# TEXT-ADBENCH: TEXT ANOMALY DETECTION BENCHMARK BASED ON LLMS EMBEDDING

**Feng Xiao**
The Chinese University of Hong Kong, Shenzhen
Shenzhen, China
fengxiao1@link.cuhk.edu.cn

**Jicong Fan**[*]
The Chinese University of Hong Kong, Shenzhen
Shenzhen, China
fanjicong@cuhk.edu.cn

July 17, 2025

## ABSTRACT

Text anomaly detection is a critical task in natural language processing (NLP), with applications spanning fraud detection, misinformation identification, spam detection and content moderation, etc. Despite significant advances in large language models (LLMs) and anomaly detection algorithms, the absence of standardized and comprehensive benchmarks for evaluating the existing anomaly detection methods on text data limits rigorous comparison and development of innovative approaches. This work performs a comprehensive empirical study and introduces a benchmark for text anomaly detection, leveraging embeddings from diverse pre-trained language models across a wide array of text datasets. Our work systematically evaluates the effectiveness of embedding-based text anomaly detection by incorporating (1) early language models (GloVe, BERT); (2) multiple LLMs (LLaMa-2, LLama-3, Mistral, OpenAI (small, ada, large)); (3) multi-domain text datasets (news, social media, scientific publications); (4) comprehensive evaluation metrics (AUROC, AUPRC). Our experiments reveal a critical empirical insight: embedding quality significantly governs anomaly detection efficacy, and deep learning-based approaches demonstrate no performance advantage over conventional shallow algorithms (e.g., KNN, Isolation Forest) when leveraging LLM-derived embeddings.In addition, we observe strongly low-rank characteristics in cross-model performance matrices, which enables an efficient strategy for rapid model evaluation (or embedding evaluation) and selection in practical applications. Furthermore, by open-sourcing our benchmark toolkit that includes all embeddings from different models and code (see `https://github.com/jicongfan/Text-Anomaly-Detection-Benchmark`), this work provides a foundation for future research in robust and scalable text anomaly detection systems.

***Keywords*** Text Anomaly Detection · LLMs · Text Representation · Benchmark

## 1 Introduction

Anomaly detection (AD) Chandola et al. [2009], Aggarwal [2016], Ruff et al. [2021], broadly defined as the task of identifying patterns that deviates significantly from the majority, is a cornerstone of data analysis and decision-marking in different fields. In past few decades, researchers achieved significant progress in medical diagnosis Kononenko [2001], Amato et al. [2013], Bakator and Radosav [2018], Richens et al. [2020], intrusion detection Mukherjee et al. [2002], Tsai et al. [2009], Liao et al. [2013], Khraisat et al. [2019] in cybersecurity, fraud detection Bolton and Hand [2002], Abdallah et al. [2016], Motie and Raahemi [2024] in finance, and fault detection Isermann [2005], Fan et al. [2017, 2022] in industry. While these AD methods commonly excel in structured data type, such as time-series Zamanzadeh Darban et al. [2024], tabular data, images Fernando et al. [2021] and videos Liu et al. [2025], text-based anomaly detection (see the formal definition given by Definition 1 and the corresponding examples) presents challenges

---

[*]Corresponding Author

due to the unstructured, high-dimensional nature of language, complex and diverse anomalies, where anomalies may manifest as subtle semantic inconsistencies, syntactic irregularities, or contextual deviations.

Due to the inherent characteristics of human language, text representation techniques have become a fundamental component of natural language processing (NLP). Early text representation techniques Harris [1954], Sparck Jones [1972], Aizawa [2003] relied heavily on handcrafted features, such as term frequency (TF) and inverse document frequency (IDF), which aimed to quantify the importance of words within documents. These methods, while simple and computationally efficient, treated text as unordered collections of words (i.e., bag-of-words models), failing to capture syntactic structure or semantic relationships. Subsequent developments introduced vector space models, such as Latent Semantic Analysis (LSA) Blei et al. [2003] and Latent Dirichlet Allocation (LDA) Dumais [2004], incorporating co-occurrence patterns and topic distributions to produce more meaningful document embeddings. However, these techniques still struggled to represent contextual meaning and polysemy, limiting their expressiveness in downstream tasks. A major breakthrough occurred with the rise of neural word embeddings, particularly Word2Vec Mikolov et al. [2013], GloVe Pennington et al. [2014a], and FastText Joulin et al. [2016], which mapped words into dense, low-dimensional vectors based on their distributional properties in large corpora. These models significantly improved the semantic fidelity of word representations, but they remained static, assigning each word a single representation regardless of its context. To address this limitation, researchers introduced contextualized word embeddings, such as ELMo Peters et al. [2018], BERT Devlin et al. [2019] and GPT Radford et al. [2018], which generated dynamic representations conditioned on surrounding text. These models marked a paradigm shift in NLP by allowing representations to reflect both syntax and semantics within a specific context. With the advent of large language models (LLMs) such as GPT-series, LLama-series and etc., leveraging transformer architectures and pretraining on large-scale text data, text representation based on LLMs exhibits remarkable performance in a wide range of language understanding and generation tasks Naveed et al. [2023], Patil and Gudivada [2024].

However, the integration of the representation technique of LLMs into text anomaly detection remains limited exploration and there is a significant lack for a comprehensive analysis and comparison for existing conventional shallow AD algorithms Schölkopf et al. [1999], Breunig et al. [2000], Liu et al. [2008], deep learning-based AD methods Ruff et al. [2018], Fu et al. [2024], and tailored text AD methods Ruff et al. [2019], Manolache et al. [2021], which has hindered the progress and development of text anomaly detection techniques. We have noticed that there are two related works Li et al. [2024a], Cao et al. [2025] recently that provided some comparison and valuable insights. However, the two works exhibit at least the following limitations: (1) insufficient comparative analysis, with a focus primarily on shallow anomaly detection algorithms; (2) the absence of discussion or comparison regarding different pooling strategies when utilizing text embeddings from large language models (LLMs), despite their critical impact on embedding effectiveness; (3) a narrow evaluation framework, relying on a single performance metric; and (4) a lack of open-source embeddings and code, limiting reproducibility and further research.

To address these gaps, this work introduces Text-ADBench, a comprehensive benchmark for text anomaly detection based on embeddings derived from LLMs. Specifically, we construct a large number of two-stage text AD methods. First, we generate the text embedding using a diverse suite of language models, including early language models (GloVe-6B, BERT), open-source LLMs (LLaMA2-7B, LLaMA3-8B, Mistral-7B) and OpenAI's text-embedding models (small, ada, large). To aggregate sequential token embeddings into a single vector representation, we utilize three pooling strategies, including "mean", "end-of-sequence (EOS) token", and "weighted mean". Second, we develop anomaly detection tasks by applying these embeddings to a range of AD methods, encompassing both shallow and deep learning-based techniques. The workflow of Text-ADBench is illustrated in Fig. 1. Additionally, we incorporate two specialized text AD methods (CVDD Ruff et al. [2018] and DATE Manolache et al. [2021]) into our comparative analysis. Our experiments are conducted on eight real-world text datasets. The main contributions of this work are summarized as follows.

- We present a comprehensive benchmark for text anomaly detection, addressing the existing gap in leveraging large language models (LLMs) embeddings for this task. Our framework incorporates diverse language models, multiple pooling strategies, a wide range of anomaly detection methods, and comprehensive evaluation metrics

- We conduct further analysis of the results and identify a low-rank property in performance. This finding has two important implications: (1) the detection performance of novel text datasets or AD methods can be reliably predicted using only a subset of performance measurements, and (2) this property enables an efficient strategy for rapid model evaluation (embedding evaluation) and selection in practical applications.

- We release Text-ADBench at `https://github.com/jicongfan/Text-Anomaly-Detection-Benchmark` that including all data, embeddings, and code used in our experiments. We provide an integrated and easy-to-use framework whatever related researchers want to reproduce results or evaluate novel embeddings and methods, which would foster the development of these fields.
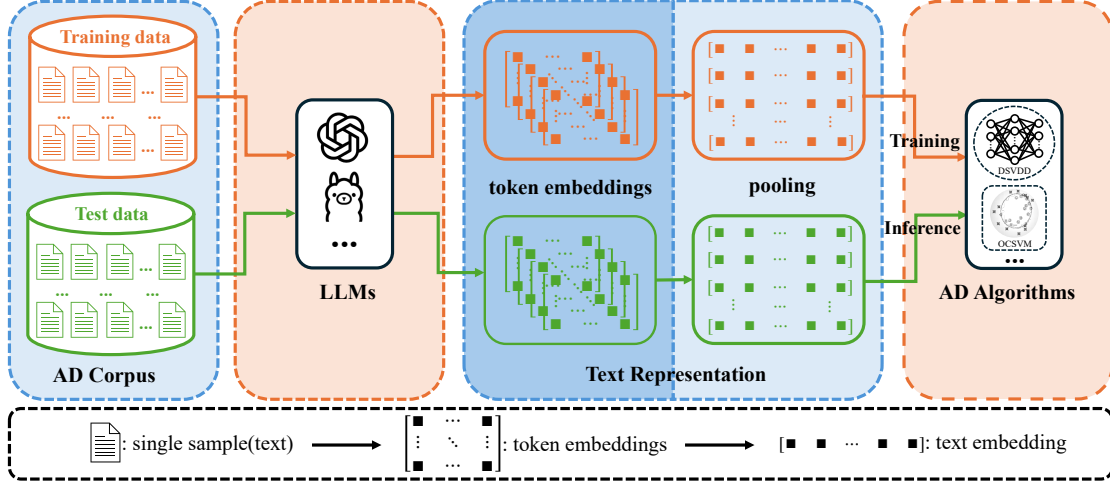
Figure 1: Flowchart of Text-ADBench.

This work serves as a foundational resource for both researchers and practitioners, advancing the text AD field through its methodological rigor and holistic evaluation. By open-sourcing our benchmark framework, precomputed embeddings, we aim to catalyze further research into hybrid pooling strategies, domain adaptation, and efficient LLM utilization for text anomaly detection.

The remainder of this paper is organized as follows: Section 2 reviews related work; Section 3 formulates the problem of text anomaly detection; Section 4 details the specific settings including datasets, LLMs, embedding, AD methods. Section 5 presents experimental results and basic analysis; Section 6 explores the low-rank property of performance; and Section 7 concludes this work and provide future directions.

## 2  Related Work

### 2.1  Unsupervised Anomaly Detection

Unsupervised Anomaly Detection (UAD) is a fundamental task in data analysis and decision-making, with broad applications across various domains, such as fault detection in industry, fraud detection in finance, medical diagnosis, and intrusion detection in cybersecurity. Consequently, a flood of unsupervised anomaly detection methods have been proposed in recent decades Schölkopf et al. [1999], Breunig et al. [2000], Liu et al. [2008], Ruff et al. [2018], Deecke et al. [2019], Qiu et al. [2021], Fu et al. [2024], Xiao et al. [2025a], Dai and Fan [2025]. Generally, these UAD methods are grounded in distinct assumptions about data distribution Aggarwal [2016], Pang et al. [2021], Ruff et al. [2021]. For instance, anomalies may reside in low-density regions or exhibit deviations in feature space compared to normal instances. The performance of these methods is often contingent upon the degree of alignment between the input data and the underlying assumptions. Roughly, these methods can be organized into two categories, including conventional machine learning-based shallow algorithms and deep learning-based methods Pang et al. [2021]. Shallow anomaly detection algorithms Schölkopf et al. [1999], Breunig et al. [2000], Liu et al. [2008], Li et al. [2022] typically have a straightforward learning process and exhibit superior interpretability. For instance, OCSVM Schölkopf et al. [1999] learns a maximum-margin hyperplane in a kernel space to separate normal data from the origin. Isolation Forest Liu et al. [2008] builds an ensemble of trees to isolate data points, where anomalies are identified by their shorter average path lengths on these trees. KNN Ramaswamy et al. [2000] defines the anomaly score of a sample as its distance to the k-th nearest neighbor, based on the assumption that anomalies reside in sparser regions of the feature space compared to normal data. Deep learning-based anomaly detection methods Schlegl et al. [2017], Zong et al. [2018], Ruff et al. [2018], Pang et al. [2019], Perera et al. [2019], Goyal et al. [2020], Fan et al. [2022], Qiu et al. [2021], Cai and Fan [2022], Xu et al. [2023a], Fu et al. [2024], Ye et al. [2025] often achieve superior detection performance, particularly in high-dimensional data settings. For instance, AutoEncoder (AE) Hinton and Salakhutdinov [2006] minimizes the reconstruction error on normal data. It assumes that abnormal data cannot be reconstructed well when an autoencoder is trained only on normal data. Naturally, reconstruction error is used as an anomaly score. Deep SVDD Ruff et al. [2018] learns a minimum-radius hypersphere to encompass all the normal data while the unseen abnormal data are expected to fall outside, where its anomaly score is defined by the distance to the hypersphere's center. The PLAD proposed by Cai

and Fan [2022] trains a perturbator to perturb normal data and a classifier to distinguish normal data and perturbed data, leading to an effective decision boundary for anomaly identification.

Note that in addition to UAD methods, there are also a few semi-supervised AD methods Hendrycks et al. [2018], Pang et al. [2018], Akcay et al. [2019], Ruff et al. [2020], Zhou et al. [2021], Pang et al. [2023], Xiao et al. [2025b] that can utilize labeled anomalies to improve the detection performance, since in some scenarios, a few known anomalous samples, though scarce, are available during the training stage.

## 2.2 Text Representation

Text representation is a fundamental technique of natural language processing (NLP) that transforms unstructured textual data into structured numerical formats, thereby enabling machine learning models to process and analyze linguistic information effectively. The evolution of text representation techniques has undergone several significant stages.

Early text representation techniques primarily relied on statistical and frequency-based approaches. The Bag-of-Words (BOW) model Harris [1954] and Term Frequency-Inverse Document Frequency (TF-IDF) Sparck Jones [1972] represented text as vectors of word frequencies, thereby only preserving the statistical lexical information but neglecting word order. As a result, the BOW model fails to capture semantic relationships between words and lacks contextual understanding. To address the limitations of the BOW model, distributed representations Mikolov et al. [2013], Pennington et al. [2014b] were developed, which encode words as dense vectors in a continuous space. For instance, Word2Vec Mikolov et al. [2013] introduced two architectures, Skip-gram and Continuous Bag-of-Words (CBOW), that capture semantic and syntactic relationships by predicting word contexts. Similarly, GloVe (Global Vectors for Word Representation) Pennington et al. [2014b] integrated global co-occurrence statistics with local context window-based learning. However, despite their advancements, Word2Vec and GloVe produce static and context-free word embeddings, which are unable to effectively handle polysemy and out-of-vocabulary (OOV) words. In 2017, the introduction of Transformer architecture Vaswani et al. [2017] revolutionized text representation by generating context-dependent embeddings. Subsequent models, such as GPT (Generative Pretrained Transformer) Radford et al. [2018] and BERT (Bidirectional Encoder Representations from Transformers) Devlin et al. [2019], leveraged the self-attention mechanism to produce dynamic word representations conditioned on surrounding text, markedly improving performance on downstream NLP tasks. With the advent of LLMs (Large Language Models), recently, the NLP community has begun adopting decoder-only LLMs for text embedding Wang et al. [2023], BehnamGhader et al. [2024a].

## 2.3 Text Anomaly Detection

In the field of natural language processing, anomaly detection tasks are employed to identify harmful content, spam reviews, and abusive language. A straightforward strategy is to construct pipeline (two-stage) methods by integrating text representation techniques with existing anomaly detection methods. Xu et al. Xu et al. [2023b] systematically estimated such combinations based on SBERT Reimers and Gurevych [2019]. Additionally, several specialized text anomaly detection algorithms Ruff et al. [2019], Manolache et al. [2021], Das et al. [2023] were proposed in recent years. For instance, CVDD Ruff et al. [2018], a self-attention-based model for unsupervised text anomaly detection method, minimizes the weighted cosine distance between the attention-weighted text embeddings and "context vectors", where it employs self-attention heads to generate "context vectors" that capture diverse semantic themes in normal data and each head produces an orthogonal attention matrix, ensuring independent feature extraction. During the inference stage, anomalies are identified by averaging the cosine distances between the new sample's embeddings and all "context vectors". DATE Manolache et al. [2021] trains Transformers using masked language modeling (MLM) and two tailored pretext tasks to capture contextual anomalies. FATE Das et al. [2023], a semi-supervised text anomaly detection method, employs a contrastive objective to adjust anomaly scores. This approach ensures that normal samples cluster around a reference score while anomalies deviate significantly. Additionally, FATE leverages attention mechanisms and multi-instance learning to capture distinctive anomaly behaviors. More recently, Pantin and Marsala Pantin and Marsala [2024] introduced ERLA, a robust ensemble-based anomaly detection method utilizing autoencoders, where each autoencoder incorporates a local robust subspace recovery projection of the original data in its encoding embedding, thereby leveraging the geometric properties of the k-nearest neighbors to optimize subspace recovery and identify anomalous patterns in textual data.

## 2.4 Anomaly Detection based on LLMs

The remarkable success of large language models (LLMs) in natural language processing has spurred their adaptation to anomaly detection (AD) tasks Su et al. [2024] for further improvement of detection accuracy via their advanced semantic understanding. Current LLM-based AD approaches can be categorized into three paradigms: ***(i) prompt-based***

***methods*** Yang et al. [2024], Dong et al. [2024] directly employ pretrained LLMs to perform few-shot or zero-shot learning by designing specialized prompt functions, and ***(ii) fine-tuning methods*** Gu et al. [2024], Li et al. [2024b] adapt general-purpose LLMs into domain-specific anomaly detectors by fine-tuning them on task-relevant datasets, thereby tailoring their capabilities to the nuances of AD, and ***(iii) embedding based methods*** Li et al. [2024a], Cao et al. [2025] regard LLMs as feature extractors, generating high-dimensional embeddings that capture high-quality semantic and contextual patterns and then leveraging existing AD algorithms (e.g., OCSVM, Isolation Forest) to detect deviations from normal pattern based on these embeddings.

## 2.5 Existing Benchmarks for Text Anomaly Detection

There are some notable works Xu et al. [2023b], Bejan et al. [2023], Li et al. [2024a], Yang et al. [2024], Cao et al. [2025] that take effort to text anomaly detection. For instance, Xu et al. Xu et al. [2023b] evaluated 22 AD algorithms on 17 text corpora and mainly analyzed the performance effects of weak supervised signals. To the best of our knowledge, the first text anomaly detection benchmark based on LLMs is NLP-ADBench Li et al. [2024a], where Li et al. evaluated 11 AD algorithms based on embeddings of BERT and OpenAI (text-embedding-3-large) across 8 text corpora. To clearly delineate the different emphases and advantages of existing benchmarks and to further compare them with Text-ADBench, we provide a statistical comparison in Table 1. In summary, this benchmark emphasizes the following three aspects.

- First, while previous works Xu et al. [2023b], Bejan et al. [2023] mainly focus on the embeddings derived from early language models, recent works Li et al. [2024a], Cao et al. [2025] have begun to incorporate embeddings from LLMs, albeit with a limited scope. Considering the rapid advancement of the representation capabilities of LLMs, we compare 33 distinct embeddings for each dataset by integrating embedding models with pooling strategies. Notably, previous works do not compare the impact of different pooling strategies.

- Second, this benchmark identifies a low-rank property in performance matrices and demonstrates that the detection performance of novel text datasets or AD methods can be effectively predicted based on historical performance measurements. This finding enables an efficient strategy for rapid model evaluation (or embedding evaluation) or selection based on the results of this benchmark.

- More importantly, we release all resources used in this benchmark including all original text datasets, embeddings and code. These open-access resources ensures researchers and practitioners can easily end efficient reproduce our results, estimate novel datasets or AD algorithms, thereby fostering progress in the field. Additionally, the released embeddings can be leveraged for other downstream tasks.

