |
| 3193 | graham's radio show | six years of Ciancis administration violent crime | 3244 |
| 3194 | Label true statement | was down in the United States. It was down in the | 3245 |
| 3195 | Test Modified Rick Scott says All Aboard | region. It was down in Rhode Island. But it was up | 3246 |
| 3196 | Florida is a 100 percent private venture. There | in Providence. </s> a debate </s> <false statement | 3247 |
| 3197 | is no state money involved. </s> a TV interview | 3.746598958969116> Donald Trump says You will | 3248 |
| 3198 | </s> <false statement 3.664231777191162> Don- | learn more about Donald Trump by going down to | 3249 |
| 3199 | ald Trump says The federal government is sending | the Federal Elections to see the financial disclo- | 3250 |
| 3200 | refugees to states with governors who are Republi- | sure form than by looking at tax returns. </s> a | 3251 |
| 3201 | cans, not to the Democrats. </s> an interview on | Presidential debate at Hofstra University | 3252 |
| 3202 | Laura Ingraham's radio show </s> <true statement | Label true statement | 3253 |
| 3203 | 3.831820011138916> Patrick Murphy says Marco | Test Modified Rand Paul says Of the roughly 15 | 3254 |
| 3204 | Rubio opposes immigration reform. Worse, Rubio | percent of Americans who dont have health insur- | 3255 |
| 3205 | supports Donald Trump. His plan would deport | ance, half of them made more than 50,000 a year. | 3256 |
| 3206 | 800,000 children, destroying families. </s> a TV | </s> an interview on Comedy Central's "The Daily | 3257 |
| 3207 | ad | Show" </s> <true statement 3.7997491359710693> | 3258 |
| 3208 | Label true statement | Bernie S says We have the highest rate of child- | 3259 |
| 3209 | Test Modified Julie Pace says The Obama ad- | hood poverty of any major country on Earth. | 3260 |
| 3210 | ministration is using as its legal justification for | </s> an interview on CNN </s> <false statement | 3261 |
| 3211 | these airstrikes (on the Islamic State), an autho- | 3.9633538722991943> Donald Trump says The | 3262 |
| 3212 | rization for military force that the president him- | federal government is sending refugees to states | 3263 |
| 3213 | self has called for repeal of. </s> a question to | with governors who are Republicans, not to the | 3264 |
| 3214 | White House Press Secretary Josh Earnest </s> | Democrats. </s> an interview on Laura Ingraham's | 3265 |
| 3215 | <false statement 3.5803754329681396> Donald | radio show | 3266 |
| 3216 | Trump says Hillary Clinton invented ISIS with | Label false statement | 3267 |
| 3217 | her stupid policies. She is responsible for ISIS. | Test Modified Barack Obama says Stimulus tax | 3268 |
| 3218 | </s> an interview on 60 Minutes </s> <true state- | cuts "began showing up in paychecks of 4.8 mil- | 3269 |
| 3219 | ment 3.869307518005371> Donald Trump says | lion Indiana households about three months ago." | 3270 |

