

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

Anonymous ACL submission

Abstract

Text classification systems have impressive capabilities but are infeasible to deploy and use reliably due to their dependence on prompting and billion-parameter language models. SetFit (Tunstall et al., 2022) is a recent, practical approach that fine-tunes a Sentence Transformer under a contrastive learning paradigm and achieves similar results to more unwieldy systems. Text classification is important for addressing the problem of domain drift in detecting harmful content, which plagues social media platforms. Here, we propose Like a Good Nearest Neighbor (LAGONN), a modification to SetFit that requires no additional parameters or hyperparameters 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 identifying harmful content and generally improves SetFit’s performance. To demonstrate LAGONN’s value, we conduct a thorough study of text classification systems in the context of content moderation under four label distributions.¹

1 Introduction

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

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

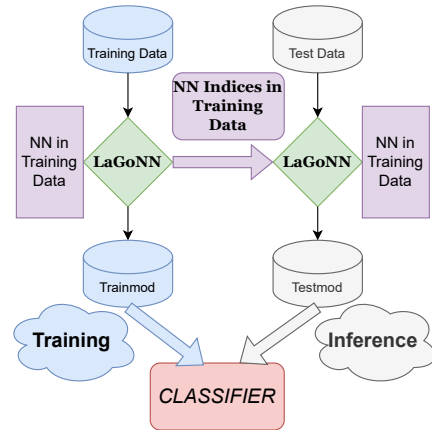


Figure 1: We embed training data, retrieve the text, gold label, and distance for each instance from its nearest neighbor and modify the original text with this information. Then we embed the modified training data and train a classifier. During inference, the NN from the training data is selected, the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

use because of their reliance on billion-parameter PLMs, pay-to-use APIs, and/or prompting. Constructing prompts is not trivial and may require domain expertise.

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

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

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

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

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

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

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

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

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

2 Related Work

There is 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 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) 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

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) |
| DIST | "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]" (?) |
| DIST | "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.5 for examples of LAGONN ALL modified text.

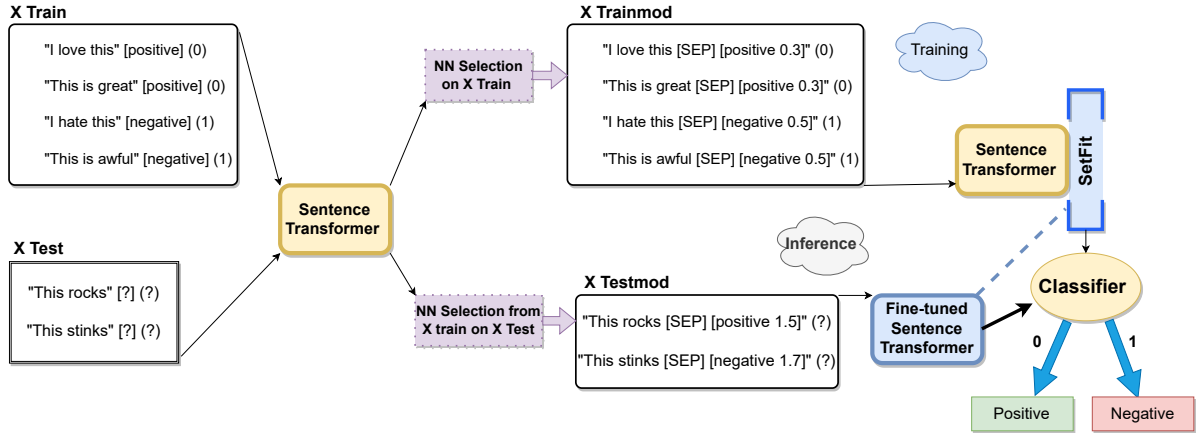


Figure 2: LAGONN LABDIST uses an ST to encode training data, performs NN lookup, appends the NN’s gold label and distance, and optionally SetFit to fine-tune the embedding model. We then embed this new instance and train a classifier. During inference, we use the embedding model to modify the test data 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.

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

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 neighbor and use it to decorate the original example. We experimented with different text that LAGONN could use. The first configuration we consider is the gold label of the NN, which we call LABEL. We then consider the Euclidean distance of the NN, which we call DIST, giving the model access to 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 2). Information from the NN is appended to the text following a separator token to indicate this instance is composed of multiple sequences. While the ALL and TEXT configurations are arguably the most interesting, we find LABDIST to result in the most performant version of LAGONN, and this is the version about which we report results. See Appendix A.4.1 for a detailed study of and comparison between 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. As LAGONN utilizes similarity and clear distinctions between classes, we believe it fitting for our use case of content moderation, where it is realistic to have few labels, harmful or neutral, for example.

4 Experiments

4.1 Data and label distributions

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

Here, we provide a summary of the five datasets we studied. LIAR was created from Politifact⁴ for fake news detection and is composed of the data fields *context*, *speaker*, and *statement*, which are labeled with varying levels of truthfulness (Wang, 2017). We used a collapsed version of this dataset where a statement can only be true or false. We did not use *speaker*, but did use *context* and *statement*, separated by a separator token. Quora Insincere Questions⁵ is composed of neutral and toxic questions, where the author is not asking in good faith. Hate Speech Offensive⁶ has three labels and is composed of tweets that can contain either neutral text, offensive language, or hate speech (Davidson et al., 2017). Amazon Counterfactual⁷ 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 Conversations⁸ is a dataset of comments where the author wrote with unintended bias⁹ (see Table 2).

We study our system by simulating growing

⁴<https://www.politifact.com/>

⁵<https://www.kaggle.com/c/quora-insincere-questions-classification>

⁶https://huggingface.co/datasets/hate_speech_offensive

⁷https://huggingface.co/datasets/SetFit/amazon_counterfactual_en

⁸https://huggingface.co/datasets/SetFit/toxic_conversations

⁹<https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification/overview>

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

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

training data over ten discrete steps sampled under four different label distributions: extreme, imbalanced, moderate, and balanced (see Table 3). 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.¹⁰ On each step, we train our method with the sampled data and evaluate on the test split. Considering growing training data has two benefits: 1) We can simulate a streaming data scenario, where new data is labeled and added for training and 2) We can investigate each method’s sensitivity to the number of training examples. We sampled over five seeds, reporting the mean and standard deviation.

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

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

4.2 Baselines

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

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

¹⁰For Hate Speech Offensive, 0 and 2 denote undesirable text and 1 denotes neither.

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.¹¹ 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 also uses logistic regression as its final classifier and refer to this method as SetFit.

***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

¹¹<https://huggingface.co/sentence-transformers/paraphrase-mpnet-base-v2>

336 it is comparable to SetFit_{exp} . On the first step, this
337 method is equivalent to LAGONN.

338 5 Results

339 Table 4 and Figure 3 show our results. In the
340 cases of the extreme and imbalanced regimes, the
341 performance of SetFit_{exp} steadily increases with
342 the number of training examples. As the label
343 distribution shifts to the balanced regime, how-
344 ever, the performance quickly saturates or even
345 degrades as the number of training examples grows.
346 LAGONN, RoBERTa_{full} , and SetFit , other fine-
347 tuned PLM classifiers, do not exhibit this behavior.
348 LAGONN_{exp} , being based on SetFit_{exp} , exhibits a
349 similar trend, but the performance degradation is
350 mitigated; on the 10th step of Amazon Counterfactual
351 in Table 4 SetFit_{exp} 's performance decreased
352 by 9.7, while LAGONN_{exp} only fell by 3.7.

353 LAGONN and LAGONN_{exp} generally outper-
354 form SetFit and SetFit_{exp} , respectively, often re-
355 sulting in a more stable model, as reflected in the
356 standard deviation. We find that LAGONN and
357 LAGONN_{exp} exhibit stronger predictive power
358 with fewer examples than RoBERTa_{full} despite
359 having fewer trainable parameters. For example,
360 on the first step of Insincere Questions under the
361 extreme setting, LAGONN's performance is more
362 than 10 points higher.

363 LAGONN_{cheap} outperforms all other methods
364 on the Insincere Questions dataset for all balance
365 regimes, despite being the third fastest (see Table
366 5) and having the second fewest trainable param-
367 eters. We attribute this result to the fact that this
368 dataset is composed of questions from Quora¹² and
369 our ST backbone was pretrained on similar data.
370 This intuition is supported by Probe, the cheapest
371 method, which despite having the fewest trainable
372 parameters, shows comparable performance.

373 5.1 SetFit for efficient many-shot learning

374 Respectively comparing SetFit to SetFit_{exp} and
375 LAGONN to LAGONN_{exp} suggests that fine-
376 tuning the ST embedding model on moderate or bal-
377 anced data hurts model performance as the number
378 of training samples grows. We therefore hypoth-
379 esize that randomly sampling a subset of training
380 data to fine-tune the encoder, freezing, embedding
381 the remaining data, and training the classifier will
382 result in a stronger model.

¹²<https://www.quora.com/>

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-
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.

Figure 4 shows 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 on
Toxic Conversations. We note very different re-
sults for both LAGONN_{lite} and SetFit_{lite} compared
to LAGONN_{exp} and SetFit_{exp} on Toxic Conversa-
tions and Amazon Counterfactual under the moder-
ate and balanced label distributions. As their expen-
sive counterparts start to plateau or degrade on the
fourth step, the predictive power of these two new
models dramatically increases, showing improved
or comparable performance to RoBERTa_{full} , de-
spite being optimized on less data; for example,
 LAGONN_{lite} reaches an average precision of ap-
proximately 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 pro-
vide a performance boost for SetFit -based classi-
fication.

5.2 LAGONN's computational expense

LAGONN is more computationally expensive than
Sentence Transformer- or SetFit -based text classi-
fication. LAGONN introduces additional inference
with the encoder, NN-lookup, and string modifi-
cation. As the computational complexity of trans-
formers increases with sequence length (Vaswani
et al., 2017), additional expense is created when
LAGONN appends textual information before in-
ference with the ST. In Table 5, we provide a speed
comparison of comparable methods computed on
the same hardware.¹³ On average, LAGONN in-
troduced 24.2 additional seconds of computation
compared to its relative counterpart.

6 Discussion

Flagging potentially dangerous text presents a chal-
lenge even for state-of-the-art approaches. 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

¹³We used a 40 GB NVIDIA A100 Tensor Core GPU.

| Method | InsincereQs | | | | AmazonCF | | | |
|-------------------------|----------------------------|----------------------------|----------------------------|----------------------------|-----------------------------|----------------------------|----------------------------|-----------------------------|
| | 1 st | 5 th | 10 th | Average | 1 st | 5 th | 10 th | Average |
| <i>Extreme</i> | | | | | | | | |
| 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} |
| <i>Balanced</i> | | | | | | | | |
| RoBERTa _{full} | 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} |
| <i>Balanced</i> | | | | | | | | |
| RoBERTa _{full} | 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 4: Average performance (average precision \times 100) on Insincere Questions and Amazon Counterfactual. The first, fifth, and tenth step are followed by the average over all ten steps. The average gives insight into the overall strongest performer by aggregating all steps. We group methods with a comparable number of trainable parameters together. The extreme label distribution results are followed by balanced (see Appendix A.2 for additional results).