Table 1: Comparison among Text-ADBench and existing benchmarks, where "DL" and "TFT" refers to "Deep Learning" and "Tailored for Text", respectively. Note that the term "# Emb." denotes the number of distinct embeddings derived from each text dataset.

| Benchmark | Coverage (numbers) | | | | | AD Algorithm Type | | | Representation | | Open Source | | Performance Matrices | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Datasets | Algo. | LLMs | Emb. | Metrics | Shallow | DL | TFT | Early LM | Pooling | Code | Embedding | Analysis | Prediction |
| Xu et al. Xu et al. [2023b] (2023) | 22 | 17 | 0 | 2 | 1 | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Bejan et al. Bejan et al. [2023] (2023) | 7 | 4 | 0 | 2 | 2 | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ |
| Li et al. Li et al. [2024a] (2024) | 8 | 11 | 1 | 2 | 1 | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ |
| Yang et al. Yang et al. [2024] (2024) | 5 | 10 | 2 | - | 2 | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ |
| Cao et al. Cao et al. [2025] (2025) | 6 | 8 | 5 | 8 | 1 | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Text-ADBench (ours) | 12 | 12 | 12 | 33 | 2 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

# 3  Problem Formulation

Text anomaly detection, considered in this work, aims to identify the text instances that deviate significantly from the majority. The following presents a formal definition.

**Definition 1** (Text anomaly detection). *Let $\mathcal{C} = \{s_1, s_2, \cdots, s_n\}$ be a corpus of $n$ textual sequences, where all or most sequences are in some unknown normal condition or pattern $P$. Text anomaly detection aims to learn a detector $f$ from $\mathcal{C}$ that can determine whether a new textual sequence $s_{new}$ is in $P$ or not.*

For instance, if $\mathcal{C}$ is composed of emails, a textual sequence containing spam information should be detected as abnormal. One more example, if $\mathcal{C}$ consists of movie reviews, a review with abusive content would also be detected as abnormal.

In this work, we split the text anomaly detection into two stages. In the first stage, we extract text embeddings based on early language models (GloVe and BERT) and LLMs. Subsequently, in the second stage, we design a general unsupervised anomaly detection task on the text embeddings. Before delving into the details, we first introduce the notation convention as follows.

1. **Text Representation based on Embedding Models**    Let $\mathbf{M}_{\text{emb}}$ be a language model (e.g., BERT and LLaMA3). We obtain the embeddings $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_n\}$ of $\mathcal{C}$ as

$$\mathbf{x}_i = Pooling\big(\mathbf{M}_{\text{emb}}(s_i)\big), \quad i = 1, 2, \cdots, n, \tag{1}$$

   where $Pooling(\cdot)$ aims to aggregate token-level embeddings of the sequence $s_i$ into a single vector $\mathbf{x}_i \in \mathbb{R}^d$, where $d$ represents the embedding dimension of the language model.

2. **Unsupervised Anomaly Detection**    We assume that the $\mathcal{X}$ is drawn from an unknown distribution $\mathcal{D}_{\mathbf{x}} \subseteq \mathbb{R}^d$. A point $\mathbf{x} \in \mathbb{R}^d$ is deemed to be anomalous if $\mathbf{x} \notin \mathcal{D}_{\mathbf{x}}$. Then, the goal of the unsupervised AD is to obtain a decision function $h_{\text{UAD}} : \mathbb{R}^d \to \{0, 1\}$ by utilizing only $\mathcal{X}$, such that $h_{\text{UAD}}(\mathbf{x}) = 0$ if $\mathbf{x} \in \mathcal{D}_{\mathbf{x}}$ and $h_{\text{UAD}}(\mathbf{x}) = 1$ if $\mathbf{x} \notin \mathcal{D}_{\mathbf{x}}$. The primary difference among AD methods lies in the design of the decision function $f(\cdot)$.

Based on the two stages above, the detector $f$ can be formulated as

$$f(s) := h_{\text{UAD}}\big(Pooling\big(\mathbf{M}_{\text{emb}}(s)\big)\big) \tag{2}$$

The value of $f(s_{\text{new}})$ can determine whether the textual sequence $s_{\text{new}}$ is normal or anomalous.

In this work, we consider the combinations of different embedding models, different pooling operations, and different UAD algorithms, leading to a comprehensive evaluation of key techniques for textual anomaly detection.

## 4    Benchmark Settings

In this section, we will introduce the detailed experimental preparations, including dataset splitting, embedding models, and the anomaly detection methods used in this benchmark.

### 4.1    Datasets

Different from anomaly detection tasks on images, tabular data, and time series, where there are standard datasets with normal and abnormal samples defined clearly, text anomaly detection presents unique challenges. The richness of natural language allows anomalies to manifest in diverse forms, such as rare topics, unusual language styles, irregular syntax or grammar, and harmful or abusive content. We observe that existing textual AD benchmarks Xu et al. [2023b], Bejan et al. [2023], Li et al. [2024a], Yang et al. [2024], Cao et al. [2025] consistently employ the NLP classification datasets. Aligning with this established methodology, our benchmark utilizes 8 text classification datasets from various NLP domains to construct 14 specialized Text-AD datasets. The statistical information of the Text-AD datasets are shown in Table 2.

Table 2: Statistical information of the Text-AD datasets.

| Text-AD Dataset | Domain | # Training Samples | # Test Samples | Anomaly Ratio |
|---|---|---|---|---|
| 20Newsgroups | News | 10,419 | 7,876 | 0.12 |
| Reuters21578 | News | 4,435 | 4,723 | 0.45 |
| IMDB | Movie Review | 10,000 | 40,000 | 0.625 |
| SST2 | Movie Review | 10,000 | 58,078 | 0.51 |
| SMS-Spam | Phone Message | 3,000 | 2,534 | 0.29 |
| Enron | Email | 10,000 | 21,924 | 0.31 |
| WOS | Paper Abstract | 28,918 | 18,065 | 0.31 |
| DBpedia14-0 | Wikipedia Term | 20,000 | 44,999 | 0.44 |
| DBpedia14-1 | Wikipedia Term | 20,000 | 44,999 | 0.44 |
| DBpedia14-2 | Wikipedia Term | 20,000 | 44,999 | 0.44 |
| DBpedia14-3 | Wikipedia Term | 20,000 | 45,000 | 0.44 |
| DBpedia14-4 | Wikipedia Term | 20,000 | 44,997 | 0.44 |

The specific splitting and descriptions of each dataset are summarized below.

- **20Newsgroups**[2] Lang [1995] is a collection of newsgroup documents.  It is organized into 20 different newsgroups, each corresponding to a different topic.  Some of the newsgroups are very closely related to each

---

[2]http://qwone.com/ jason/20Newsgroups/

other (e.g., "comp.sys.ibm.pc.hardware / comp.sys.mac.hardware"), while others are highly unrelated (e.g "misc.forsale / soc.religion.christian"). For constructing the text-AD dataset, class "misc.forsale" is regarded as an abnormal class, and the remaining classes are regarded as normal.

- **Reuters21578**[3] Lewis [1997] is a collection of documents that appeared on Reuters newswire in 1987. It is organized into 59 different groups with 21578 samples in the original version. In this work, we use the ApteMod version provided in NLTK Corpora[4] and classes "earn" and "acq" are regarded as normal, and the remaining are regarded as abnormal.

- **IMDB**[5] Maas et al. [2011] is a dataset for binary sentiment ("pos" and "neg") classification containing substantially more data than previous benchmark datasets. For constructing the text-AD dataset, class "pos" is regarded as normal, and "neg" is regarded as the abnormal class.

- **SST2**[6] Socher et al. [2013] is a corpus with fully labeled parse trees that allows for a complete analysis of the compositional effects of sentiment in language. The corpus is based on the dataset introduced by Pang and Lee (2005) and consists of 11,855 single sentences extracted from movie reviews. It is a binary sentiment classification dataset. In this benchmark, class "1-(positive)" is regarded as the normal class and class "0-(negative)" is regarded as the abnormal class.

- **SMS-Spam**[7] Almeida et al. [2011] is a public set of SMS labeled messages that have been collected for mobile phone spam research. It has one collection composed by 5,574 English, real and non-encoded messages, tagged according being legitimate or spam. Naturally, legitimate messages are regarded the normal samples, and spam messages are regarded as abnormal samples.

- **Enron**[8] contains a total of about 0.5M messages from about 150 users, mostly senior management of Enron. We counted the entire dataset and found the two email accounts ("kay.mann@enron.com", "vince.kaminski@enron.com") with the most emails under them, both greater than 10,000, as well as many accounts with just 1 email under them. Based on this observation, we use the emails from accounts ("kay.mann@enron.com", "vince.kaminski@enron.com") as normal samples and the emails from the accounts with only 1 email as abnormal samples.

- **WOS**[9] Kowsari et al. [2017] is Web of Science Dataset (WOS-46985). This dataset contains the abstracts of 46,985 published papers that have 7 parent categories, including "Computer Science (CS), Electrical Engineering (ECE), Psychology, Mechanical Engineering (MAE), Civil Engineering (Civil), Medical Science (Medical), and Biochemistry". For constructing the text-AD dataset, class "Psychology" is regarded as the normal class, and the remaining classes are regarded as abnormal.

- **DBpedia14**[10] Zhang et al. [2015] is constructed by picking 14 non-overlapping classes from DBpedia 2014. Each class contains 40,000 training samples and 5,000 testing samples. We construct five text-AD datasets (DBpedia14-0, DBpedia14-1, DBpedia14-2, DBpedia14-3, DBpedia14-4) based on the original dataset. Specifically, for DBpedia14-0, we use class 0 as the normal class, and the samples from the test set of classes [1, 2, 3, 4] are regarded as abnormal samples. For DBpedia14-1, we use class 1 as the normal class, and the samples from the test set of classes [0, 2, 3, 4] are regarded as abnormal samples. The text-AD datasets (DBpedia14-2, DBpedia14-3, DBpedia14-4) have the same construction process.

## 4.2 Embedding Models

Text embedding is widely recognized as the foundational and essential step in NLP tasks, as it transforms the semantic content of natural language into a structured vector representation. Over the past several years, the emergence of innovative embedding techniques has consistently propelled the field of NLP forward, significantly advancing the development and applications of NLP techniques.

For instance, researchers at Stanford University advanced the field of word embeddings with the introduction of GloVe (Global Vectors for Word Representation) Pennington et al. [2014a]. GloVe enhanced the capabilities of Word2Vec by incorporating global statistical information across the entire corpus to generate word vectors, which enabled a

---

[3]https://raw.githubusercontent.com/nltk/nltk_data/gh-pages/packages/corpora/reuters.zip

[4]https://www.nltk.org/nltk_data/

[5]http://ai.stanford.edu/ amaas/data/sentiment/

[6]https://huggingface.co/datasets/stanfordnlp/sst2

[7]https://huggingface.co/datasets/ucirvine/sms_spam

[8]https://huggingface.co/datasets/Hellisotherpeople/enron_emails_parsed

[9]https://huggingface.co/datasets/river-martin/web-of-science-with-label-texts

[10]https://huggingface.co/datasets/fancyzhx/dbpedia_14

more nuanced semantic understanding by integrating both local context windows and global corpus statistics. The landscape of NLP was further revolutionized in 2017 with the introduction of the transformer architecture Vaswani et al. [2017], which introduced the self-attention mechanism. Building on this breakthrough, BERT(Bidirectional Encoder Representation from Transformers) Devlin et al. [2019], released in 2018, pioneered context-dependent word embeddings. Unlike earlier models such as Word2Vec and GloVe, which produced statistic, context-free embeddings, BERT leveraged a transformer-based architecture to create dynamic representations that capture the meaning of words based on their surrounding context, considering both preceding and succeeding words within a sentence. With the advent of LLMs (Large Language models), only recently, the NLP community started to adopt decoder-only LLMs for text embedding BehnamGhader et al. [2024a].

In this work, we primarily employ LLMs as embedding models and also consider two classical language models including GloVe and BERT. Detailed information about the embedding models used in this work are summarized in Table 3. Building upon Llama-2-7B-chat Touvron et al. [2023], Mistral-7B-Instruct-v0.2 Jiang et al. [2023] and Llama-3-8B-Instruct AI@Meta [2024], we employ their fine-tuned versions tailored for text embedding from LLM2Vec [11] BehnamGhader et al. [2024b].Additionally, three specialized text embedding models [12] (text-embedding-3-small, text-embedding-ada-002, text-embedding-3-large) from OpenAI are also evaluated in this work. For brevity, we refer to these models as "small, ada, large" in the following sections. Our text representation pipeline involves two sequential stages: (1) extracting token-level embeddings via the embedding models, followed by (2) applying pooling strategies to aggregate the token embeddings into final text representations. As detailed in Table 3, three distinct pooling strategies are evaluated in this benchmark. For each AD dataset, we derive 33 distinct text representations, each generated by varying combinations of embedding models and pooling strategies. Notably, the terms "mntp"[13], "mntp-unsup-simcse"[14] and "mntp-supervised"[15] denote the different fine-tuning techniques. More details on these fine-tuning techniques can be found in LLM2Vec BehnamGhader et al. [2024b].

Table 3: Basic information of embedding models.

| Embedding Models | Pooling | Max Tokens | Dimensions | Parameters | Open Source |
|---|---|---|---|---|---|
| GloVe (glove-6B) | Mean | 512 | 300 | 120M | ✓ |
| BERT (large-uncased) | Mean, CLS | 512 | 1024 | 340M | ✓ |
| LLaMA-2(mntp) | Mean, Weighted Mean, EOS | 8191 | 4096 | 7B | ✓ |
| LLaMA-2(mntp-unsup-simcse) | Mean, Weighted Mean, EOS | 8191 | 4096 | 7B | ✓ |
| LLaMA-2(mntp-supervised) | Mean, Weighted Mean, EOS | 8191 | 4096 | 7B | ✓ |
| LLaMA-3(mntp) | Mean, Weighted Mean, EOS | 8191 | 4096 | 8B | ✓ |
| LLaMA-3(mntp-unsup-simcse) | Mean, Weighted Mean, EOS | 8191 | 4096 | 8B | ✓ |
| LLaMA-3(mntp-supervised) | Mean, Weighted Mean, EOS | 8191 | 4096 | 8B | ✓ |
| Mistral(mntp) | Mean, Weighted Mean, EOS | 32768 | 4096 | 7B | ✓ |
| Mistral(mntp-unsup-simcse) | Mean, Weighted Mean, EOS | 32768 | 4096 | 7B | ✓ |
| Mistral(mntp-supervised) | Mean, Weighted Mean, EOS | 32768 | 4096 | 7B | ✓ |
| OpenAI (embedding-3-small) | - | 8192 | 1536 | - | ✗ |
| OpenAI (embedding-ada-002) | - | 8192 | 1536 | - | ✗ |
| OpenAI (embedding-3-large) | - | 8192 | 3072 | - | ✗ |

## 4.3 Anomaly Detection Methods

For comprehensive benchmarking, we consider both shallow machine learning-based AD algorithms, deep learning based AD methods and specialized textual AD methods. The details are described below.

**Shallow Machine Learning-based AD Algorithms:**

- **One-Class SVM (OCSVM)** Schölkopf et al. [1999] OCSVM maps input data to a high-dimensional space by kernel methods and finds a hyperplane that separate normal data from the origin with maximum margin. RBF kernel was used in our experiments.

- **Isolation Forest (IForest)** Liu et al. [2008] IForest isolates instances by building an ensemble of Trees, then anomalies are those instances which have short average path lengths on the Trees.

---

[11]https://github.com/McGill-NLP/llm2vec/tree/main

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

[13]mntp: masked next token prediction

[14]mntp-unsup-simcse: unsupervised contrastive learning

[15]mntp-supervised: supervised contrastive learning

- **Local Outlier Factor (LOF)** Breunig et al. [2000] LOF measures the local deviation of the density of each instance with respect to its neighbors. We set the hyperparameter "n_neighbors"=30 in our experiments.

- **Principal Component Analysis (PCA)** Shyu et al. [2003] PCA is a linear dimensionality reduction method that employs matrix decomposition to transform data into a new orthogonal coordinate system defined by its principal components. These principal components are ordered such that the first component captures the maximal variance in the data, with subsequent components representing decreasingly smaller variations. When applied in anomaly detection, PCA projects data onto a lower-dimensional subspace spanned by the top-k eigenvectors. The anomaly score for a new sample is computed as the reconstruction error.

- **K-Nearest Neighbors (KNN)** Ramaswamy et al. [2000] KNN views the anomaly score of a new sample as the distance to its $k$-th nearest neighbor. We set $k = 3$ in our experiments.

- **Kernel Density Estimation (KDE)** Kim and Scott [2012] KDE is a non-parametric method to estimate the probability density function (PDF) of a dataset. It smooths observed data points using a kernel function to approximate the underlying distribution. When applied in anomaly detection, anomalies are identified as points in low-density regions of the estimated PDF. We use the Gaussian kernel in our experiments.

- **Empirical-Cumulative-distributed-based Outlier Detection (ECOD)** Li et al. [2022] ECOD is a hyperparameter-free outlier detection algorithm based on the empirical cumulative distribution function (ECDF). It employs ECDF to estimate the density of each feature independently and assumes that outliers are located in the tails of the distribution.

**Deep Learning-based AD Methods:**

- **AutoEncoder (AE)** Hinton and Salakhutdinov [2006] An AutoEncoder is a type of neural network that consists of two parts including an encoder and a decoder. Encoder learns to compress input data into a lower-dimensional representation space and then decoder reconstructs them back to the original data space. When applied in anomaly detection, reconstruction error is a natural selection as the anomaly score.

- **Deep Support Vector Data Description (SVDD)** Ruff et al. [2018] Deep SVDD is a deep one-class learning method that learns a minimal enclosing hypersphere in a latent space to characterize normal data. The anomaly score for a sample is given by the squared Euclidean distance to the hypersphere center in latent space.

- **Dense Projection for Anomaly Detection (DPAD)** Fu et al. [2024] DPAD is a density based method that learns to obtain locally dense low-dimensional representation for training data (normal data). The anomaly score is the KNN score in the low-dimensional space.

**Specialized Textual AD Methods:**

- **Context Vector Data Description (CVDD)** Ruff et al. [2019] CVDD is a method that builds upon pretrained word embeddings to learn multiple sentence representations (context vector) that capture multiple semantic contexts via the self-attention mechanism and to project the word representations near these context vector by minimizing the cosine distance between them.

- **Detecting Anomalies in Text using ELECTRA (DATE)** Manolache et al. [2021] DATE an end-to-end approach for the discrete text domain that combines a powerful representation learner for Text (ELECTRA) Clark et al. [2020] and an anomaly score is tailed for sequential data.

To evaluate these anomaly detection approaches, for seven shallow and three deep learning based AD methods, we combine them with all 33 text embeddings, which results in 330 distinct two-stage anomaly detection configurations (10 AD methods × 33 embeddings) per dataset. Regarding the specialized text AD methods, CVDD utilizes word embeddings from GloVe and BERT, and DATE operates directly on raw text inputs without requiring pretrained embeddings.

## 4.4 Evaluation Metrics

For a holistic evaluation, this work utilizes both AUROC (Area Under the Receiver Operating Characteristic curve) and AUPRC (Area Under the Precision-Recall curve) to measure detection performance. AUROC is more suitable for scenarios where the number of positive and negative samples is relatively balanced or where high classification accuracy is required for both types of samples. It can comprehensively reflect the performance of the model. AUPRC is more suitable for scenarios with imbalanced positive and negative samples, especially where high precision and recall are required for the identification of positive samples. It can better reflect the model's predictive ability for the minority class. The experiments on DBpedia14 were conducted on Intel(R) Xeon(R) Gold 5117 CPU @ 2.00GHz with 4×

NVIDIA GeForce RTX 4090 GPUs, Python 3.8, and all other experiments all were conducted Intel(R) Xeon(R) CPU E5-2678 v3 @ 2.50GHz with 4× NVIDIA GeForce RTX 3090 GPUs, Python 3.8. We report the average results over five runs.

