TAD-Bench: A Comprehensive Benchmark for Embedding-Based Text Anomaly Detection

Anonymous ACL submission

Abstract

Text anomaly detection is crucial for identifying spam, misinformation, and offensive language in natural language processing tasks. Despite the growing adoption of embeddingbased methods, their effectiveness and generalizability across diverse application scenarios remain insufficiently explored. To address this, we present TAD-Bench, a comprehensive benchmark designed to systematically evaluate embedding-based approaches for text anomaly detection. TAD-Bench integrates mul-011 tiple datasets spanning different domains, com-012 bining state-of-the-art embeddings from large language models with a variety of anomaly detection algorithms. Through extensive experiments, we analyze the interplay between 017 embeddings and detection methods, uncovering their strengths, weaknesses, and applicability to different tasks. These findings offer 019 new perspectives on building more robust, efficient, and generalizable anomaly detection systems for real-world applications. All the code are available at https://anonymous.4open. science/r/TAD-Bench-B4C6/.

1 Introduction

037

041

Anomaly detection (AD) is a critical task in machine learning, widely applied in fraud detection and content moderation to user behavior analysis (Pang et al., 2021). Within natural language processing (NLP), anomaly detection has become increasingly relevant for identifying outliers such as harmful content, phishing attempts, and spam reviews. However, while AD tasks in structured data (e.g., tabular, time series, graphs) (Steinbuss and Böhm, 2021; Blázquez-García et al., 2021; Qiao et al., 2024) have been extensively studied, anomaly detection in the unstructured and highdimensional domain of text remains underexplored. The inherent complexity of textual data, driven by its diverse syntactic, semantic, and pragmatic structures, presents significant challenges for robust and

reliable anomaly detection.

The rise of deep learning and transformer-based models has revolutionized NLP, enabling the development of contextualized embeddings that encode rich semantic and syntactic information. Techniques such as BERT (Devlin et al., 2019) and OpenAI's text-embedding models (OpenAI, 2024) have demonstrated remarkable success across a wide range of NLP tasks, offering dense, highdimensional representations that effectively capture linguistic nuances. These embeddings have become a cornerstone for many downstream tasks, providing powerful tools for applications such as text classification (da Costa et al., 2023) and retrieval (Zhu et al., 2023). Their ability to generalize across tasks and domains positions them as a promising foundation for complex challenges, including anomaly detection.

042

043

044

047

048

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

078

079

081

082

In recent years, embedding-based methods have gained significant attention in anomaly detection tasks due to their ability to capture semantic and contextual nuances in data (Wang et al., 2024). These methods typically involve two key stages: 1) extracting high-dimensional representations from textual data using pre-trained language models, which encode rich contextual and semantic features. 2) Applying specialized algorithms to identify anomalies based on these embeddings. The embeddings serve as a compact and expressive feature space, enabling downstream algorithms to efficiently identify deviations or outliers. Figure 1 shows the steps involved in embedding-based anomaly detection.

However, existing studies often lack systematic evaluations of how different embeddings perform across diverse anomaly types, raising questions about their generalization capabilities in complex, real-world scenarios such as multilingual settings or domain-specific anomalies. Recent efforts, such as AD-NLP (Bejan et al., 2023) and NLP-ADBench (Li et al., 2024), have significantly



Figure 1: Illustration of the embedding-based anomaly detection pipeline, encompassing embedding extraction and anomaly scoring.

advanced anomaly detection in NLP. AD-NLP provides valuable insights into different types of anomalies, while NLP-ADBench expands evaluations to a wide range of algorithms and datasets. However, AD-NLP evaluates few detection algorithms while NLP-ADBench considers only a few embedding methods, respectively. Our work aims to move beyond simply filling these gaps, by systematically exploring the following questions:

086

097

101

102

104

105

106

107

109

110

111

112

113

114

115

116

117

118

- What types of tasks are LLMs (Large Language Models) embeddings paired with anomaly detectors most suitable for, and where do they face limitations?
 - Which embedding methods consistently excel across different anomaly detection tasks?
 - Which anomaly detection algorithms perform robustly across various embeddings and tasks?

In this work, we introduce TAD-Bench, a novel benchmark specifically designed for text anomaly detection. Our objective is to enable a more comprehensive and systematic evaluation of state-ofthe-art embeddings, anomaly detection techniques, and their various combinations, offering valuable insights for a broad spectrum of NLP applications. By incorporating a diverse range of embedding models and rigorously evaluating an extensive suite of anomaly detection methods, TAD-Bench facilitates an in-depth understanding of their effectiveness on static datasets, with a strong emphasis on robustness, adaptability, and real-world applicability. The main contributions of this work are summarized as follows:

> • We propose TAD-Bench, a benchmark integrating diverse datasets for text anomaly detection across domains such as spam, fake news, and offensive language.

• We conduct a systematic evaluation of LLMbased embeddings and anomaly detection algorithms, revealing their relative strengths and weaknesses. 119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

148

149

150

151

152

153

154

• We provide insights into effective embeddingdetector configurations for improving robustness and generalizability in NLP anomaly detection tasks.

2 **Problem Definitions**

In the context of NLP, an anomaly refers to a text instance that exhibits characteristics or patterns that deviate significantly from the majority of the dataset. Such anomalies can manifest in various ways, including rare or niche topics, unusual or complex syntax and semantics, domain-specific jargon, or even intentionally manipulated language, such as spam, fake news, deceptive reviews, or offensive and harmful content. Detecting these anomalies is crucial for numerous real-world applications, such as content moderation, fraud detection, cybersecurity threat analysis, and identifying novel or emerging patterns in large-scale text corpora to enhance decision-making and knowledge discovery.

Formally, let $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ represent a corpus consisting of N textual instances, where each instance $x_i \in \mathcal{X}$ is represented as a sequence of tokens:

$$x_i = [t_1, t_2, \dots, t_{L_i}], \tag{14}$$

where L_i denotes the sequence length of x_i . The goal of text anomaly detection is to identify a subset of instances $\mathcal{D}_{anomaly} \subset \mathcal{D}$, such that $\mathcal{D}_{anomaly}$ contains samples that deviate significantly from the majority of the dataset $\mathcal{D}_{normal} = \mathcal{D} \setminus \mathcal{D}_{anomaly}$.

To achieve this, an anomaly detection algorithm $g: \mathbb{R}^d \to \mathbb{R}$ is applied to the representations of

the textual instances to identify potential anomalies. (1) Each text instance x_i is first mapped to a fixed-dimensional vector $\mathbf{z}_i \in \mathbb{R}^d$ using an embedding model $\phi : \mathcal{X} \to \mathbb{R}^d$, such that $\mathbf{z}_i = \phi(x_i)$. (2) The anomaly detection algorithm then assigns an anomaly score $s_i = g(\mathbf{z}_i)$ to each instance, $s_i \in [0, 1]$. Based on a predefined threshold τ , an instance x_i is classified as anomalous if:

$$x_i \in \mathcal{D}_{\text{anomaly}} \iff s_i \ge \tau.$$

The objective of text anomaly detection is to ensure that g effectively distinguishes between normal and anomalous instances, even in the absence of labeled data, while being robust to the inherent variability and high dimensionality of textual data.