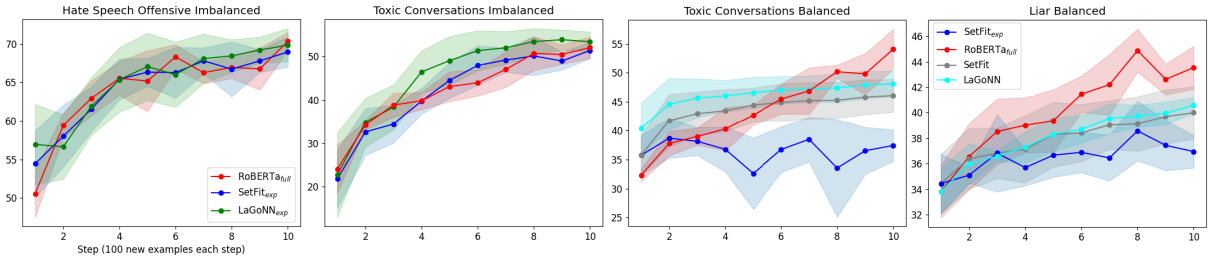


Figure 3: Average performance in the imbalanced and balanced regimes relative to comparable methods. We include RoBERTa_{full} results for reference. The metric is macro-F1 for Hate Speech Offensive, average precision elsewhere.

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

of hate speech, toxicity, and fake news. LAGONN exploits semantic similarity and clear boundaries between labels, which we believe is reflected in scenarios with fewer classes, such as quickly filtering out harmful content.

Our results suggest that LAGONN_{exp} or SetFit_{exp}, relatively expensive techniques, can detect harmful content when dealing with imbalanced label distributions, 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

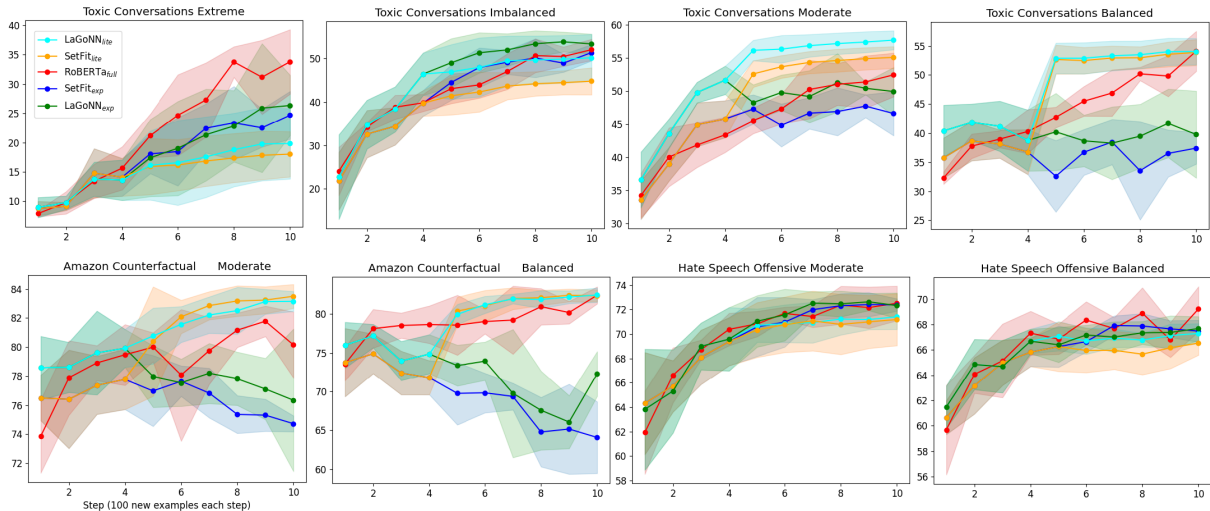


Figure 4: Average performance for all sampling regimes on Toxic Conversations and the moderate and balanced regimes for Amazon Counterfactual and Hate Speech Offensive. More expensive models, such as LAGoNN_{exp} , SetFit_{exp} , and RoBERTa_{full} perform best when the label distribution is imbalanced.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

LAGoNN_{cheap} excelled in the extreme and balanced settings. This highlights the fact that we can inexpensively extract pretrained knowledge if PLMs are chosen with care for related tasks.

Standard fine-tuning with SetFit does not help performance on more balanced datasets that are not few-shot. SetFit was developed for few-shot learning, but we have observed that it should not be applied "out of the box" to balanced, non-few-shot data. This can be detrimental to performance, directly affecting our own approach. However, we have observed that LAGoNN can stabilize SetFit's predictions and reduce its performance drop. Figures 3 and 4 show that when the label distribution is moderate or balanced (see Table 3), SetFit_{exp} plateaus, yet cheaper systems, such as LAGoNN, continue to learn. We believe this is due to SetFit's fine-tuning objective, which optimizes an ST using cosine similarity loss to separate examples belonging to different labels in feature space, assuming independence between labels. This may be too strong an assumption as we optimize with more examples, which is counter-intuitive for data-hungry transformers; RoBERTa_{full} , optimized with cross-entropy loss, generally showed improved performance as we added training data.

When dealing with balanced data, it is sufficient to fine-tune the Sentence Transformer via SetFit with 50 to 100 examples per label, while 150 to 200 instances appear to be sufficient when the training data are moderately balanced. The encoder can

then be frozen and all available data embedded to train a classifier. This improves performance and is more efficient than full-model fine-tuning. LAGoNN is directly applicable to this case, boosting the performance of SetFit_{lite} without introducing trainable parameters. In this setup, all models fine-tuned on Hate Speech Offensive exhibited similar, upward-trending learning curves, but we note the speed of LAGoNN relative to RoBERTa_{full} or SetFit_{exp} (see Figure 4 and Table 5).

7 Conclusion

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

8 Limitations

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

9 Ethics Statement

It is our sincere goal that our work contributes to the social good in multiple ways. We first hope to have furthered research on text classification that can be feasibly applied to combat undesirable content, such as misinformation, on the Internet, which could potentially cause someone harm. To this end, we have tried to describe our approach as accurately as possible and released our code and data, such that our work is transparent and can be easily reproduced and expanded upon. We hope that we have also created a useful but efficient system which reduces the need to expend energy in the form expensive computation. For example, LAGONN does

not rely on billion-parameter language models that demand thousand-dollar GPUs to use. LAGONN makes use of GPUs no more than SetFit, despite being more computationally expensive. We have additionally proposed a simple method to make SetFit, an already relatively inexpensive method, even more efficient.

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, Matiusz 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.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of](#)

| | | | |
|-----|---|--|--|
| 611 | deep bidirectional transformers for language understanding. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics. | <i>Empirical Methods in Natural Language Processing</i> , pages 7092–7108, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. | 667 668 669 670 |
| 612 | | | |
| 613 | | | |
| 614 | | | |
| 615 | | | |
| 616 | | | |
| 617 | | | |
| 618 | Ulrich K. H. Ecker, Stephan Lewandowsky, John Cook, Philipp Schmid, Lisa K. Fazio, Nadia Brashier, Panayiota Kendeou, Emily K. Vraga, and Michelle A. Amazeen. 2022. The psychological drivers of misinformation belief and its resistance to correction. <i>Nature Reviews Psychology</i> , 1(1):13–29. | Christian S. Perone, Roberto Pereira Silveira, and Thomas S. Paula. 2018. Evaluation of sentence embeddings in downstream and linguistic probing tasks. <i>arXiv preprint arXiv:1806.06259</i> . | 671 672 673 674 |
| 619 | | | |
| 620 | | | |
| 621 | | | |
| 622 | | | |
| 623 | | | |
| 624 | Rabeeh Karimi Mahabadi, Luke Zettlemoyer, James Henderson, Lambert Mathias, Marzieh Saeidi, Veselin Stoyanov, and Majid Yazdani. 2022. Prompt-free and efficient few-shot learning with language models. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3638–3652, Dublin, Ireland. Association for Computational Linguistics. | Guangyuan Piao. 2021. Scholarly text classification with sentence bert and entity embeddings. In <i>Trends and Applications in Knowledge Discovery and Data Mining</i> , pages 79–87, Cham. Springer International Publishing. | 675 676 677 678 679 |
| 625 | | | |
| 626 | | | |
| 627 | | | |
| 628 | | | |
| 629 | | | |
| 630 | | | |
| 631 | | | |
| 632 | Gregory Koch, Richard Zemel, Ruslan Salakhutdinov, et al. 2015. Siamese neural networks for one-shot image recognition. In <i>ICML Deep Learning Workshop</i> , volume 2, page 0. Lille. | Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics. | 680 681 682 683 684 685 686 687 |
| 633 | | | |
| 634 | | | |
| 635 | | | |
| 636 | Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohata, Tenghao Huang, Mohit Bansal, and Colin Raffel. 2022a. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. <i>arXiv preprint arXiv:2205.05638</i> . | Sarah T. Roberts. 2017. <i>Content Moderation</i> , pages 1–4. Springer International Publishing, Cham. | 688 689 |
| 637 | | | |
| 638 | | | |
| 639 | | | |
| 640 | | | |
| 641 | Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022b. What makes good in-context examples for GPT-3? In <i>Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures</i> , pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics. | Stephen Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: Bm25 and beyond. <i>Found. Trends Inf. Retr.</i> , 3(4):333–389. | 690 691 692 |
| 642 | | | |
| 643 | | | |
| 644 | | | |
| 645 | | | |
| 646 | | | |
| 647 | | | |
| 648 | | | |
| 649 | Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> . | Timo Schick and Hinrich Schütze. 2021a. Exploiting cloze-questions for few-shot text classification and natural language inference. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 255–269, Online. Association for Computational Linguistics. | 693 694 695 696 697 698 699 |
| 650 | | | |
| 651 | | | |
| 652 | | | |
| 653 | | | |
| 654 | Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. <i>Introduction to Information Retrieval</i> . Cambridge University Press, USA. | Timo Schick and Hinrich Schütze. 2021b. It’s not just size that matters: Small language models are also few-shot learners. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2339–2352, Online. Association for Computational Linguistics. | 700 701 702 703 704 705 706 |
| 655 | | | |
| 656 | | | |
| 657 | Todor Markov, Chong Zhang, Sandhini Agarwal, Tyna Eloundou, Teddy Lee, Steven Adler, Angela Jiang, and Lilian Weng. 2022. A holistic approach to undesired content detection. <i>arXiv preprint arXiv:2208.03274</i> . | Timo Schick and Hinrich Schütze. 2022. True few-shot learning with Prompts—A real-world perspective. <i>Transactions of the Association for Computational Linguistics</i> , 10:716–731. | 707 708 709 710 |
| 658 | | | |
| 659 | | | |
| 660 | | | |
| 661 | | | |
| 662 | James O’Neill, Polina Rozenshtein, Ryuichi Kiryo, Motoko Kubota, and Danushka Bollegala. 2021. I wish I would have loved this one, but I didn’t – a multilingual dataset for counterfactual detection in product review. In <i>Proceedings of the 2021 Conference on</i> | Yusuke Shido, Hsien-Chi Liu, and Keisuke Umezawa. 2022. Textual content moderation in C2C marketplace. In <i>Proceedings of the Fifth Workshop on e-Commerce and NLP (ECNLP 5)</i> , pages 58–62, Dublin, Ireland. Association for Computational Linguistics. | 711 712 713 714 715 716 |
| 663 | | | |
| 664 | | | |
| 665 | | | |
| 666 | | | |
| | | Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tiejian Liu. 2020. MpNet: Masked and permuted pre-training for language understanding. In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 16857–16867. Curran Associates, Inc. | 717 718 719 720 721 |