## 5 Experimental results

Table 4: Average AUROC(%) of all two-stage methods. The top two results on each dataset are marked in **bold** (top, second).

| Embedding | Pooling | 20News. | Reuters. | IMDB | SST2 | SMS | Enron | WOS | DBp.0 | DBp.1 | DBp.2 | DBp.3 | DBp.4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 56.26 | 76.26 | 46.97 | 48.95 | 40.26 | 65.22 | 45.94 | 76.89 | 87.70 | 81.03 | 85.07 | 82.46 |
| BERT | CLS | 55.73 | 60.05 | 45.27 | 53.99 | 41.84 | 60.11 | 47.62 | 85.41 | 85.96 | 77.30 | 74.33 | 74.19 |
| | Mean | 54.89 | 80.33 | 45.40 | 54.25 | 51.77 | 59.88 | 48.19 | 83.41 | 85.92 | 83.66 | 82.48 | 82.00 |
| LLaMA-2 (mntp) | EOS | 47.40 | 64.82 | 48.90 | 53.26 | 62.18 | 57.73 | 48.49 | 66.25 | 72.31 | 65.09 | 72.45 | 69.68 |
| | Mean | 51.95 | 61.16 | 48.03 | 52.29 | 56.07 | 60.36 | 47.15 | 66.95 | 74.54 | 72.93 | 80.45 | 80.20 |
| | Weighted Mean | 51.15 | 61.10 | 48.21 | 52.18 | 60.52 | 59.18 | 46.61 | 65.45 | 73.26 | 71.01 | 78.46 | 78.25 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 59.30 | **92.49** | 48.19 | 64.65 | 79.80 | 82.89 | 59.22 | **92.66** | **98.00** | 88.81 | **97.01** | 94.72 |
| | Mean | 52.26 | 86.72 | 49.41 | 62.91 | 78.64 | 74.76 | 44.81 | 78.22 | 91.18 | 81.15 | 94.99 | 91.85 |
| | Weighted Mean | 53.02 | 85.27 | 49.93 | 63.11 | 84.83 | 72.49 | 44.87 | 77.34 | 90.86 | 79.05 | 94.11 | 90.97 |
| LLaMA-2 (mntp-supervised) | EOS | 57.53 | 90.41 | 51.74 | 68.61 | **91.40** | 82.74 | **67.18** | 88.16 | 97.59 | 84.21 | 94.88 | 95.22 |
| | Mean | 56.81 | 89.84 | 53.38 | **71.35** | 78.93 | 80.53 | 49.95 | 83.96 | 97.20 | 87.67 | 93.44 | 93.50 |
| | Weighted Mean | 57.13 | 88.55 | 53.13 | 69.40 | 84.67 | 80.11 | 50.95 | 82.39 | 96.49 | 86.69 | 92.69 | 92.91 |
| LLaMA-3 (mntp) | EOS | 62.57 | 82.40 | 45.53 | 56.90 | **90.86** | 74.07 | 55.80 | 76.92 | 89.09 | 87.06 | **95.20** | 90.57 |
| | Mean | 56.82 | 85.16 | 47.40 | 58.54 | 72.42 | 70.57 | 45.89 | 66.82 | 82.35 | 83.03 | 93.84 | 89.98 |
| | Weighted Mean | 54.59 | 81.82 | 47.72 | 59.11 | 80.20 | 67.15 | 46.30 | 67.44 | 82.70 | 81.68 | 93.16 | 90.05 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 54.35 | 77.95 | 47.36 | 58.47 | 89.98 | 81.19 | 55.76 | 63.48 | 78.27 | 64.78 | 75.43 | 72.86 |
| | Mean | 48.68 | 86.53 | 49.26 | 60.56 | 78.63 | 72.96 | 42.38 | 64.76 | 86.10 | 80.40 | 94.79 | 91.02 |
| | Weighted Mean | 46.56 | 84.57 | 49.42 | 60.17 | 86.78 | 70.13 | 43.35 | 66.28 | 85.92 | 79.62 | 94.03 | 90.56 |
| LLaMA-3 (mntp-supervised) | EOS | 62.17 | 89.04 | 52.98 | 67.22 | 90.70 | 80.01 | 56.24 | 87.97 | 97.31 | 82.84 | 91.48 | 92.61 |
| | Mean | 60.06 | 92.21 | 50.63 | 69.56 | 86.64 | 79.87 | **63.59** | 87.30 | **98.16** | **90.85** | 95.11 | **96.04** |
| | Weighted Mean | 57.91 | 90.95 | 50.58 | 68.84 | 86.28 | 79.98 | 62.76 | 87.37 | 97.88 | **89.90** | 94.57 | **95.42** |
| Mistral (mntp) | EOS | **71.99** | 76.96 | 43.67 | 58.10 | 87.77 | 80.89 | 59.08 | 79.32 | 90.56 | 86.33 | 91.47 | 88.10 |
| | Mean | **64.55** | 74.79 | 50.63 | 55.39 | 78.32 | 70.94 | 46.56 | 69.48 | 83.72 | 81.16 | 92.50 | 88.31 |
| | Weighted Mean | 62.19 | 71.77 | 50.60 | 54.86 | 80.42 | 68.24 | 46.53 | 67.95 | 83.10 | 80.12 | 91.79 | 88.33 |
| Mistral (mntp-unsup-simcse) | EOS | 56.89 | 85.30 | 52.79 | 64.78 | 84.88 | **89.53** | 63.80 | 79.68 | 94.45 | 75.24 | 86.40 | 82.07 |
| | Mean | 51.04 | 78.63 | 53.92 | 64.58 | 76.65 | 75.98 | 39.71 | 65.20 | 81.42 | 72.29 | 87.38 | 82.09 |
| | Weighted Mean | 50.97 | 76.25 | 54.13 | 64.70 | 84.21 | 73.95 | 40.38 | 63.85 | 80.40 | 69.78 | 85.89 | 81.87 |
| Mistral (mntp-supervised) | EOS | 53.07 | 88.51 | **56.46** | 71.12 | 83.70 | 81.25 | 60.27 | 87.57 | 95.57 | 88.31 | 94.46 | 91.52 |
| | Mean | 50.88 | 91.10 | 50.67 | **71.98** | 65.36 | **84.37** | 52.86 | 83.67 | 94.65 | 87.77 | 92.47 | 90.17 |
| | Weighted Mean | 50.75 | 89.24 | 51.07 | 71.17 | 70.84 | 83.99 | 52.75 | 81.09 | 94.41 | 87.15 | 91.68 | 88.42 |
| OpenAI-3-small | - | 55.61 | 91.10 | 51.69 | 66.50 | 57.88 | 74.08 | 58.52 | 91.20 | 94.34 | 88.10 | 93.08 | 91.64 |
| OpenAI-ada-002 | - | 61.51 | 88.66 | 51.60 | 65.49 | 77.54 | 78.12 | 60.28 | 86.05 | 92.45 | 87.58 | 93.05 | 92.54 |
| OpenAI-3-large | - | 51.84 | **92.10** | **54.78** | 67.60 | 65.41 | 77.55 | 63.31 | **92.52** | 94.15 | 88.03 | 93.43 | 91.58 |

Table 4 and Table 5 report average AUROC (%) and AUPRC (%) on all two-stage detectors across different embeddings and datasets, respectively. Depending on the results of Table 4 and Table 5, we have the following observations and analysis:

- Across all datasets, the top two performing results originate from the detectors utilizing LLM-derived embeddings. This indicates that LLM-based embeddings boost the performance for two-stage anomaly detection frameworks and exhibit marked advantages over traditional embedding methods such as GloVe and BERT for text anomaly detection tasks.

- No single LLM-derived embedding universally outperforms others across all datasets. This suggests that the optimal choice of embedding model may depend on specific dataset characteristics or task requirements.

- The results of AUPRC (%) in Table 5 reveals that the embedding models from OpenAI perform much better than other embedding models, which suggests that the models from OpenAI achieve both high precision and high recall, a critical advantage in the scenarios where false negatives are costly.

To observe the effect of different pooling strategies across all datasets, we compute the row-wise averages based on Table 4 to get the mean performance of each embedding-pooling combination across all datasets. We visualize these aggregated mean results for the LLaMA-2, LLaMA-3, and Mistral in Figure 2, enabling direct comparison for the overall effectiveness of different pooling strategies, where "EOS" exhibits significant advantages over "Mean" and "Weighted Mean" in most cases. In addition, we also notice that "Mean" and "Weighted Mean" have quite similar performance in

Table 5: Average AUPRC(%) of all two-stage methods. The top two results on each dataset are marked in **bold** (top, second).

| Embedding | Pooling | 20News. | Reuters. | IMDB | SST2 | SMS | Enron | WOS | DBp.0 | DBp.1 | DBp.2 | DBp.3 | DBp.4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 13.62 | 83.03 | 59.97 | 51.77 | 26.40 | 39.46 | 29.55 | 71.45 | 82.95 | 75.84 | 81.68 | 77.87 |
| BERT | CLS | 14.51 | 66.85 | 59.59 | 54.69 | 25.07 | 38.22 | 30.71 | 82.73 | 81.00 | 72.35 | 68.03 | 66.32 |
| | Mean | 14.05 | 83.73 | 58.99 | 54.50 | 31.27 | 37.28 | 32.26 | 78.99 | 80.30 | 79.26 | 77.77 | 74.87 |
| LLaMA-2 (mntp) | EOS | 10.99 | 75.29 | 61.57 | 54.04 | 38.91 | 39.29 | 29.66 | 56.95 | 64.90 | 58.91 | 67.42 | 63.20 |
| | Mean | 12.31 | 72.68 | 60.86 | 53.47 | 38.03 | 37.47 | 28.99 | 56.74 | 66.76 | 65.35 | 77.06 | 73.65 |
| | Weighted Mean | 12.15 | 72.37 | 61.02 | 53.25 | 43.37 | 37.19 | 28.65 | 56.02 | 66.11 | 63.95 | 75.18 | 71.15 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 14.69 | 93.86 | 60.98 | 62.97 | 57.42 | 64.46 | 38.03 | 90.79 | 97.47 | 86.86 | 96.33 | 92.88 |
| | Mean | 12.78 | 88.45 | 61.47 | 62.98 | 59.27 | 52.26 | 28.30 | 71.08 | 88.06 | 77.33 | 93.84 | 88.99 |
| | Weighted Mean | 13.01 | 87.46 | 61.76 | 62.86 | 68.98 | 51.36 | 28.27 | 70.44 | 88.16 | 74.64 | 92.81 | 87.80 |
| LLaMA-2 (mntp-supervised) | EOS | 16.44 | 93.34 | 63.17 | 66.15 | 77.74 | 64.81 | 45.06 | 84.99 | 96.63 | 83.42 | 94.39 | 93.76 |
| | Mean | 16.00 | 93.05 | 64.39 | 69.89 | 57.70 | 59.51 | 32.04 | 80.70 | 96.40 | 85.42 | 92.26 | 92.09 |
| | Weighted Mean | 16.67 | 92.28 | 64.26 | 67.85 | 66.43 | 59.85 | 32.63 | 79.16 | 95.58 | 84.22 | 91.06 | 91.29 |
| LLaMA-3 (mntp) | EOS | 18.43 | 85.76 | 59.26 | 57.55 | 78.56 | 56.45 | 35.15 | 67.60 | 82.73 | 82.22 | 92.94 | 86.25 |
| | Mean | 14.53 | 87.69 | 60.23 | 59.10 | 46.72 | 47.12 | 29.02 | 56.89 | 75.39 | 77.99 | 92.22 | 85.68 |
| | Weighted Mean | 13.97 | 85.18 | 60.37 | 59.30 | 56.77 | 45.10 | 28.88 | 57.53 | 76.42 | 76.86 | 91.41 | 85.93 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 13.28 | 82.85 | 60.27 | 57.92 | 77.25 | 66.06 | 36.67 | 57.35 | 71.57 | 58.27 | 68.15 | 65.09 |
| | Mean | 12.07 | 88.82 | 61.28 | 60.82 | 58.63 | 51.45 | 27.78 | 57.52 | 82.19 | 78.73 | 93.93 | 89.08 |
| | Weighted Mean | 11.53 | 87.49 | 61.26 | 60.21 | 70.42 | 50.34 | 27.96 | 58.80 | 82.46 | 76.84 | 92.98 | 88.22 |
| LLaMA-3 (mntp-supervised) | EOS | 18.38 | 93.08 | 64.30 | 66.75 | 81.73 | 63.40 | 35.55 | 86.77 | 96.76 | 77.69 | 88.92 | 90.01 |
| | Mean | 18.90 | 94.93 | 63.07 | 68.39 | 75.57 | 63.20 | 42.54 | 84.23 | 97.62 | 89.40 | 94.20 | 95.11 |
| | Weighted Mean | 17.24 | 94.09 | 63.04 | 67.86 | 74.50 | 63.32 | 41.61 | 84.61 | 97.37 | 88.10 | 93.43 | 94.43 |
| Mistral (mntp) | EOS | 23.38 | 84.77 | 57.89 | 57.71 | 72.76 | 65.53 | 38.72 | 73.90 | 86.92 | 84.52 | 90.61 | 85.90 |
| | Mean | 16.84 | 79.46 | 62.53 | 55.88 | 54.97 | 46.45 | 29.20 | 59.40 | 74.81 | 74.21 | 89.55 | 81.57 |
| | Weighted Mean | 16.08 | 77.65 | 62.55 | 55.52 | 58.02 | 44.84 | 29.25 | 57.82 | 74.88 | 72.74 | 88.81 | 81.92 |
| Mistral (mntp-unsup-simcse) | EOS | 15.65 | 89.32 | 63.84 | 62.75 | 68.35 | 78.10 | 41.85 | 74.75 | 91.78 | 72.03 | 84.83 | 79.31 |
| | Mean | 12.37 | 83.75 | 64.48 | 63.99 | 57.76 | 57.44 | 26.60 | 57.14 | 76.79 | 66.83 | 84.69 | 78.56 |
| | Weighted Mean | 12.32 | 82.29 | 64.60 | 63.76 | 69.51 | 56.85 | 26.66 | 56.55 | 76.26 | 63.28 | 82.76 | 77.12 |
| Mistral (mntp-supervised) | EOS | 14.20 | 92.42 | 67.12 | 68.78 | 67.82 | 67.41 | 38.98 | 84.36 | 94.50 | 86.62 | 93.57 | 90.07 |
| | Mean | 12.55 | 94.15 | 62.69 | 70.39 | 42.49 | 66.87 | 33.17 | 79.29 | 93.36 | 84.55 | 90.88 | 88.17 |
| | Weighted Mean | 12.53 | 92.84 | 62.93 | 69.54 | 48.78 | 66.84 | 33.05 | 77.46 | 93.25 | 83.42 | 89.67 | 86.33 |
| OpenAI-3-small | - | 19.04 | 95.35 | 65.29 | 67.93 | 38.38 | 52.64 | 40.09 | 92.24 | 96.27 | 89.35 | 94.94 | 92.76 |
| OpenAI-ada-002 | - | 21.83 | 93.08 | 64.43 | 66.24 | 54.38 | 59.43 | 42.37 | 83.67 | 92.01 | 88.47 | 94.99 | 93.16 |
| OpenAI-3-large | - | 17.58 | 96.10 | 66.96 | 69.12 | 42.00 | 56.88 | 45.75 | 91.64 | 93.02 | 83.55 | 93.88 | 91.17 |



Figure 2: The average AUROC(%) across all datasets. The "unsup." and "super." correspond to the embedding models with "mntp-unsup-simcse" and "mntp-supervised", respectively.

all cases. Figure 3 presents the average AUROC ranking across different LLM-derived embeddings. The visualization reveals that embeddings fine-tuned using the "mntp-supervised" approach (denoted as -3-) consistently achieve superior ranking positions compared to other embeddings. This finding aligns with the results reported in BehnamGhader et al. [2024b], where representations fine-tuned by "mntp-supervised" achieved state-of-the-art performance on the Massive Text Embedding Benchmark (MTEB).

Table 6 summarizes detection performance (AUROC%) of different AD methods across all datasets. Notably, CVDD, a specialized text anomaly detection method, operates directly on token-level embedding rather than aggregated

(a) Average Rank (AUROC) under EOS     (b) Average Rank (AUROC) under Mean     (c) Average Rank (AUROC) under Weighted Mean

Figure 3: Performance comparison of different embeddings from LLMs. Note that "llama2-1, llama3-1, mistral-1" denote those corresponding embedding models with "mntp", and "llama2-2, llama3-2, mistral-2" denote those corresponding embedding models with "mntp-unsup-simcse" and "llama2-3, llama3-3, mistral-3" denote those corresponding embedding models with "mntp-supervised". The plots (a), (b), (c) compare the average AUROC from all the embedding models under EOS, Mean, and Weighted Mean, respectively.



Figure 4: Performance comparison of different AD methods. In plot (a), a pair of methods with p-value > 0.05 indicates no statistically significant difference in their detection performance at the 95% confidence level.

sentence-level representations. Due to computational constraints associated with processing all token embeddings from LLMs, the experiments on CVDD were limited to conventional embeddings, specifically GloVe and BERT. Futhermore, DATE, another text-specific AD method, requires access to raw text data rather than precomputed embeddings. Figure 4 presents a comprehensive comparison of AD methods based on the results from Table 6 and two-stage AD methods only use the "best" results. In Figure 4, plot (a), a heatmap visualization of p-values derived from the Nemenyi post-hoc test, and plot (b), an average rank diagram where lower rank indicates better performance. Depending on the results of Table 6 and Figure 4, we have the following observations:

- Analyzing the "best" results of all two-stage AD methods, deep learning based detectors (AE, DSVDD, DPAD) exhibit no advantage over conventional "shallow" algorithms(OCSVM, IForest, LOF, KNN, KDE) when using LLM-derived embeddings. This finding suggests that (1) the high-quality representations from LLMs effectively encode textual nuances, enabling conventional algorithms to achieve competitive even better detection performance directly in the input space; and (2) the added complexity of deep anomaly detectors may be unnecessary when leveraging such high-quality text representation.

- When utilizing the embedding from GloVe and BERT, CVDD exhibits a slight advantage over two-stage AD methods including conventional shallow and deep learning based methods.

- No single method universally outperforms others across all datasets. Across all datasets, the average performance of KNN outperforms all others methods.

