

# Self-Supervised Losses for One-Class Textual Anomaly Detection

Anonymous ACL submission

## Abstract

Current deep learning methods for anomaly detection in text rely on supervisory signals in inliers that may be unobtainable or bespoke architectures that are difficult to tune. We study a simpler alternative: fine-tuning Transformers on the inlier data with self-supervised objectives and using the losses as an anomaly score. Overall, the self-supervision approach outperforms other methods under various anomaly detection scenarios, improving the AUROC score on semantic anomalies by 11.6% and on syntactic anomalies by 22.8% on average. Additionally, the optimal objective and resultant learnt representation depend on the type of downstream anomaly. The separability of anomalies and inliers signals that a representation is more effective for detecting semantic anomalies, whilst the presence of narrow feature directions signals a representation that is effective for detecting syntactic anomalies.

## 1 Introduction

Anomaly detection is the task of identifying unusual samples relative to an exemplar inlier distribution. It has numerous applications in natural language processing (NLP), including fake news detection (Lee et al., 2021), spam detection (Crawford et al., 2015), and flagging atypical reviews (Ruff et al., 2019).

The difficulty of anomaly detection depends on the magnitude of difference between an anomalous representation and the distribution of inlier representations. Existing works in NLP focus on the far out-of-distribution (OOD) setting (Winkens et al., 2020) in which the anomalies are derived from a distinct dataset (Hendrycks et al., 2020; Arora et al., 2021; Li et al., 2021; Podolskiy et al., 2021; Zhou et al., 2021). For example, a model is trained on a sentiment classification dataset, and then that model is used to identify news articles as anomalies. These approaches also often assume the model is trained to classify the distinct inlier

sub-classes. The anomaly scoring mechanisms typically leverage these supervisory signals by fitting a Mahalanobis distance (Lee et al., 2018) to each sub-class or by obtaining the highest probability in the softmax layer (Hendrycks and Gimpel, 2017). However, these supervisory signals may not always be available.

As an alternative configuration, we analyse the one-class anomaly detection setting on more challenging near-OOD anomalies. One-class anomaly detection assumes only inlier data are available at training time and only have one label. Instead of supervisory signals, we study the performance of fine-tuning a Transformer on the inlier data using various self-supervised objectives, and we use the loss as the anomaly score. We examine anomaly detection performance on two near-OOD anomaly types: semantic anomalies, which are created by partitioning a single dataset by class label, and syntactic anomalies, which are created by randomly shuffling inlier sentences. We find that fine-tuning on a pre-trained Transformer outperforms existing and more complex methods, boosting AUROC score on semantic anomalies by 11.6% and on syntactic anomalies by 22.8% on average.

Our findings also suggest that the separation of anomalies and inlier classes in the learnt representation space of the detectors is a strong signal for detecting semantic anomalies, whilst adversarially brittle features are a better indicator of performance in the syntactic anomaly detection setting. Overall, our results indicate the fine-tuning paradigm is a simple baseline that can achieve good results, and the self-supervised objectives used for fine-tuning exploit different cues to identify anomalies.

## 2 Approach

### 2.1 Models

Using the loss of a fine-tuned Transformer for anomaly detection is analogous to using an autoen-

042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080

081 coder’s reconstruction error as an anomaly score in  
082 vision (Sakurada and Yairi, 2014). We anticipate  
083 that the fine-tuned models can learn the underlying  
084 characteristics of inlier data but not those of anom-  
085 alies. Hence, the loss is used as the anomaly score  
086 as it should be higher for anomalous instances.

087 We analyse three self-supervised objectives in  
088 our experiments. To minimise the influence of archi-  
089 tectural differences, we use the encoder from  
090 a pre-trained uncased BERT<sub>BASE</sub> (Devlin et al.,  
091 2019) and append different heads depending on  
092 the objective. We fine-tune each model for a max-  
093 imum of 30,000 steps on inlier data, employing  
094 early stopping based on the inlier validation set’s  
095 loss.

096 **Masked language modelling (MLM).** We re-  
097 tain the default configuration for BERT<sub>BASE</sub> and  
098 randomly mask 15% of tokens. At inference time,  
099 we mask the same proportion of tokens in the test  
100 sentences and use the error between the predicted  
101 and true tokens as the anomaly score.

102 **Causal language modelling (CLM).** We fine-  
103 tune the model to predict the next token given pre-  
104 vious tokens in the sequence and use perplexity  
105 as the anomaly score. Perplexity has been used  
106 to evaluate evidence-supported fact-checking (Lee  
107 et al., 2021) and far-OOD detection (Arora et al.,  
108 2021). Our work differs as it uses perplexity to eval-  
109 uate more difficult anomalies and does not require  
110 auxiliary data.

111 **Contrastive loss (SimCSE).** Previous works in  
112 vision suggest a contrastive loss can help discrim-  
113 inate anomalies from inliers (Tack et al., 2020;  
114 Schwag et al., 2021). However, these methods  
115 require data augmentations that are not directly  
116 transferrable to NLP.