| | | | |
|------|---|--|------|
| 3271 | </s> a speech in Wakarusa, Ind. </s> <true state- | Rubio supports Donald Trump. His plan would | 3322 |
| 3272 | ment 3.8199117183685303> Jorge Elorza says In | deport 800,000 children, destroying families. </s> | 3323 |
| 3273 | the last six years of Ciancis administration vio- | a TV ad | 3324 |
| 3274 | lent crime was down in the United States. It was | Label false statement | 3325 |
| 3275 | down in the region. It was down in Rhode Is- | LIAR (collapsed) Step 5 | 3326 |
| 3276 | land. But it was up in Providence. </s> a debate | Test Modified Afscmc says In labor negotiations | 3327 |
| 3277 | </s> <false statement 3.916092872619629> Don- | with city employees, Milwaukee Mayor Tom Bar- | 3328 |
| 3278 | ald Trump says The federal government is sending | rett demanded concessions that went beyond those | 3329 |
| 3279 | refugees to states with governors who are Republi- | mandated by Gov. Scott Walkers collective bargain- | 3330 |
| 3280 | cans, not to the Democrats. </s> an interview on | ing law </s> a letter to members </s> <false state- | 3331 |
| 3281 | Laura Ingraham's radio show | ment 3.131746292114258> Tom Barrett says Gov. | 3332 |
| 3282 | Label false statement | Scott Walker said no to equal pay for equal work | 3333 |
| 3283 | Test Modified Allen West says If you look | for women. </s> a TV ad </s> <true statement | 3334 |
| 3284 | at the application for a security clearance, I | 3.1800825595855713> Scott Walker says If public | 3335 |
| 3285 | have a clearance that even the president of the | employees dont pay more for benefits starting April | 3336 |
| 3286 | United States cannot obtain because of my back- | 1, 2011, the equivalent is 1,500 state employee lay- | 3337 |
| 3287 | ground. </s> a candidate forum </s> <false | offs by June 30, 2011 and 10,000 to 12,000 state | 3338 |
| 3288 | statement 3.760773181915283> Rush Limbaugh | and local government employee layoffs in the next | 3339 |
| 3289 | says 11 straight years of no major hurricanes | two years. </s> a news conference | 3340 |
| 3290 | striking land in the United States bores a hole | Label true statement | 3341 |
| 3291 | right through the whole climate change argument. | Test Modified Rick Scott says All Aboard | 3342 |
| 3292 | </s> a radio show broadcast </s> <true statement | Florida is a 100 percent private venture. There | 3343 |
| 3293 | 3.77760648727417> Arizona Citizens Defense | is no state money involved. </s> a TV interview | 3344 |
| 3294 | League says a gun bill before the Senate would | </s> <true statement 3.0582425594329834> Char- | 3345 |
| 3295 | make it a federal felony to leave town for more | lie Crist says All Aboard Florida is receiving mil- | 3346 |
| 3296 | than seven days, and leave someone else at home | lions in Florida taxpayer dollars. </s> a fundraising | 3347 |
| 3297 | with your firearms. </s> an email to supporters | email </s> <false statement 3.1522974967956543> | 3348 |
| 3298 | Label false statement | Corey Lewandowski says Mr. Trump is self- | 3349 |
| 3299 | Test Modified Bernie S says We now work | financing his campaign, so we dont have any | 3350 |
| 3300 | the longest hours of any people around the world. | donors. </s> a radio interview. | 3351 |
| 3301 | </s> a C-SPAN interview </s> <true statement | Label true statement | 3352 |
| 3302 | 3.7155606746673584> Bernie S says We have the | Test Modified Julie Pace says The Obama ad- | 3353 |
| 3303 | highest rate of childhood poverty of any major | ministration is using as its legal justification for | 3354 |
| 3304 | country on Earth. </s> an interview on CNN </s> | these airstrikes (on the Islamic State), an authoriza- | 3355 |
| 3305 | <false statement 4.0561442375183105> Rush Lim- | tion for military force that the president himself | 3356 |
| 3306 | baugh says 11 straight years of no major hurricanes | has called for repeal of. </s> a question to White | 3357 |
| 3307 | striking land in the United States bores a hole right | House Press Secretary Josh Earnest </s> <true | 3358 |
| 3308 | through the whole climate change argument. </s> | statement 2.9627556800842285> Martha Raddatz | 3359 |
| 3309 | a radio show broadcast | says The Obama administration originally wanted | 3360 |
| 3310 | Label false statement | 10,000 troops to remain in Iraq – not combat troops, | 3361 |
| 3311 | Test Modified Sarah Palin says Donald Trumps | but military advisers, special operations forces, | 3362 |
| 3312 | conversion to pro-life beliefs are akin to Justin | to watch the counterterrorism effort. </s> com- | 3363 |
| 3313 | Biebers, who said in the past that abortion was | ments on ABC's "This Week" </s> <false statement | 3364 |
| 3314 | no big deal to him. </s> an interview on CNN | 3.246009588241577> Rick Perry says Obama has | 3365 |
| 3315 | </s> <false statement 3.7367687225341797> Don- | chosen to deny the vicious anti-Semitic motivation | 3366 |
| 3316 | ald Trump says The federal government is sending | of the attack on a kosher Jewish grocery in Paris. | 3367 |
| 3317 | refugees to states with governors who are Republi- | </s> a statement | 3368 |
| 3318 | cans, not to the Democrats. </s> an interview | Label true statement | 3369 |
| 3319 | on Laura Ingraham's radio show </s> <true state- | Test Modified John Kasich says We are now | 3370 |
| 3320 | ment 3.7425951957702637> Patrick Murphy says | eighth in the nation in job creation . . . we are | 3371 |
| 3321 | Marco Rubio opposes immigration reform. Worse, | No. 1 in the Midwest. </s> a news conference </s> | 3372 |