To further evaluate the performance of two-stage AD methods on different text embeddings, we visualize the heatmaps of AUROC(%) across all datasets in Figure 5. This visualization clearly demonstrates that LLM-derived embeddings

Table 6: Comparison among different AD methods. The top two results for each dataset are marked in **bold** (top, second). Note that "mean" denotes the average performance on all embeddings and "best" refers to the best performance across all embeddings.

| Methods | 20News. | Reuters. | IMDB | SST2 | SMS | Enron | WOS | DBp.0 | DBp.1 | DBp.2 | DBp.3 | DBp.4 | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| OCSVM (GloVe) | 55.88 | 76.85 | 46.35 | 46.38 | 48.02 | 63.12 | 52.49 | 82.33 | 91.59 | 85.27 | 87.97 | 86.64 | 68.57 |
| OCSVM (BERT) | 52.80 | 80.74 | 45.11 | 51.94 | 54.87 | 53.26 | 53.09 | 81.63 | 82.06 | 80.26 | 75.88 | 76.02 | 65.64 |
| OCSVM (mean) | 55.69 | 85.33 | 49.47 | 63.04 | 81.94 | 73.90 | 51.52 | 79.87 | 90.66 | 83.74 | 91.14 | 89.98 | 74.69 |
| OCSVM (best) | 72.74 | **99.22** | 62.97 | **85.22** | 94.66 | **99.24** | 69.13 | 98.70 | 99.87 | 95.83 | 99.46 | 98.89 | **89.66** |
| IForest (GloVe) | 56.09 | 78.21 | 46.76 | 44.69 | 33.43 | 62.45 | 38.61 | 77.03 | 88.37 | 83.18 | 82.11 | 82.12 | 64.42 |
| IForest (BERT) | 55.17 | 80.79 | 45.34 | 50.73 | 46.89 | 54.83 | 50.35 | 78.96 | 77.26 | 76.38 | 72.88 | 72.27 | 63.49 |
| IForest (mean) | 54.42 | 80.39 | 48.20 | 55.86 | 72.32 | 64.73 | 48.01 | 70.05 | 83.63 | 73.82 | 82.84 | 81.24 | 67.96 |
| IForest (best) | 71.87 | 94.08 | 54.68 | 66.90 | 92.83 | 75.01 | 62.69 | 84.70 | 95.89 | 84.84 | 93.48 | 89.15 | 80.51 |
| LOF (GloVe) | 60.16 | 84.54 | 49.31 | 58.73 | 44.65 | 67.86 | 64.93 | 86.25 | 96.57 | 87.29 | 97.15 | 92.22 | 74.14 |
| LOF (BERT) | 58.86 | 89.45 | 47.75 | 58.66 | 56.05 | 76.18 | 59.02 | 91.35 | 95.04 | 92.97 | 96.42 | 93.65 | 76.28 |
| LOF (mean) | 54.87 | 80.66 | 54.30 | 64.83 | 77.35 | 81.31 | 62.12 | 87.77 | 94.73 | 90.42 | 96.73 | 94.26 | 78.28 |
| LOF (best) | 70.79 | 93.86 | **64.66** | 71.96 | 93.62 | 94.52 | 73.27 | **99.30** | **99.89** | **97.91** | **99.81** | 99.22 | 88.23 |
| PCA (GloVe) | 56.36 | 77.77 | 46.51 | 44.79 | 34.60 | 62.64 | 39.30 | 78.19 | 90.11 | 83.39 | 83.69 | 83.09 | 65.04 |
| PCA (BERT) | 54.21 | 80.45 | 44.59 | 51.45 | 50.44 | 55.01 | 49.05 | 80.39 | 80.88 | 78.62 | 70.94 | 72.75 | 64.06 |
| PCA (mean) | 55.66 | 83.99 | 48.02 | 57.57 | 78.04 | 65.66 | 48.11 | 77.04 | 88.63 | 80.97 | 88.73 | 87.43 | 71.65 |
| PCA (best) | 73.68 | 99.13 | 56.58 | 69.76 | 95.01 | 79.23 | 67.61 | 98.50 | 99.83 | 94.18 | 99.02 | 98.42 | 85.91 |
| KNN (GloVe) | 59.90 | 82.75 | 46.74 | 48.68 | 35.43 | 73.63 | 33.74 | 79.83 | 89.22 | 83.86 | 92.02 | 85.92 | 67.64 |
| KNN (BERT) | 55.61 | 83.61 | 47.65 | 59.00 | 33.36 | 61.86 | 37.77 | 87.41 | 92.76 | 89.61 | 95.96 | 93.23 | 69.82 |
| KNN (mean) | 59.40 | 85.34 | 56.82 | 68.89 | 74.57 | 78.89 | 56.09 | 84.32 | 93.50 | 86.97 | 95.10 | 93.02 | 77.74 |
| KNN (best) | **74.56** | 96.58 | **73.44** | 80.11 | 94.61 | 92.78 | **80.91** | **99.19** | **99.90** | **97.54** | 99.64 | **99.24** | **90.71** |
| KDE (GloVe) | 55.99 | 75.96 | 46.37 | 45.23 | 41.38 | 63.57 | 47.93 | 81.04 | 90.89 | 84.03 | 85.17 | 84.82 | 66.86 |
| KDE (BERT) | 53.90 | 78.06 | 44.29 | 51.32 | 51.41 | 54.88 | 49.52 | 80.11 | 79.44 | 77.28 | 69.22 | 71.38 | 63.40 |
| KDE (mean) | 56.04 | 84.09 | 48.84 | 62.06 | 81.43 | 71.38 | 48.67 | 78.68 | 89.24 | 81.94 | 89.55 | 88.47 | 73.37 |
| KDE (best) | 73.18 | 99.04 | 61.74 | **84.67** | **95.38** | **99.11** | 68.70 | 98.59 | 99.84 | 94.03 | 98.99 | 98.40 | 89.31 |
| ECOD (GloVe) | 55.42 | 64.57 | 42.95 | 43.83 | 32.81 | 58.58 | 36.11 | 65.91 | 76.42 | 74.83 | 69.89 | 70.75 | 57.67 |
| ECOD (BERT) | 52.41 | 66.28 | 39.24 | 48.39 | 32.27 | 51.93 | 43.80 | 71.79 | 70.92 | 71.73 | 56.43 | 63.34 | 55.71 |
| ECOD (mean) | 53.05 | 71.35 | 40.10 | 51.79 | 59.89 | 61.85 | 43.24 | 64.61 | 79.13 | 72.60 | 79.69 | 79.91 | 63.10 |
| ECOD (best) | 71.97 | 93.74 | 48.10 | 57.75 | 89.95 | 75.51 | 63.07 | 91.73 | 99.57 | 87.13 | 96.97 | 95.86 | 80.95 |
| AE (GloVe) | 55.60 | 80.93 | 48.66 | 49.60 | 31.25 | 69.41 | 40.26 | 79.30 | 89.48 | 83.76 | 91.06 | 86.25 | 67.13 |
| AE (BERT) | 54.84 | 83.64 | 46.87 | 57.39 | 42.61 | 59.28 | 41.25 | 90.65 | 94.50 | 93.38 | 96.35 | 94.07 | 71.24 |
| AE (mean) | 57.28 | 87.62 | 53.86 | 71.12 | 75.08 | 87.45 | 56.29 | 85.17 | 94.24 | 88.99 | 95.77 | 93.68 | 78.88 |
| AE (best) | 71.51 | 98.75 | 59.17 | 80.95 | 94.37 | 98.72 | **75.89** | 98.32 | 99.75 | 97.14 | **99.71** | 99.15 | 89.45 |
| DSVDD (GloVe) | 49.77 | 57.92 | 47.72 | 53.47 | 64.28 | 56.91 | 56.32 | 56.61 | 68.34 | 57.31 | 66.49 | 61.66 | 58.07 |
| DSVDD (BERT) | 54.49 | 71.45 | 46.19 | 51.35 | 89.43 | 54.60 | 49.32 | 81.18 | 91.28 | 84.91 | 93.83 | 88.91 | 71.41 |
| DSVDD (mean) | 54.32 | 80.11 | 49.49 | 58.62 | 76.25 | 70.74 | 52.04 | 70.89 | 84.70 | 74.30 | 87.78 | 81.38 | 70.05 |
| DSVDD (best) | 68.32 | **99.12** | 54.06 | 67.80 | 92.98 | 85.81 | 59.92 | 98.27 | 99.87 | 95.31 | 99.69 | **99.26** | 85.03 |
| DPAD (GloVe) | 57.43 | 83.08 | 48.31 | 54.09 | 36.78 | 73.99 | 49.69 | 82.42 | 96.02 | 87.38 | 95.19 | 91.17 | 71.30 |
| DPAD (BERT) | 56.64 | 88.87 | 46.96 | 62.26 | 60.38 | 79.74 | 48.70 | 90.67 | 92.02 | 91.41 | 96.92 | 94.38 | 75.75 |
| DPAD (mean) | 56.36 | 82.95 | 51.32 | 67.58 | 76.57 | 86.76 | 54.23 | 78.27 | 88.81 | 81.59 | 90.13 | 86.90 | 75.12 |
| DPAD (best) | **73.41** | 94.56 | 57.02 | 78.36 | **95.02** | 97.43 | 66.74 | 93.53 | 99.25 | 93.59 | 99.62 | 98.61 | 87.26 |
| CVDD (GloVe) | 62.60 | 86.48 | 46.66 | 51.71 | 49.00 | 69.38 | 50.51 | 91.26 | 97.13 | 74.83 | 92.03 | 82.14 | 71.14 |
| CVDD (BERT) | 54.98 | 51.76 | 47.64 | 46.83 | 44.80 | 52.86 | 40.23 | 57.24 | 61.44 | 56.21 | 55.69 | 52.47 | 51.85 |
| DATE | 51.47 | - | 48.56 | 56.96 | 72.89 | 71.67 | 56.36 | 79.61 | 85.92 | 77.18 | 84.14 | 78.31 | 69.37 |

outperform conventional embedding techniques (GloVe and BERT) in most cases, with particularly pronounced improvements in sematic anomaly detection tasks.

# 6 Low-rank analysis and prediction

To systematically analyze the performance properties across datasets, we conduct singular value decomposition (SVD) on the performance matrices (AUROC on each dataset) (Table 8 to Table 19) where each performance matrix (or table) is a holistic evaluation of detection performance on one dataset across all text embeddings (rows) and AD algorithms (columns). Let $\mathbf{P} \in \mathbb{R}^{m \times n}$ be the performance matrix and $\sigma_i$ denotes the $i$-th singular value, where $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_n > 0$, $m$ denotes the number of embeddings, $n$ denotes the number of AD algorithms, and $n \leq m$. Figure 6 presents the cumulative contribution ratio of the singular values across all datasets, where the cumulative contribution ratio (*ccr*) is defined by $ccr(j) = (\sum_{i=1}^{j} \sigma_i)/(\sum_{i=1}^{n} \sigma_i)$.

Figure 5: The heatmap of AUROC(%) of each two-stage AD method with different text embeddings. Note that "llama2-1, llama3-1, mistral-1" denote those corresponding "mntp" models, "llama2-2, llama3-2, mistral-2" denote those corresponding "mntp-unsup-simcse" models and "llama2-3, llama3-3, mistral-3" denote those corresponding "mntp-supervised" models.

As evidenced by Figure 6, the cumulative contribution ratio of the first two singular values ($ccr(2)$) exceeds 0.90 on most datasets. This demonstrates that the AUROC matrices possess strongly low-rank characteristics. This finding has two important implications: (1) the detection performance of novel text datasets or anomaly detection methods can be reliably predicted using only a subset of performance measurements, and (2) this property enables an efficient strategy for rapid model evaluation (embedding evaluation) and selection in practical applications.

To further verify the low-rank characteristics of the AUROC matrices (Table 8 to Table 19), we construct the matrix completion task using the performance matrices. Specifically, for every AUROC matrix, we construct matrix completion Candes and Recht [2012], Fan et al. [2019] tasks by MCAR (missing completely at random) mechanism with missing rate $\in \{0.5, 0.6, 0.7\}$. We use matrix factorization and non-convex optimization techniques Chi [2018] to recover the missed entries of the performance matrices. Formally, we consider a rank-constrained least-squares problem:

$$\min_{\mathbf{\Phi} \in \mathbb{R}^{m \times n}} \|\mathcal{P}_\Omega(\mathbf{\Phi} - \mathbf{P})\|_F^2, \quad \text{s.t. } \text{rank}(\mathbf{\Phi}) \leq r, \tag{3}$$

where $\mathbf{P} \in \mathbb{R}^{m \times n}$ denotes the performance matrix with $m$ rows (text embeddings) and $n$ columns (AD algorithms), $\Omega$ consists of the locations of observed entries and the observation operator $\mathcal{P}_\Omega : \mathbb{R}^{m \times n} \to \mathbb{R}^{m \times n}$ as

$$[\mathcal{P}_\Omega(\mathbf{P})]_{ij} = \begin{cases} \mathbf{P}_{ij}, & (i,j) \in \Omega \\ 0, & \text{otherwise} \end{cases}. \tag{4}$$

Figure 6: The cumulation contribution ratio of singular value of detection performance (AUROC) across all datasets.

Invoking the low-rank factorization $\mathbf{\Phi} = \mathbf{U}\mathbf{V}^T$, where $\mathbf{U} \in \mathbb{R}^{m \times r}$ and $\mathbf{V} \in \mathbb{R}^{n \times r}$, we can rewrite (3) as an unconstrained optimization problem:

$$\min_{\mathbf{U},\mathbf{V}} \|\mathcal{P}_\Omega(\mathbf{U}\mathbf{V}^T - \mathbf{P})\|_F^2. \tag{5}$$

In our experiments, we set $r = 1$ and utilize the singular value decomposition technique to initialize the $\mathbf{U} = \mathbf{U}_0\mathbf{\Sigma}_0^{1/2}$ and $\mathbf{V} = \mathbf{V}_0\mathbf{\Sigma}_0^{1/2}$, where $\mathbf{U}_0\mathbf{\Sigma}_0\mathbf{V}_0^T$ is the best rank-r approximation of $\mathcal{P}_\Omega(\mathbf{P})$. The normalized Mean Absolute Percentage Error (MAPE), as defined in Equation (6), is employed to quantify the accuracy of the prediction.

$$\text{MAPE} = \frac{1}{k} \sum_{i=1}^{k} \left| \frac{p_i - \hat{p}_i}{p_i} \right|, \tag{6}$$

where $p_i, \hat{p}_i \in \mathbb{R}$ denote the missing (or unknown) values of detection performance and the corresponding recovered (or predicted) values, respectively. The $k$ denotes the number of unknown values, and the missing rate is mr $= k/mn$.

The recovery results are reported in Table 7, where most of the recovery errors are below 0.1. The results verify the feasibility of a rapid model evaluation (embedding evaluation) and selection in practical applications based on the empirical evaluations in this benchmark. Take the 20Newsgroups dataset as an example, we use 50% entries of the AUROC performance matrix (see Table 8) to predict other values and show the result of six AD methods on six embeddings in Figure 7. We can see that the predictions are close to the real performance values. Thant means, a rapid and relatively accurate prediction of the performance of the new pairs (text embedding, AD method) can be obtained based on some observed results from this benchmark, which guides an efficient and reliable model evaluation (embedding evaluation) and selection.

## 7 Conclusion

In this work, we present Text-ADBench, a comprehensive benchmark for text anomaly detection, which both considers shallow and deep AD approaches, constructs abundant two-stage text AD methods by concatenating LLMs with different pooling strategies and estimates the detection accuracy across eight datasets. Experimental results reveal that *(i)* LLM-based embeddings boost the detection performance for these two-stage AD methods and "EOS" pooling

**Observed matrix** (Observed / Unknown):

| | OCSVM | IForest | KNN | AE | DSVDD | DPAD |
|---|---|---|---|---|---|---|
| GloVe | 55.88 | 56.09 | 59.90 | 55.60 | 49.77 | Unknown |
| BERT | 53.33 | Unknown | Unknown | 59.47 | 47.29 | Unknown |
| LLaMA-2 | Unknown | 54.45 | 57.85 | 61.05 | 58.22 | Unknown |
| LLaMA-3 | Unknown | Unknown | Unknown | 66.69 | 55.05 | 62.81 |
| Mistral | Unknown | Unknown | Unknown | Unknown | Unknown | 73.41 |
| OpenAI-Large | Unknown | 49.55 | 62.51 | Unknown | Unknown | Unknown |

**Real Performance:**

| OCSVM | IForest | KNN | AE | DSVDD | DPAD |
|---|---|---|---|---|---|
| 55.88 | 56.09 | 59.90 | 55.60 | 49.77 | **57.43** |
| 53.33 | **55.26** | **56.63** | 59.47 | 47.29 | **56.73** |
| **58.76** | 54.45 | 57.85 | 61.05 | 58.22 | **59.04** |
| **62.89** | **60.91** | **65.48** | 66.69 | 55.05 | 62.81 |
| **72.74** | **71.87** | **72.40** | **71.51** | **68.32** | 73.41 |
| **52.96** | 49.55 | 62.51 | **56.33** | **51.18** | **50.00** |

**Predicted Performance:**

| OCSVM | IForest | KNN | AE | DSVDD | DPAD |
|---|---|---|---|---|---|
| 55.88 | 56.09 | 59.90 | 55.60 | 49.77 | **56.76** |
| 53.33 | **50.82** | **57.19** | 59.47 | 47.29 | **55.00** |
| **57.75** | 54.45 | 57.85 | 61.05 | 58.22 | **59.15** |
| **61.35** | **58.06** | **65.34** | 66.69 | 55.05 | 62.81 |
| **71.55** | **67.72** | **76.21** | **76.11** | **65.95** | 73.41 |
| **55.83** | 49.55 | 62.51 | **59.41** | **51.47** | **57.21** |

Figure 7: The visualization of AUROC performance prediction on the 20Newsgroups dataset, where the missing rate is 0.5. The prediction process is completed in 1.88 seconds, while the actual performance evaluation requires 5099.50 seconds.

Table 7: Recovery performance (MAPE) across all datasets in the setting of MCAR. Note that 'mr' denotes missing rate.

| | 20News. | IMDB | Enton | Reuters. | DBp.0 | DBp.1 | DBp.2 | DBp.3 | DBp.4 | SMS-SPAM | SST2 | WOS |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| mr=0.5 | 0.0348 | 0.0254 | 0.0419 | 0.0380 | 0.0357 | 0.0375 | 0.0301 | 0.0279 | 0.0286 | 0.0698 | 0.0336 | 0.0570 |
| mr=0.6 | 0.0431 | 0.0460 | 0.0624 | 0.0404 | 0.0433 | 0.0452 | 0.0438 | 0.0344 | 0.0433 | 0.0828 | 0.0482 | 0.0749 |
| mr=0.7 | 0.0746 | 0.0535 | 0.1233 | 0.0684 | 0.0787 | 0.0795 | 0.0661 | 0.0589 | 0.0753 | 0.1534 | 0.0771 | 0.1282 |

exhibits significant advantages over "Mean" and "Weighted Mean" in most cases, and *(ii)* deep learning based detectors (AE, DSVDD, DPAD) exhibit no advantage over conventional "shallow" algorithms (OCSVM, IForest, KNN, LOF, KDE) when using LLM-derived embeddings, and *(iii)* the average performance of KNN outperforms all other methods, although no single method universally outperforms other across all datasets. Furthermore, we analyze the performance matrices across datasets and find strongly low-rank characteristics, which enables an efficient strategy for rapid model evaluation (embedding evaluation) and selection in practical applications by using the historical estimation from Text-ADBench.

# References

Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM Comput. Surv.*, 41(3), July 2009. ISSN 0360-0300. doi:10.1145/1541880.1541882. URL https://doi.org/10.1145/1541880.1541882.

Charu C Aggarwal. An introduction to outlier analysis. In *Outlier analysis*, pages 1–34. Springer, 2016.

Lukas Ruff, Jacob R Kauffmann, Robert A Vandermeulen, Grégoire Montavon, Wojciech Samek, Marius Kloft, Thomas G Dietterich, and Klaus-Robert Müller. A unifying review of deep and shallow anomaly detection. *Proceedings of the IEEE*, 109(5):756–795, 2021.

Igor Kononenko. Machine learning for medical diagnosis: history, state of the art and perspective. *Artificial Intelligence in medicine*, 23(1):89–109, 2001.

Filippo Amato, Alberto López, Eladia María Peña-Méndez, Petr Vaňhara, Aleš Hampl, and Josef Havel. Artificial neural networks in medical diagnosis, 2013.

Mihalj Bakator and Dragica Radosav. Deep learning and medical diagnosis: A review of literature. *Multimodal Technologies and Interaction*, 2(3):47, 2018.

Jonathan G Richens, Ciarán M Lee, and Saurabh Johri. Improving the accuracy of medical diagnosis with causal machine learning. *Nature communications*, 11(1):3923, 2020.

Biswanath Mukherjee, L Todd Heberlein, and Karl N Levitt. Network intrusion detection. *IEEE network*, 8(3):26–41, 2002.

Chih-Fong Tsai, Yu-Feng Hsu, Chia-Ying Lin, and Wei-Yang Lin. Intrusion detection by machine learning: A review. *expert systems with applications*, 36(10):11994–12000, 2009.

Hung-Jen Liao, Chun-Hung Richard Lin, Ying-Chih Lin, and Kuang-Yuan Tung. Intrusion detection system: A comprehensive review. *Journal of network and computer applications*, 36(1):16–24, 2013.

Ansam Khraisat, Iqbal Gondal, Peter Vamplew, and Joarder Kamruzzaman. Survey of intrusion detection systems: techniques, datasets and challenges. *Cybersecurity*, 2(1):1–22, 2019.

Richard J Bolton and David J Hand. Statistical fraud detection: A review. *Statistical science*, 17(3):235–255, 2002.

Aisha Abdallah, Mohd Aizaini Maarof, and Anazida Zainal. Fraud detection system: A survey. *Journal of Network and Computer Applications*, 68:90–113, 2016.

Soroor Motie and Bijan Raahemi. Financial fraud detection using graph neural networks: A systematic review. *Expert Systems with Applications*, 240:122156, 2024.