117 SimCSE (Gao et al., 2021) resolves the data  
118 augmentation issue by applying different dropout  
119 masks to sentences and trains the model to select  
120 the same sentence from a minibatch of other sen-  
121 tence pairs. We fine-tune the model using the de-  
122 fault dropout probability ( $p = 0.1$ ) and tempera-  
123 ture ( $\tau = 0.05$ ) described in SimCSE and evalu-  
124 ate anomalies using the NT-Xent loss (Chen et al.,  
125 2020).

126 We compare the three fine-tuned models to four  
127 baselines:

128 **Pre-trained BERT (Pre-trained).** We evaluate  
129 MLM loss on BERT<sub>BASE</sub> without any fine-tuning.  
130 This configuration can be compared to MLM to  
131 examine the incremental benefit of fine-tuning. We

132 disregard the auxiliary next-sentence prediction ob-  
133 jective as we do not use sentence pairs for anomaly  
134 detection.

135 **Other attention-based anomaly detectors.** We  
136 compare our approach to two state-of-the-art meth-  
137 ods which use attention. CVDD (Ruff et al., 2019)  
138 learns a set of compact context vectors to describe  
139 the inlier data using a multi-head self-attention  
140 mechanism. It evaluates a sentence through the  
141 average cosine distance of the sentence’s contex-  
142 tual embedding to the context vectors.

143 DATE (Manolache et al., 2021) adapts ELEC-  
144 TRA (Clark et al., 2020) for the anomaly detection  
145 task. DATE includes an additional objective to  
146 predict which pre-defined pattern was used by the  
147 generator to mask the input tokens. At inference  
148 time, the input text is fed into the discriminator di-  
149 rectly. The average probability of each token being  
150 uncorrupted serves as the anomaly score.

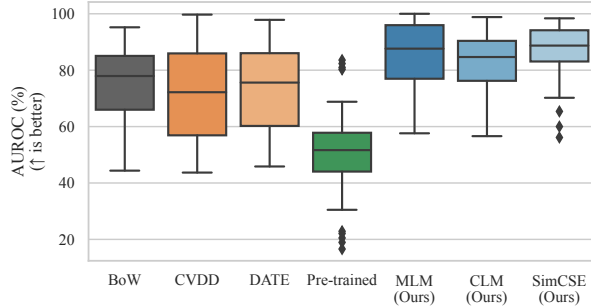
151 **Bag-of-words models (BoW).** We follow the ap-  
152 proach in CVDD and compute the mean over word  
153 embeddings extracted from FastText (Bojanowski  
154 et al., 2017) to create a sentence embedding for  
155 each datum. We use these sentence embeddings to  
156 train linear OC-SVMs, which worked better than  
157 using  $k$ -NNs or Mahalanobis distances in our ex-  
158 periments.

## 159 2.2 Datasets and anomaly detection setup

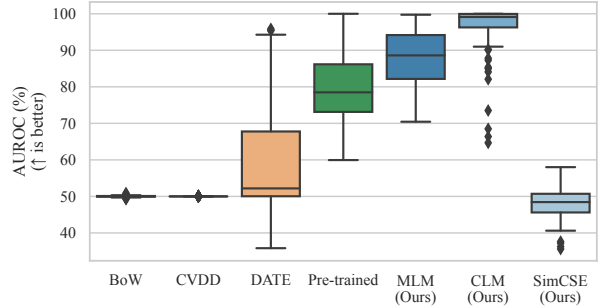
160 To allow comparison with the baseline methods,  
161 we evaluate anomaly detection performance on 20  
162 Newsgroups (Lang, 1995), Reuters-21578 (Lewis,  
163 1997), AG News (Zhang et al., 2015) and IMDb  
164 Movie Reviews (Maas et al., 2011). We also per-  
165 form experiments on Snopes (Vo and Lee, 2020) (a  
166 fact-checking dataset) and the Enron Spam Dataset  
167 (Metsis et al., 2006) to simulate more realistic  
168 anomaly detection applications. We pre-process  
169 each dataset by lowercasing text, stripping punc-  
170 tuation and removing stopwords as per Ruff et al.  
171 (2019).

172 We use the datasets’ class labels to construct  
173 two setups for the inlier training data. This allows  
174 us to examine anomaly detection performance in  
175 the settings where the inliers are narrow and more  
176 diverse. For a dataset with  $m$  class labels:

- 177 • Unimodal normality: We construct the inliers  
178 using data from a single label.
- 179 • Multimodal normality: We construct the in-  
180 liers using data from  $m - 1$  labels.



(a) Semantic anomaly results.



(b) Syntactic anomaly results encompassing all  $n$ -grams.

Figure 1: Anomaly detection results aggregated by model.

Class	Sentence					
Inlier	voip	gaining	ground	despite	cost	concerns
Anomaly	concerns	voip	despite	cost	ground	gaining

Table 1: Example of a syntactic anomaly derived from the AG News dataset. We look at  $n$ -grams ( $n \in \{1, 2, 3, 4\}$ ) and shuffle them until each  $n$ -gram is no longer in its original position.