| | | | |
|-----|---|--|-----|
| 722 | Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. | Meng Ye, Karan Sikka, Katherine Atwell, Sabit Hassan, | 778 |
| 723 | Conceptnet 5.5: An open multilingual graph of gen- | Ajay Divakaran, and Malihe Alikhani. 2023. Multi- | 779 |
| 724 | eral knowledge . <i>Proceedings of the AAAI Conference</i> | lingual content moderation: A case study on reddit . | 780 |
| 725 | on Artificial Intelligence , 31(1). | <i>arXiv preprint arXiv:2302.09618</i> . | 781 |
| 726 | Lewis Tunstall, Nils Reimers, Unso Eun Seo Jo, Luke | A Appendix | 782 |
| 727 | Bates, Daniel Korat, Moshe Wasserblat, and Oren | A.1 Observations about LAGONN | 783 |
| 728 | Pereg. 2022. Efficient few-shot learning without | Here, at the suggestion of an anonymous reviewer, | 784 |
| 729 | prompts . <i>arXiv preprint arXiv:2209.11055</i> . | we include a little background on LAGONN. We | 785 |
| 730 | Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob | originally attempted to use Sentence Transform- | 786 |
| 731 | Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz | ers/SetFit as a retrieval model that would modify | 787 |
| 732 | Kaiser, and Illia Polosukhin. 2017. Attention is all | input text and then pass this input to a Transformer- | 788 |
| 733 | you need . In <i>Advances in Neural Information Pro-</i> | based classifier, such as RoBERTa, instead of back | 789 |
| 734 | cessing Systems , volume 30. Curran Associates, Inc. | into the ST as in LaGoNN. We experimented with | 790 |
| 735 | Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. | different ST retrieval models and Transformer clas- | 791 |
| 736 | The spread of true and false news online . <i>Science</i> , | sifiers, but this system was often beaten by base- | 792 |
| 737 | 359(6380):1146–1151. | lines, and performant versions were too expensive | 793 |
| 738 | Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman- | to justify their use. The failure of this system is | 794 |
| 739 | preet Singh, Julian Michael, Felix Hill, Omer Levy, | what ultimately inspired LAGONN. We had hoped | 795 |
| 740 | and Samuel Bowman. 2019. Superglue: A stickier | to construct a system that did not need to be up- | 796 |
| 741 | benchmark for general-purpose language understand- | dated after step one and could simply perform infer- | 797 |
| 742 | ing systems . In <i>Advances in Neural Information</i> | ence on subsequent steps, an active learning setup. | 798 |
| 743 | Processing Systems , volume 32. Curran Associates, | While the performance of this version of LAGONN | 799 |
| 744 | Inc. | did not degrade, it also did not appear to learn any- | 800 |
| 745 | Shuohang Wang, Yichong Xu, Yuwei Fang, Yang Liu, | thing and we found it necessary to update parame- | 801 |
| 746 | Siqi Sun, Ruochen Xu, Chenguang Zhu, and Michael | ters on each step. We additionally tried fine-tuning | 802 |
| 747 | Zeng. 2022. Training data is more valuable than you | the embedding model via SetFit first before mod- | 803 |
| 748 | think: A simple and effective method by retrieving | ifying data, however, this hurt performance in all | 804 |
| 749 | from training data . In <i>Proceedings of the 60th Annual</i> | cases. We include this information for transparency | 805 |
| 750 | Meeting of the Association for Computational Lin- | and because we find it interesting. | 806 |
| 751 | guistics (Volume 1: Long Papers) , pages 3170–3179, | A.2 Additional results for initial experiments | 807 |
| 752 | Dublin, Ireland. Association for Computational Lin- | Here we provide additional results from our ini- | 808 |
| 753 | guistics. | tial experimental setup that, due to space limita- | 809 |
| 754 | William Yang Wang. 2017. “Liar, liar pants on fire”: | could not be included in the main text. We | 810 |
| 755 | A new benchmark dataset for fake news detection . | note that a version of LAGONN outperforms or | 811 |
| 756 | In <i>Proceedings of the 55th Annual Meeting of the</i> | has the same performance of all methods, includ- | 812 |
| 757 | Association for Computational Linguistics (Volume 2: | ing our upper bound RoBERTa _{full} , on 54% of all | 813 |
| 758 | Short Papers) , pages 422–426, Vancouver, Canada. | displayed results, and is the best performer rela- | 814 |
| 759 | Association for Computational Linguistics. | tive to Sentence Transformer-based methods on | 815 |
| 760 | Thomas Wolf, Lysandre Debut, Victor Sanh, Julien | 72%. This excludes LAGONN _{cheap} . This method | 816 |
| 761 | Chaumond, Clement Delangue, Anthony Moi, Pier- | showed strong performance on the Insincere Ques- | 817 |
| 762 | ric Cistac, Tim Rault, Remi Louf, Morgan Funtow- | tions dataset, but hurts performance in other cases. | 818 |
| 763 | icz, Joe Davison, Sam Shleifer, Patrick von Platen, | In cases when SetFit-based methods do outper- | 819 |
| 764 | Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, | form our system, the performances are compara- | 820 |
| 765 | Teven Le Scao, Sylvain Gugger, Mariama Drame, | ble, usually within a point, yet they can be quite | 821 |
| 766 | Quentin Lhoest, and Alexander Rush. 2020. Trans- | dramatic when LAGONN-based methods are the | 822 |
| 767 | formers: State-of-the-art natural language processing . | strongest. Below, we report the mean average pre- | 823 |
| 768 | In <i>Proceedings of the 2020 Conference on Empirical</i> | cision $\times 100$ for all methods over five seeds with | 824 |
| 769 | Methods in Natural Language Processing: System | the standard deviation, except in the case of Hate | 825 |
| 770 | Demonstrations , pages 38–45, Online. Association | Speech Offensive, where the evaluation metric is | 826 |
| 771 | for Computational Linguistics. | | |
| 772 | Yichong Xu, Chenguang Zhu, Shuohang Wang, Siqi | | |
| 773 | Sun, Hao Cheng, Xiaodong Liu, Jianfeng Gao, | | |
| 774 | Pengcheng He, Michael Zeng, and Xuedong Huang. | | |
| 775 | 2021. Human parity on commonsenseqa: Aug- | | |
| 776 | menting self-attention with external attention . <i>arXiv</i> | | |
| 777 | preprint arXiv:2112.03254 , abs/2112.03254. | | |

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 | | | | |
|---------------------------|---------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| | <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} |
| 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 6: 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 | | | | |
|---------------------------|---------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| | <i>Moderate</i> | 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 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 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 | | | | |
|---------------------------|-----------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| | <i>Imbalanced</i> | 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 8: 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 | | | | |
|---------------------------|-----------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| | <i>Moderate</i> | 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 9: LAGONN_{exp} and LAGONN are the strongest performers on the first step, but LAGONN is strongest classifier on subsequent steps and is also the overall strongest performer based on the average over all steps.

| Method | Toxic Conversations | | | | |
|---------------------------|---------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| | <i>Extreme</i> | 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} |
| 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 10: 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 | | | | |
|---------------------------|---------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| | <i>Imbalanced</i> | 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} |
| 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 11: 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 | | | | |
|---------------------------|---------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| | <i>Moderate</i> | 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} |
| 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 12: 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 | | | | |
|---------------------------|---------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| | <i>Balanced</i> | 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} |
| 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 13: 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 14: 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 15: 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 16: 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 17: 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 18: 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 19: SetFit, SetFit_{exp}, LAGONN, and LAGONN_{exp} start out as the strongest performers. On the 5th step, SetFit is overtaken 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 20: 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 21: 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.3 Additional results for secondary experiments

Here, we provide additional results from our second set of experiments that, due to space limitations, could not be included in the main text. We note that a version of LAGONN outperforms or has the same performance of all methods, including our upper bound RoBERTa_{full}, on 60% of all displayed results, and is the best performer relative to Sentence Transformer-based methods on 65%. This excludes LAGONN_{cheap}. This method showed strong performance on the Insincere Questions dataset, but hurts performance in other cases.

In cases when SetFit-based methods do outperform our system, the performances are comparable, usually within one point, yet they can be quite

different when LAGONN-based methods are the strongest. Below, we report the mean average precision $\times 100$ for all methods over five seeds with the standard deviation, except in the case of Hate Speech Offensive, where the evaluation metric is the macro-F1. Each table shows the results for a given dataset and a given label-balance distribution on the first, fifth, and tenth step followed by the average for all ten steps. In the table caption we provide a summary/interpretation of the results for a given setting. Liar appears to be the most difficult dataset for all methods. This is expected because it likely does not include enough context to determine the truth of a statement.

| Method | Insincere Questions | | | | |
|---------------------------|----------------------------|----------------------------|----------------------------|----------------------------|---------|
| | <i>Extreme</i> | 1 st | 5 th | 10 th | Average |
| RoBERTa _{full} | 19.9 _{8.4} | 30.9 _{7.9} | 42.0 _{7.4} | 33.5 _{6.7} | |
| SetFit _{exp} | 24.1 _{6.3} | 29.2 _{6.7} | 36.7 _{7.3} | 31.7 _{3.4} | |
| LAGONN _{exp} | 30.7 _{8.9} | 37.6 _{6.1} | 39.0 _{6.1} | 36.1 _{2.3} | |
| SetFit _{lite} | 24.1 _{6.3} | 38.1 _{6.3} | 41.1 _{6.5} | 35.6 _{5.5} | |
| LAGONN _{lite} | 30.7 _{8.9} | 41.8 _{8.3} | 43.4 _{8.5} | 39.3 _{4.4} | |
| RoBERTa _{freeze} | 19.9 _{8.4} | 34.1 _{5.4} | 37.9 _{5.2} | 32.5 _{5.4} | |
| kNN | 6.8 _{0.4} | 15.9 _{3.4} | 16.9 _{4.3} | 14.4 _{3.0} | |
| SetFit | 24.1 _{6.3} | 31.7 _{4.9} | 36.1 _{5.4} | 31.8 _{3.6} | |
| LAGONN | 30.7 _{8.9} | 39.3 _{4.9} | 41.2 _{4.7} | 38.4 _{3.0} | |
| Probe | 24.3 _{8.4} | 39.8 _{5.6} | 44.8 _{4.2} | 38.3 _{6.2} | |
| LAGONN _{cheap} | 23.6 _{7.8} | 40.7 _{5.9} | 45.3 _{4.4} | 38.6 _{6.6} | |

Table 22: LAGONN, LAGONN_{lite}, and LAGONN_{exp} start out as the strongest models, but LAGONN_{lite} remains the most performant by the 10th step. It is also the overall strongest performer based on the average. We note the strength of LAGONN_{cheap} relative to far more expensive methods.