Rolf Isermann. Model-based fault-detection and diagnosis–status and applications. *Annual Reviews in control*, 29(1): 71–85, 2005.

Jicong Fan, Wei Wang, and Haijun Zhang. Autoencoder based high-dimensional data fault detection system. In *2017 ieee 15th international conference on industrial informatics (indin)*, pages 1001–1006. IEEE, 2017.

Jicong Fan, Tommy W. S. Chow, and S. Joe Qin. Kernel-based statistical process monitoring and fault detection in the presence of missing data. *IEEE Transactions on Industrial Informatics*, 18(7):4477–4487, 2022.

Zahra Zamanzadeh Darban, Geoffrey I. Webb, Shirui Pan, Charu Aggarwal, and Mahsa Salehi. Deep learning for time series anomaly detection: A survey. *ACM Comput. Surv.*, 57(1), October 2024. ISSN 0360-0300. doi:10.1145/3691338. URL https://doi.org/10.1145/3691338.

Tharindu Fernando, Harshala Gammulle, Simon Denman, Sridha Sridharan, and Clinton Fookes. Deep learning for medical anomaly detection – a survey. *ACM Comput. Surv.*, 54(7), July 2021. ISSN 0360-0300. doi:10.1145/3464423. URL https://doi.org/10.1145/3464423.

Jing Liu, Yang Liu, Jieyu Lin, Jielin Li, Liang Cao, Peng Sun, Bo Hu, Liang Song, Azzedine Boukerche, and Victor C.M. Leung. Networking systems for video anomaly detection: A tutorial and survey. *ACM Comput. Surv.*, 57(10), May 2025. ISSN 0360-0300. doi:10.1145/3729222. URL https://doi.org/10.1145/3729222.

Zellig S Harris. Distributional structure. *Word*, 10(2-3):146–162, 1954.

Karen Sparck Jones. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28(1):11–21, 1972.

Akiko Aizawa. An information-theoretic perspective of tf–idf measures. *Information Processing & Management*, 39(1): 45–65, 2003.

David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.

Susan T Dumais. Latent semantic analysis. *Annual Review of Information Science and Technology (ARIST)*, 38: 189–230, 2004.

Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans, editors, *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, October 2014a. Association for Computational Linguistics. doi:10.3115/v1/D14-1162. URL https://aclanthology.org/D14-1162/.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*, 2016.

Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations, 2018. URL https://arxiv.org/abs/1802.05365.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 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, 2019.

Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.

Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language models. *arXiv preprint arXiv:2307.06435*, 2023.

Rajvardhan Patil and Venkat Gudivada. A review of current trends, techniques, and challenges in large language models (llms). *Applied Sciences*, 14(5):2074, 2024.

Bernhard Schölkopf, Robert C Williamson, Alex Smola, John Shawe-Taylor, and John Platt. Support vector method for novelty detection. *Advances in neural information processing systems*, 12, 1999.

Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. Lof: identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pages 93–104, 2000.

Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In *2008 eighth ieee international conference on data mining*, pages 413–422. IEEE, 2008.

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

Dazhi Fu, Zhao Zhang, and Jicong Fan. Dense projection for anomaly detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 8398–8408, 2024.

Lukas Ruff, Yury Zemlyanskiy, Robert Vandermeulen, Thomas Schnake, and Marius Kloft. Self-attentive, multi-context one-class classification for unsupervised anomaly detection on text. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4061–4071, 2019.

Andrei Manolache, Florin Brad, and Elena Burceanu. Date: Detecting anomalies in text via self-supervision of transformers. *arXiv preprint arXiv:2104.05591*, 2021.

Yuangang Li, Jiaqi Li, Zhuo Xiao, Tiankai Yang, Yi Nian, Xiyang Hu, and Yue Zhao. Nlp-adbench: Nlp anomaly detection benchmark. *arXiv preprint arXiv:2412.04784*, 2024a.

Yang Cao, Sikun Yang, Chen Li, Haolong Xiang, Lianyong Qi, Bo Liu, Rongsheng Li, and Ming Liu. Tad-bench: A comprehensive benchmark for embedding-based text anomaly detection. *arXiv preprint arXiv:2501.11960*, 2025.

Lucas Deecke, Robert Vandermeulen, Lukas Ruff, Stephan Mandt, and Marius Kloft. Image anomaly detection with generative adversarial networks. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2018, Dublin, Ireland, September 10–14, 2018, Proceedings, Part I 18*, pages 3–17. Springer, 2019.

Chen Qiu, Timo Pfrommer, Marius Kloft, Stephan Mandt, and Maja Rudolph. Neural transformation learning for deep anomaly detection beyond images. In *International conference on machine learning*, pages 8703–8714. PMLR, 2021.

Feng Xiao, Jianfeng Zhou, Kunpeng Han, Haoyuan Hu, and Jicong Fan. Unsupervised anomaly detection using inverse generative adversarial networks. *Information Sciences*, 689:121435, 2025a.

Wei Dai and Jicong Fan. AutoUAD: Hyper-parameter optimization for unsupervised anomaly detection. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=ErQPdaD5wJ.

Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. Deep learning for anomaly detection: A review. *ACM computing surveys (CSUR)*, 54(2):1–38, 2021.

Zheng Li, Yue Zhao, Xiyang Hu, Nicola Botta, Cezar Ionescu, and George H Chen. Ecod: Unsupervised outlier detection using empirical cumulative distribution functions. *IEEE Transactions on Knowledge and Data Engineering*, 35(12):12181–12193, 2022.

Sridhar Ramaswamy, Rajeev Rastogi, and Kyuseok Shim. 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, 2000.

Thomas Schlegl, Philipp Seeböck, Sebastian M Waldstein, Ursula Schmidt-Erfurth, and Georg Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *International conference on information processing in medical imaging*, pages 146–157. Springer, 2017.

Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Daeki Cho, and Haifeng Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *International conference on learning representations*, 2018.

Guansong Pang, Chunhua Shen, and Anton Van Den Hengel. Deep anomaly detection with deviation networks. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 353–362, 2019.

Pramuditha Perera, Ramesh Nallapati, and Bing Xiang. Ocgan: One-class novelty detection using gans with constrained latent representations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2898–2906, 2019.

Sachin Goyal, Aditi Raghunathan, Moksh Jain, Harsha Vardhan Simhadri, and Prateek Jain. Drocc: Deep robust one-class classification. In *International conference on machine learning*, pages 3711–3721. PMLR, 2020.

Jinyu Cai and Jicong Fan. Perturbation learning based anomaly detection. *Advances in Neural Information Processing Systems*, 35:14317–14330, 2022.

Hongzuo Xu, Guansong Pang, Yijie Wang, and Yongjun Wang. Deep isolation forest for anomaly detection. *IEEE Transactions on Knowledge and Data Engineering*, 35(12):12591–12604, 2023a.

Hangting Ye, He Zhao, Wei Fan, Mingyuan Zhou, Dan dan Guo, and Yi Chang. Drl: Decomposed representation learning for tabular anomaly detection. In *The Thirteenth International Conference on Learning Representations*, 2025.

Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006.

Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. *arXiv preprint arXiv:1812.04606*, 2018.

Guansong Pang, Longbing Cao, Ling Chen, and Huan Liu. Learning representations of ultrahigh-dimensional data for random distance-based outlier detection. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 2041–2050, 2018.

Samet Akcay, Amir Atapour-Abarghouei, and Toby P Breckon. Ganomaly: Semi-supervised anomaly detection via adversarial training. In *Computer Vision–ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part III 14*, pages 622–637. Springer, 2019.

Lukas Ruff, Robert A Vandermeulen, Nico Görnitz, Alexander Binder, Emmanuel Müller, Klaus-Robert Müller, and Marius Kloft. Deep semi-supervised anomaly detection. In *Proceedings of the International Conference on Learning Representations*, 2020.

Yingjie Zhou, Xucheng Song, Yanru Zhang, Fanxing Liu, Ce Zhu, and Lingqiao Liu. Feature encoding with autoencoders for weakly supervised anomaly detection. *IEEE Transactions on Neural Networks and Learning Systems*, 33(6): 2454–2465, 2021.

Guansong Pang, Chunhua Shen, Huidong Jin, and Anton van den Hengel. Deep weakly-supervised anomaly detection. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1795–1807, 2023.

Feng Xiao, Youqing Wang, S Joe Qin, and Jicong Fan. Semi-supervised anomaly detection using restricted distribution transformation. *IEEE Transactions on Neural Networks and Learning Systems*, 2025b.

Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014b.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Improving text embeddings with large language models. *arXiv preprint arXiv:2401.00368*, 2023.

Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados, and Siva Reddy. Llm2vec: Large language models are secretly powerful text encoders. *arXiv preprint arXiv:2404.05961*, 2024a.

Yizhou Xu, Jérôme Milleret, and Frédérique Segond. Comparative analysis of anomaly detection algorithms in text data. In *Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing*, pages 1234–1245, 2023b.

Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.

Anindya Sundar Das, Aravind Ajay, Sriparna Saha, and Monowar Bhuyan. Few-shot anomaly detection in text with deviation learning. In *International Conference on Neural Information Processing*, pages 425–438. Springer, 2023.

Jeremie Pantin and Christophe Marsala. A robust autoencoder ensemble-based approach for anomaly detection in text. *arXiv preprint arXiv:2405.13031*, 2024.

Jing Su, Chufeng Jiang, Xin Jin, Yuxin Qiao, Tingsong Xiao, Hongda Ma, Rong Wei, Zhi Jing, Jiajun Xu, and Junhong Lin. Large language models for forecasting and anomaly detection: A systematic literature review. *arXiv preprint arXiv:2402.10350*, 2024.

Tiankai Yang, Yi Nian, Shawn Li, Ruiyao Xu, Yuangang Li, Jiaqi Li, Zhuo Xiao, Xiyang Hu, Ryan Rossi, Kaize Ding, et al. Ad-llm: Benchmarking large language models for anomaly detection. *arXiv preprint arXiv:2412.11142*, 2024.

Manqing Dong, Hao Huang, and Longbing Cao. Can llms serve as time series anomaly detectors? *arXiv preprint arXiv:2408.03475*, 2024.

Zhaopeng Gu, Bingke Zhu, Guibo Zhu, Yingying Chen, Ming Tang, and Jinqiao Wang. Anomalygpt: Detecting industrial anomalies using large vision-language models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pages 1932–1940, 2024.

Aodong Li, Yunhan Zhao, Chen Qiu, Marius Kloft, Padhraic Smyth, Maja Rudolph, and Stephan Mandt. Anomaly detection of tabular data using llms. *arXiv preprint arXiv:2406.16308*, 2024b.

Matei Bejan, Andrei Manolache, and Marius Popescu. 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, 2023.

Ken Lang. Newsweeder: Learning to filter netnews. In *Proceedings of the Twelfth International Conference on Machine Learning*, pages 331–339, 1995.

David D. Lewis. Reuters-21578 text categorization collection data set, 1997. URL `https://kdd.ics.uci.edu/databases/reuters21578/reuters21578.html`.

Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL `http://www.aclweb.org/anthology/P11-1015`.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics. URL `https://www.aclweb.org/anthology/D13-1170`.

Tiago A. Almeida, Jose Maria Gomez Hidalgo, and Akebo Yamakami. Contributions to the study of sms spam filtering: New collection and results. In *Proceedings of the 2011 ACM Symposium on Document Engineering (DOCENG'11)*, 2011.

Kamran Kowsari, Donald E Brown, Mojtaba Heidarysafa, Kiana Jafari Meimandi, , Matthew S Gerber, and Laura E Barnes. Hdltex: Hierarchical deep learning for text classification. In *Machine Learning and Applications (ICMLA), 2017 16th IEEE International Conference on*. IEEE, 2017.

Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL `https://proceedings.neurips.cc/paper_files/paper/2015/file/250cf8b51c773f3f8dc8b4be867a9a02-Paper.pdf`.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL `https://arxiv.org/abs/2310.06825`.

AI@Meta. Llama 3 model card. 2024. URL `https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md`.

Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados, and Siva Reddy. LLM2Vec: Large language models are secretly powerful text encoders. In *First Conference on Language Modeling*, 2024b. URL `https://openreview.net/forum?id=IW1PR7vEBf`.

Mei-Ling Shyu, Shu-Ching Chen, Kanoksri Sarinnapakorn, and LiWu Chang. A novel anomaly detection scheme based on principal component classifier. In *Proceedings of the IEEE foundations and new directions of data mining workshop*, pages 172–179. IEEE Press Piscataway, NJ, USA, 2003.

JooSeuk Kim and Clayton D Scott. Robust kernel density estimation. *The Journal of Machine Learning Research*, 13 (1):2529–2565, 2012.

Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*, 2020.

Emmanuel Candes and Benjamin Recht. Exact matrix completion via convex optimization. *Communications of the ACM*, 55(6):111–119, 2012.

Jicong Fan, Lijun Ding, Yudong Chen, and Madeleine Udell. Factor group-sparse regularization for efficient low-rank matrix recovery. *Advances in neural information processing Systems*, 32, 2019.

Yuejie Chi. Low-rank matrix completion [lecture notes]. *IEEE Signal Processing Magazine*, 35(5):178–181, 2018.

Table 8: Average AUROC(%) on 20Newsgroups. The best results are marked in **bold**.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 55.88 | 56.09 | 60.16 | 56.36 | 59.90 | 55.99 | 55.42 | 55.60 | 49.77 | 57.43 | 56.26 |
| BERT | CLS | 53.33 | 55.26 | 64.69 | 55.50 | 56.63 | 54.03 | 54.34 | 59.47 | 47.29 | 56.73 | 55.73 |
| | Mean | 52.80 | 55.17 | 58.86 | 54.21 | 55.61 | 53.90 | 52.41 | 54.84 | 54.49 | 56.64 | 54.89 |
| LLaMA-2 (mntp) | EOS | 50.34 | 48.51 | 44.65 | 47.76 | 49.60 | 47.07 | 47.84 | 46.56 | 45.08 | 46.57 | 47.40 |
| | Mean | 57.73 | 51.73 | 51.06 | 53.31 | 49.92 | 54.96 | 52.76 | 49.34 | 49.99 | 48.66 | 51.95 |
| | Weighted Mean | 57.49 | 51.05 | 50.27 | 52.93 | 49.36 | 54.23 | 52.41 | 47.51 | 48.49 | 47.79 | 51.15 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 62.87 | 61.58 | 58.53 | 62.84 | 58.16 | 62.38 | 58.74 | 55.69 | 55.60 | 56.63 | 59.30 |
| | Mean | 48.36 | 49.17 | 49.70 | 50.27 | 59.76 | 51.34 | 48.28 | 54.71 | 54.54 | 56.44 | 52.26 |
| | Weighted Mean | 49.07 | 53.43 | 48.44 | 51.19 | 59.91 | 52.20 | 49.22 | 55.03 | 55.02 | 56.74 | 53.03 |
| LLaMA-2 (mntp-supervised) | EOS | 58.76 | 54.45 | 59.81 | 56.97 | 57.85 | 57.09 | 52.09 | 61.05 | 58.22 | 59.04 | 57.53 |
| | Mean | 52.37 | 50.39 | 56.44 | 54.85 | 68.25 | 55.07 | 51.98 | 64.36 | 56.01 | 58.41 | 56.81 |
| | Weighted Mean | 52.81 | 52.91 | 56.70 | 54.90 | 68.61 | 55.04 | 51.96 | 64.62 | 55.60 | 58.11 | 57.13 |
| LLaMA-3 (mntp) | EOS | 64.87 | 61.85 | 62.98 | 64.39 | 61.12 | 65.39 | 61.66 | 60.73 | 59.64 | 63.07 | 62.57 |
| | Mean | 56.78 | 54.21 | 54.59 | 56.13 | 59.78 | 58.09 | 53.46 | 56.84 | 58.48 | 59.89 | 56.83 |
| | Weighted Mean | 55.01 | 52.59 | 52.73 | 54.18 | 56.27 | 56.51 | 51.75 | 53.52 | 57.07 | 56.29 | 54.59 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 56.19 | 54.69 | 53.46 | 55.46 | 53.86 | 54.12 | 54.23 | 53.16 | 54.98 | 53.31 | 54.35 |
| | Mean | 45.40 | 48.45 | 46.00 | 45.76 | 52.27 | 46.84 | 43.63 | 49.65 | 54.34 | 54.42 | 48.68 |
| | Weighted Mean | 43.75 | 40.98 | 45.34 | 43.63 | 49.93 | 44.73 | 41.58 | 47.93 | 54.45 | 53.31 | 46.56 |
| LLaMA-3 (mntp-supervised) | EOS | 62.89 | 60.91 | 64.37 | 62.32 | 65.48 | 62.76 | 58.42 | 66.69 | 55.05 | 62.81 | 62.17 |
| | Mean | 56.77 | 57.46 | 64.29 | 57.19 | 68.83 | 57.62 | 54.25 | 66.26 | 56.08 | 61.88 | 60.06 |
| | Weighted Mean | 54.98 | 52.65 | 62.84 | 55.14 | 67.38 | 55.57 | 52.30 | 64.81 | 53.43 | 59.95 | 57.91 |
| Mistral (mntp) | EOS | **72.74** | **71.87** | **70.79** | **73.68** | 72.40 | **73.18** | **71.97** | **71.51** | **68.32** | **73.41** | **71.99** |
| | Mean | 66.07 | 66.54 | 59.40 | 66.17 | 67.03 | 69.00 | 63.94 | 62.30 | 61.46 | 63.56 | 64.55 |
| | Weighted Mean | 64.60 | 63.53 | 57.32 | 64.29 | 64.03 | 67.08 | 62.18 | 59.96 | 59.78 | 59.18 | 62.20 |
| Mistral (mntp-unsup-simcse) | EOS | 50.49 | 60.63 | 49.75 | 72.13 | 47.56 | 51.66 | 67.61 | 51.58 | 56.42 | 61.03 | 56.89 |
| | Mean | 50.49 | 49.25 | 48.71 | 48.63 | 54.34 | 50.22 | 47.54 | 52.94 | 54.12 | 54.19 | 51.04 |
| | Weighted Mean | 50.68 | 49.99 | 49.06 | 48.68 | 53.77 | 50.29 | 47.57 | 52.65 | 53.06 | 53.91 | 50.97 |
| Mistral (mntp-supervised) | EOS | 53.66 | 52.78 | 50.09 | 53.35 | 52.03 | 53.61 | 47.48 | 57.82 | 52.10 | 57.78 | 53.07 |
| | Mean | 55.51 | 46.15 | 46.68 | 47.31 | 59.08 | 53.51 | 44.75 | 55.35 | 48.74 | 51.67 | 50.88 |
| | Weighted Mean | 55.18 | 48.75 | 45.39 | 46.87 | 58.30 | 53.17 | 44.34 | 54.91 | 49.47 | 51.07 | 50.75 |
| OpenAI-3-small | - | 56.21 | 54.65 | 52.08 | 57.52 | 66.17 | 57.85 | 51.91 | 59.08 | 50.59 | 50.00 | 55.61 |
| OpenAI-ada-002 | - | 60.79 | 58.61 | 58.62 | 65.32 | **74.56** | 64.94 | 60.85 | 67.56 | 53.84 | 50.00 | 61.51 |
| OpenAI-3-large | - | 52.96 | 49.55 | 56.77 | 47.58 | 62.51 | 49.73 | 41.77 | 56.33 | 51.18 | 50.00 | 51.84 |