| | | | |
|------|---|---|------|
| 3373 | <false statement 2.610369920730591> Ted Strick- | convention | 3424 |
| 3374 | land says Gov. John Kasich incorrectly claimed | Label false statement | 3425 |
| 3375 | Ohios economy was 38th in the nation when he | Test Modified Allen West says If you look at | 3426 |
| 3376 | took office. We were sixth in the nation in terms of | the application for a security clearance, I have a | 3427 |
| 3377 | economic job growth. </s> an interview on CNN | clearance that even the president of the United | 3428 |
| 3378 | </s> <true statement 3.028876543045044> Terry | States cannot obtain because of my background. | 3429 |
| 3379 | Mcauliffe says If you take the population growth | </s> a candidate forum </s> <false statement | 3430 |
| 3380 | here in Virginia, we are net zero on job creation | 3.050549268722534> Ted Cruz says One of the | 3431 |
| 3381 | since (Bob McDonnell) became governor. </s> a | most troubling aspects of the Rubio-Schumer Gang | 3432 |
| 3382 | speech. | of Eight bill was that it gave President Obama | 3433 |
| 3383 | Label true statement | blanket authority to admit refugees, including Syr- | 3434 |
| 3384 | Test Modified Mike Pence says It was Hillary | ian refugees, without mandating any background | 3435 |
| 3385 | Clinton who left Americans in harms way in Beng- | checks whatsoever. </s> a Republican presiden- | 3436 |
| 3386 | hazi and after four Americans fell said, What | tial debate in Las Vegas </s> <true statement | 3437 |
| 3387 | difference at this point does it make? </s> the | 3.196129560470581> David Shuster says Said former | 3438 |
| 3388 | Republican national convention </s> <true state- | U.S. Ambassador to Kenya Scott Gration was | 3439 |
| 3389 | ment 2.5875017642974854> Hillary Clinton says | forced to resign two years ago because of his per- | 3440 |
| 3390 | When terrorists killed more than 250 Americans | sonal use of emails. </s> a Hillary Clinton press | 3441 |
| 3391 | in Lebanon under Ronald Reagan, the Democrats | conference | 3442 |
| 3392 | didnt make that a partisan issue. </s> a CNN town | Label false statement | 3443 |
| 3393 | hall </s> <false statement 2.9331557750701904> | Test Modified Bernie S says We now work | 3444 |
| 3394 | Facebook Posts says Hillary Clinton refuses to tes- | the longest hours of any people around the world. | 3445 |
| 3395 | tify before Congress about the 2012 attack in Beng- | </s> a C-SPAN interview </s> <true statement | 3446 |
| 3396 | hazi. </s> a meme on social media | 3.08957576751709> Jim Sensenbrenner says We | 3447 |
| 3397 | Label true statement | have the highest corporate tax rate in the world. Its | 3448 |
| 3398 | Test Modified Rand Paul says Of the roughly 15 | 35 percent. </s> an interview </s> <false statement | 3449 |
| 3399 | percent of Americans who dont have health insur- | 3.3488667011260986> Mitt Romney says Today | 3450 |
| 3400 | ance, half of them made more than 50,000 a year. | there are more men and women out of work in | 3451 |
| 3401 | </s> an interview on Comedy Central's "The Daily | America than there are people working in Canada. | 3452 |
| 3402 | Show" </s> <true statement 2.932455062866211> | </s> a speech to the Conservative Political Action | 3453 |
| 3403 | Joe Biden says Among the money spent on health | Conference | 3454 |
| 3404 | care in the United States, "46 cents on every dollar | Label false statement | 3455 |
| 3405 | spent is through Medicare and Medicaid." </s> an | Test Modified Sarah Palin says Donald Trumps | 3456 |
| 3406 | interview on NBC's 'Meet the Press' </s> <false | conversion to pro-life beliefs are akin to Justin | 3457 |
| 3407 | statement 3.02447247505188> Trent Franks says | Biebers, who said in the past that abortion was | 3458 |
| 3408 | The top 1 percent pay over half of the entire revenue | no big deal to him. </s> an interview on CNN </s> | 3459 |
| 3409 | for this country. </s> an interview on MSNBC's | <false statement 3.1018259525299072> Herman | 3460 |
| 3410 | 'The Dylan Ratigan Show' | Cain says Said Planned Parenthoods early objective | 3461 |
| 3411 | Label false statement | was to help kill black babies before they came into | 3462 |
| 3412 | Test Modified Barack Obama says Stimulus tax | the world. </s> a talk at a conservative think tank | 3463 |
| 3413 | cuts "began showing up in paychecks of 4.8 million | </s> <true statement 3.1297004222869873> Greg | 3464 |
| 3414 | Indiana households about three months ago." </s> | Abbott says After Texas defunded Planned Parent- | 3465 |
| 3415 | a speech in Wakarusa, Ind. </s> <false statement | hood, both the unintended pregnancy and abortion | 3466 |
| 3416 | 2.8908281326293945> Paul Broun says Stimulus | rates dropped. </s> a tweet | 3467 |
| 3417 | money funded a government board that made rec- | Label false statement | 3468 |
| 3418 | ommendations that would cost 378,000 jobs and | LIAR (collapsed) Step 10 | 3469 |
| 3419 | 28.3 billion in sales. </s> a tweet </s> <true state- | Test Modified Afscome says In labor negotia- | 3470 |
| 3420 | ment 2.9225375652313232> Sarah Palin says "One | tions with city employees, Milwaukee Mayor Tom | 3471 |
| 3421 | state even spent a million bucks to put up signs that | Barrett demanded concessions that went beyond | 3472 |
| 3422 | advertise that they were spending on the federal | those mandated by Gov. Scott Walkers collective | 3473 |
| 3423 | stimulus projects." </s> an address at the Tea Party | bargaining law </s> a letter to members </s> <false | 3474 |