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

Table 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 | | | |
|---------------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| | 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 24: 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 25: 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 26: 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} | 77.1 _{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 27: 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 |
| <i>Moderate</i> | | | | |
| RoBERTa _{full} | 73.9 _{2.5} | 80.0 _{1.0} | 80.1 _{2.3} | 79.1 _{2.1} |
| SetFit _{exp} | 76.5 _{1.6} | 77.0 _{2.4} | 74.7 _{0.5} | 76.5 _{1.0} |
| LAGONN _{exp} | 78.6 _{2.2} | 78.0 _{2.1} | 76.3 _{4.9} | 78.2 _{1.0} |
| 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 28: 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 |
| <i>Balanced</i> | | | | |
| 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 29: 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 |
| <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} |
| 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} |
| <i>Freeze</i> | | | | |
| 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 30: 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 |
| <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} |
| 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} |
| <i>Freeze</i> | | | | |
| 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 31: 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 |
| <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} |
| 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} |
| <i>Freeze</i> | | | | |
| 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 32: 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 |
| <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} |
| 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} |
| <i>Freeze</i> | | | | |
| 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 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 | 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} |
| <i>Freeze</i> | | | | |
| 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 34: 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} |
| <i>Freeze</i> | | | | |
| 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 35: 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 36: 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 37: 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 38: 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>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 40: 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>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 39: 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>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 41: 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.

866 A.4 Ablations

867 In this Appendix section, we perform ablation stud-
868 ies with LAGONN to support our findings in the
869 main text.

870 A.4.1 Ablation: LAGONN configurations

871 Here, we provide an in-depth comparison between
872 all LAGONN configurations, LABEL, DIST,
873 LABDIST, TEXT, and ALL (see Table 1) for all
874 datasets, balances, and levels of expense. The eval-
875 uation metric is the mean average precision ($\times 100$)
876 over five seeds in all cases except for Hate Speech
877 Offensive where the metric is the macro-F1.

878 Below, Figures 5 through 9 are the results for
879 the LAGONN_{cheap} training strategy, Figures 10
880 through 14 are the results for LAGONN, Figures
881 15 through 19 are the results for LAGONN_{lite},
882 and Figures 20 through 24 are the results for
883 LAGONN_{exp}. We place the figures on a new page
884 for ease of viewing.

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

901 Overall, we believe that LABDIST is the most
902 performant/stable configuration of LAGONN, and
903 it is about this version that we present results in the
904 main text. We note that we could have presented the
905 best performer for each evaluation scenario, how-
906 ever, this is not in the spirit of our work as it adds
907 yet another hyperparameter to configure, standing
908 in the way of practical usage and convoluting our
909 analysis. However, in our codebase, we hope that
910 we have made it easy for one to change these con-
911 figurations for their own usage, be it scientific or
912 otherwise.

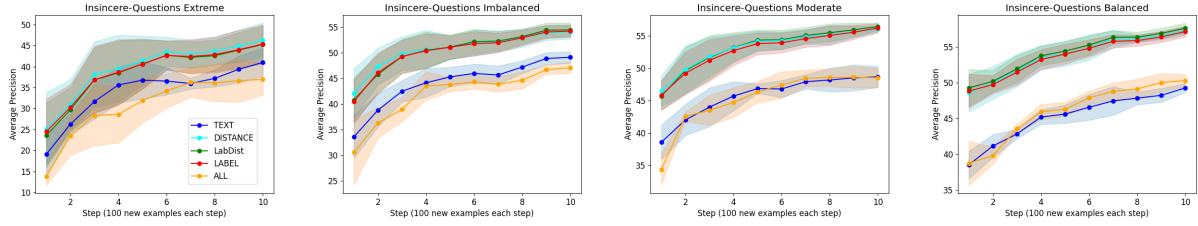


Figure 5: LAGONN_{cheap} performance for all configurations and balance regimes on the Insincere Questions dataset. The relevant balance is in the title of each panel.

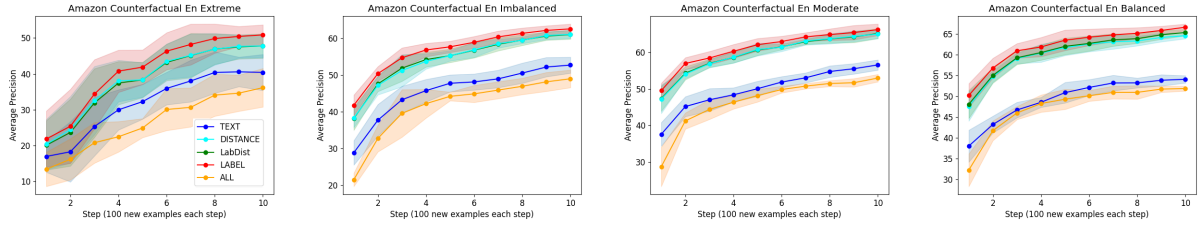


Figure 6: LAGONN_{cheap} performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

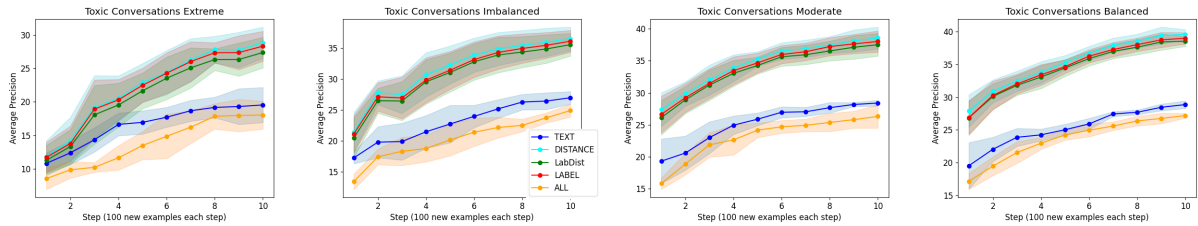


Figure 7: LAGONN_{cheap} performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

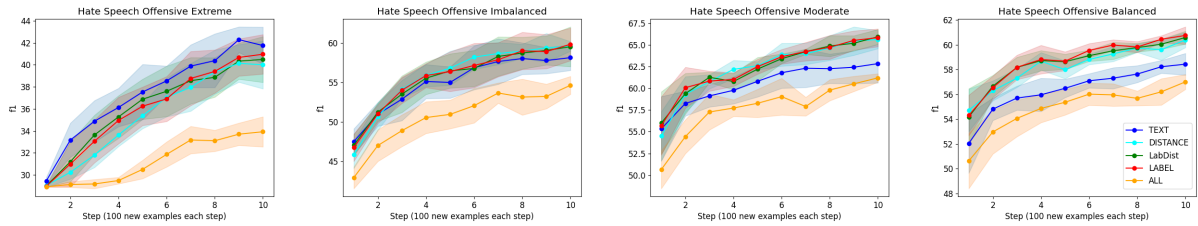


Figure 8: LAGONN_{cheap} performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

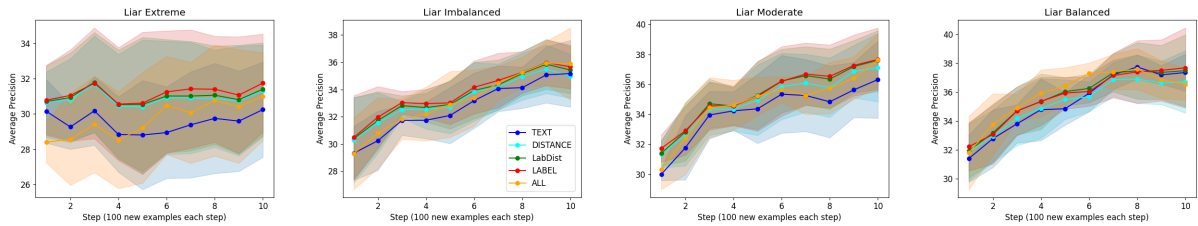


Figure 9: LAGONN_{cheap} performance for all configurations and balance regimes on the Liar dataset. The relevant balance is in the title of each panel.

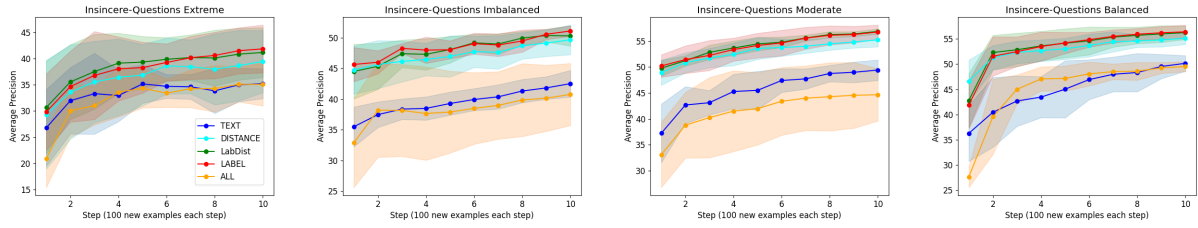


Figure 10: LAGoNN performance for all configurations and balance regimes on the InSincere Questions dataset. The relevant balance is in the title of each panel.