Table 9: Average AUROC(%) on Enron.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 63.12 | 62.45 | 67.86 | 62.64 | 73.63 | 63.57 | 58.58 | 69.41 | 56.91 | 73.99 | 65.22 |
| BERT | CLS | 53.74 | 53.72 | 76.81 | 54.28 | 68.02 | 53.49 | 52.78 | 63.68 | 44.82 | 79.74 | 60.11 |
| | Mean | 53.26 | 54.83 | 76.18 | 55.01 | 61.86 | 54.88 | 51.93 | 59.28 | 54.60 | 76.96 | 59.88 |
| LLaMA-2 (mntp) | EOS | 49.45 | 51.32 | 63.05 | 50.81 | 62.33 | 51.07 | 49.14 | 62.76 | 50.53 | 86.86 | 57.73 |
| | Mean | 53.25 | 57.01 | 75.37 | 54.82 | 65.48 | 53.28 | 52.89 | 63.47 | 52.50 | 75.49 | 60.36 |
| | Weighted Mean | 52.72 | 55.04 | 73.45 | 53.73 | 62.56 | 52.70 | 52.25 | 61.04 | 53.65 | 74.64 | 59.18 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 79.37 | 70.29 | 93.97 | 73.96 | 90.00 | 73.96 | 69.56 | **98.72** | 82.24 | 96.87 | 82.89 |
| | Mean | 66.36 | 64.54 | 75.49 | 65.92 | 79.29 | 66.77 | 60.60 | 96.50 | 78.66 | 93.51 | 73.96 |
| | Weighted Mean | 62.76 | 62.23 | 73.83 | 62.02 | 77.20 | 62.67 | 56.72 | 96.07 | 77.46 | 93.92 | 72.29 |
| LLaMA-2 (mntp-supervised) | EOS | 77.73 | 73.46 | 90.54 | 74.64 | 89.13 | 73.96 | 71.92 | 98.34 | 80.39 | 97.27 | 82.63 |
| | Mean | 72.64 | 71.80 | 89.83 | 72.30 | 84.22 | 71.79 | 69.74 | 98.10 | 78.61 | 96.26 | 80.53 |
| | Weighted Mean | 72.18 | 70.57 | 89.84 | 71.67 | 85.30 | 71.04 | 68.78 | 98.27 | 76.99 | 96.43 | 80.15 |
| LLaMA-3 (mntp) | EOS | 65.56 | 64.01 | 74.72 | 64.44 | 84.38 | 65.51 | 60.31 | 94.98 | 72.05 | 94.72 | 75.67 |
| | Mean | 62.76 | 62.63 | 76.23 | 63.19 | 75.43 | 63.50 | 59.38 | 87.32 | 66.59 | 88.66 | 69.63 |
| | Weighted Mean | 58.53 | 59.51 | 72.36 | 58.66 | 71.95 | 58.71 | 55.18 | 85.99 | 63.28 | 87.29 | 67.34 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 83.39 | 68.32 | 88.18 | 68.99 | 90.19 | 73.90 | 64.50 | 97.69 | 80.84 | 95.88 | 81.19 |
| | Mean | 64.98 | 62.56 | 75.92 | 63.12 | 77.73 | 64.95 | 56.47 | 94.95 | 76.24 | 92.65 | 73.86 |
| | Weighted Mean | 61.01 | 58.90 | 74.23 | 58.55 | 74.75 | 59.90 | 52.26 | 94.54 | 74.84 | 92.29 | 70.13 |
| LLaMA-3 (mntp-supervised) | EOS | 76.73 | 71.47 | 87.94 | 72.70 | 80.99 | 72.37 | 70.04 | 95.00 | 78.12 | 94.70 | 79.91 |
| | Mean | 76.73 | 71.13 | 87.94 | 72.70 | 80.99 | 72.37 | 70.04 | 95.00 | 76.83 | 94.94 | 79.62 |
| | Weighted Mean | 76.73 | 71.66 | 87.94 | 72.70 | 80.99 | 72.37 | 70.04 | 95.00 | 77.57 | 94.83 | 79.89 |
| Mistral (mntp) | EOS | 98.59 | 64.01 | 79.51 | 66.56 | 84.17 | 98.50 | 60.53 | 91.67 | 71.18 | 94.22 | 80.89 |
| | Mean | 68.90 | 63,21 | 76.57 | 63.60 | 74.84 | 66.84 | 59.02 | 83.59 | 67.56 | 85.28 | 70.94 |
| | Weighted Mean | 66.16 | 59.70 | 74.99 | 60.55 | 71.16 | 62.31 | 57.39 | 81.00 | 64.72 | 84.42 | 69.34 |
| Mistral (mntp-unsup-simcse) | EOS | **99.24** | **75.01** | 94.14 | **79.23** | **92.78** | **99.11** | 74.83 | 97.76 | **85.81** | **97.43** | **89.53** |
| | Mean | 97.55 | 57.20 | 76.67 | 57.84 | 77.89 | 78.38 | 51.44 | 94.24 | 74.39 | 94.25 | 75.75 |
| | Weighted Mean | 97.59 | 53.20 | 75.89 | 53.52 | 75.20 | 76.75 | 48.21 | 93.33 | 72.42 | 93.38 | 73.67 |
| Mistral (mntp-supervised) | EOS | 98.47 | 63.79 | 83.80 | 65.75 | 78.17 | 98.38 | 61.87 | 89.47 | 78.20 | 94.57 | 81.26 |
| | Mean | 98.39 | 68.29 | 89.60 | 71.18 | 83.33 | 95.51 | 67.86 | 96.65 | 77.04 | 95.88 | 84.77 |
| | Weighted Mean | 98.43 | 68.67 | 88.62 | 69.93 | 83.02 | 96.88 | 66.40 | 96.45 | 75.92 | 95.60 | 84.14 |
| OpenAI-3-small | - | 74.64 | 73.79 | 86.20 | 74.50 | 86.77 | 74.01 | 70.52 | 85.30 | 65.05 | 50.00 | 73.98 |
| OpenAI-ada-002 | - | 78.75 | 76.39 | **94.52** | 78.84 | 90.82 | 78.13 | **75.51** | 84.88 | 73.35 | 50.00 | 78.03 |
| OpenAI-3-large | - | 79.01 | 75.35 | 90.90 | 78.50 | 88.76 | 77.98 | 74.41 | 85.51 | 75.09 | 50.00 | 77.55 |

Table 10: Average AUROC(%) on IMDB.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 46.35 | 46.76 | 49.31 | 46.51 | 46.74 | 46.37 | 42.95 | 48.66 | 47.72 | 48.31 | 46.97 |
| BERT | CLS | 44.33 | 45.77 | 45.66 | 44.46 | 46.31 | 44.40 | 40.95 | 45.12 | 48.91 | 46.76 | 45.27 |
| | Mean | 45.11 | 45.34 | 47.75 | 44.59 | 47.65 | 44.29 | 39.24 | 46.87 | 46.19 | 46.96 | 45.40 |
| LLaMA-2 (mntp) | EOS | 48.36 | 48.33 | 50.26 | 48.33 | 49.32 | 48.33 | **48.10** | 49.20 | 48.90 | 49.83 | 48.90 |
| | Mean | 46.58 | 47.22 | 49.44 | 46.70 | 49.39 | 47.03 | 46.07 | 49.45 | 48.34 | 50.11 | 48.03 |
| | Weighted Mean | 46.78 | 47.37 | 50.07 | 46.81 | 49.53 | 47.15 | 46.22 | 49.41 | 48.55 | 50.16 | 48.21 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 46.81 | 45.85 | 54.24 | 45.19 | 54.25 | 45.31 | 36.69 | 55.03 | 47.54 | 50.96 | 48.19 |
| | Mean | 47.02 | 47.93 | 54.04 | 47.11 | 56.24 | 47.05 | 40.53 | 55.06 | 48.73 | 50.39 | 49.51 |
| | Weighted Mean | 47.95 | 47.23 | 54.95 | 48.11 | 55.54 | 48.04 | 41.90 | 55.10 | 49.29 | 51.17 | 49.93 |
| LLaMA-2 (mntp-supervised) | EOS | 52.96 | 49.81 | 61.47 | 50.41 | 57.19 | 50.04 | 38.27 | 54.25 | 50.49 | 52.53 | 51.74 |
| | Mean | 50.70 | 50.66 | 57.43 | 49.84 | 70.00 | 49.87 | 36.92 | 62.78 | 51.28 | 54.36 | 53.38 |
| | Weighted Mean | 50.62 | 49.70 | 57.68 | 49.71 | 69.58 | 49.76 | 36.86 | 62.01 | 50.82 | 54.53 | 53.13 |
| LLaMA-3 (mntp) | EOS | 44.69 | 43.88 | 48.10 | 43.96 | 46.39 | 44.07 | 40.66 | 48.87 | 46.64 | 48.05 | 45.53 |
| | Mean | 45.31 | 45.56 | 50.65 | 45.51 | 50.82 | 45.57 | 40.46 | 51.42 | 48.06 | 50.65 | 47.40 |
| | Weighted Mean | 45.82 | 45.68 | 51.25 | 46.07 | 50.43 | 46.07 | 41.47 | 51.50 | 48.34 | 50.59 | 47.72 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 47.15 | 47.01 | 51.15 | 46.52 | 48.09 | 46.23 | 42.05 | 48.92 | 48.02 | 48.45 | 47.36 |
| | Mean | 47.21 | 46.62 | 53.39 | 47.43 | 57.03 | 47.62 | 39.07 | 54.22 | 49.16 | 50.80 | 49.26 |
| | Weighted Mean | 47.57 | 48.13 | 53.43 | 47.81 | 55.85 | 47.99 | 39.76 | 53.70 | 49.40 | 50.61 | 49.43 |
| LLaMA-3 (mntp-supervised) | EOS | 53.22 | 52.45 | 58.45 | 52.58 | 59.94 | 52.09 | 39.02 | 55.60 | 53.36 | 53.12 | 52.98 |
| | Mean | 45.68 | 46.96 | 55.31 | 45.91 | 68.50 | 45.94 | 35.48 | 60.64 | 49.38 | 52.50 | 50.63 |
| | Weighted Mean | 45.83 | 46.36 | 55.18 | 46.07 | 68.54 | 46.10 | 35.20 | 60.44 | 49.03 | 53.02 | 50.58 |
| Mistral (mntp) | EOS | 43.55 | 42.30 | 44.95 | 42.96 | 43.61 | 43.44 | 39.85 | 44.56 | 46.37 | 45.07 | 43.67 |
| | Mean | 49.30 | 49.24 | 53.55 | 49.12 | 53.51 | 49.59 | 44.37 | 54.63 | 49.89 | 53.13 | 50.63 |
| | Weighted Mean | 49.54 | 48.71 | 53.65 | 49.36 | 53.00 | 49.78 | 44.95 | 54.36 | 50.07 | 52.57 | 50.60 |
| Mistral (mntp-unsup-simcse) | EOS | 60.00 | 49.21 | 56.65 | 48.48 | 58.41 | 57.97 | 37.61 | 55.18 | 51.41 | 53.01 | 52.79 |
| | Mean | 53.44 | 51.29 | 58.19 | 51.71 | 59.78 | 52.94 | 46.47 | 58.36 | 52.54 | 54.46 | 53.92 |
| | Weighted Mean | 53.75 | 52.36 | 58.93 | 52.08 | 58.87 | 53.21 | 46.92 | 57.98 | 52.53 | 54.63 | 54.13 |
| Mistral (mntp-supervised) | EOS | **62.97** | **54.68** | 62.49 | **56.58** | 59.78 | **61.74** | 39.73 | 55.53 | **54.06** | **57.02** | **56.46** |
| | Mean | 52.80 | 47.19 | 54.34 | 44.50 | 63.85 | 49.12 | 34.83 | 57.37 | 49.56 | 53.15 | 50.67 |
| | Weighted Mean | 53.44 | 48.84 | 54.47 | 45.03 | 63.61 | 49.85 | 35.07 | 57.01 | 50.21 | 53.14 | 51.07 |
| OpenAI-3-small | - | 51.72 | 49.98 | 62.25 | 50.70 | 63.05 | 50.55 | 35.90 | 53.48 | 49.23 | 50.00 | 51.69 |
| OpenAI-ada-002 | - | 49.98 | 50.40 | 58.47 | 49.49 | 66.97 | 49.26 | 38.19 | 51.47 | 48.14 | 53.66 | 51.60 |
| OpenAI-3-large | - | 56.02 | 51.94 | **64.66** | 55.10 | **73.44** | 54.99 | 31.40 | **59.17** | 51.13 | 50.00 | 54.79 |

Table 11: Average AUROC(%) on Reuters21578.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 76.85 | 78.21 | 84.54 | 77.77 | 82.75 | 75.96 | 64.57 | 80.93 | 57.92 | 83.08 | 76.26 |
| BERT | CLS | 65.90 | 58.08 | 76.69 | 58.43 | 52.54 | 60.90 | 49.07 | 63.65 | 54.05 | 61.18 | 60.05 |
| | Mean | 80.74 | 80.79 | 89.45 | 80.45 | 83.61 | 78.06 | 66.28 | 83.64 | 71.45 | 88.87 | 80.33 |
| LLaMA-2 (mntp) | EOS | 62.17 | 64.19 | 59.43 | 65.07 | 67.87 | 63.69 | 59.71 | 69.37 | 63.61 | 73.05 | 64.82 |
| | Mean | 56.63 | 64.05 | 61.33 | 60.15 | 65.36 | 55.25 | 51.87 | 72.80 | 57.44 | 66.75 | 61.16 |
| | Weighted Mean | 56.80 | 64.30 | 61.29 | 60.16 | 65.15 | 56.23 | 53.06 | 67.89 | 58.62 | 67.48 | 61.10 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 95.38 | 91.17 | 86.94 | 95.37 | 95.48 | 95.03 | 87.49 | 96.42 | 87.06 | **94.56** | **92.49** |
| | Mean | 88.38 | 85.75 | 82.98 | 86.88 | 88.28 | 87.17 | 69.20 | 91.01 | 93.85 | 93.73 | 86.72 |
| | Weighted Mean | 86.85 | 82.27 | 80.99 | 85.55 | 87.56 | 86.01 | 67.09 | 90.01 | 93.00 | 93.40 | 85.27 |
| LLaMA-2 (mntp-supervised) | EOS | 94.91 | 87.13 | 78.69 | 94.63 | 95.00 | 94.57 | 89.01 | 95.48 | 81.75 | 92.90 | 90.41 |
| | Mean | 93.01 | 86.39 | 80.14 | 93.74 | 93.45 | 93.97 | 80.80 | 95.82 | 87.85 | 93.22 | 89.84 |
| | Weighted Mean | 91.79 | 84.43 | 78.29 | 92.80 | 92.61 | 93.17 | 78.40 | 94.96 | 86.85 | 92.19 | 88.55 |
| LLaMA-3 (mntp) | EOS | 85.68 | 81.94 | 83.20 | 85.15 | 91.29 | 84.94 | 69.60 | 89.73 | 65.39 | 87.03 | 82.40 |
| | Mean | 87.78 | 87.23 | 88.00 | 88.29 | 90.27 | 85.52 | 70.45 | 91.42 | 71.72 | 90.90 | 85.16 |
| | Weighted Mean | 83.23 | 84.13 | 86.00 | 83.58 | 88.53 | 79.73 | 65.95 | 89.88 | 68.13 | 89.07 | 81.82 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 82.44 | 75.30 | 75.92 | 80.25 | 81.38 | 80.89 | 71.27 | 81.72 | 72.18 | 78.12 | 77.95 |
| | Mean | 88.81 | 86.12 | 82.93 | 87.32 | 86.69 | 87.51 | 72.88 | 89.19 | 92.31 | 91.53 | 86.53 |
| | Weighted Mean | 86.65 | 84.48 | 81.57 | 84.96 | 85.76 | 85.77 | 68.92 | 87.91 | 89.82 | 89.89 | 84.57 |
| LLaMA-3 (mntp-supervised) | EOS | 92.05 | 85.86 | 90.48 | 91.02 | 87.55 | 90.04 | 83.79 | 91.49 | 88.21 | 89.91 | 89.04 |
| | Mean | 95.15 | 89.31 | 85.14 | 95.43 | 94.43 | 95.37 | 86.48 | 96.14 | 91.00 | 93.68 | 92.21 |
| | Weighted Mean | 94.20 | 87.61 | 83.18 | 94.64 | 93.21 | 94.61 | 85.08 | 95.44 | 88.59 | 92.97 | 90.95 |
| Mistral (mntp) | EOS | 76.36 | 82.49 | 72.82 | 86.65 | 74.52 | 76.81 | 68.80 | 77.98 | 68.81 | 84.37 | 76.96 |
| | Mean | 83.02 | 73.65 | 85.07 | 71.73 | 78.16 | 74.15 | 51.18 | 80.69 | 70.46 | 79.78 | 74.79 |
| | Weighted Mean | 79.11 | 67.31 | 84.11 | 67.05 | 75.84 | 68.29 | 51.52 | 79.00 | 68.48 | 77.02 | 71.77 |
| Mistral (mntp-unsup-simcse) | EOS | 91.31 | 77.16 | 78.51 | 88.63 | 91.46 | 91.46 | 76.21 | 91.32 | 79.13 | 87.85 | 85.30 |
| | Mean | 82.74 | 69.52 | 77.99 | 72.18 | 82.54 | 82.56 | 52.60 | 85.73 | 90.93 | 89.49 | 78.63 |
| | Weighted Mean | 80.27 | 67.28 | 76.98 | 68.18 | 80.85 | 80.14 | 51.16 | 84.49 | 85.46 | 87.65 | 76.25 |
| Mistral (mntp-supervised) | EOS | 92.27 | 85.17 | 75.34 | 92.81 | 91.92 | 92.50 | 83.30 | 93.17 | 87.50 | 91.08 | 88.51 |
| | Mean | 96.03 | 85.75 | 83.32 | 94.96 | 94.16 | 96.27 | 81.86 | 96.07 | 88.45 | 94.15 | 91.10 |
| | Weighted Mean | 94.63 | 80.60 | 80.77 | 93.60 | 92.50 | 94.99 | 78.34 | 94.93 | 89.51 | 92.52 | 89.24 |
| OpenAI-3-small | - | 98.28 | **94.08** | 89.56 | 98.19 | 95.15 | 97.98 | 91.06 | **98.75** | 97.96 | 50.00 | 91.10 |
| OpenAI-ada-002 | - | 97.36 | 89.57 | 86.28 | 96.76 | 93.65 | 96.40 | 83.91 | 95.60 | 97.05 | 50.00 | 88.66 |
| OpenAI-3-large | - | **99.22** | 91.51 | **93.86** | **99.13** | **96.58** | **99.04** | **93.74** | 98.81 | **99.12** | 50.00 | 92.10 |