We use the test splits of each dataset to formulate two types of near-OOD anomalies:

- Semantic anomalies: Data belonging to the same original class label(s) as the training data are categorised as inliers whilst the remainder are categorised as anomalies.
- Syntactic anomalies: Inlier and anomaly data are derived from the same class of data used to construct the training set. Inlier data are unchanged; anomalies have shuffled word order. To create the anomalies, we implement the seeded random function algorithm in [Sinha et al. \(2021\)](#). This setup allows us to measure the anomaly detectors’ sensitivity to the underlying syntactic information whilst fixing the word frequency statistics. We illustrate an example of a syntactic anomaly in Table 1.

### 3 Results

Figure 1 shows the overall anomaly detection results for both types of anomalies. The full results split by dataset and normality are in Appendix A.

**Fine-tuning a pre-trained Transformer boosts anomaly detection performance.** In the case of semantic anomalies, although the BoW performance suggests anomaly detection can be performed through analysing word frequency statis-

tics, fine-tuning helps to give additional information about the nature of inliers. This observation aligns with observations in vision ([Fort et al., 2021](#)). Our approach also outperforms CVDD and DATE, particularly in the multimodal normality setting.

Fine-tuning also improves syntactic anomaly detection, where frequency statistics are insufficient for discrimination. SimCSE is an exception, and we attribute this to the NT-Xent loss considering the entire sentence representation at inference.

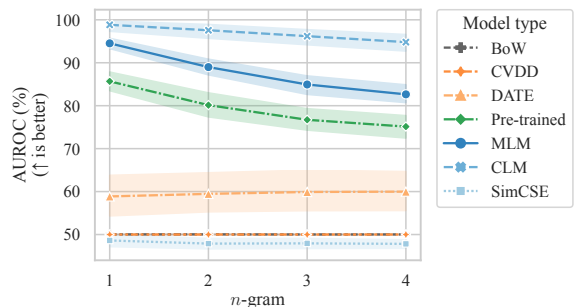


Figure 2: Mean AUROC across datasets on syntactic anomalies by  $n$ -gram level. Larger  $n$ -grams are more challenging to differentiate from inliers as fewer individual tokens are shuffled.

**Density models are much better at detecting syntactic anomalies.** We conducted an ablation study of performance under different permutation strengths. CLM is more stable under more challenging anomaly detection conditions (Figure 2), experiencing a decline of only 4% between 1-grams and 4-grams. Pre-trained and fine-tuned MLM experience similar drops (11%), which indicates the choice of objective for anomaly scoring is a core component for performance. As CLM calculates its score at the token level, it is more sensitive to syntactic changes compared to MLM, which considers spans of tokens through its masking mechanism.

In the following experiments, we extracted the embeddings at the last hidden BERT layer and mean-pooled over the positions to analyse the characteristics of the learnt embeddings.

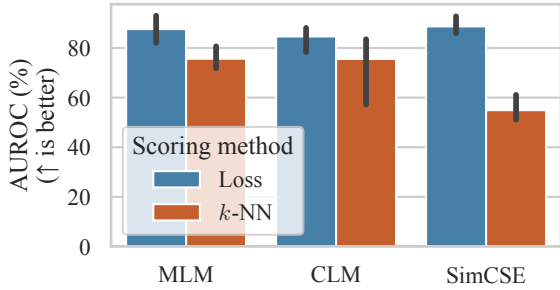


Figure 3: Comparison between using the loss as an anomaly score and  $k$ -NNs for semantic anomaly detection. We also experimented with OC-SVMs and Mahalanobis distances, but  $k$ -NNs performed best overall.

**Using the loss combined with the embedding is better than using the embeddings as a feature extractor.** Figure 3 shows the median semantic anomaly detection AUROC score when using the models end-to-end compared to extracting the embedding to train a  $k$ -NN. Although the raw embeddings are generally capable of performing anomaly detection, end-to-end use of the methods is more discriminative.

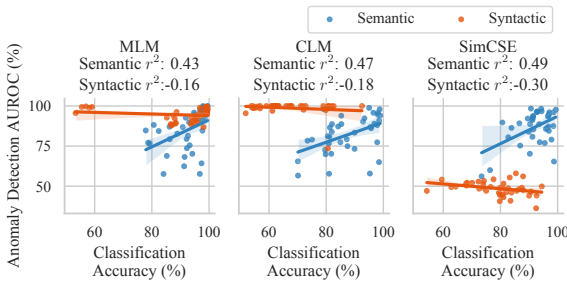


Figure 4: Scatter plot comparing classification accuracy of test inliers versus anomalies to anomaly detection performance across datasets.

**Separability of inliers and anomalies is a stronger signal for better semantic anomaly detection.** To examine the separability of embeddings for each learnt representation, we extracted both inlier and anomalous embeddings at the last hidden state and trained a logistic classifier. The correlation between classification accuracy and anomaly detection is more apparent for semantic anomalies (Figure 4), suggesting separability is a good indicator for better representations in this case, whereas there is no such relationship for syntactic anomalies.