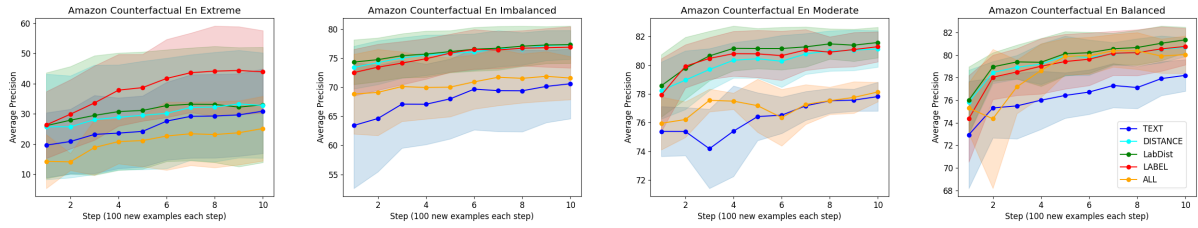


Figure 11: LAGoNN performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

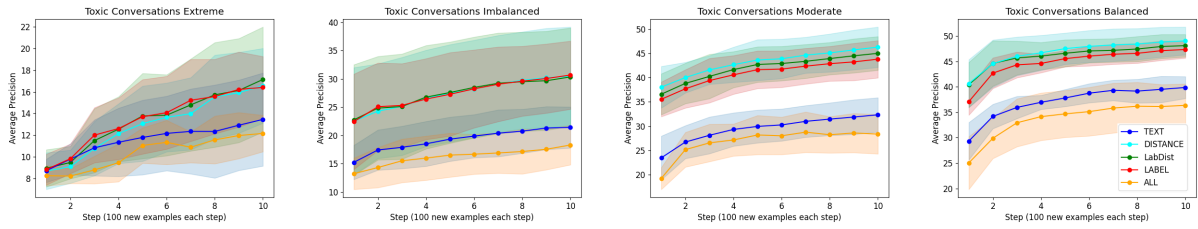


Figure 12: LAGoNN performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

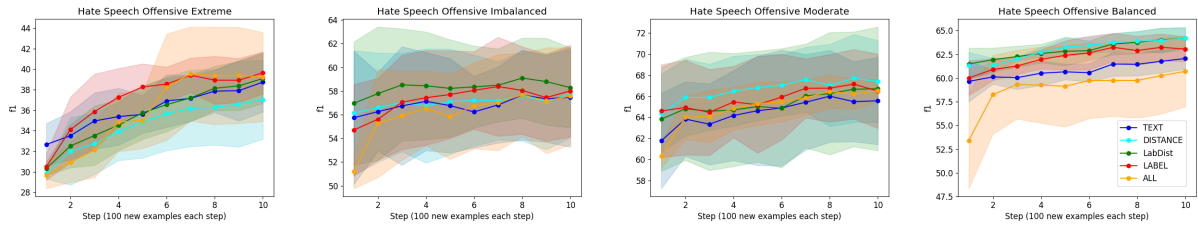


Figure 13: LAGoNN performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

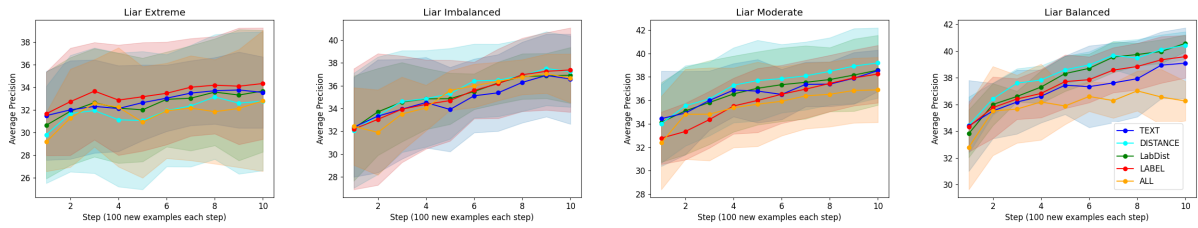


Figure 14: LAGoNN performance for all configurations and balance regimes on the Liar dataset. The relevant balance is in the title of each panel.

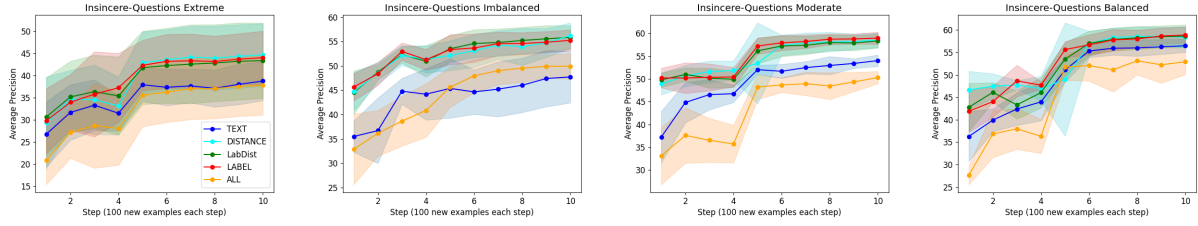


Figure 15: LAGONN_{lite} performance for all configurations and balance regimes on the Insincere Questions dataset. The relevant balance is in the title of each panel.

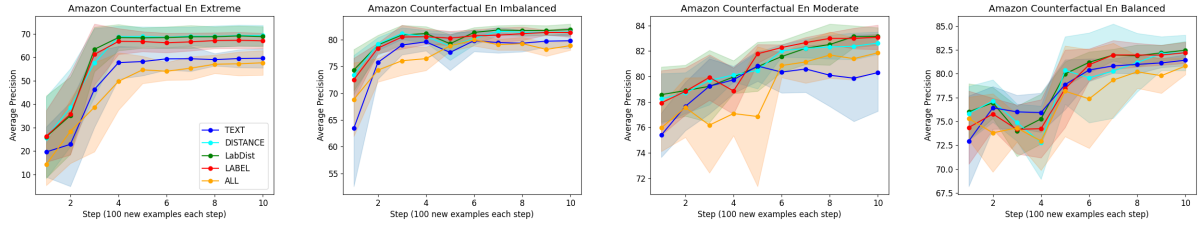


Figure 16: LAGONN_{lite} performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

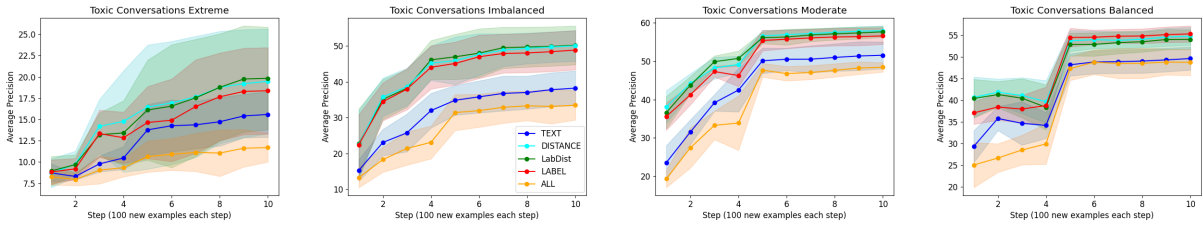


Figure 17: LAGONN_{lite} performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

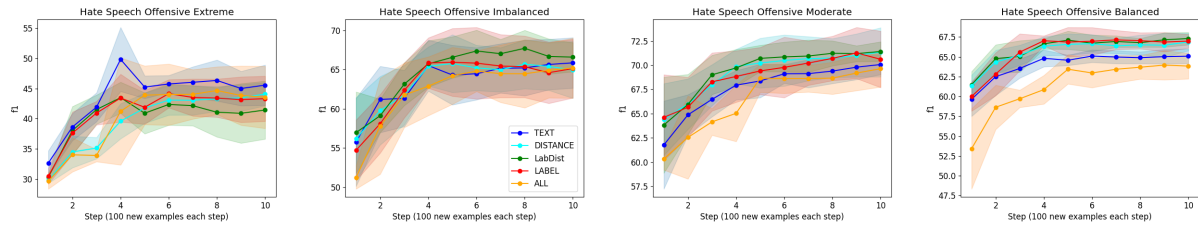


Figure 18: LAGONN_{lite} performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

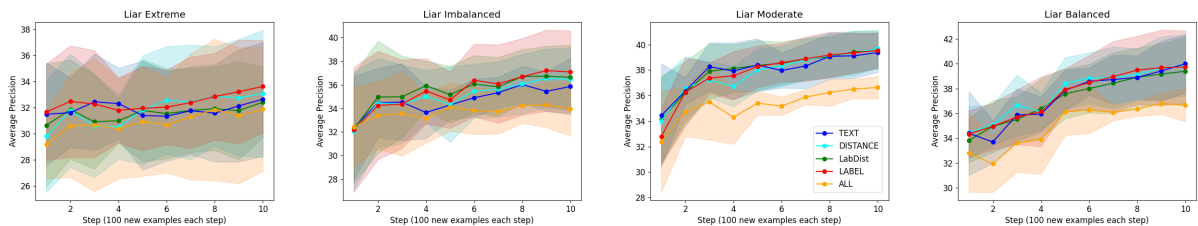


Figure 19: LAGONN_{lite} performance for all configurations and balance regimes on the Liar dataset. The relevant balance is in the title of each panel.

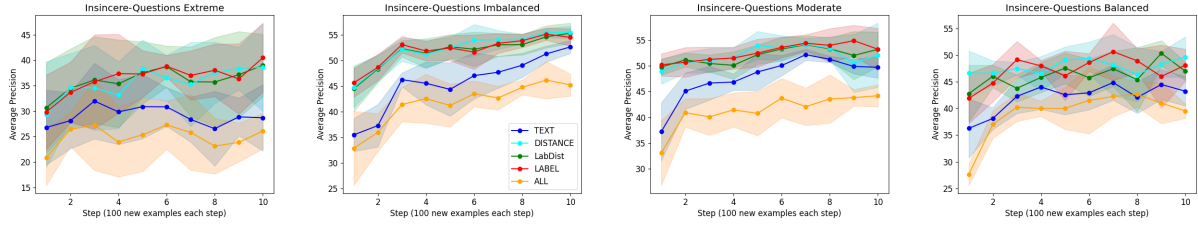


Figure 20: LAGONN_{exp} performance for all configurations and balance regimes on the Insincere Questions dataset. The relevant balance is in the title of each panel.

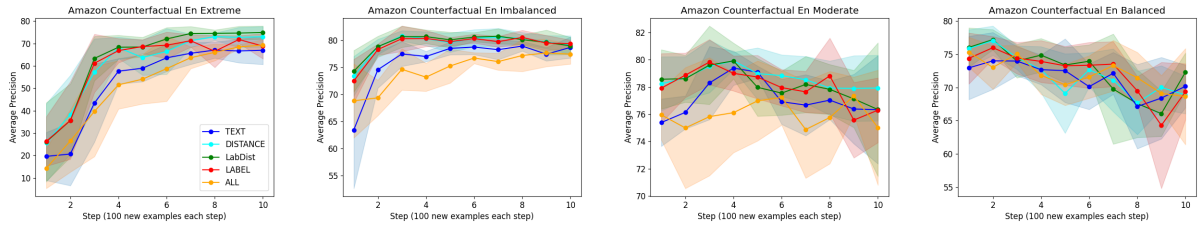


Figure 21: LAGONN_{exp} performance for all configurations and balance regimes on the Amazon Counterfactual dataset. The relevant balance is in the title of each panel.