| | | | |
|------|--|---|------|
| 3475 | statement 3.131746292114258> Tom Barrett says | Clinton who left Americans in harms way in Beng- | 3526 |
| 3476 | Gov. Scott Walker said no to equal pay for equal | hazi and after four Americans fell said, What | 3527 |
| 3477 | work for women. </s> a TV ad </s> <true state- | difference at this point does it make? </s> the | 3528 |
| 3478 | ment 3.1403446197509766> Portland Association | Republican national convention </s> <true state- | 3529 |
| 3479 | Teachers says Did you know that if you accepted | ment 2.5874826908111572> Hillary Clinton says | 3530 |
| 3480 | the Districts proposal today you would have NO | When terrorists killed more than 250 Americans | 3531 |
| 3481 | pay increase for 4 years? Seven years of frozen | in Lebanon under Ronald Reagan, the Democrats | 3532 |
| 3482 | wages = Disrespect. </s> a newsletter | didn't make that a partisan issue. </s> a CNN town | 3533 |
| 3483 | Label true statement | hall </s> <false statement 2.849807024002075> | 3534 |
| 3484 | Test Modified Rick Scott says All Aboard | Donald Trump says Sidney Blumenthal wrote | 3535 |
| 3485 | Florida is a 100 percent private venture. There | that the Benghazi attack was almost certainly pre- | 3536 |
| 3486 | is no state money involved. </s> a TV interview | ventable. Clinton was in charge of the State Depart- | 3537 |
| 3487 | </s> <true statement 3.0582022666931152> Char- | ment, and it failed to protect U.S. personnel and | 3538 |
| 3488 | lie Crist says All Aboard Florida is receiving mil- | an American consulate in Libya. </s> a rally in | 3539 |
| 3489 | lions in Florida taxpayer dollars. </s> a fundrais- | Wilkes-Barre, Pa. | 3540 |
| 3490 | ing email </s> <false statement 3.152191162109375> | Label true statement | 3541 |
| 3491 | Corey Lewandowski says Mr. Trump is self- | Test Modified Rand Paul says Of the roughly | 3542 |
| 3492 | financing his campaign, so we don't have any | 15 percent of Americans who don't have health | 3543 |
| 3493 | donors. </s> a radio interview. | insurance, half of them made more than 50,000 | 3544 |
| 3494 | Label true statement | a year. </s> an interview on Comedy Cen- | 3545 |
| 3495 | Test Modified Julie Pace says The Obama ad- | tral's "The Daily Show" </s> <false statement | 3546 |
| 3496 | ministration is using as its legal justification for | 2.9004263877868652> Rand Paul says Over half | 3547 |
| 3497 | these airstrikes (on the Islamic State), an autho- | of the young people in medical, dental and law | 3548 |
| 3498 | rization for military force that the president him- | schools are women. </s> an interview with CNN | 3549 |
| 3499 | self has called for repeal of. </s> a question to | </s> <true statement 2.932455062866211> Joe | 3550 |
| 3500 | White House Press Secretary Josh Earnest </s> | Biden says Among the money spent on health care | 3551 |
| 3501 | <true statement 2.962770462036133> Martha Rad- | in the United States, "46 cents on every dollar spent | 3552 |
| 3502 | datz says The Obama administration originally | is through Medicare and Medicaid." </s> an inter- | 3553 |
| 3503 | wanted 10,000 troops to remain in Iraq – not com- | view on NBC's 'Meet the Press' | 3554 |
| 3504 | bat troops, but military advisers, special opera- | Label false statement | 3555 |
| 3505 | tions forces, to watch the counterterrorism effort. </s> | Test Modified Barack Obama says Stimulus tax | 3556 |
| 3506 | comments on ABC's "This Week" </s> <false state- | cuts "began showing up in paychecks of 4.8 million | 3557 |
| 3507 | ment 2.9962246417999268> Rand Paul says The | Indiana households about three months ago." </s> | 3558 |
| 3508 | president is advocating a drone strike program in | a speech in Wakarusa, Ind. </s> <false statement | 3559 |
| 3509 | America. </s> a tweet | 2.8908281326293945> Paul Broun says Stimulus | 3560 |
| 3510 | Label true statement | money funded a government board that made rec- | 3561 |
| 3511 | Test Modified John Kasich says We are now | ommendations that would cost 378,000 jobs and | 3562 |
| 3512 | eighth in the nation in job creation . . . we are | 28.3 billion in sales. </s> a tweet </s> <true state- | 3563 |
| 3513 | No. 1 in the Midwest. </s> a news conference </s> | ment 2.898074150085449> Chain Email says Hav- | 3564 |
| 3514 | <false statement 2.610369920730591> Ted Strick- | ing an entirely Democrat congressional delegation | 3565 |
| 3515 | land says Gov. John Kasich incorrectly claimed | in 2009, when the [federal stimulus] bill passed, | 3566 |
| 3516 | Ohio's economy was 38th in the nation when he | increases the per capita stimulus dollars that the | 3567 |
| 3517 | took office. We were sixth in the nation in terms of | state receives per person by 460. </s> a message | 3568 |
| 3518 | economic job growth. </s> an interview on CNN | via the Internet | 3569 |
| 3519 | </s> <true statement 2.896986246109009> John | Label false statement | 3570 |
| 3520 | Kasich says We are in the bottom 10 in dollars in | Test Modified Allen West says If you look at | 3571 |
| 3521 | the classroom and the top 10 in dollars in the bu- | the application for a security clearance, I have a | 3572 |
| 3522 | reaucracy and red tape. </s> an interview on Fox | clearance that even the president of the United | 3573 |
| 3523 | News | States cannot obtain because of my background. | 3574 |
| 3524 | Label true statement | </s> a candidate forum </s> <false statement | 3575 |
| 3525 | Test Modified Mike Pence says It was Hillary | 3.02140736579895> Steve Southerland says 92 | 3576 |