3 Related Work

155

156

157

158

160

161

162

165

166

167

169

170

171

172

174

175 176

177

178

179

180

181

182

183

186

187

188

190

192

193

195

196

198 199

203

3.1 Text representations

Advancements in text representation extraction techniques have been instrumental in driving significant progress in the field of natural language processing. Early methods like TF-IDF (Term Frequency-Inverse Document Frequency) (Salton and Buckley, 1988) represented text in sparse vector spaces by measuring word importance relative to a corpus. While interpretable and computationally efficient, TF-IDF could not capture semantic relationships between words. Later, dense embeddings such as Word2Vec (Word to Vector) (Mikolov, 2013) and GloVe (Global Vectors for Word Representation) (Pennington et al., 2014) addressed this limitation by mapping words into continuous vector spaces based on their co-occurrence patterns in large corpora. However, these embeddings were static, assigning the same vector to a word regardless of its context.

To overcome the limitations of static embeddings, contextualized embeddings were introduced, with models like ELMo (Embeddings from Language Models) (Peters et al., 2018) producing word representations that vary based on context. This innovation was further advanced by transformer-based models like BERT (Bidirectional Encoder Representations from Transformers) (Devlin et al., 2019), which used bidirectional attention mechanisms to simultaneously capture left and right context. BERT set new benchmarks in NLP and inspired numerous improvements, including RoBERTa (Robustly Optimized BERT Approach) (Zhuang et al., 2021) and ALBERT (A Lite BERT) (Lan et al., 2020). More recently, large language models such as GPT (Generative Pre-trained Transformer (Brown et al., 2020) have significantly advanced the capabilities of embedding methods. These models, trained on massive and diverse datasets, generate highly expressive embeddings that capture both deep semantic relationships and rich generative properties of text. LLMs have exhibited unprecedented performance across a broad spectrum of NLP tasks, solidifying their role as dominant tools for text representation in numerous applications, including anomaly detection, information retrieval, and text generation. 204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

245

246

247

248

249

250

251

252

253

254

3.2 Anomaly Detection

Existing anomaly detection methods can be broadly categorized into 6 classes: distance, density, isolation, statistical, projection and deep learning-based approaches. Each category offers distinct advantages and is suited for different types of data distributions and anomaly patterns.

Distance-based methods, such as kNN (k-Nearest Neighbors) (Ramaswamy et al., 2000), identify anomalies by measuring the distance of a given data point to its nearest neighbors. Points that are far from their neighbors are considered anomalous. These methods are intuitive and straightforward but suffer from the curse of dimensionality in high-dimensional spaces, where distances lose their discriminative power, reducing their effectiveness.

Density-based methods identify points with significantly lower local density compared to their surroundings as anomaly. LOF (Local Outlier Factor) (Breunig et al., 2000) measures the local density of a point relative to its neighbors. HBOS (Histogram-Based Outlier Score) (Goldstein and Dengel, 2012) estimates densities using histograms for individual features.

Isolation-based methods assume anomalies are rare and different, iForest (Isolation Forest) (Liu et al., 2008, 2012), detect anomalies by recursively partitioning the feature space where anomalies require fewer partitions than normal points. Improved techniques, such as iNNE (Isolation-based Nearest Neighbor Ensembles) (Bandaragoda et al., 2018), use hypersphere to partition data space and assigns larger hyperspheres to anomalies, improving robustness in detecting local anomalies.

Probabilistic and statistical methods identify anomalies based on deviations from the data distribution. These approaches assume that normal instances follow a certain statistical pattern,

336

337

338

339

340

341

342

344

345

Table 1: Dataset description. Nor. and Ano. stand for Normal and Anomaly.

Dataset	# Samples	# Nor.	# Ano.	% Ano.
Email Spam	3578	3432	146	4.0805
SMS Spam	4969	4825	144	2.8980
COVID-Fake	1173	1120	53	4.5183
LIAR2	2130	2068	62	2.9108
OLID	641	620	21	3.2761
Hate Speech	4287	4163	124	2.8925

and anomalies appear as outliers that do not conform to this pattern. ECOD (Empirical Cumulative Distribution Function-based Outlier Detection) (Li et al., 2022) uses cumulative distribution functions for efficient anomaly scoring, while COPOD (Copula-Based Outlier Detection) (Li et al., 2020) leverages copulas to model feature dependencies, handling multivariate data effectively. Projection-based methods, such as OCSVM (One-Class SVM) (Schölkopf et al., 2001), separate normal and anomalous data by learning a decision boundary in a high-dimensional feature space. While effective for complex distributions.

Deep learning-based methods train on normal data to learn representations, identifying anomalies as deviations. Approaches like Deep SVDD (Deep Support Vector Data Description) (Ruff et al., 2018) and LUNAR (Unifying Local Outlier Detection Methods via Graph Neural Networks) (Goodge et al., 2022) capture nonlinear patterns but require substantial data and computational resources.

4 Benchmark Settings

4.1 Datasets

255

256

257

258

261

262

263

264

269

270

271

274

275

276

277

278

279

290

291

294

The scarcity of dedicated datasets poses a challenge to the development and evaluation of effective anomaly detection methods in NLP. To address this gap, we curated and transformed 6 existing classification datasets from three common NLP domains: spam detection, fake news detection, and offensive language detection. By incorporating datasets from diverse domains, our benchmark facilitates a comprehensive evaluation of embeddingbased anomaly detection methods across various NLP tasks.

Anomalies, as defined in our problem, are inherently rare. However, due to the lack of dedicated datasets for text anomaly detection, we adapted classification datasets by designating specific classes as anomalies and down-sampling them to simulate realistic anomaly rates. For each dataset, the anomaly rate was set to approximately 3%, reflecting the typical rarity of anomalies in real-world scenarios.

While some studies treat anomaly detection as novelty detection-assuming only normal instances in training (e.g., NLP-ADBench (Li et al., 2024)). TAD-Bench removes this constraint and directly utilizes all available data for anomaly detection. Additionally, we retain the original text without extra pre-processing, as any token, word, or symbol may carry critical information indicative of an anomaly. This approach preserves linguistic, structural, and contextual features essential for detecting anomalies. Table 1 presents the statistics of the six pre-processed datasets used in this benchmark, including Email-Spam(Metsis et al., 2006), SMS-Spam(Almeida et al., 2011), COVID-Fake(Das et al., 2021), LIAR2(Xu and Kechadi, 2024), OLID(Zampieri et al., 2019a), and Hate-Speech(Davidson et al., 2017).