This pattern suggests there is another factor that influences syntactic anomaly detection.

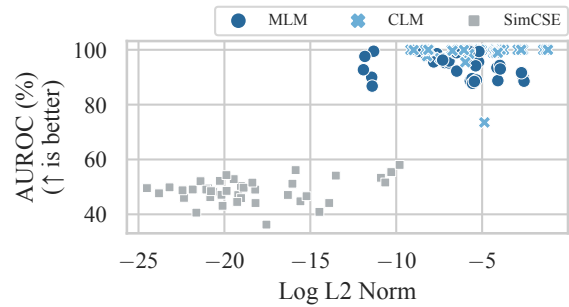


Figure 5: Scatter plot comparing average log L2 norms of the training inlier data to 1-gram syntactic anomaly detection performance. Higher norms are more brittle. The pattern is similar across different  $n$ -gram levels.

**Syntactic anomaly detection performance is more correlated to brittle features.** We hypothesise that a narrower<sup>1</sup> inlier representation is a better signal for syntactic anomaly performance as it provides more directions for anomalies to manifest.

We adopt the procedure in Mai et al. (2021) and calculate the average L2 gradient norms divided by the trace of the covariance matrix with respect to the training data. We observe similar behaviour across all datasets (summarised in Figure 5), whereby higher gradient norms clearly correspond to better anomaly detection performance.

Among the methods, CLM-based embeddings tend to be the most brittle and SimCSE the least. This corresponds with previous literature which states that autoregressive models like GPT (Radford et al., 2018) are highly anisotropic (Cai et al., 2021), and models such as SimCSE which are trained on contrastive objectives are more isotropic (Wang and Isola, 2020; Gao et al., 2021).

## 4 Conclusion

We studied the performance of fine-tuned Transformers using three self-supervised losses through a range of datasets and anomaly detection tasks. We show that this approach outperforms more complex methods, and employing the loss as an anomaly detector is better than using the learnt embeddings as a feature extractor. The best self-supervised loss depends on the nature of the anomalies, which suggests there is scope for analysing ensemble models or outlier exposure in future work.

<sup>1</sup>Narrow and brittle features refer to non-robust features as defined in adversarial machine learning literature (Ilyas et al., 2019).

287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
  
300  
  
301  
302  
303  
304  
305  
306  
  
307  
308  
309  
310  
  
311  
312  
313  
314  
  
315  
316  
317  
318  
319  
320  
321  
  
322  
323  
324  
325  
326  
  
327  
328  
329  
330  
331  
  
332  
333  
334  
335  
336  
337  
338

## Ethical considerations

Anomaly detectors are practical tools for indicating whether a system is working as intended and for flagging potential hazards (Hendrycks et al., 2021). An adversary may learn how to bypass systems by leveraging anomaly detection research. We restrict this by manually curating inliers and anomalies from publicly available datasets (as described in Section 2.2). By construction, our experiments are limited to the English language and may not represent features in other languages. We encourage extending our work to other domains and languages to investigate these differences.

## References

Udit Arora, William Huang, and He He. 2021. [Types of out-of-distribution texts and how to detect them](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10687–10701, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. [Enriching word vectors with subword information](#). *Transactions of the Association for Computational Linguistics*, 5:135–146.

Xingyu Cai, Jiayi Huang, Yuchen Bian, and Kenneth Church. 2021. [Isotropy in the contextual embedding space: Clusters and manifolds](#). In *International Conference on Learning Representations*.

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

Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. [Electra: Pre-training text encoders as discriminators rather than generators](#). In *International Conference on Learning Representations*.

Michael Crawford, Taghi M Khoshgoftaar, Joseph D Prusa, Aaron N Richter, and Hamzah Al Najada. 2015. Survey of review spam detection using machine learning techniques. *Journal of Big Data*, 2(1):1–24.

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

4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics. 339  
340

Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. 2021. [Exploring the limits of out-of-distribution detection](#). In *Advances in Neural Information Processing Systems*. 341  
342  
343  
344

Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. [SimCSE: Simple contrastive learning of sentence embeddings](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 345  
346  
347  
348  
349  
350  
351

Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. 2021. [Unsolved problems in ml safety](#). *arXiv preprint arXiv:2109.13916*. 352  
353  
354

Dan Hendrycks and Kevin Gimpel. 2017. [A baseline for detecting misclassified and out-of-distribution examples in neural networks](#). In *International Conference on Learning Representations*. 355  
356  
357  
358

Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song. 2020. [Pretrained transformers improve out-of-distribution robustness](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2744–2751, Online. Association for Computational Linguistics. 359  
360  
361  
362  
363  
364  
365

Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. 2019. [Adversarial examples are not bugs, they are features](#). In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc. 366  
367  
368  
369  
370