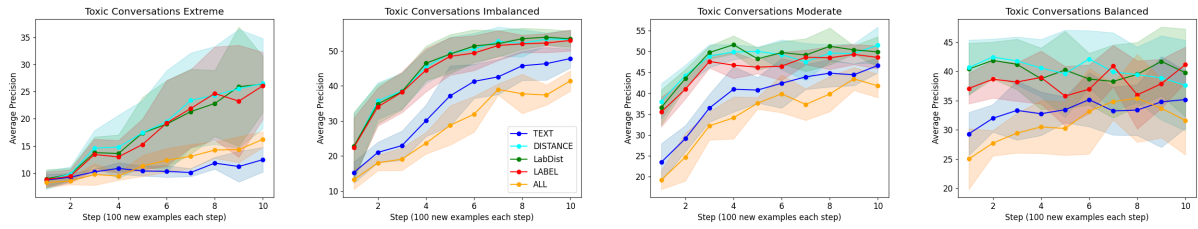


Figure 22: LAGONN_{exp} performance for all configurations and balance regimes on the Toxic Conversations dataset. The relevant balance is in the title of each panel.

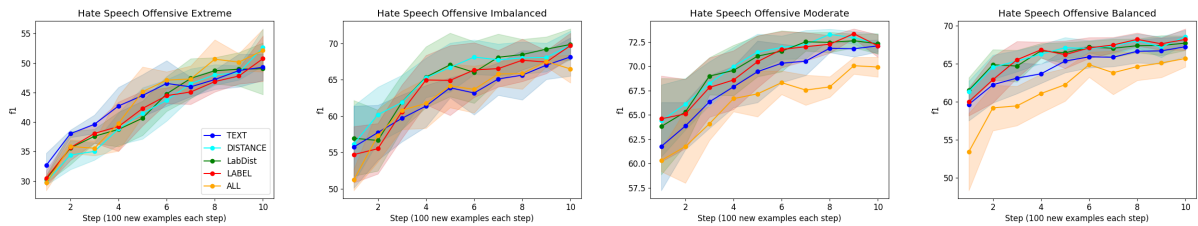


Figure 23: LAGONN_{exp} performance for all configurations and balance regimes on the Hate Speech Offensive dataset. The relevant balance is in the title of each panel.

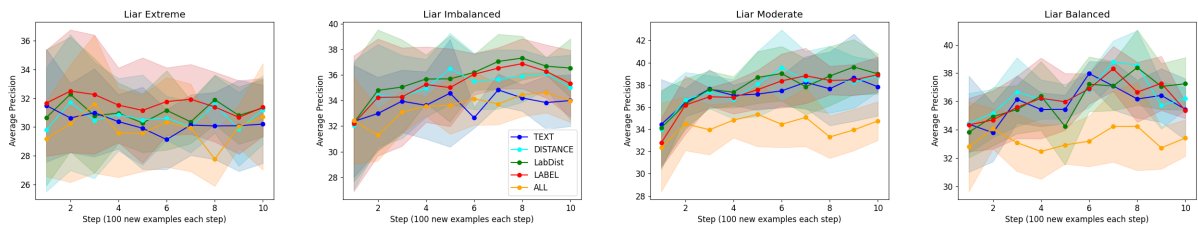


Figure 24: LAGONN_{exp} performance for all configurations and balance regimes on the Liar dataset. The relevant balance is in the title of each panel.

913 **A.4.2 Ablation: the effect of encoding distance**

914 Here, at the suggestion of an anonymous reviewer,
915 we present ablation results and analysis of how en-
916 coding distance affects LAGONN, because PLMs
917 often struggle to understand numbers. Note that
918 during our development stage, we ensured that our
919 tokenizer was capable of encoding floats with trail-
920 ing digits. To examine the effect of trailing digits
921 on LAGONN, we consider the DIST configuration
922 (see Table 1), where we append only the Euclidean
923 distance to the input text. In this ablation, however,
924 we round to different levels of precision. For exam-
925 ple, if the distance were a float of 0.123456789, we
926 round it to the nearest whole number, 0.0, single
927 digit float, 0.1, three digit float, 0.123, six digit
928 float, 0.123457, and finally keep it unrounded, that
929 is, the original DIST configuration, 0.123456789.
930 The below results are only for the LAGONN_{lite}
931 training strategy. We chose LAGONN_{lite} for this
932 ablation because it provides insight into both how
933 distance affects full-model fine-tuning and only re-
934 fitting the classification head. The results can be
935 seen below in Figures 25 through 29. We place the
936 figures on a new page for ease of viewing.

937 Interestingly, we tend to observe very similar per-
938 formance curves for all rounding precisions. The
939 exceptions to this would perhaps be Amazon Coun-
940 terfactual and Hate Speech Offensive in the bal-
941 anced regime where DIST and rounding to the
942 third trailing digit respectively exhibit large insta-
943 bility.

944 Although not always the case, it appears that
945 providing the model with the distance rounded to
946 the nearest whole number tends to result in the
947 strongest and stablest performer, however, we em-
948 phasize that in general there does not seem to a
949 dramatic difference between the rounding preci-
950 sions we considered. Longer digits slightly worsen
951 model performance and the model might learn the
952 most from simpler or abbreviated representations
953 of distance. This finding motivated us to consider
954 the ablation in Appendix A.4.3.

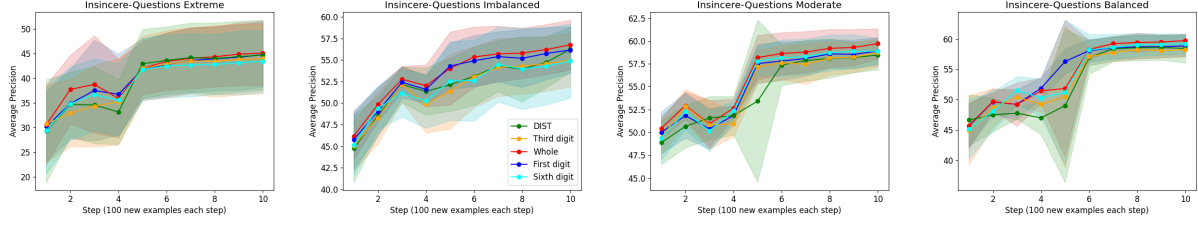


Figure 25: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the Insincere Questions dataset and the relevant balance is in the title of each panel.

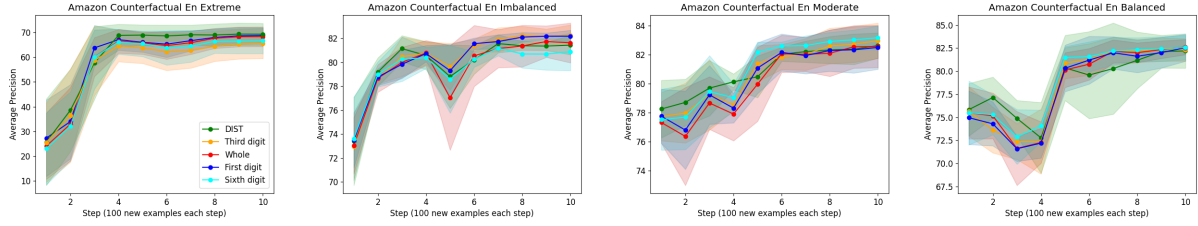


Figure 26: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the Amazon Counterfactual dataset and the relevant balance is in the title of each panel.

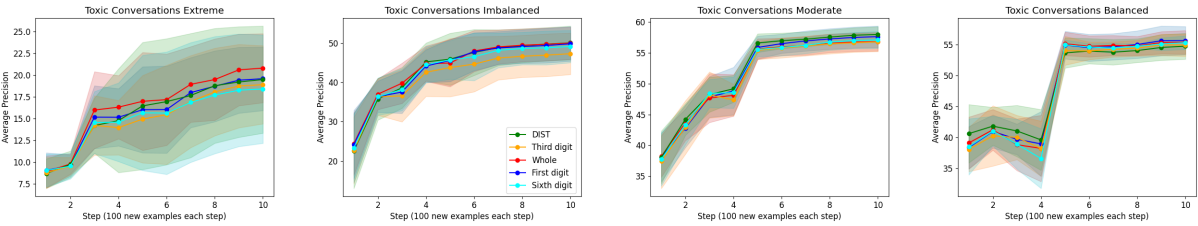


Figure 27: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the Toxic Conversations dataset and the relevant balance is in the title of each panel.

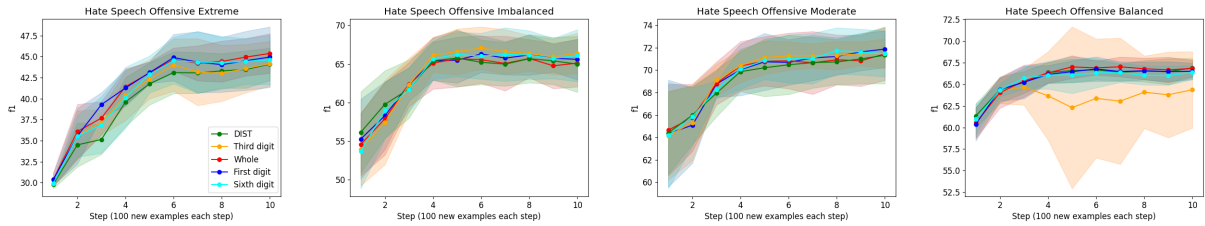


Figure 28: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the Hate Speech Offensive dataset and the relevant balance is in the title of each panel.

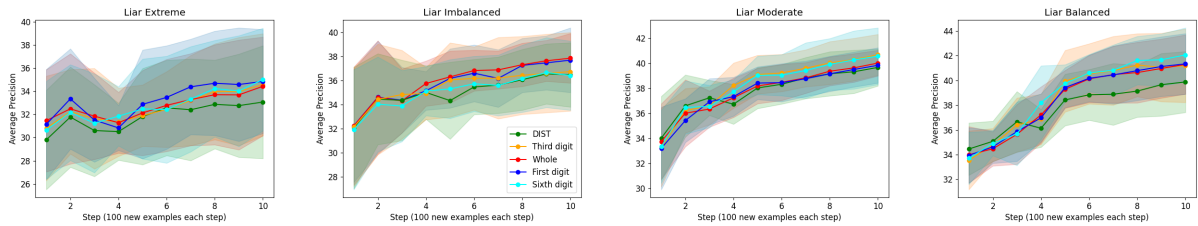


Figure 29: LAGONN_{lite} performance when considering different rounding precisions for the Euclidean distance before appending it to a modified instance. We consider all balance regimes on the Liar dataset and the relevant balance is in the title of each panel.