Table 12: Average AUROC(%) on SMS-SPAM.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 48.02 | 33.43 | 44.65 | 34.60 | 35.43 | 41.38 | 32.81 | 31.25 | 64.28 | 36.78 | 40.26 |
| BERT | CLS | 61.94 | 33.87 | 43.46 | 37.87 | 24.09 | 53.32 | 30.65 | 32.28 | 64.43 | 36.44 | 41.84 |
| | Mean | 54.87 | 46.89 | 56.05 | 50.44 | 33.36 | 51.41 | 32.27 | 42.61 | 89.43 | 60.38 | 51.77 |
| LLaMA-2 (mntp) | EOS | 77.78 | 72.10 | 69.49 | 74.03 | 37.45 | 73.46 | 53.77 | 53.72 | 46.29 | 63.70 | 62.18 |
| | Mean | 83.29 | 78.76 | 57.48 | 78.45 | 24.03 | 79.76 | 55.17 | 23.46 | 49.16 | 31.14 | 56.07 |
| | Weighted Mean | 86.29 | 82.84 | 73.59 | 82.85 | 24.68 | 84.22 | 59.15 | 30.42 | 42.61 | 38.52 | 60.52 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 84.47 | 73.75 | 85.16 | 84.62 | 82.78 | 83.84 | 63.55 | 85.56 | 70.46 | 83.84 | 79.80 |
| | Mean | 81.30 | 69.30 | 78.38 | 81.35 | 85.65 | 83.29 | 51.53 | 85.92 | 79.59 | 90.06 | 78.64 |
| | Weighted Mean | 90.15 | 77.71 | 82.54 | 90.03 | 89.54 | 91.15 | 67.09 | 90.14 | 77.15 | 92.77 | 84.83 |
| LLaMA-2 (mntp-supervised) | EOS | **94.66** | **92.83** | 79.26 | **95.01** | **94.61** | **95.38** | **89.95** | 94.04 | 83.25 | **95.02** | **91.40** |
| | Mean | 81.31 | 75.32 | 71.73 | 83.60 | 86.17 | 84.04 | 63.34 | 88.26 | 68.52 | 87.06 | 78.94 |
| | Weighted Mean | 88.52 | 82.91 | 73.14 | 90.55 | 91.40 | 90.92 | 75.79 | 92.04 | 70.33 | 91.06 | 84.67 |
| LLaMA-3 (mntp) | EOS | 92.22 | 89.65 | 91.47 | 91.78 | 93.96 | 92.14 | 85.26 | 92.63 | 86.98 | 92.53 | 90.86 |
| | Mean | 76.28 | 69.36 | 71.61 | 71.36 | 76.02 | 80.55 | 50.00 | 67.66 | 85.48 | 75.86 | 72.42 |
| | Weighted Mean | 85.66 | 78.23 | 74.60 | 81.86 | 84.55 | 88.40 | 62.24 | 76.60 | 85.39 | 84.51 | 80.20 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 91.73 | 89.18 | 88.12 | 91.52 | 92.30 | 91.61 | 87.39 | 93.21 | 82.83 | 91.93 | 89.98 |
| | Mean | 78.70 | 70.04 | 79.80 | 78.66 | 83.08 | 81.52 | 54.77 | 80.07 | 90.05 | 89.63 | 78.63 |
| | Weighted Mean | 89.12 | 82.52 | 79.30 | 88.98 | 90.44 | 90.84 | 71.52 | 88.28 | **92.98** | 93.87 | 86.79 |
| LLaMA-3 (mntp-supervised) | EOS | 92.98 | 89.27 | 88.63 | 92.30 | 92.43 | 92.57 | 86.64 | **94.37** | 83.90 | 93.86 | 90.70 |
| | Mean | 88.38 | 83.92 | 83.94 | 88.12 | 90.71 | 88.51 | 76.90 | 91.60 | 82.33 | 92.01 | 86.64 |
| | Weighted Mean | 88.49 | 82.51 | 80.21 | 88.33 | 91.05 | 88.73 | 77.01 | 91.87 | 81.95 | 92.64 | 86.28 |
| Mistral (mntp) | EOS | 91.19 | 82.91 | **93.62** | 85.86 | 90.31 | 91.28 | 77.19 | 90.52 | 84.68 | 90.14 | 87.77 |
| | Mean | 87.73 | 76.24 | 89.00 | 80.09 | 77.92 | 85.22 | 59.60 | 69.45 | 79.48 | 78.45 | 78.32 |
| | Weighted Mean | 87.05 | 78.91 | 89.90 | 81.99 | 76.55 | 84.70 | 63.97 | 72.57 | 84.97 | 83.58 | 80.42 |
| Mistral (mntp-unsup-simcse) | EOS | 90.57 | 74.40 | 91.78 | 86.37 | 89.20 | 90.86 | 72.87 | 90.42 | 74.88 | 87.41 | 84.88 |
| | Mean | 86.56 | 62.12 | 77.74 | 69.28 | 80.45 | 85.71 | 44.16 | 82.17 | 90.03 | 88.29 | 76.65 |
| | Weighted Mean | 93.31 | 74.03 | 81.37 | 83.82 | 88.78 | 92.96 | 59.97 | 89.13 | 90.83 | 87.94 | 84.21 |
| Mistral (mntp-supervised) | EOS | 88.10 | 75.49 | 80.83 | 85.18 | 85.82 | 88.77 | 66.90 | 87.21 | 88.33 | 90.42 | 83.71 |
| | Mean | 67.30 | 56.11 | 68.55 | 60.83 | 65.23 | 65.35 | 38.87 | 73.20 | 77.94 | 80.23 | 65.36 |
| | Weighted Mean | 72.84 | 69.16 | 68.22 | 71.76 | 69.69 | 71.85 | 47.77 | 76.31 | 76.61 | 84.21 | 70.84 |
| OpenAI-3-small | - | 61.90 | 51.99 | 83.73 | 58.73 | 70.73 | 61.47 | 24.67 | 60.85 | 54.70 | 50.00 | 57.88 |
| OpenAI-ada-002 | - | 87.52 | 72.89 | 89.00 | 87.78 | 84.46 | 88.38 | 66.59 | 87.85 | 63.97 | 46.97 | 77.54 |
| OpenAI-3-large | - | 73.78 | 58.06 | 86.32 | 67.39 | 77.96 | 73.45 | 26.93 | 72.12 | 72.48 | 45.65 | 65.41 |

Table 13: Average AUROC(%) on SST2.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 46.38 | 44.69 | 58.73 | 44.79 | 48.68 | 45.23 | 43.83 | 49.60 | 53.47 | 54.09 | 48.95 |
| BERT | CLS | 52.93 | 51.93 | 58.17 | 51.82 | 55.93 | 52.40 | 49.58 | 55.04 | 51.98 | 60.13 | 53.99 |
|  | Mean | 51.94 | 50.73 | 58.66 | 51.45 | 59.00 | 51.32 | 48.39 | 57.39 | 51.35 | 62.26 | 54.25 |
| LLaMA-2 (mntp) | EOS | 49.80 | 50.25 | 53.66 | 49.80 | 57.31 | 49.37 | 48.97 | 58.31 | 50.41 | 64.76 | 53.26 |
|  | Mean | 49.30 | 48.99 | 54.87 | 48.44 | 55.45 | 48.39 | 47.83 | 56.74 | 50.29 | 62.60 | 52.29 |
|  | Weighted Mean | 48.75 | 48.94 | 54.63 | 48.18 | 55.68 | 48.04 | 47.49 | 56.41 | 49.73 | 63.92 | 52.18 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 59.54 | 55.85 | 69.28 | 60.01 | 75.26 | 60.30 | 52.98 | 79.73 | 59.35 | 73.34 | 64.56 |
|  | Mean | 56.08 | 54.24 | 69.40 | 56.10 | 74.60 | 56.22 | 50.87 | 80.42 | 57.77 | 73.43 | 62.91 |
|  | Weighted Mean | 56.33 | 55.41 | 68.69 | 56.75 | 74.28 | 57.03 | 51.18 | 80.11 | 58.95 | 72.38 | 63.11 |
| LLaMA-2 (mntp-supervised) | EOS | 68.26 | 65.75 | 69.25 | 67.43 | 73.06 | 66.84 | 59.87 | 75.91 | 65.34 | 74.43 | 68.61 |
|  | Mean | 69.42 | **66.90** | **71.96** | **69.76** | 79.30 | 69.87 | 59.56 | **80.95** | **67.80** | 78.00 | 71.35 |
|  | Weighted Mean | 67.63 | 65.33 | 69.90 | 67.63 | 76.72 | 67.58 | 58.33 | 79.01 | 65.51 | 76.34 | 69.40 |
| LLaMA-3 (mntp) | EOS | 53.44 | 53.01 | 58.52 | 53.01 | 61.97 | 53.29 | 50.71 | 66.14 | 54.67 | 64.20 | 56.90 |
|  | Mean | 53.34 | 52.50 | 61.56 | 52.86 | 65.73 | 53.18 | 49.74 | 71.71 | 56.56 | 68.17 | 58.54 |
|  | Weighted Mean | 54.01 | 53.22 | 60.36 | 53.98 | 66.88 | 54.55 | 50.38 | 72.34 | 56.09 | 69.29 | 59.11 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 55.67 | 54.47 | 61.22 | 55.28 | 61.80 | 55.53 | 53.08 | 66.40 | 56.37 | 64.89 | 58.47 |
|  | Mean | 54.18 | 52.57 | 64.17 | 53.40 | 71.01 | 54.25 | 48.57 | 78.33 | 57.43 | 71.68 | 60.56 |
|  | Weighted Mean | 54.26 | 53.52 | 63.34 | 54.02 | 68.19 | 54.59 | 50.35 | 76.12 | 57.32 | 69.94 | 60.17 |
| LLaMA-3 (mntp-supervised) | EOS | 67.37 | 64.56 | 66.66 | 65.65 | 72.67 | 65.23 | 57.07 | 75.38 | 63.72 | 73.93 | 67.22 |
|  | Mean | 68.40 | 65.00 | 69.25 | 68.05 | 77.93 | 68.46 | 56.30 | 79.29 | 65.78 | 77.17 | 69.56 |
|  | Weighted Mean | 67.93 | 63.80 | 68.28 | 67.27 | 76.67 | 67.53 | 56.26 | 78.47 | 65.22 | 76.93 | 68.84 |
| Mistral (mntp) | EOS | 68.11 | 49.06 | 59.58 | 49.75 | 60.00 | 67.15 | 46.48 | 63.11 | 54.85 | 62.90 | 58.10 |
|  | Mean | 57.58 | 49.75 | 61.80 | 49.23 | 59.01 | 50.04 | 47.23 | 64.01 | 53.35 | 61.93 | 55.39 |
|  | Weighted Mean | 58.85 | 48.24 | 60.76 | 47.80 | 58.34 | 50.29 | 45.61 | 63.02 | 53.43 | 62.23 | 54.86 |
| Mistral (mntp-unsup-simcse) | EOS | 78.33 | 55.41 | 64.16 | 55.57 | 68.95 | 77.53 | 48.89 | 71.74 | 58.47 | 68.78 | 64.78 |
|  | Mean | 76.40 | 51.94 | 65.86 | 52.15 | 72.60 | 73.37 | 47.95 | 77.75 | 56.58 | 71.19 | 64.58 |
|  | Weighted Mean | 76.84 | 51.12 | 65.90 | 52.85 | 72.66 | 74.02 | 47.98 | 76.92 | 57.85 | 70.89 | 64.70 |
| Mistral (mntp-supervised) | EOS | 82.86 | 62.85 | 67.00 | 65.44 | 75.32 | 82.57 | **57.75** | 75.92 | 66.17 | 75.35 | 71.12 |
|  | Mean | **85.22** | 60.37 | 70.55 | 63.38 | 79.46 | **84.67** | 53.42 | 79.77 | 64.61 | **78.36** | **71.98** |
|  | Weighted Mean | 84.58 | 58.70 | 69.98 | 62.82 | 78.02 | 83.96 | 53.60 | 78.47 | 65.03 | 76.54 | 71.17 |
| OpenAI-3-small | - | 69.22 | 61.05 | 74.36 | 68.21 | **80.11** | 68.29 | 55.22 | 76.10 | 62.47 | 50.00 | 66.50 |
| OpenAI-ada-002 | - | 66.86 | 61.98 | 72.88 | 67.42 | 79.33 | 67.43 | 56.94 | 69.61 | 62.43 | 50.00 | 65.49 |
| OpenAI-3-large | - | 70.49 | 61.21 | 76.14 | 69.56 | 81.38 | 69.48 | 56.76 | 76.90 | 64.13 | 50.00 | 67.61 |

Table 14: Average AUROC(%) on WOS.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 52.49 | 38.61 | 64.93 | 39.30 | 33.74 | 47.93 | 36.11 | 40.26 | 56.32 | 49.69 | 45.94 |
| BERT | CLS | 47.15 | 47.21 | 58.02 | 45.24 | 46.23 | 45.26 | 42.84 | 46.73 | 47.94 | 49.57 | 47.62 |
|  | Mean | 53.09 | 50.35 | 59.02 | 49.05 | 37.77 | 49.52 | 43.80 | 41.25 | 49.32 | 48.70 | 48.19 |
| LLaMA-2 (mntp) | EOS | 47.46 | 46.17 | 49.46 | 46.92 | 50.92 | 45.99 | 46.68 | 51.51 | 48.03 | 51.74 | 48.49 |
|  | Mean | 51.40 | 48.20 | 49.87 | 47.89 | 44.30 | 49.80 | 45.72 | 41.10 | 46.66 | 46.57 | 47.15 |
|  | Weighted Mean | 50.81 | 46.20 | 49.01 | 47.56 | 43.84 | 49.00 | 45.80 | 40.44 | 47.12 | 46.34 | 46.61 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 59.93 | 55.46 | 60.92 | 59.98 | 63.82 | 60.33 | 52.78 | 62.55 | 56.07 | 60.37 | 59.22 |
|  | Mean | 43.21 | 40.96 | 63.18 | 37.96 | 44.43 | 37.21 | 32.79 | 47.23 | 48.75 | 52.38 | 44.81 |
|  | Weighted Mean | 43.55 | 40.89 | 61.41 | 38.99 | 43.46 | 38.17 | 34.00 | 46.54 | 49.30 | 52.43 | 44.87 |
| LLaMA-2 (mntp-supervised) | EOS | 63.75 | **62.69** | 63.58 | **67.61** | 80.63 | **68.70** | **63.07** | 75.20 | 59.79 | **66.74** | **67.18** |
|  | Mean | 42.33 | 44.13 | 71.32 | 38.35 | 62.02 | 37.53 | 33.12 | 61.07 | 51.96 | 57.71 | 49.95 |
|  | Weighted Mean | 43.48 | 45.23 | 71.22 | 40.30 | 62.89 | 39.56 | 34.94 | 61.45 | 52.09 | 58.34 | 50.95 |
| LLaMA-3 (mntp) | EOS | 57.53 | 56.04 | 58.65 | 56.68 | 56.21 | 55.15 | 51.04 | 56.67 | 54.18 | 55.83 | 55.80 |
|  | Mean | 46.21 | 40.74 | 64.35 | 39.34 | 41.85 | 38.43 | 32.59 | 51.39 | 51.10 | 52.91 | 45.89 |
|  | Weighted Mean | 46.47 | 43.61 | 61.57 | 40.70 | 41.70 | 39.55 | 34.43 | 51.18 | 51.50 | 52.25 | 46.30 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 55.43 | 54.32 | 57.35 | 54.56 | 59.02 | 55.75 | 52.04 | 58.71 | 54.23 | 56.24 | 55.77 |
|  | Mean | 39.86 | 35.45 | 67.20 | 32.38 | 41.82 | 31.03 | 28.28 | 47.85 | 48.64 | 51.26 | 42.38 |
|  | Weighted Mean | 41.18 | 38.75 | 64.85 | 34.76 | 41.69 | 33.31 | 30.64 | 47.64 | 48.72 | 51.93 | 43.35 |
| LLaMA-3 (mntp-supervised) | EOS | 52.83 | 50.45 | 63.67 | 50.43 | 69.22 | 51.45 | 46.91 | 66.72 | 52.78 | 57.94 | 56.24 |
|  | Mean | 57.58 | 56.96 | 68.03 | 58.56 | **80.91** | 60.32 | 52.28 | **75.89** | **59.92** | 65.48 | 63.59 |
|  | Weighted Mean | 57.09 | 53.65 | 67.86 | 57.88 | 80.45 | 59.60 | 51.82 | 75.11 | 58.97 | 65.12 | 62.76 |
| Mistral (mntp) | EOS | 62.76 | 56.11 | 56.25 | 58.24 | 63.09 | 62.29 | 55.07 | 62.98 | 54.97 | 59.00 | 59.08 |
|  | Mean | 45.85 | 43.45 | 58.23 | 45.42 | 43.69 | 39.01 | 39.31 | 50.99 | 48.41 | 51.22 | 46.56 |
|  | Weighted Mean | 46.08 | 47.59 | 56.45 | 45.30 | 42.62 | 39.85 | 39.57 | 50.21 | 47.35 | 50.24 | 46.53 |
| Mistral (mntp-unsup-simcse) | EOS | **69.13** | 61.76 | 55.32 | 63.37 | 70.04 | 68.34 | 61.41 | 68.10 | 57.69 | 62.82 | 63.80 |
|  | Mean | 35.26 | 32.69 | 66.85 | 29.61 | 39.89 | 25.44 | 27.66 | 46.58 | 45.95 | 47.14 | 39.71 |
|  | Weighted Mean | 35.84 | 35.69 | 65.07 | 31.09 | 39.67 | 26.82 | 29.11 | 46.36 | 46.95 | 47.20 | 40.38 |
| Mistral (mntp-supervised) | EOS | 67.36 | 51.96 | 63.23 | 56.02 | 67.46 | 66.76 | 50.75 | 65.73 | 54.28 | 59.12 | 60.27 |
|  | Mean | 53.03 | 44.95 | 66.33 | 42.45 | 63.01 | 48.78 | 37.68 | 62.63 | 52.41 | 57.38 | 52.87 |
|  | Weighted Mean | 53.01 | 46.27 | 66.04 | 42.53 | 62.26 | 48.95 | 37.85 | 61.62 | 51.993 | 56.95 | 52.75 |
| OpenAI-3-small | - | 58.28 | 54.04 | 65.73 | 60.26 | 77.54 | 58.95 | 50.42 | 58.66 | 52.45 | 48.87 | 58.52 |
| OpenAI-ada-002 | - | 57.94 | 56.71 | 61.64 | 64.10 | 75.75 | 63.28 | 52.17 | 65.50 | 55.66 | 50.00 | 60.28 |
| OpenAI-3-large | - | 62.88 | 57.12 | **73.27** | 64.90 | 78.96 | 64.10 | 54.21 | 71.66 | 55.96 | 49.99 | 63.31 |

Table 15: Average AUROC(%) on DBpedia-0.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 82.33 | 77.03 | 86.25 | 78.19 | 79.83 | 81.04 | 65.91 | 79.30 | 56.61 | 82.42 | 76.89 |
| BERT | CLS | 84.94 | 83.48 | 94.62 | 84.91 | 88.83 | 83.56 | 78.23 | 91.72 | 73.72 | 90.05 | 85.41 |
| | Mean | 81.63 | 78.96 | 91.35 | 80.39 | 87.41 | 80.11 | 71.79 | 90.65 | 81.18 | 90.67 | 83.41 |
| LLaMA-2 (mntp) | EOS | 66.76 | 65.36 | 73.55 | 66.98 | 66.16 | 65.87 | 62.14 | 69.01 | 59.94 | 66.76 | 66.25 |
| | Mean | 71.12 | 67.53 | 75.38 | 68.81 | 64.65 | 68.76 | 60.57 | 64.77 | 64.68 | 63.27 | 66.95 |
| | Weighted Mean | 68.45 | 65.79 | 76.16 | 66.89 | 64.00 | 66.50 | 60.16 | 61.86 | 61.99 | 62.72 | 65.45 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 95.81 | **84.70** | 96.94 | 95.58 | 97.19 | 95.72 | 86.18 | 96.90 | 84.03 | **93.53** | **92.66** |
| | Mean | 74.23 | 68.99 | 89.01 | 73.36 | 90.10 | 74.03 | 58.97 | 91.77 | 68.61 | 93.10 | 78.22 |
| | Weighted Mean | 73.66 | 68.05 | 87.25 | 72.96 | 88.96 | 73.46 | 60.43 | 90.74 | 66.59 | 91.35 | 77.35 |
| LLaMA-2 (mntp-supervised) | EOS | 92.60 | 76.40 | 97.15 | 91.94 | 95.01 | 92.53 | 71.28 | 94.77 | 80.94 | 88.98 | 88.16 |
| | Mean | 85.12 | 70.96 | 96.43 | 84.64 | 95.57 | 85.43 | 67.99 | 95.08 | 69.66 | 88.76 | 83.96 |
| | Weighted Mean | 82.66 | 70.59 | 95.45 | 82.01 | 94.07 | 82.87 | 65.51 | 93.65 | 70.64 | 86.41 | 82.39 |
| LLaMA-3 (mntp) | EOS | 79.66 | 75.02 | 85.66 | 79.01 | 82.40 | 77.65 | 62.42 | 79.92 | 72.24 | 75.25 | 76.92 |
| | Mean | 67.77 | 61.99 | 78.51 | 64.46 | 72.96 | 64.46 | 49.11 | 75.34 | 59.41 | 74.14 | 66.82 |
| | Weighted Mean | 67.23 | 63.40 | 79.63 | 64.99 | 72.87 | 64.01 | 52.51 | 75.60 | 60.46 | 73.72 | 67.44 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 64.99 | 59.96 | 78.33 | 62.28 | 64.22 | 60.07 | 53.30 | 69.45 | 59.06 | 63.10 | 63.48 |
| | Mean | 57.96 | 52.51 | 77.22 | 55.97 | 72.86 | 56.43 | 39.12 | 79.26 | 68.53 | 87.71 | 64.76 |
| | Weighted Mean | 60.54 | 54.98 | 78.71 | 58.64 | 73.21 | 58.85 | 44.32 | 79.66 | 67.49 | 86.36 | 66.28 |
| LLaMA-3 (mntp-supervised) | EOS | 91.28 | 81.33 | 96.52 | 89.50 | 91.92 | 89.30 | 79.18 | 92.78 | 79.51 | 88.42 | 87.97 |
| | Mean | 89.58 | 75.19 | 97.70 | 88.85 | 97.10 | 89.57 | 72.66 | 95.77 | 76.45 | 90.12 | 87.30 |
| | Weighted Mean | 90.05 | 74.91 | 97.39 | 89.40 | 96.54 | 90.10 | 75.55 | 95.25 | 75.27 | 89.26 | 87.37 |
| Mistral (mntp) | EOS | 84.28 | 74.80 | 87.38 | 79.19 | 84.40 | 84.09 | 60.87 | 82.58 | 74.01 | 81.65 | 79.33 |
| | Mean | 73.53 | 64.10 | 78.53 | 66.68 | 78.63 | 71.00 | 53.56 | 75.60 | 62.50 | 70.67 | 69.48 |
| | Weighted Mean | 72.03 | 63.49 | 77.92 | 64.06 | 76.14 | 68.74 | 53.46 | 73.52 | 62.46 | 67.64 | 67.95 |
| Mistral (mntp-unsup-simcse) | EOS | 84.17 | 69.58 | 87.92 | 83.64 | 83.72 | 84.88 | 69.46 | 88.11 | 67.04 | 78.28 | 79.68 |
| | Mean | 62.80 | 57.52 | 75.63 | 56.29 | 73.78 | 61.24 | 49.01 | 77.50 | 63.44 | 74.81 | 65.20 |
| | Weighted Mean | 61.06 | 55.33 | 74.31 | 55.30 | 71.51 | 59.71 | 49.57 | 75.98 | 62.69 | 73.02 | 63.85 |
| Mistral (mntp-supervised) | EOS | 95.09 | 73.39 | 96.62 | 89.64 | 94.79 | 95.24 | 72.87 | 95.15 | 74.24 | 88.68 | 87.57 |
| | Mean | 92.75 | 69.24 | 96.86 | 79.56 | 95.32 | 91.00 | 65.31 | 96.01 | 63.39 | 87.23 | 83.67 |
| | Weighted Mean | 91.19 | 67.42 | 96.15 | 77.83 | 93.97 | 89.39 | 64.72 | 95.04 | 63.79 | 71.38 | 81.09 |
| OpenAI-3-small | - | 98.23 | 83.36 | 98.86 | 98.18 | 98.86 | 98.12 | **91.73** | 97.52 | 97.12 | 50.00 | 91.20 |
| OpenAI-ada-002 | - | 93.55 | 76.89 | 97.77 | 93.76 | 96.55 | 94.05 | 73.06 | 92.03 | 93.44 | 49.35 | 86.05 |
| OpenAI-3-large | - | **98.70** | 79.29 | **99.30** | **98.50** | **99.19** | **98.59** | 91.03 | **98.32** | **98.27** | 63.99 | 92.52 |