Ken Lang. 1995. [Newsweeder: Learning to filter news](#). In *Machine Learning, Proceedings of the Twelfth International Conference on Machine Learning, Tahoe City, California, USA, July 9-12, 1995*, pages 331–339. Morgan Kaufmann. 371  
372  
373  
374  
375

Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. 2018. [A simple unified framework for detecting out-of-distribution samples and adversarial attacks](#). In *Advances in Neural Information Processing Systems*, pages 7167–7177. 376  
377  
378  
379  
380

Nayeon Lee, Yejin Bang, Andrea Madotto, and Pascale Fung. 2021. [Towards few-shot fact-checking via perplexity](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1971–1981, Online. Association for Computational Linguistics. 381  
382  
383  
384  
385  
386  
387

David D. Lewis. 1997. Reuters-21578 text categorization test collection, distribution 1.0. 388  
389

Xiaoya Li, Jiwei Li, Xiaofei Sun, Chun Fan, Tianwei Zhang, Fei Wu, Yuxian Meng, and Jun Zhang. 2021. [kFolden: k-fold ensemble for out-of-distribution detection](#). In *Proceedings of the 2021 Conference on* 390  
391  
392  
393

394				
395				
396				
397				
398	Andrew L. Maas, Raymond E. Daly, Peter T. Pham,			
399	Dan Huang, Andrew Y. Ng, and Christopher Potts.			
400	2011. <a href="#">Learning word vectors for sentiment analysis</a> .			
401	In <i>Proceedings of the 49th Annual Meeting of the</i>			
402	<i>Association for Computational Linguistics: Human</i>			
403	<i>Language Technologies</i> , pages 142–150, Portland,			
404	Oregon, USA. Association for Computational Lin-			
405	guistics.			
406	Kimberly T. Mai, Toby Davies, and Lewis D. Grif-			
407	fin. 2021. Brittle features may help anomaly de-			
408	tection. <i>Women in Computer Vision Workshop at</i>			
409	<i>the IEEE/CVF Conference on Computer Vision and</i>			
410	<i>Pattern Recognition</i> .			
411	Andrei Manolache, Florin Brad, and Elena Burceanu.			
412	2021. <a href="#">DATE: Detecting anomalies in text via self-</a>			
413	<a href="#">supervision of transformers</a> . In <i>Proceedings of the</i>			
414	<i>2021 Conference of the North American Chapter of</i>			
415	<i>the Association for Computational Linguistics: Hu-</i>			
416	<i>man Language Technologies</i> , pages 267–277, Online.			
417	Association for Computational Linguistics.			
418	Vangelis Metsis, Ion Androutsopoulos, and Georgios			
419	Paliouras. 2006. <a href="#">Spam filtering with naive bayes -</a>			
420	<a href="#">which naive bayes?</a> In <i>CEAS 2006 - The Third Con-</i>			
421	<i>ference on Email and Anti-Spam, July 27-28, 2006,</i>			
422	<i>Mountain View, California, USA</i> .			
423	Alexander Podolskiy, Dmitry Lipin, Andrey Bout, Eka-			
424	terina Artemova, and Irina Piontkovskaya. 2021. <a href="#">Re-</a>			
425	<a href="#">visiting mahalanobis distance for transformer-based</a>			
426	<a href="#">out-of-domain detection</a> . In <i>Proceedings of the AAAI</i>			
427	<i>Conference on Artificial Intelligence</i> , volume 35,			
428	pages 13675–13682.			
429	Alec Radford, Karthik Narasimhan, Tim Salimans, and			
430	Ilya Sutskever. 2018. Improving language under-			
431	standing by generative pre-training.			
432	Lukas Ruff, Yury Zemlyanskiy, Robert Vandermeulen,			
433	Thomas Schnake, and Marius Kloft. 2019. <a href="#">Self-</a>			
434	<a href="#">attentive, multi-context one-class classification for</a>			
435	<a href="#">unsupervised anomaly detection on text</a> . In <i>Proceed-</i>			
436	<i>ings of the 57th Annual Meeting of the Association for</i>			
437	<i>Computational Linguistics</i> , pages 4061–4071, Flo-			
438	rence, Italy. Association for Computational Linguis-			
439	tics.			
440	Mayu Sakurada and Takehisa Yairi. 2014. Anomaly			
441	detection using autoencoders with nonlinear dimen-			
442	sionality reduction. In <i>Proceedings of the MLSDA</i>			
443	<i>2014 2nd workshop on machine learning for sensory</i>			
444	<i>data analysis</i> , pages 4–11.			
445	Vikash Sehwal, Mung Chiang, and Prateek Mittal. 2021.			
446	<a href="#">SSD: A unified framework for self-supervised out-</a>			
447	<a href="#">lier detection</a> . In <i>9th International Conference on</i>			
448	<i>Learning Representations</i> .			
	Koustuv Sinha, Robin Jia, Dieuwke Hupkes, Joelle			449
	Pineau, Adina Williams, and Douwe Kiela. 2021.			450
	<a href="#">Masked language modeling and the distributional hy-</a>			451
	<a href="#">pothesis: Order word matters pre-training for little</a> .			452
	In <i>Proceedings of the 2021 Conference on Empiri-</i>			453
	<i>cal Methods in Natural Language Processing</i> , pages			454
	2888–2913, Online and Punta Cana, Dominican Re-			455
	public. Association for Computational Linguistics.			456
	Jihoon Tack, Sangwoo Mo, Jongheon Jeong, and Jinwoo			457
	Shin. 2020. <a href="#">Csi: Novelty detection via contrastive</a>			458
	<a href="#">learning on distributionally shifted instances</a> . <i>Ad-</i>			459
	<i>vances in Neural Information Processing Systems</i> ,			460
	33:11839–11852.			461
	Nguyen Vo and Kyumin Lee. 2020. <a href="#">Where are the</a>			462
	<a href="#">facts? searching for fact-checked information to al-</a>			463
	<a href="#">leviate the spread of fake news</a> . In <i>Proceedings of</i>			464
	<i>the 2020 Conference on Empirical Methods in Natu-</i>			465
	<i>ral Language Processing</i> , pages 7717–7731, Online.			466
	Association for Computational Linguistics.			467
	Tongzhou Wang and Phillip Isola. 2020. <a href="#">Understanding</a>			468
	<a href="#">contrastive representation learning through alignment</a>			469
	<a href="#">and uniformity on the hypersphere</a> . In <i>Proceedings of</i>			470
	<i>the 37th International Conference on Machine</i>			471
	<i>Learning</i> , volume 119 of <i>Proceedings of Machine</i>			472
	<i>Learning Research</i> , pages 9929–9939. PMLR.			473
	Jim Winkens, Rudy Bunel, Abhijit Guha Roy, Robert			474
	Stanforth, Vivek Natarajan, Joseph R Ledsam, Patri-			475
	cia MacWilliams, Pushmeet Kohli, Alan Karthike-			476
	salingam, Simon Kohl, et al. 2020. Contrastive train-			477
	ing for improved out-of-distribution detection. <i>arXiv</i>			478
	<i>preprint arXiv:2007.05566</i> .			479
	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.			480
	<a href="#">Character-level convolutional networks for text clas-</a>			481
	<a href="#">sification</a> . In <i>Advances in Neural Information Pro-</i>			482
	<i>cessing Systems</i> , volume 28. Curran Associates, Inc.			483
	Wenxuan Zhou, Fangyu Liu, and Muhao Chen. 2021.			484
	<a href="#">Contrastive out-of-distribution detection for pre-</a>			485
	<a href="#">trained transformers</a> . In <i>Proceedings of the 2021</i>			486
	<i>Conference on Empirical Methods in Natural Lan-</i>			487
	<i>guage Processing</i> , pages 1100–1111, Online and			488
	Punta Cana, Dominican Republic. Association for			489
	Computational Linguistics.			490