955 **A.4.3 Ablation: support for LABDIST**

956 The results from the ablation in Appendix A.4.2
957 suggest that rounding the distance to the nearest
958 whole number results in a stronger classifier than
959 appending the unrounded distance. Thus far, we
960 have asserted that LABDIST, where we append
961 both the gold label of the NN and unrounded dis-
962 tance is the most performant version of LAGONN
963 (see Table 1). To demonstrate that this is reason-
964 able, in this ablation study, we compare the orig-
965 inal LABDIST configuration against three mod-
966 els, namely the LABEL configuration, distance
967 rounded to near whole number (Whole), and finally
968 a new configuration similar to LABDIST, but where
969 we append the gold label and distance rounded to a
970 whole number, which we refer to as LABROUND.
971 As in Appendix A.4.2, in this ablation we consider
972 only the LAGONN_{lite} fine-tuning strategy. We
973 chose for this ablation because it provides insight
974 into both how the different configurations affect
975 full-model fine-tuning and only re-fitting the clas-
976 sification head. The results can be seen below in
977 Figures 30 through 34. We place the figures on a
978 new page for ease of viewing.

979 In general, we note very similar performance
980 curves for these four models. In the case of Insin-
981 cere Questions, appending the distance after round-
982 ing it to the nearest whole number (Whole, the red
983 curve), is a strong model, except in the balanced
984 regime where we note large instability. The results
985 for Amazon Counterfactual tell a different story,
986 where rounding the Euclidean distance to the near-
987 est whole number causes large instability and even
988 degrades performance on the fifth step.

989 For the other evaluation scenarios, it is unclear
990 what is the strongest method as sometimes LAB-
991 DIST is the best performer and sometimes it is
992 Whole (the red curve). However, we believe that in
993 general LABDIST is the most stable model while
994 also often being the most performant. We therefore
995 choose it as our default LAGONN configuration as
996 a compromise between strength and stability. It is
997 about this configuration which we report results in
998 the main text. Our interpretation of this is that pass-
999 ing the model both a discrete prediction (the gold
1000 label of the NN) and a truly continuous measure
1001 of similarity (the unrounded Euclidean distance)
1002 gives it the most consistent and dependable reason-
1003 ing ability.

1004 We note, as we did in Appendix A.4.1, that we
1005 could have presented the best performer for each

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

1006
1007
1008
1009
1010

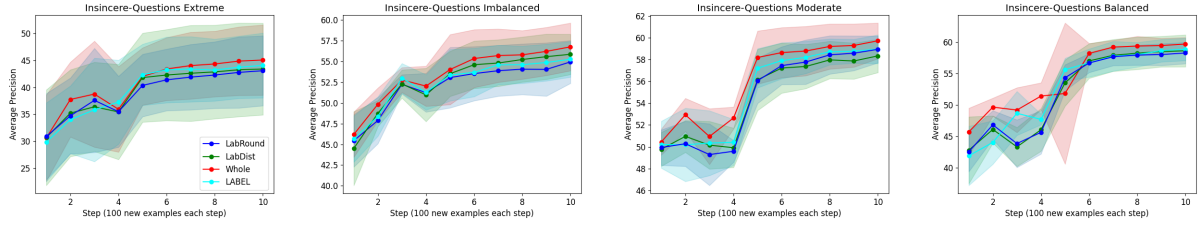


Figure 30: LAGoNN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Insincere Questions dataset and the relevant balance is in the title of each panel.

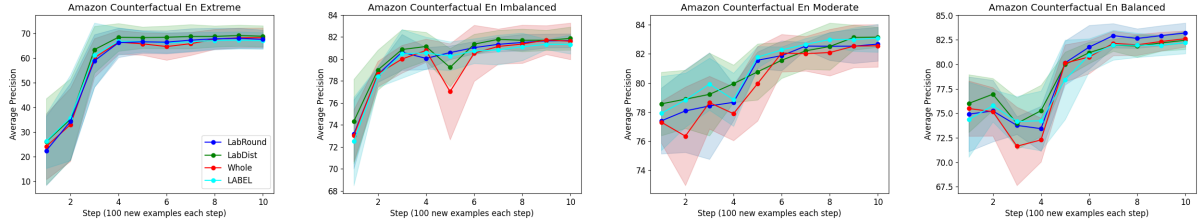


Figure 31: LAGoNN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Amazon Counterfactual dataset and the relevant balance is in the title of each panel.

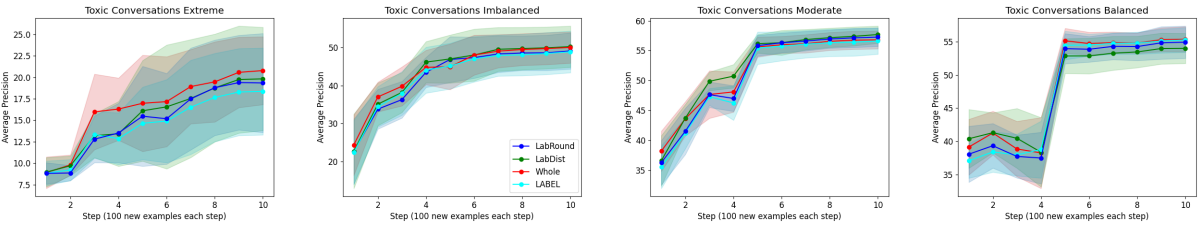


Figure 32: LAGoNN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Toxic Conversations dataset and the relevant balance is in the title of each panel.

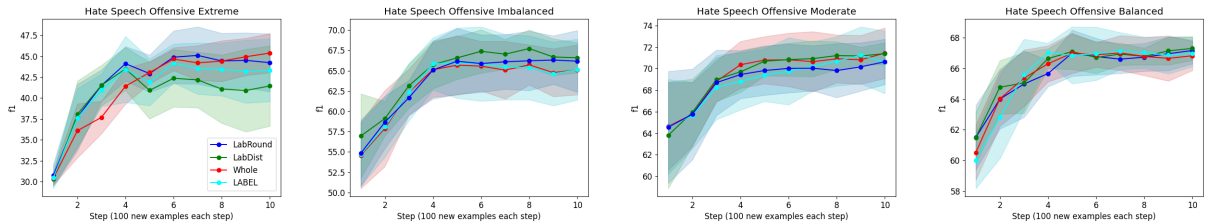


Figure 33: LAGoNN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Hate Speech Offensive dataset and the relevant balance is in the title of each panel.

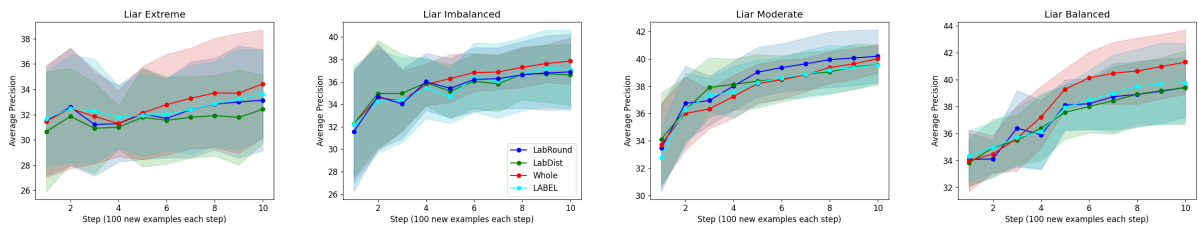


Figure 34: LAGoNN_{lite} performance where we compare the LABDIST against LABEL, LABROUND, and rounding the distance to the nearest whole number. We consider all balance regimes on the Liar dataset and the relevant balance is in the title of each panel.

1011 **A.5 Examples of LAGONN modified text**

1012 **WARNING:** Some of the examples below are of
1013 an offensive nature. Please view with caution.

1014 In this section, we provide examples of how
1015 LAGONN_{exp} modifies test text from the datasets
1016 we studied under the ALL configuration. We
1017 choose this configuration because the informa-
1018 tion it appends from a NN in the training data
1019 to a test instance encapsulates all configurations.
1020 LAGONN_{exp} was trained under a balanced dis-
1021 tribution and five examples per label were chosen
1022 randomly on the first, fifth, and tenth step to demon-
1023 strate how the same test instance might be deco-
1024 rated with different training examples as the train-
1025 ing data grow. We recognize that some the images
1026 below are difficult to see and have made the .csv
1027 files available with our code and data files. Note
1028 that MPNET’s separator token is </s>, not [SEP].

| Text Modified | Gold Label |
|---|--------------------|
| Cling to the wall, doesn't flop around when a bag is pulled out, the mess of bags falling out is gone. </> </not-counterfactual 3.161405993965967> And the dvd cases were tightly packed to ensure they didn't move around. </> </counterfactual 3.289605140686035> If I had to come up with anything negative, I would say that the attachments don't seem to stay on the vacuum cleaner when not in use - but that could be me not putting them on properly! | not-counterfactual |
| I like these jeans they sit low without being inappropriate when you sit or bend over. </> </not-counterfactual 2.447404623031616> These shorts fit really well and look good too. </> </counterfactual 2.550638198852539> The top fits great just the bottoms fit too. | not-counterfactual |
| He was very professional and with all transactions I make through Amazon were his good. </> </not-counterfactual 3.2391127681732178> The new speaker was just what the doctor ordered and I couldn't be more pleased. </> </counterfactual 3.3897111418685037> But the author alleviated my concerns quickly with a few well-timed comments about how it was the man could have known that the arrangement was something Jack wanted. | not-counterfactual |
| Well written with a twist I didn't expect. </> </not-counterfactual 2.557440620290205> "A bit workmanlike, not up to Lind's high standard of 'A Night to Remember'" but well-detailed, and a story that made me now know. </> </counterfactual 2.792485475540161> Wow I am really glad I didn't read these reviews BEFORE I read this book because I would have passed on the book and missed a really great start to a series that captured my attention and made me laugh all the while using my imagination and painting a clear picture of the author's world she was building for us. | not-counterfactual |
| Doesn't feel like the quality level I am used to. </> </counterfactual 2.5612920641296387> I was hoping the pants would be thicker but being that it's not too expensive it's understandable. </> </not-counterfactual 2.572395086288452> But I don't have a lining like the last couple models I bought. | not-counterfactual |
| I wish I had studied, I believe the enclosed hardware would have been sufficient. </> </not-counterfactual 2.663814544677344> It 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. </> </counterfactual 2.771895206451416> Wish it had a little more padding, otherwise just as advertised. | counterfactual |
| if this ever turns into a film, I hope they do it justice. </> </not-counterfactual 2.71574354171753> I read this book because of the motion picture that is coming out soon. </> </counterfactual 3.141676187515259> Wish this story would have been longer and turned into a book, with some gut wrenching action, love/hate lovers quarrels scenes, with a happy ending at the end. | counterfactual |
| If you don't want a prominent display this rack to large for most bed or living rooms, it is wider and taller than my tall Brooklyn wardrobe style dresser which was the largest piece in the room until this shoe rack. </> </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. </> </not-counterfactual 2.873617172241121> Good for under the bed shoe storage, if the wife wants to use it. | counterfactual |
| I wish I could have seen all of the places he recommends! </> </counterfactual 2.7999041080474854> I wish I had had him as an instructor at college. </> </not-counterfactual 3.2804622840881348> I wanted to order him a few affordable hats I wouldn't mind him losing | counterfactual |
| I wish I could replace just that small stupid piece, since there's nothing wrong with the rest of the hose assembly. </> </counterfactual 2.474032402038574> I wish I could just hook up a hose rather than connecting and routing an tube internally. </> </not-counterfactual 2.920058199157715> 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) | counterfactual |