4.2 Embedding Models

Table 2 summarizes the embedding models employed in this paper. To extract highquality embeddings from the datasets, 8 embedding models were utilized. These include BERT (bert-base-uncased) (Devlin et al., 2019), MiniLM (all-MiniLM-L6-v2) (Wang et al., 2020), LLAMA (Llama-3.2-1B), stella (stella_en_400M_v5) (Zhang et al., 2024), and **Qwen** (*Qwen2.5-1.5B*) (Yang et al., 2024a; Team, 2024) from the HuggingFace platform, as well as OpenAI-provided models: O-ada (text-embeddingada-002), O-small (text-embedding-3-small), and O-large (text-embedding-3-large) (OpenAI, 2024). All these models are based on the Transformer architecture, which has become the standard for representation learning in NLP tasks. The OpenAI models (O-ada, O-small, O-large) are specifically designed for embedding generation, offering embeddings with varying levels of granularity. On the other hand, LLAMA and Qwen are primarily auto-regressive language models optimized for text generation. In this paper, we repurposed these models for embedding extraction by computing the attention-weighted mean of their last hidden states, ensuring that only valid tokens contribute to the final sentence embeddings.

Notably, LLAMA and Qwen were constrained to a maximum token length of 512 tokens, same as BERT, due to computational resource limitations. Other models, such as MiniLM, Stella, and the

OpenAI embeddings, utilized automatic truncation to process longer input sequences. This limitation 347 may restrict LLAMA and Qwen's ability to fully 348 leverage their extended context capabilities, particularly for datasets with longer text instances, such as LIAR2 and Hate-Speech. However, this unified 351 token length ensures a fair comparison of runtime 352 efficiency across models under consistent experimental conditions. It also highlights the trade-offs between computational cost and embedding quality, particularly when resource constraints are a factor in model deployment. 357

Table 2: Embedding Models Overview. M and B are for million and billion, respectively.

Models	Max Tokens	# Dimensions	# Parameters
BERT	512	768	110 M
MINILM	512	384	22.7 M
O-ada	8191	1536	-
O-small	8191	1536	-
O-large	8191	3072	-
LLAMA	4096	2048	1.24 B
stella	2048	1024	435 M
Qwen	8192	1536	1.54 B

4.3 Anomaly Detectors

The embeddings derived from these models were subsequently used as input features for anomaly detection algorithms. To identify anomalous instances, we employed 8 anomaly detection methods sourced from the *PyOD* library¹ (Zhao et al., 2019). These algorithms include **KNN**, **LOF**, **OCSVM**, **iForest**, **INNE**, **ECOD**, **HBOS** and **COPOD**. These algorithms were selected to capture diverse anomaly detection paradigms, ensuring robust detection across datasets with varying characteristics, structures, and distributions.

For reproducibility and consistency, the default hyperparameter settings provided in the respective algorithm implementations and original papers were adopted. This approach minimizes user bias and allows for a fair comparison of algorithm performance when applied to the embeddings generated by the different models. The combination of diverse embedding models and anomaly detection algorithms ensures a comprehensive evaluation of text anomaly detection in terms of both computational efficiency and detection effectiveness.

4.4 Evaluation Criteria and Trials

Performance was evaluated using the Area Under the Receiver Operating Characteristic Curve (AU-ROC), a widely adopted metric in anomaly detection tasks for measuring the trade-off between true positive and false positive rates. To ensure the reliability and robustness of the results, each experiment was repeated 5 times, and the average AUROC score was reported. 381

383

384

386

387

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

5 Experiments

5.1 Applicability of LLM Embeddings with Anomaly Detectors

Table 2 summarizes the performance of various anomaly detectors combined with LLM-derived embeddings across different datasets, while Figure 1 highlights their strong performance in specific tasks, particularly spam detection. In both the email spam and SMS spam detection tasks (Figure 2a and Figure 2d, many embedding-detector combinations achieve high AUC scores, with several exceeding 0.8. This strong performance can be attributed to the explicit nature of spam-related features, such as the presence of URLs, nonsensical text, or repetitive patterns. An example from the Email Spam dataset is shown below:

Subject: oxycconttin no script needeeed your place to ggo too for all ur prreexxxiscrlpt 10 n pi sx , paaaaain killerzxss noeoo presscippt http : // hyyydroccodeeeine vicccodinne / vic geeet reeeliefff noowee http : // offfmeebabyy

These features are effectively represented in the semantic spaces created by general-purpose embeddings, enabling anomaly detectors to distinguish spam messages from legitimate ones. Additionally, the relatively small variance in detection performance across embeddings suggests that spam detection primarily relies on surface-level linguistic patterns, which are effectively captured by the embeddings employed in this study.

For fake news detection, the results indicate a more mixed performance across datasets. On the Covid Fake News dataset (Figure 2b, multiple embedding-detector combinations achieve AUC scores close to or exceeding 0.8, suggesting that these methods are capable of identifying subtle stylistic and linguistic differences between fake and