A.1 Semantic anomaly detection results

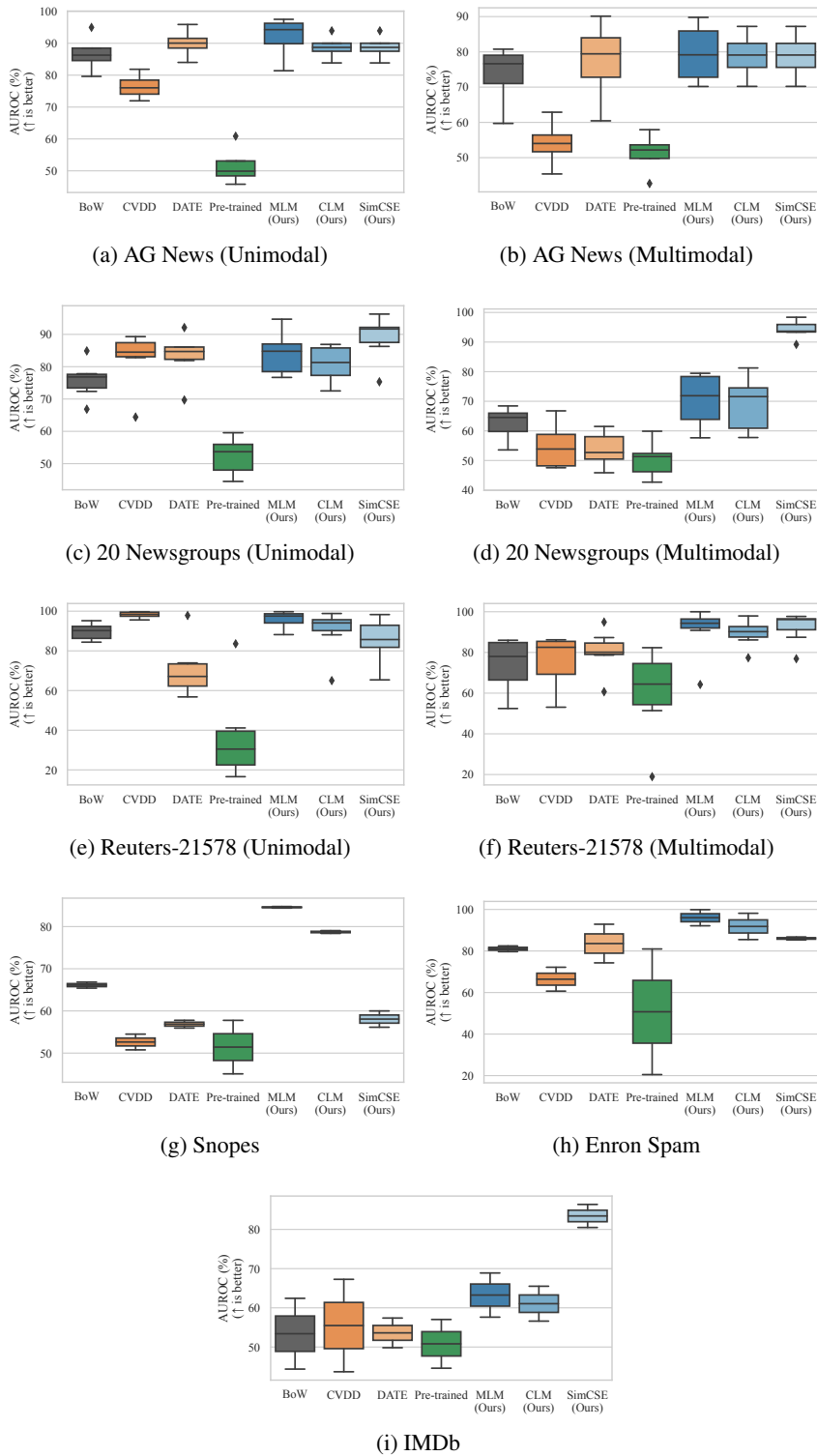


Figure 6: Semantic anomaly detection results split by dataset.