3577 percent of President Barack Obamas administra-
3578 tion has never worked outside government. </s>
3579 comments at the Liberty County Chamber of
3580 Commerce annual dinner. </s> <true statement
3581 3.1747167110443115> John McCain says "The fact
3582 is it's not amnesty." </s> a debate in Manchester,
3583 N.H.

3584 **Label** false statement

3585 **Test Modified** Bernie S says We now work
3586 the longest hours of any people around the world.
3587 </s> a C-SPAN interview </s> <false statement
3588 3.0254147052764893> Bernie S says We spend
3589 twice as much per capita on health care as any
3590 other nation on Earth. </s> an appearance on
3591 the Rachel Maddow Show </s> <true statement
3592 3.08957576751709> Jim Sensenbrenner says We
3593 have the highest corporate tax rate in the world. Its
3594 35 percent. </s> an interview

3595 **Label** false statement

3596 **Test Modified** Sarah Palin says Donald Trumps
3597 conversion to pro-life beliefs are akin to Justin
3598 Biebers, who said in the past that abortion was
3599 no big deal to him. </s> an interview on CNN
3600 </s> <false statement 2.7887768745422363> Don-
3601 ald Trump says Public support for abortion is actu-
3602 ally going down a little bit, polls show. </s> com-
3603 ments on CNN's "State of the Union" </s> <true
3604 statement 3.1297004222869873> Greg Abbott says
3605 After Texas defunded Planned Parenthood, both the
3606 unintended pregnancy and abortion rates dropped.
3607 </s> a tweet

3608 **Label** false statement