373

¹PyOD: https://pyod.readthedocs.io/en/latest/ index.html

Table 3: Evaluation across 6 datasets in terms of AU-ROO
--

Embeddings	Detectors	Email-Spam	SMS-Spam	COVID-Fake	LIAR2	Hate-Speech	OLID	Average
	kNN	0.7625	0.4484	0.8467	0.6594	0.5033	0.5137	0.6223
	OCSVM	0.7362	0.6323	0.7867	0.6237	0.4866	0.4866	0.6254
	IForest	0.7152	0.6164	0.7701	0.6051	0.4925	0.4783	0.6129
BERT	LOF	0.6786	0.3230	0.8713	0.6717	0.4632	0.4970	0.5841
	ECOD	0.7309	0.6235	0.7722	0.6175	0.4889	0.4933	0.6211
	INNE	0.7732	0.6497	0.8012	0.6362	0.4850	0.4740	0.6366
	HBOS	0.7145	0.6251	0.7698	0.6190	0.4935	0.5002	0.5317
	COPOD	0.6454	0.5929	0.7714	0.6242	0.4971	0.5189	0.5214
	kNN	0.9414	0.3180	0.8413	0.7249	0.5804	0.5063	0.6520
	OCSVM	0.9626	0.5915	0.7843	0.6470	0.4062	0.4520	0.6406
	IForest	0.9078	0.5472	0.7455	0.5936	0.4697	0.4531	0.6195
MINILM	LOF	0.5587	0.5024	0.7433	0.6804	0.5078	0.5422	0.5891
	ECOD	0.9525	0.5934	0.7581	0.6532	0.3786	0.4208	0.6261
	INNE	0.9526	0.5737	0.8035	0.6601	0.4223	0.4824	0.6491
	HBOS	0.9478	0.6137	0.7441	0.6447	0.3888	0.4316	0.5387
	COPOD	0.9453	0.6317	0.7416	0.6695	0.3710	0.4037	0.5375
	kNN	0.8865	0.3212	0.9094	0.7921	0.6341	0.5243	0.6779
	OCSVM	0.9310	0.8221	0.8143	0.7169	0.4807	0.5048	0.7116
	IForest	0.8872	0.7376	0.7432	0.6421	0.4632	0.4891	0.6604
O-ada	LOF	0.3808	0.5033	0.7316	0.7541	0.4328	0.5376	0.5567
	ECOD	0.9380	0.8822	0.8150	0.7200	0.4610	0.4986	0.7191
	INNE	0.8507	0.8031	0.8533	0.7378	0.4820	0.5102	0.7062
	HBOS	0.9433	0.8813	0.8164	0.7186	0.4583	0.5098	0.6182
	COPOD	0.9502	0.8759	0.8153	0.7201	0.4513	0.4811	0.6134
	KNN	0.8921	0.2290	0.9400	0.7756	0.6416	0.5587	0.6728
	OCSVM	0.9475	0.5755	0.8932	0.7024	0.4577	0.5547	0.6885
0 1	lforest	0.9058	0.6177	0.8085	0.5973	0.5025	0.5580	0.6650
O-small	LOF	0.3863	0.5257	0.7809	0.7489	0.4139	0.5581	0.5690
	ECOD	0.9481	0.6301	0.8808	0.7022	0.4249	0.5295	0.0859
	INNE	0.8673	0.6080	0.9185	0.7198	0.4491	0.5382	0.0835
	COPOD	0.9322	0.0275	0.8719	0.7008	0.4243	0.3137	0.3840
		0.9005	0.3722	0.004	0.0974	0.4017	0.4903	0.5700
	OCSVM	0.8292	0.1098	0.9337	0.6621	0.0291	0.3497	0.6535
	IForest	0.9403	0.5050	0.8924	0.0021	0.4200	0.4971	0.0055
O -large	LOF	0.0777	0.3277	0.8233	0.7356	0.3833	0.5000	0.5559
0-laige	ECOD	0.9487	0.4712	0.8255	0.7550	0.3059	0.3107	0.5555
	INNE	0.8230	0.5970	0.9261	0.6876	0.4197	0.4207	0.6617
	HBOS	0.9538	0.6525	0.8849	0.6404	0.3835	0.4989	0.5734
	COPOD	0.9639	0.6798	0.8854	0.6536	0.3537	0.4980	0.5763
	kNN	0.8715	0.3655	0.8668	0.7229	0.4991	0.4081	0.6223
	OCSVM	0.9023	0.7379	0.8132	0.6892	0.4774	0.4057	0.6710
	IForest	0.8962	0.7275	0.7833	0.6860	0.4647	0.4082	0.6610
Llama	LOF	0.6056	0.4053	0.8673	0.7274	0.4376	0.3972	0.5734
	ECOD	0.8844	0.7573	0.7819	0.6989	0.4643	0.3998	0.6644
	INNE	0.9122	0.7065	0.8160	0.6935	0.4702	0.3917	0.6650
	HBOS	0.9017	0.7895	0.7758	0.7064	0.4580	0.3898	0.5745
	COPOD	0.9153	0.8163	0.7584	0.7291	0.4435	0.3526	0.5736
	kNN	0.8654	0.3212	0.9034	0.6884	0.4746	0.5016	0.6258
	OCSVM	0.8922	0.7165	0.8063	0.5103	0.3729	0.4439	0.6237
	IForest	0.8862	0.7377	0.7738	0.4999	0.3545	0.4325	0.6141
stella	LOF	0.3931	0.4733	0.7129	0.6549	0.4036	0.5285	0.5277
	ECOD	0.9075	0.7894	0.8115	0.5023	0.3421	0.4395	0.6321
	INNE	0.8271	0.6926	0.8366	0.5330	0.3325	0.4532	0.6125
	HBOS	0.9178	0.8017	0.8086	0.4952	0.3355	0.4252	0.5406
	COPOD	0.9300	0.8589	0.8167	0.4936	0.3018	0.3797	0.5401
	kNN	0.8618	0.2110	0.8438	0.6626	0.5163	0.4602	0.5926
	OCSVM	0.8804	0.6229	0.7868	0.6216	0.4916	0.4882	0.6486
0	lForest	0.8829	0.6195	0.7686	0.6155	0.4825	0.4869	0.6427
Qwen	LOF	0.6043	0.3600	0.8555	0.6894	0.4600	0.4518	0.5702
	ECOD	0.8678	0.6648	0.7680	0.6172	0.4852	0.4773	0.6467
	INNE	0.8839	0.5940	0.7833	0.6339	0.4902	0.4693	0.6424
	HBOS	0.8854	0.6877	0.7638	0.6170	0.4847	0.4685	0.5582
	COPOD	0.9044	0.7393	0.7463	0.6291	0.4794	0.4336	0.5617



Figure 2: Boxplot of AUCROC scores for anomaly detectors on different embeddings across six datasets.

real news. These differences may include devia-423 424 tions in tone, phrasing, or structural composition of the text. However, on the LIAR2 dataset (Figure 2e, 425 the AUC scores exhibit much greater variability 426 across different combinations of embeddings and 427 detectors. This variability likely stems from the 428 429 greater factual complexity of the LIAR2 dataset, where detecting anomalies may require external 430 knowledge or sophisticated reasoning that is not in-431 herently encoded within the embeddings. Despite 432 this variability, the relatively strong performance 433 434 on the Covid Fake News dataset underscores the potential of embedding-based approaches for fake 435 news detection, particularly when the anomalies 436 are stylistic or linguistic in nature. 437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

In contrast, the performance on hate speech and offensive language detection tasks (Figure 2c and Figure 2f) is consistently weaker, with AUC scores rarely exceeding 0.6 across embeddingdetector combinations. This suggests that the embeddings struggle to capture the nuanced and context-dependent features necessary for these tasks. For instance, hate speech often relies on implicit cues such as sarcasm, cultural references, or subtle forms of hostility, which may not be fully captured by standard embeddings. Similarly, offensive language detection, as observed in the OLID dataset, requires identifying fine-grained differences in tone, intent, and subjectivity, such as distinguishing between neutral, offensive, and sarcastic expressions. These distinctions often depend

on broader contextual information, such as the discourse or dialogue in which the language appears. 454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

For example, without additional context, such as the speaker's intent or the conversational background, the following statement from OLID dataset remains ambiguous whether this statement qualifies as hate speech:

@USER #metoo are all racist!