Figure 40: Amazon Counterfactual, step 10.

| Text Modified | Gold Label |
|---|------------|
| Homos demand that you accept their fudge packing. But none of us ever will. </> </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? </> </toxic 4.095604249572734> I would imagine it is a miserable existence. Being a whining sore loser is no way to go through life. | not toxic |
| 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 1 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. </> </not toxic 3.234884929569582> Christ never said he would give grace, mercy, and acceptance to those who determinedly violate Scripture. In fact, he often spoke of hell. | not toxic |
| Also intellectually bankrupt... </> </not toxic 3.850763455344045> 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? </> </toxic 3.3897111418685037> Rather than call you a liar, I'll just post a link to HART's helicopter footage of this "most densely populated" ag land. | not toxic |
| https://www.youtube.com/watch?v=opizRWVWdWw | |
| 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. </> </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? </> </toxic 3.3897111418685037> Get the "real news" from Breitbart and Prager Report? Ahaahahaha! Are you for real? </> </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? </> </toxic 4.239322663353516> I would imagine it is a miserable existence. Being a whining sore loser is no way to go through life. | not toxic |
| Also intellectually bankrupt... </> </not toxic 3.850763455344045> 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? </> </toxic 3.3897111418685037> I wish I could have seen all of the places he recommends! </> </counterfactual 2.7999041080474854> I wish I had had him as an instructor at college. </> </not-counterfactual 3.2804622840881348> I wanted to order him a few affordable hats I wouldn't mind him losing | not toxic |
| Also intellectually bankrupt... </> </not toxic 3.850763455344045> 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? </> </toxic 3.3897111418685037> I wish I could have seen all of the places he recommends! </> </counterfactual 2.7999041080474854> I wish I had had him as an instructor at college. </> </not-counterfactual 3.2804622840881348> I wanted to order him a few affordable hats I wouldn't mind him losing | toxic |
| 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. </> </toxic 3.2624108791351313> Calling everybody that disagrees with you a racist and anti-Muslim is pretty uncivil. | toxic |
| Also intellectually bankrupt... </> </not toxic 3.8916428089141846> It's always important to remember what can happen when you have sociopaths as leaders and also have compliant followers. Some of the younger posters on this site might want to Google "Jim Jones and Jonestown"... There were no "checks and balances" in Jonestown, fear these are none in North Korea... and it can only hope in our country are firmly in place and functioning. Gary Crum | toxic |
| 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 price for her stupid stunt. </> </toxic 3.0406124591827393> Calling everybody that disagrees with you a racist and anti-Muslim is pretty uncivil. | toxic |
| Also intellectually bankrupt... </> </not toxic 3.094666481018066> Christ never said he would give grace, mercy, and acceptance to those who determinedly violate Scripture. In fact, he often spoke of hell. | toxic |
| no one cares what a paid lobbyist telling back like you believe lunatic. </> </not toxic 3.811786556243896> Calling everybody that disagrees with you a racist and anti-Muslim is pretty uncivil. | toxic |
| Also intellectually bankrupt... </> </not toxic 3.850763455344045> 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? </> </toxic 3.3897111418685037> Rather than call you a liar, I'll just post a link to HART's helicopter footage of this "most densely populated" ag land. | toxic |
| Ok all you GOP "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 dictatorship could ever be sanctioned. </> </toxic 3.361635120391846> Calling everybody that disagrees with you a racist and anti-Muslim is pretty uncivil. | toxic |
| Also intellectually bankrupt... </> </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 student's 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 some admin and union union). But the majority of teachers here instead of doing their jobs they band together wear purple shirts and mob the government for a better contract, and no accountability in the form of testing students. | toxic |
| And as for my post being "speculation" which part: that the liberals are the party in power, or that this involves money? </> </not toxic 3.042002267868684> As for me not knowing what is going on, you are correct. I am not a member of the Liberal party insider class, as you apparently are. | toxic |

Figure 41: Toxic Conversations, step 1.

| Text Modified | Gold Label |
|--|------------|
| Homos demand that you accept their fudge packing. But none of us ever will. </> </not toxic 3.138349330810547> So you admit you would exterminate inferior humans. </> </not toxic 3.2354952716827393> Mark Macdonald and the interests he work for would like us to 'get used to it, because they don't want to do anything practical to stop it. | not toxic |
| 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 1 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. </> </not toxic 3.234884929569582> Christ never said he would give grace, mercy, and acceptance to those who determinedly violate Scripture. In fact, he often spoke of hell. | not toxic |
| Also intellectually bankrupt... </> </not toxic 3.850763455344045> 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? </> </toxic 3.3897111418685037> Rather than call you a liar, I'll just post a link to HART's helicopter footage of this "most densely populated" ag land. | not toxic |
| https://www.youtube.com/watch?v=opizRWVWdWw | |
| 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. </> </not toxic 3.0351966421819665> So you are saying should build the road. Lots of down and groom but I'm missing your point. </> </toxic 3.146400279386226> Isn't that terrible? Caldwel and Inouye do not know the amount a 10-year extension would raise yet they are proposing it. Caldwell asked for two years and the get it. Two years later he is now asking for ten years, which is only weeks later. Must correct is coming in 10 years or so. This scenario shows our mayor and heart do not know the scope of this project. Pathetic. | not toxic |
| Get the "real news" from Breitbart and Prager Report? Ahaahahaha! Are you for real? </> </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? </> </toxic 4.239322663353516> I would imagine it is a miserable existence. Being a whining sore loser is no way to go through life. | not toxic |
| Well, just stupid, and it would be easy to give the behavior a pass if he wasn't stupid. </> </not toxic 3.2488777366203> With drama, I noticed that ADN had primed yet another column from Rob Knutson's column but after reading it I got the editors choice as the feature article on the Trump election. Knutson is also a psychiatrist so he analyzes of Trump and a man who had near emotionally, intellectual, developed beyond adolescent hobbies no doubt. But what does it say about Trump's supporters that he primed them to see through the boreson, conformational results of the man? How can so many Americans have devolved into such anger, fear and irrationality that they would/could find redemption in Trump after he has raped his own people's trust and honor? I'm sure she's awfully hysterical over the thought that such a pathetic dictatorship could ever be sanctioned. </> </toxic 3.361635120391846> Calling everybody that disagrees with you a racist and anti-Muslim is pretty uncivil. | not toxic |
| Wow! Lu to see one of those American Boats with all the planes on it at the bottom of the sea. </> </not toxic 3.390105489941377> I bet Regent Seven Seas will never offer Mr Hammond another trip. | not toxic |
| Wow, what a snappy article. He makes a strong, some wild points about food and atmosphere. However, after dipping the treats available on his "massive deck" the "blew off" his remaining restaurant reservations, don't his comely antibodies and ordered. | not toxic |
| He was certainly not an ideal passenger and, for one Reason on a Peak, he is a total legend. </> </not toxic 3.462364060406060> How respectful to the lady (a 140 driver) Using a preexisting military Operating Area Get over!!!!!! | not toxic |
| Wow, 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. </> </not toxic 3.78887806682111> I was informed that my friend (a PhD) had been contacted the evening of the day when I had inserted a profer note in the mailing list in Jerusalem that said God had my daughter conceive after a year of supposed sterility. Maybe Elizabeth did the same thing after all. She was in the same neighborhood. </> </not toxic 3.2888888888888888> Christ never said he would give grace, mercy, and acceptance to those who determinedly violate Scripture. In fact, he often spoke of hell. | not toxic |
| Be aware that of useless grumbling, and keep your tongue from slander, because no secret word is without result, and a lying mouth destroys the soul (4/8) (Romans 1:31) | toxic |
| If that is the case, then Trump's soul was utterly destroyed decades ago | toxic |
| 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. </> </toxic 3.203758887206473> That's what happens when you betray the people of your country for foreign ba. Let's go Le Pen, Geert Wilders. If the media refuses to mention the muslim crisis the total accountability, uneducated muslim makes swarming countries and turning them into misogynistic fundamentalist religious areas then we need these people to save us from YOU! </> </not toxic 3.2354952716827393> Your first mistake is believing what a politician says because generally, it has nothing to do with what he do. | toxic |
| The US will be happy to let this disease metastasize as well as very poor american gain our immigration dishonor. That is the advice of the elite prospects so significantly for the US! I have no doubts that PM Butts will ram through banned Ballot | toxic |
| 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 price for her stupid stunt. </> </toxic 3.0406124591827393> Calling everybody that disagrees with you a racist and anti-Muslim is pretty uncivil. | toxic |
| Also intellectually bankrupt... </> </not toxic 3.850763455344045> 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? </> </toxic 3.3897111418685037> Rather than call you a liar, I'll just post a link to HART's helicopter footage of this "most densely populated" ag land. | toxic |
| Ok all you GOP "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 dictatorship could ever be sanctioned. </> </toxic 3.361635120391846> Calling everybody that disagrees with you a racist and anti-Muslim is pretty uncivil. | toxic |
| Also intellectually bankrupt... </> </not toxic 3.850763455344045> 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? </> </toxic 3.3897111418685037> Rather than call you a liar, I'll just post a link to HART's helicopter footage of this "most densely populated" ag land. | toxic |
| You are completely ignorant of the shortage of nurses in this country. In some cases, critical shortages, and why would anyone want to be a nurse when they are disrespected by a former state and federal prosecutor such as you. </> </not toxic 3.02802267868684> As for me not knowing what is going on, you are correct. I am not a member of the Liberal party insider class, as you apparently are. | toxic |

Figure 42: Toxic Conversations, step 5.