## A.2 Syntactic anomaly detection results

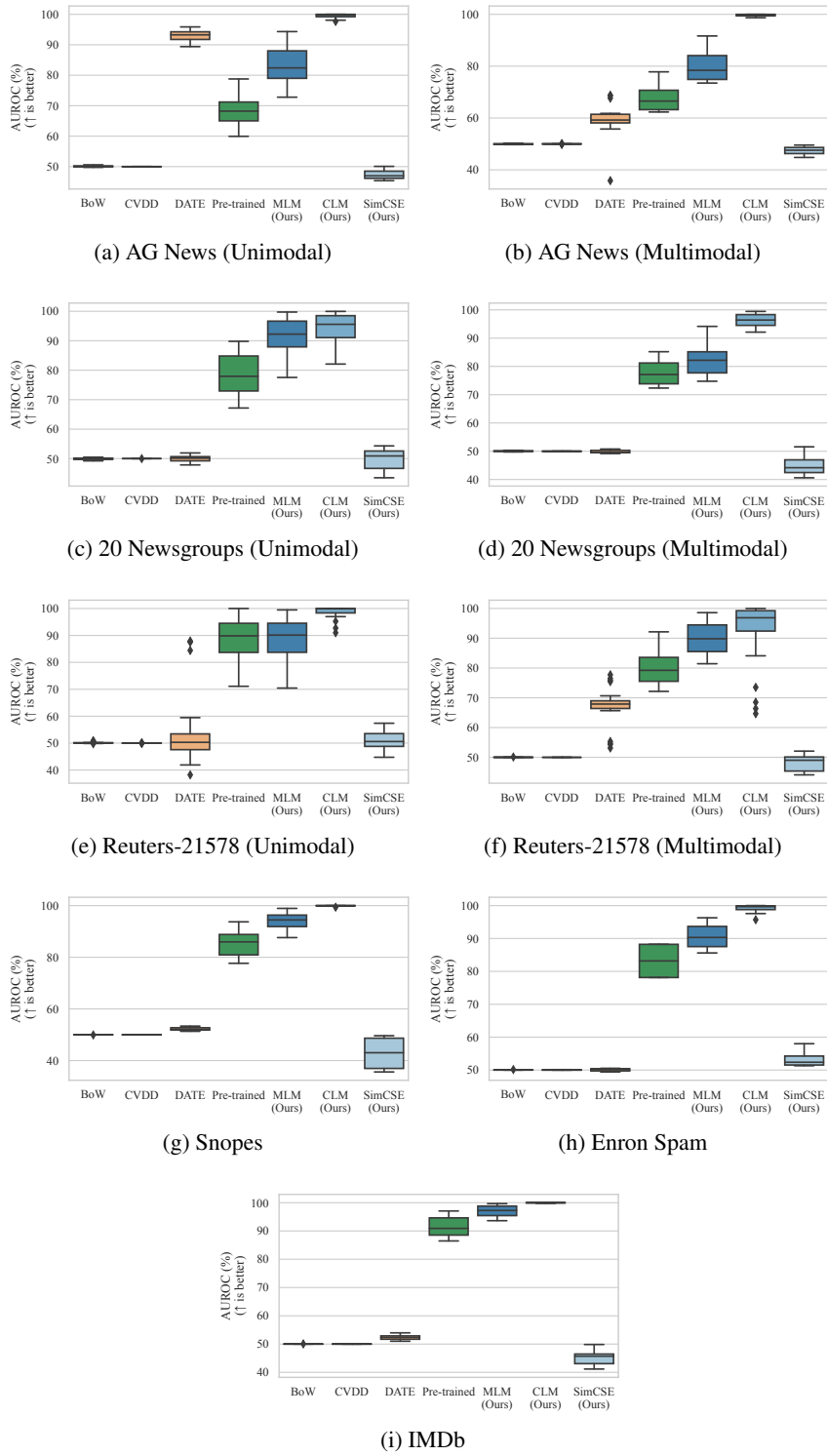


Figure 7: Syntactic anomaly detection results split by dataset. The figures include all  $n$ -gram runs.



### A.3 Contamination results

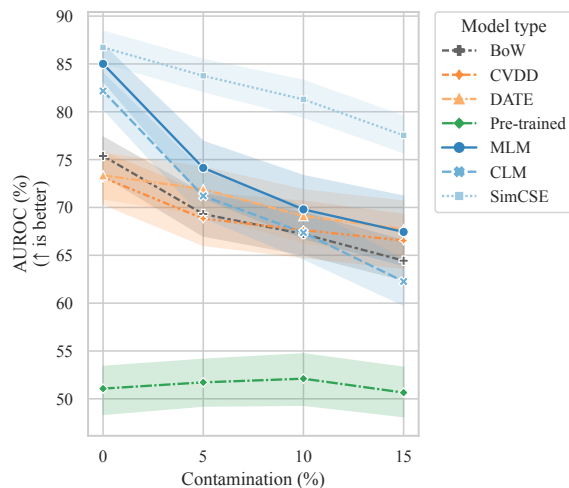


Figure 8: Mean AUROC scores across datasets by contamination percentage. Experiments are conducted using semantic anomalies.

We simulate a purely unsupervised anomaly detection setup by incorporating a set percentage of semantic anomalies  $\{5\%, 10\%, 15\%\}$  into the training data. The self-supervised losses on average elicit higher AUROC scores compared to the other model types, and SimCSE appears to be the most robust approach.

### A.4 Implementation details

We used an NVIDIA RTX Titan X and NVIDIA Tesla V100s to run our experiments depending on availability.

**Model implementation.** We used Huggingface’s<sup>2</sup> implementation of BERT<sub>BASE</sub> and Sentence-Transformers<sup>3</sup> for our Transformer experiments. In addition, we used nltk<sup>4</sup> for pre-processing, spaCy<sup>5</sup> for encoding the bag-of-words models, Faiss<sup>6</sup> to train the  $k$ -NNs, and sci-kit learn<sup>7</sup> for constructing OC-SVMs.

**Dataset details.** All of the datasets used in our paper are publicly available.

- 20 Newsgroups (Lang, 1995) is a collection of 20,000 newsgroup documents split across

<sup>2</sup><https://huggingface.co>

<sup>3</sup><https://sbnet.net>

<sup>4</sup><https://nltk.org>

<sup>5</sup><https://spacy.io>

<sup>6</sup><https://faiss.ai>

<sup>7</sup><https://scikit-learn.org>