5.2 Comparative Effectiveness of Embeddings in Anomaly Detection

The results in Table 3 demonstrate the remarkable capabilities of the OpenAI family of embeddings (O-ada, O-small, and O-large), consistently outperforming other embeddings across a variety of anomaly detection tasks. Specifically, O-ada achieves the highest average AUC scores with the ECOD detector (0.8822) on the SMS Spam dataset and with kNN (0.7921) on the LIAR2 dataset. Similarly, O-small demonstrates outstanding performance, achieving the highest AUC scores with kNN on the Hate Speech (0.6416) and OLID (0.5587) datasets. Additionally, O-large secures top AUC scores with COPOD (0.9639) on the Email Spam dataset and with kNN (0.9537) on the COVID Fake News dataset.

In comparison, other embeddings, such as MINILM, exhibit strong performance in specific tasks but lack consistency across more complex datasets. For instance, MINILM achieves excep-

553

554

555

556

557

558

559

560

561

562

563

564

515

tional AUC scores of 0.9526 and 0.9626 on the 483 Email Spam datasets when paired with INNE and 484 OCSVM, respectively. However, its performance 485 declines significantly on datasets like OLID and 486 LIAR2, suggesting limitations in capturing deeper 487 contextual or domain-specific cues essential for 488 these tasks. Similarly, embeddings such as stella 489 and Qwen exhibit moderate performance, excelling 490 in a limited subset of tasks but failing to match 491 the versatility of OpenAI embeddings. Their in-492 consistent performance across datasets indicates 493 that while they may effectively capture certain lin-494 guistic patterns, they struggle with tasks requiring 495 a broader understanding of context, intent, or nu-496 anced semantics. 497

498

499

501

505

506

507

508

510

511

512

513

514

These observations suggest that OpenAI embeddings, deliver the most robust and consistent performance across a diverse set of tasks. Their ability to effectively capture both explicit textual features (e.g., in spam detection) and nuanced contextual variations (e.g., in Covid Fake News and OLID) highlights their versatility. This underscores their suitability for anomaly detection scenarios that demand both surface-level pattern recognition and deeper linguistic comprehension, making them well-equipped for handling a wide range of textbased anomalies.

5.3 Performance Across Anomaly Detectors



Figure 3: Average rank (lower the better) of 3 different OpenAI embeddings-based methods on AUCROC across 6 datasets.

To evaluate the robustness of anomaly detection algorithms across various embeddings and tasks, we analyze their average rankings using OpenAI embeddings (O-ada, O-small, and O-large) as representative examples (Figure 3). These embeddings were selected based on their strong and consistent performance across datasets, as demonstrated in Section 5.2. The rankings provide insight into which detection algorithms perform reliably regardless of the embedding or task.

Across all three embeddings, kNN and INNE consistently rank as the top-performing algorithms. This indicates their robustness and adaptability to the semantic structures of LLM-derived embeddings. kNN, in particular, excels due to its ability to effectively model local density variations in feature space, making it well-suited for both explicit-pattern tasks like spam detection and nuanced tasks like fake news and hate speech detection. INNE, with its efficiency and strong generalization capabilities, complements kNN as a reliable alternative in diverse anomaly detection scenarios.

ECOD also ranks highly, consistently appearing among the top three detectors across embeddings. Its lightweight design and ability to estimate density-based anomalies make it a strong candidate for scenarios where computational efficiency is critical. On the other hand, methods like LOF, COPOD, and iForest consistently rank lower, highlighting their limitations in high-dimensional and semantically complex embedding spaces. These methods struggle with noise, data sparsity, and the nuanced patterns encoded in LLM embeddings, which limits their effectiveness across diverse tasks. Overall, kNN, INNE, and ECOD perform well, while LOF, COPOD, and iForest struggle with highdimensional embeddings.

6 Conclusion

In this study, we present a comprehensive benchmark for embedding-based anomaly detection in NLP, systematically evaluating the interplay between LLM embeddings and classical anomaly detection algorithms across three diverse domains: spam detection, fake news detection, and offensive language detection. Our results reveal both the strengths and limitations of embedding-based anomaly detection methods, demonstrating their effectiveness in tasks with explicit and well-defined patterns while highlighting challenges in capturing implicit, context-dependent anomalies that require broader contextual cues. These findings emphasize the need for more adaptive embeddings and hybrid detection strategies that integrate external knowledge and contextual reasoning.

568

571

574

576

577

580

582

584

585

589

591

595

596

597

610

611

612

614

Limitations

TAD-Bench evaluates anomaly detection across three domains: spam detection, fake news detection, and offensive language detection. While these tasks provide diverse and relevant benchmarks, they do not fully capture the complexity of realworld applications. Strong performance in spam detection highlights the ability of LLM embeddings to capture explicit patterns, while mixed results in fake news detection and poor performance in offensive language detection reveal their limitations in modeling implicit, context-sensitive cues. Expanding to domains like medical, financial, or legal texts that involve unique challenges, and exploring datasets with more implicit anomalies, could better evaluate the adaptability and robustness of these methods.

In addition, our work uses pre-trained LLM embeddings and default hyperparameters for anomaly detectors, ensuring consistency but potentially underestimating their best-case performance. Finetuning LLMs on domain-specific data could improve embedding quality, while systematic hyperparameter optimization might unlock the full potential of anomaly detectors. Future research should explore these directions, leveraging techniques like AutoML to streamline both embedding fine-tuning and parameter tuning, thereby achieving more competitive performance.

Moreover, TAD-Bench focuses solely on embedding-based methods, excluding end-to-end approaches that directly process raw text. While embeddings offer modularity and efficiency, end-toend models like autoencoders or transformer-based methods may capture richer contextual information and handle more complex anomalies. Future work should incorporate end-to-end models and explore hybrid approaches that combine the strengths of both paradigms, providing a more comprehensive evaluation of anomaly detection methods in NLP.

Furthermore, due to computational resource constraints, we primarily evaluate small-scale LLMs rather than larger, more powerful models. While this allows for a fair comparison across different embedding methods, it may not fully reflect the capabilities of state-of-the-art LLMs in anomaly detection. Future studies could benefit from leveraging larger models with extended context windows and more sophisticated representations to assess their impact on anomaly detection performance.

Ethic Statement

This study adheres to ethical research practices and considerations in the development and evaluation of text anomaly detection methods.

Use of Potentially Offensive Language. Some examples in this paper may contain offensive, harmful, or misleading language. These examples are used purely for illustrative purposes to demonstrate the challenges of text anomaly detection in realworld scenarios. They do not reflect the opinions, beliefs, or endorsements of the authors.

Data Sources and Usage. All datasets used in this study are sourced from publicly available research datasets that have been previously used in NLP and anomaly detection research. Proper citations and references to the original datasets are provided in the paper. No private, proprietary, or personally identifiable information was used in this study.

Risks and Responsible Use. Because anomaly detection models can be misused for purposes such as censorship, surveillance, or unfair content moderation. We strongly emphasize that our benchmark is intended for research and academic purposes only and should be used responsibly with consideration of ethical and societal implications.