20 different newsgroups. We use the six top-level subjects (*computer, recreation, science, miscellaneous, politics, religion*) to partition the classes. Partitioning by class label, there are 577-2859 training samples and 382-1909 test samples.

- Reuters-21578 (Lewis, 1997) is a collection of 10,788 news articles split across 90 topics. We only use a subset of data that have only one label (*earn, acq, crude, trade, money-fx, interest, ship*). Partitioning by class label, there are 108-2,840 training samples and 36-1,083 testing samples.
- AG News (Zhang et al., 2015) is a topic classification dataset gathered from more than 2,000 news sources over one year of activity. It contains four classes (*business, sci, sports, world*), each with 30,000 samples for training and 1,900 for testing.
- IMDb (Maas et al., 2011) is a sentiment classification dataset consisting of film reviews. It contains two classes (*pos, neg*), each with 25,000 samples for training and 25,000 for testing.
- Snopes (Vo and Lee, 2020) is a fact-checking dataset containing paired examples of tweets and a fact-checking article from *snopes.com*. There are four classes (*true, mostly true, mostly false, false*). We only use *true* (7,363) and *false* (21,256) tweets in our experiments and do not use the articles. We randomly partition 80% of this smaller dataset for training and use the remaining 20% for testing.
- The Enron Spam Dataset (Metsis et al., 2006) is derived from the Enron Email Dataset. There are two classes, *ham* (16,458) and *spam* (17,171) emails. We randomly partition 80% of the dataset for training and the remaining 20% for testing.