Use of AI Assistance We acknowledge the use of AI-based writing assistants for grammar refinement, spelling correction, and improving the clarity of our manuscript. However, all intellectual contributions, experimental designs, analyses, and conclusions in this paper are solely the work of the authors.

References

- Tiago A Almeida, José María G Hidalgo, and Akebo Yamakami. 2011. Contributions to the study of sms spam filtering: new collection and results. In *Proceedings of the 11th ACM symposium on Document engineering*, pages 259–262.
- Tharindu R Bandaragoda, Kai Ming Ting, David Albrecht, Fei Tony Liu, Ye Zhu, and Jonathan R Wells. 2018. Isolation-based anomaly detection using nearest-neighbor ensembles. *Computational Intelligence*, 34(4):968–998.
- Matei Bejan, Andrei Manolache, and Marius Popescu. 2023. Ad-nlp: A benchmark for anomaly detection in natural language processing. In *Proceedings of the* 2023 Conference on Empirical Methods in Natural Language Processing, pages 10766–10778.
- Ane Blázquez-García, Angel Conde, Usue Mori, and Jose A Lozano. 2021. A review on outlier/anomaly

616 617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

664

665

detection in time series data. *ACM computing surveys* (*CSUR*), 54(3):1–33. Markus M Breunig, Hans-Peter Kriegel, Raymond T

666

667

673

674

675

677

679

696

701

702

704

710

711

713

714

716

717

718

719

721

- Ng, and Jörg Sander. 2000. Lof: identifying densitybased local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference on Management* of data, pages 93–104.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Liliane Soares da Costa, Italo L Oliveira, and Renato Fileto. 2023. Text classification using embeddings: a survey. *Knowledge and Information Systems*, 65(7):2761–2803.
- Sourya Dipta Das, Ayan Basak, and Saikat Dutta. 2021. A heuristic-driven ensemble framework for covid-19 fake news detection. In *Combating Online Hostile Posts in Regional Languages during Emergency Situation: First International Workshop, CONSTRAINT* 2021, Collocated with AAAI 2021, Virtual Event, February 8, 2021, Revised Selected Papers 1, pages 164–176. Springer.
- Thomas Davidson, Dana Warmsley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In *Proceedings of the international AAAI conference on web and social media*, volume 11, pages 512–515.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Markus Goldstein and Andreas Dengel. 2012. Histogram-based outlier score (hbos): A fast unsupervised anomaly detection algorithm. *KI-2012: poster and demo track*, 1:59–63.
- Adam Goodge, Bryan Hooi, See-Kiong Ng, and Wee Siong Ng. 2022. Lunar: Unifying local outlier detection methods via graph neural networks. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 36, pages 6737–6745.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. Albert: A lite bert for self-supervised learning of language representations. *Preprint*, arXiv:1909.11942.
- Yuangang Li, Jiaqi Li, Zhuo Xiao, Tiankai Yang, Yi Nian, Xiyang Hu, and Yue Zhao. 2024. Nlpadbench: Nlp anomaly detection benchmark. *arXiv preprint arXiv:2412.04784*.

- Zheng Li, Yue Zhao, Nicola Botta, Cezar Ionescu, and Xiyang Hu. 2020. Copod: copula-based outlier detection. In 2020 IEEE international conference on data mining (ICDM), pages 1118–1123. IEEE.
- Zheng Li, Yue Zhao, Xiyang Hu, Nicola Botta, Cezar Ionescu, and George H Chen. 2022. Ecod: Unsupervised outlier detection using empirical cumulative distribution functions. *IEEE Transactions on Knowledge and Data Engineering*, 35(12):12181–12193.
- Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2008. Isolation forest. In 2008 eighth ieee international conference on data mining, pages 413–422. IEEE.
- Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2012. Isolation-based anomaly detection. ACM Transactions on Knowledge Discovery from Data (TKDD), 6(1):1–39.
- Vangelis Metsis, Ion Androutsopoulos, and Georgios Paliouras. 2006. Spam filtering with naive bayeswhich naive bayes? In *CEAS*, volume 17, pages 28–69. Mountain View, CA.
- Tomas Mikolov. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 3781.
- OpenAI. 2024. New embedding models and api updates.
- Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. 2021. Deep learning for anomaly detection: A review. ACM computing surveys (CSUR), 54(2):1–38.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Marco AF Pimentel, David A Clifton, Lei Clifton, and Lionel Tarassenko. 2014. A review of novelty detection. *Signal processing*, 99:215–249.
- Hezhe Qiao, Hanghang Tong, Bo An, Irwin King, Charu Aggarwal, and Guansong Pang. 2024. Deep graph anomaly detection: A survey and new perspectives. *arXiv preprint arXiv:2409.09957*.
- Sridhar Ramaswamy, Rajeev Rastogi, and Kyuseok Shim. 2000. Efficient algorithms for mining outliers from large data sets. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pages 427–438.

Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. 2018. Deep one-class classification. In *International conference on machine learning*, pages 4393–4402. PMLR.

778

786 787

792

793

794

801

803 804

805

810

811

812

813

815

818

819

820

821

822

824

826

827

829

833 834

- Gerard Salton and Christopher Buckley. 1988. Termweighting approaches in automatic text retrieval. *Information processing & management*, 24(5):513– 523.
- Bernhard Schölkopf, John C Platt, John Shawe-Taylor, Alex J Smola, and Robert C Williamson. 2001. Estimating the support of a high-dimensional distribution. *Neural computation*, 13(7):1443–1471.
- Georg Steinbuss and Klemens Böhm. 2021. Benchmarking unsupervised outlier detection with realistic synthetic data. ACM Transactions on Knowledge Discovery from Data (TKDD), 15(4):1–20.
- Qwen Team. 2024. Qwen2.5: A party of foundation models.
- Caihong Wang, Du Xu, and Zonghang Li. 2024. Log2graphs: An unsupervised framework for log anomaly detection with efficient feature extraction. *arXiv preprint arXiv:2409.11890*.
- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep selfattention distillation for task-agnostic compression of pre-trained transformers. Advances in Neural Information Processing Systems, 33:5776–5788.
- Cheng Xu and M-Tahar Kechadi. 2024. An enhanced fake news detection system with fuzzy deep learning. *IEEE Access*, 12:88006–88021.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024a. Qwen2 technical report. arXiv preprint arXiv:2407.10671.
- Tiankai Yang, Yi Nian, Shawn Li, Ruiyao Xu, Yuangang Li, Jiaqi Li, Zhuo Xiao, Xiyang Hu, Ryan Rossi, Kaize Ding, et al. 2024b. Ad-llm: Benchmarking large language models for anomaly detection. arXiv preprint arXiv:2412.11142.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. 2019a. Predicting the Type and Target of Offensive Posts in Social Media. In *Proceedings of NAACL*.

Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. 2019b. Predicting the type and target of offensive posts in social media. 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 1415–1420, Minneapolis, Minnesota. Association for Computational Linguistics. 835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

- Dun Zhang, Jiacheng Li, Ziyang Zeng, and Fulong Wang. 2024. Jasper and stella: distillation of sota embedding models. *arXiv preprint arXiv:2412.19048*.
- Yue Zhao, Zain Nasrullah, and Zheng Li. 2019. Pyod: A python toolbox for scalable outlier detection. *Journal of Machine Learning Research*, 20(96):1–7.
- Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Haonan Chen, Zheng Liu, Zhicheng Dou, and Ji-Rong Wen. 2023. Large language models for information retrieval: A survey. *arXiv preprint arXiv:2308.07107*.
- Liu Zhuang, Lin Wayne, Shi Ya, and Zhao Jun. 2021. A robustly optimized BERT pre-training approach with post-training. In *Proceedings of the 20th Chinese National Conference on Computational Linguistics*, pages 1218–1227, Huhhot, China. Chinese Information Processing Society of China.

933

934

935

936

937

905

906

A Clarification Between Anomaly and Novelty Detection

861

862

863

866

871

872

875

876

877

879

884

891

895

900

901

903

904

Text Anomaly Detection (TAD), as defined in Section 2, focuses on identifying instances that deviate significantly from the majority of a dataset, regardless of whether anomalies are present during training. While some prior studies (e.g., AD-NLP (Bejan et al., 2023), NLP-ADBench (Li et al., 2024) and AD-LLM (Yang et al., 2024b)) assume training data contains only normal instances and testing data includes both normal and anomalous samples, this setup aligns more closely with novelty detection (Pimentel et al., 2014). Novelty detection specifically targets never-before-seen anomalies that are absent from the training phase, often treating anomalies as entirely novel classes.

In contrast, our benchmark evaluates a broader spectrum of anomaly detection scenarios. We do not restrict the training data to purely normal instances, allowing for potential partial supervision or contaminated training sets (e.g., realistic scenarios where anomalies may unintentionally exist in training data). This setup reflects real-world applications where anomaly types are not always fully known a prior, and detection systems must generalize across domains and anomaly types.

This distinction underscores our goal of advancing generalizable anomaly detection systems for real-world NLP applications, where anomalies may exhibit both explicit and context-dependent patterns.

B Datasets

Email-Spam² (Metsis et al., 2006) contains 5,171 emails labeled as spam or ham, with spam treated as the anomaly class. We utilized the preprocessed version provided in (Li et al., 2024).

SMS-Spam³ (Almeida et al., 2011) comprises 5,574 SMS messages originally labeled as spam or ham. Spam messages are designated as the anomaly.

COVID-Fake⁴ (Das et al., 2021) comprises posts collected from social media platforms and fact-checking websites. Real news items were sourced from verified outlets providing accurate COVID-19 information, while fake news was gathered from tweets, posts, and articles containing misinformation about COVID-19. Fake news is treated as the anomaly class.

LIAR2 ⁵ (Xu and Kechadi, 2024) consists of approximately 23,000 statements manually labeled by professional fact-checkers for fake news detection tasks. The "True" class, representing accurate statements, is considered the normal class, while the "Pants on Fire" class, representing highly misleading statements, is treated as the anomaly.

OLID⁶ (Zampieri et al., 2019b) (Zampieri et al., 2019a) contains 14,200 annotated English tweets, categorized using a three-level annotation model. For this benchmark, only the Level A (Offensive Language Detection) annotations are used, where tweets labeled as offensive are considered as anomalies, and non-offensive tweets are considered as normal.

Hate-Speech⁷ (Davidson et al., 2017) contains tweets annotated by CrowdFlower users. The tweet content is used as data, with "hate speech" treated as anomalies.

C Embedding Models

To effectively represent textual data, we use various pre-trained embedding models that transform text into dense vector representations. These embeddings serve as feature inputs for anomaly detection models, enabling them to capture semantic similarities and deviations in text. We selected a diverse set of embedding models, balancing between model size, token length limits, and computational efficiency. The models used in this study are:

- BERT ⁸ (*bert-base-uncased*) 938
- MINILM ⁹ (all-MiniLM-L6-v2) 939
- O-ada ¹⁰ (*text-embedding-ada-002*) 940
- O-small ¹¹ (text-embedding-3-small) 941

```
<sup>5</sup>https://github.com/chengxuphd/liar2?tab=
readme-ov-file
<sup>6</sup>https://sites.google.com/site/
offensevalsharedtask/olid
<sup>7</sup>https://github.com/t-davidson/
hate-speech-and-offensive-language
<sup>8</sup>https://huggingface.co/google-bert/
bert-base-uncased
<sup>9</sup>https://huggingface.co/sentence-transformers/
all-MiniLM-L6-v2
<sup>10</sup>https://platform.openai.com/docs/guides/
embeddings/
<sup>11</sup>https://platform.openai.com/docs/guides/
embeddings/
```

²https://huggingface.co/datasets/kendx/ NLP-ADBench/tree/main/datasets/email_spam

³https://archive.ics.uci.edu/dataset/228/sms+ spam+collection

⁴https://github.com/diptamath/covid_fake_news? tab=readme-ov-file

Embeddings	Email-Spam	SMS-Spam	COVID-Fake	LIAR2	Hate-Speech	OLID
BERT	64.95 s	10.87 s	4.48 s	3.7 s	9.79 s	1.81 s
MINILM	3.58 s	1.48 s	0.64 s	0.52 s	1.50 s	0.45 s
O-ada	154.08 s	33.76 s	17.35 s	16.67 s	35.57 s	8.9 s
O-small	166.73 s	34.16	18.32	16.10	34.09 s	9.0 s
O-large	206.07 s	41.78 s	20.58 s	29.97 s	44.20 s	10.72 s
Llama	545.51 s	129.16 s	38.15 s	28.93 s	71.49 s	18.02 s
stella	99.75 s	19.38 s	12.41 s	9.95 s	20.18 s	5.72 s
Qwen	745.85 s	129.16 s	58.19 s	40.29 s	183.24 s	20.49 s

Table 4: Embedding time of 6 datasets in seconds.

• O-large ¹² (*text-embedding-3-large*)

• LLAMA ¹³ (*Llama-3.2-1B*)

942

943

944

947

951

955

957

960

961

962

963

964

965

966

967

968

969

970

• stella ¹⁴ (*stella_en_400M_v5*)

• Qwen ¹⁵ (*Qwen2.5-1.5B*)

Beyond model size and token limits, computational efficiency is a key factor in selecting embedding models, particularly for real-world applications where inference speed is critical. Table 4 presents the embedding time (in seconds) required to process six datasets using each embedding model.

From the Table 4, we observe a significant variation in embedding extraction time. MINILM is the fastest across all datasets, taking only a few seconds, making it ideal for applications requiring real-time embedding generation. BERT offers a moderate trade-off, with embedding times significantly lower than larger models but higher than MINILM. OpenAI's embeddings (O-ada, O-small, O-large) are relatively slow, likely due to their highdimensional output and extended token support. Llama and Qwen models require the most computation, with Qwen taking up to 745.85 seconds on the Email-Spam dataset, reflecting the high computational cost of large autoregressive models.

D Comparative Analysis of Anomaly Detection Algorithms

Anomaly detection algorithms vary in their underlying assumptions, computational efficiency, and effectiveness across different types of data distributions. In this section, we provide a comparative analysis of the eight anomaly detection methods used in this study: kNN, OCSVM, iForest, LOF, HBOS, ECOD, INNE and COPOD.

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1002

1003

1004

1005

1006

1007

1009

Distance-based methods, such as kNN, define anomalies based on their relative distance to surrounding points. kNN anomaly detection computes the distance between a data point and its kth nearest neighbor, with larger distances indicating potential anomalies. This method is conceptually simple and effective in low-dimensional spaces with clear separation between normal and anomalous points. However, its primary drawback is the curse of dimensionality, where distance metrics lose discriminative power as dimensionality increases. Additionally, kNN is computationally expensive, with a worst-case complexity of $O(n^2)$, making it impractical for large datasets without optimizations such as approximate nearest neighbor search.

Density-based approaches assume that anomalies reside in low-density regions relative to normal points. LOF estimates the local density of a point by comparing it with the densities of its neighbors. It is highly effective in detecting anomalies in datasets with non-uniform density distributions, where global models may fail. However, LOF is computationally expensive complexity $O(n^2)$ in the worst case and sensitive to the choice of neighborhood size, requiring careful hyperparameter tuning.

A more efficient density estimation approach is HBOS, which models feature distributions independently using histograms. This makes it computationally extremely fast O(n) and scalable to large datasets. However, HBOS assumes feature independence, limiting its effectiveness when strong feature correlations exist. In such cases, its effectiveness diminishes as it fails to capture intricate

¹²https://platform.openai.com/docs/guides/ embeddings/

¹³https://huggingface.co/meta-llama/Llama-3. 2-1B

¹⁴https://huggingface.co/NovaSearch/stella_en_ 400M_v5

¹⁵https://huggingface.co/Qwen/Qwen2.5-1.5B

1107

1108

1060

1061

1062

1063

1064

1065

1066

1067

relationships between features, potentially leading to suboptimal anomaly detection performance.

1010

1011

1012

1013

1014

1015

1018

1019

1020

1023

1024

1025

1026

1027

1028

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1043

1044

1045

1046

1047

1048

1049

1051

1052

1053

1055

1056

1057

1059

Isolation-based approaches, such as iForest, take a different perspective by recursively partitioning the feature space. Since anomalies are typically isolated with fewer splits, iForest identifies them based on the depth required to isolate each point. iForest is computationally efficient O(nlogn) and performs well in high-dimensional spaces compared to distance-based methods, but it is struggle with local anomalies. An extension of iForest, INNE, replaces axis-aligned splits with hypersphere-based partitions. This enhances robustness in detecting anomalies in complex distributions, particularly local anomalies.

Statistical approaches model the underlying distribution of data and identify anomalies as points that significantly deviate from expected behavior. ECOD estimates anomaly scores based on the empirical cumulative distribution function (ECDF) for each feature independently. It is parameterfree and computationally efficient O(n), making it highly scalable. However, like HBOS, ECOD assumes feature independence, which can limit its effectiveness in multivariate settings. COPOD improves upon ECOD by leveraging copula functions to model dependencies between features, making it more effective for detecting anomalies in correlated data. However, this comes at the cost of increased computational complexity, making COPOD less scalable for very large datasets.

E Embedding Analysis

To better understand how different embedding models encode normal and anomalous instances, we visualize their embedding spaces using t-SNE projections across 6 datasets. Figure 5 presents the t-SNE plots for embeddings extracted from 8 embedding models, blue points represent normal instances, while red points denote anomalies.

Separation of Normal and Anomalous Instances. As defined in Section 2, anomalies should ideally exhibit significant deviation from normal instances in the embedding space. The extent to which embeddings separate anomalies from normal data is a crucial factor in determining their effectiveness for anomaly detection.

Most embedding models exhibits clear separation, particularly in the Email Spam dataset, where anomalous points form distinct regions away from the normal distribution. BERT struggles with clear separation, with many anomalies still embedded within normal clusters. This indicates that these models may not encode sufficient discriminative features for anomaly detection tasks.

Dataset-Specific Challenges. The effectiveness of embedding-based anomaly detection varies significantly across datasets, highlighting the influence of domain characteristics:

- Spam Detection (Email Spam, SMS Spam): most embedding models perform well, reflecting their ability to capture explicit spam patterns (e.g., domain-specific keywords, unusual syntax). In contrast, BERT shows more overlap between spam and normal messages, leading to weaker anomaly separation.
- Fake News Detection (COVID-Fake, LIAR2): The separation of anomalies is less pronounced across most embeddings, likely due to the subtle and nuanced nature of misinformation. This suggests that effective detection may require external knowledge or factual reasoning beyond what standard embeddings can provide.
- Hate Speech and Offensive Language (Hate Speech, OLID): All embeddings perform poorly, with anomalies scattered among normal instances. This suggests that hate speech and offensive language often depend on implicit contextual cues rather than explicit linguistic differences, making them harder to distinguish using standard embeddings.

Clustered Anomalies in Spam Detection. For both Email Spam and SMS Spam datasets, the anomalies tend to form compact clusters rather than being scattered as isolated points. This behavior contrasts with other datasets, where anomalies are often more dispersed.

Unlike anomalies in misinformation or hate speech detection, which can manifest in subtle linguistic variations, spam messages tend to exhibit repetitive patterns, including URLs, phone numbers, irregular word spacing and excessive punctuation. Since these patterns are highly distinct but internally consistent, embeddings may cluster them into a well-defined anomaly group rather than spreading them across the feature space.

F Experiments Environment

The entire pipeline, including embedding extraction and anomaly detection, was implemented in

Embeddings	Email-Spam	SMS-Spam	COVID-Fake	LIAR2	Hate-Speech	OLID
BERT						
MINILM						
O-ada						
O-small						
O-large						
Llama						
stella						
Qwen						

Table 5: t-SNE visualization of embeddings from 8 models across 6 datasets. Blue points represent normal instances, while red points denote anomalies.

1109Python 3.9. Experiments were executed on a com-
putational setup equipped with a Ryzen 9 5900X111112-core CPU for data preprocessing and model or-
chestration, and an Nvidia RTX 3060 GPU with111312GB of memory for model inference and embed-
ding generation.