Table 16: Average AUROC(%) on DBpedia-1.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 91.59 | 88.37 | 96.57 | 90.11 | 89.22 | 90.89 | 76.42 | 89.48 | 68.34 | 96.02 | 87.70 |
| BERT | CLS | 85.12 | 82.64 | 96.96 | 84.38 | 90.72 | 82.31 | 76.57 | 93.64 | 75.22 | 92.02 | 85.69 |
| | Mean | 82.06 | 77.26 | 95.04 | 80.88 | 92.76 | 79.44 | 70.92 | 94.50 | 91.28 | 95.08 | 85.92 |
| LLaMA-2 (mntp) | EOS | 72.70 | 70.78 | 77.46 | 72.85 | 73.51 | 71.59 | 67.31 | 77.05 | 64.28 | 75.58 | 72.31 |
| | Mean | 78.29 | 77.44 | 84.46 | 77.84 | 73.07 | 75.90 | 67.14 | 71.41 | 69.54 | 70.32 | 74.54 |
| | Weighted Mean | 76.40 | 74.59 | 84.16 | 75.83 | 72.22 | 74.01 | 67.42 | 71.15 | 66.28 | 70.59 | 73.26 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 99.36 | 94.23 | 98.87 | 99.29 | 99.57 | 99.17 | 95.50 | 99.68 | 95.37 | 98.98 | **98.00** |
| | Mean | 91.62 | 85.06 | 96.51 | 91.27 | 97.89 | 90.92 | 75.64 | 98.86 | 84.71 | 99.35 | 91.18 |
| | Weighted Mean | 91.55 | 84.60 | 95.31 | 91.00 | 97.60 | 90.76 | 77.20 | 98.57 | 83.08 | 98.94 | 90.86 |
| LLaMA-2 (mntp-supervised) | EOS | 99.24 | 93.39 | 99.26 | 99.06 | 99.41 | 99.09 | 94.00 | 99.65 | 94.31 | 98.47 | 97.59 |
| | Mean | 99.02 | 92.20 | 99.34 | 98.94 | 99.47 | 98.97 | 93.46 | 99.74 | 91.65 | 99.18 | 97.20 |
| | Weighted Mean | 98.59 | 90.42 | 99.00 | 98.48 | 99.34 | 98.54 | 91.94 | 99.64 | 90.16 | 98.81 | 96.49 |
| LLaMA-3 (mntp) | EOS | 91.01 | 85.87 | 92.29 | 90.23 | 95.34 | 89.26 | 77.18 | 94.66 | 85.20 | 89.82 | 89.09 |
| | Mean | 80.39 | 76.92 | 91.91 | 79.08 | 92.66 | 77.83 | 62.68 | 94.02 | 77.41 | 90.59 | 82.35 |
| | Weighted Mean | 81.00 | 77.90 | 92.25 | 79.69 | 92.03 | 78.23 | 65.85 | 93.47 | 76.69 | 89.88 | 82.70 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 81.52 | 72.13 | 89.55 | 76.58 | 81.08 | 73.44 | 66.64 | 87.63 | 73.42 | 80.69 | 78.27 |
| | Mean | 84.90 | 79.02 | 93.02 | 83.65 | 93.64 | 84.31 | 63.48 | 96.40 | 84.63 | 97.93 | 86.10 |
| | Weighted Mean | 86.15 | 78.86 | 91.60 | 84.76 | 93.30 | 85.33 | 67.97 | 96.06 | 78.25 | 96.88 | 85.92 |
| LLaMA-3 (mntp-supervised) | EOS | 98.90 | 93.60 | 98.73 | 98.33 | 98.28 | 98.07 | 94.34 | 99.15 | 95.22 | 98.46 | 97.31 |
| | Mean | 99.45 | 94.13 | 99.64 | 99.37 | 99.78 | 99.38 | 95.29 | **99.75** | 95.51 | **99.25** | 98.16 |
| | Weighted Mean | 99.36 | 93.42 | 99.43 | 99.29 | 99.69 | 99.31 | 95.44 | 99.68 | 94.20 | 99.00 | 97.88 |
| Mistral (mntp) | EOS | 96.03 | 81.95 | 97.14 | 88.69 | 96.01 | 95.67 | 72.63 | 96.80 | 86.03 | 94.69 | 90.56 |
| | Mean | 85.54 | 77.75 | 93.35 | 78.44 | 92.88 | 80.27 | 66.97 | 92.60 | 82.60 | 86.82 | 83.72 |
| | Weighted Mean | 85.26 | 74.54 | 93.26 | 77.45 | 92.01 | 80.11 | 67.72 | 91.90 | 82.20 | 86.57 | 83.10 |
| Mistral (mntp-unsup-simcse) | EOS | 96.91 | 89.10 | 96.90 | 95.84 | 96.64 | 97.11 | 90.26 | 97.84 | 88.17 | 95.70 | 94.45 |
| | Mean | 82.61 | 68.12 | 89.88 | 72.22 | 90.87 | 80.24 | 59.96 | 93.49 | 82.63 | 94.13 | 81.42 |
| | Weighted Mean | 81.05 | 67.34 | 88.65 | 71.63 | 89.72 | 78.88 | 61.55 | 92.83 | 79.31 | 93.01 | 80.40 |
| Mistral (mntp-supervised) | EOS | 98.70 | 85.69 | 97.86 | 97.90 | 98.61 | 98.76 | 89.27 | 99.17 | 91.36 | 98.35 | 95.57 |
| | Mean | 99.36 | 84.70 | 99.36 | 96.78 | 99.44 | 99.22 | 84.83 | 99.66 | 84.52 | 98.59 | 94.65 |
| | Weighted Mean | 99.21 | 84.99 | 98.93 | 96.10 | 99.29 | 99.05 | 84.15 | 99.56 | 84.45 | 98.41 | 94.41 |
| OpenAI-3-small | - | 99.81 | **95.89** | 99.85 | 99.78 | 99.86 | 99.77 | 99.09 | 99.55 | 99.81 | 50.00 | 94.34 |
| OpenAI-ada-002 | - | 99.33 | 92.45 | 99.69 | 99.35 | 99.63 | 99.31 | 92.75 | 92.57 | 99.48 | 49.94 | 92.45 |
| OpenAI-3-large | - | **99.87** | 94.48 | **99.89** | **99.83** | **99.90** | **99.84** | **99.57** | 99.67 | **99.87** | 48.56 | 94.15 |

30

Table 17: Average AUROC(%) on DBpedia-2.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 85.27 | 83.18 | 87.29 | 83.39 | 83.86 | 84.03 | 74.83 | 83.76 | 57.31 | 87.38 | 81.03 |
| BERT | CLS | 74.38 | 72.29 | 93.54 | 73.46 | 81.22 | 71.48 | 67.47 | 86.61 | 69.06 | 83.49 | 77.30 |
|  | Mean | 80.26 | 76.38 | 92.97 | 78.62 | 89.61 | 77.28 | 71.73 | 93.38 | 84.91 | 91.41 | 83.66 |
| LLaMA-2 (mntp) | EOS | 66.84 | 63.97 | 73.12 | 66.84 | 65.10 | 65.68 | 60.01 | 66.17 | 61.51 | 61.69 | 65.09 |
|  | Mean | 77.94 | 76.14 | 84.42 | 76.70 | 72.00 | 74.73 | 67.44 | 65.56 | 67.48 | 66.92 | 72.93 |
|  | Weighted Mean | 75.68 | 74.90 | 83.16 | 74.45 | 70.01 | 72.54 | 65.22 | 64.08 | 65.20 | 64.85 | 71.01 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 91.55 | 76.95 | 95.02 | 89.97 | 93.93 | 89.72 | 80.26 | 96.38 | 82.01 | 92.33 | 88.81 |
|  | Mean | 78.48 | 72.24 | 91.36 | 78.01 | 89.49 | 77.70 | 67.63 | 94.41 | 68.62 | **93.59** | 81.15 |
|  | Weighted Mean | 75.82 | 68.91 | 91.12 | 75.03 | 89.09 | 74.74 | 63.90 | 93.55 | 66.04 | 92.27 | 79.05 |
| LLaMA-2 (mntp-supervised) | EOS | 84.92 | 74.84 | 93.14 | 84.19 | 91.30 | 84.45 | 75.18 | 91.81 | 76.86 | 85.37 | 84.21 |
|  | Mean | 89.62 | 75.58 | 97.45 | 88.93 | 94.80 | 89.24 | 78.65 | 95.88 | 75.91 | 90.64 | 87.67 |
|  | Weighted Mean | 88.65 | 73.82 | 97.05 | 87.88 | 94.80 | 88.25 | 77.30 | 95.17 | 73.50 | 90.45 | 86.69 |
| LLaMA-3 (mntp) | EOS | 88.90 | **84.84** | 90.38 | 88.47 | 92.36 | 88.14 | 79.83 | 89.97 | 81.23 | 86.44 | 87.06 |
|  | Mean | 81.14 | 79.00 | 90.45 | 81.48 | 89.25 | 79.20 | 73.15 | 91.79 | 75.44 | 89.44 | 83.03 |
|  | Weighted Mean | 80.62 | 78.41 | 90.48 | 79.85 | 88.48 | 78.50 | 69.85 | 90.37 | 74.53 | 85.76 | 81.68 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 66.82 | 59.22 | 77.74 | 61.74 | 64.43 | 60.04 | 57.90 | 72.31 | 60.26 | 67.38 | 64.78 |
|  | Mean | 80.14 | 74.68 | 87.21 | 78.89 | 84.98 | 78.95 | 71.00 | 90.39 | 67.80 | 89.99 | 80.40 |
|  | Weighted Mean | 79.17 | 71.97 | 87.73 | 77.33 | 85.45 | 77.52 | 68.29 | 90.46 | 67.73 | 90.59 | 79.62 |
| LLaMA-3 (mntp-supervised) | EOS | 87.34 | 68.89 | 94.21 | 83.58 | 87.64 | 83.65 | 68.01 | 91.19 | 77.15 | 86.73 | 82.84 |
|  | Mean | 92.93 | 80.54 | **97.91** | 91.96 | **97.54** | 92.46 | 81.73 | 96.58 | 85.96 | 90.91 | **90.85** |
|  | Weighted Mean | 92.43 | 76.96 | 97.62 | 91.41 | 97.12 | 92.01 | 80.08 | 96.15 | 84.63 | 90.57 | 89.90 |
| Mistral (mntp) | EOS | 90.46 | 78.65 | 91.07 | 84.02 | 90.74 | 89.99 | 78.01 | 92.84 | 79.47 | 88.08 | 86.33 |
|  | Mean | 83.52 | 75.35 | 91.04 | 76.09 | 90.05 | 78.12 | 67.50 | 91.18 | 74.23 | 84.52 | 81.16 |
|  | Weighted Mean | 83.98 | 73.31 | 91.23 | 74.66 | 88.97 | 78.16 | 65.20 | 89.69 | 73.24 | 82.73 | 80.12 |
| Mistral (mntp-unsup-simcse) | EOS | 78.43 | 65.80 | 78.17 | 76.89 | 78.23 | 78.61 | 69.70 | 82.59 | 65.59 | 78.43 | 75.24 |
|  | Mean | 72.22 | 62.69 | 79.31 | 65.42 | 76.27 | 70.72 | 60.51 | 85.28 | 68.39 | 82.13 | 72.29 |
|  | Weighted Mean | 68.43 | 60.36 | 79.05 | 61.66 | 74.03 | 66.83 | 57.06 | 83.96 | 66.63 | 79.75 | 69.78 |
| Mistral (mntp-supervised) | EOS | 94.16 | 75.20 | 95.98 | 87.90 | 94.24 | **94.03** | 80.71 | 95.07 | 76.60 | 89.21 | 88.31 |
|  | Mean | 94.30 | 74.90 | 96.87 | 86.51 | 95.24 | 93.38 | 78.98 | **97.14** | 70.80 | 89.60 | 87.77 |
|  | Weighted Mean | 93.93 | 73.10 | 96.87 | 85.31 | 95.54 | 92.86 | 77.72 | 96.99 | 69.49 | 89.69 | 87.15 |
| OpenAI-3-small | - | **95.83** | 77.10 | 96.96 | **94.18** | 95.41 | 93.99 | 86.62 | 95.64 | **95.31** | 50.00 | 88.10 |
| OpenAI-ada-002 | - | 94.20 | 80.73 | 96.05 | 94.16 | 92.06 | 93.92 | 87.07 | 93.26 | 94.32 | 50.00 | 87.58 |
| OpenAI-3-large | - | 95.15 | 75.31 | 97.79 | 92.95 | 96.64 | 93.15 | **87.13** | 97.05 | 94.84 | 50.25 | 88.03 |

Table 18: Average AUROC(%) on DBpedia-3.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 87.97 | 82.11 | 97.15 | 83.69 | 92.02 | 85.17 | 69.89 | 91.06 | 66.49 | 95.19 | 85.07 |
| BERT | CLS | 67.87 | 63.44 | 94.62 | 64.36 | 85.11 | 63.13 | 51.13 | 87.72 | 79.72 | 86.22 | 74.33 |
| | Mean | 75.88 | 72.88 | 96.42 | 70.94 | 95.96 | 69.22 | 56.43 | 96.35 | 93.83 | 96.92 | 82.48 |
| LLaMA-2 (mntp) | EOS | 69.63 | 68.89 | 80.81 | 70.94 | 79.30 | 69.28 | 61.81 | 77.33 | 68.06 | 78.44 | 72.45 |
| | Mean | 81.68 | 82.60 | 90.70 | 82.64 | 82.90 | 80.83 | 71.28 | 78.73 | 75.17 | 77.98 | 80.45 |
| | Weighted Mean | 78.74 | 81.06 | 88.54 | 79.83 | 81.85 | 78.22 | 68.30 | 77.42 | 73.07 | 77.54 | 78.46 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 98.52 | 91.07 | 99.52 | 98.02 | 99.28 | 97.76 | 93.74 | **99.71** | 93.79 | 98.69 | **97.01** |
| | Mean | 95.15 | 88.32 | 99.48 | 95.08 | 99.21 | 94.92 | 86.74 | 99.68 | 91.68 | **99.62** | 94.99 |
| | Weighted Mean | 94.43 | 88.46 | 99.32 | 94.13 | 99.38 | 94.04 | 83.96 | 99.62 | 88.23 | 99.51 | 94.11 |
| LLaMA-2 (mntp-supervised) | EOS | 95.87 | 87.82 | 98.84 | 95.72 | 98.02 | 95.80 | 90.23 | 98.99 | 90.89 | 96.66 | 94.88 |
| | Mean | 94.54 | 82.80 | 99.55 | 94.11 | 98.78 | 94.25 | 83.03 | 99.60 | 90.02 | 97.75 | 93.44 |
| | Weighted Mean | 93.61 | 81.93 | 99.47 | 93.12 | 98.96 | 93.29 | 81.69 | 99.53 | 87.74 | 97.59 | 92.69 |
| LLaMA-3 (mntp) | EOS | 96.56 | **93.48** | 96.65 | 95.99 | 98.33 | 96.13 | 87.10 | 97.68 | 94.09 | 95.98 | 95.20 |
| | Mean | 93.31 | 91.62 | 98.98 | 93.31 | 99.03 | 91.35 | 84.11 | 99.28 | 89.18 | 98.19 | 93.84 |
| | Weighted Mean | 93.59 | 89.09 | 98.71 | 92.64 | 98.96 | 91.69 | 80.99 | 99.02 | 89.18 | 97.74 | 93.16 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 76.36 | 67.48 | 90.21 | 71.02 | 78.37 | 69.11 | 65.69 | 84.28 | 74.08 | 77.67 | 75.43 |
| | Mean | 94.99 | 90.25 | 98.76 | 94.43 | 98.39 | 94.44 | 86.34 | 99.07 | 90.37 | 98.84 | 94.79 |
| | Weighted Mean | 94.79 | 87.39 | 98.57 | 93.80 | 98.69 | 94.02 | 86.13 | 99.10 | 88.85 | 98.97 | 94.03 |
| LLaMA-3 (mntp-supervised) | EOS | 93.77 | 80.49 | 98.11 | 91.52 | 95.54 | 91.50 | 79.30 | 98.02 | 91.44 | 95.16 | 91.48 |
| | Mean | 96.56 | 85.35 | 99.78 | 95.87 | **99.64** | 95.96 | 85.80 | 99.71 | 94.20 | 98.21 | 95.11 |
| | Weighted Mean | 96.17 | 84.75 | 99.73 | 95.24 | **99.64** | 95.35 | 83.28 | 99.67 | 94.88 | 97.02 | 94.57 |
| Mistral (mntp) | EOS | 94.69 | 85.16 | 96.14 | 89.28 | 94.75 | 94.38 | 80.82 | 96.20 | 88.45 | 94.81 | 91.47 |
| | Mean | 94.67 | 87.76 | 98.53 | 90.32 | 98.79 | 89.94 | 79.43 | 98.97 | 89.64 | 96.95 | 92.50 |
| | Weighted Mean | 95.02 | 86.60 | 98.31 | 89.23 | 98.49 | 91.01 | 76.39 | 98.48 | 88.46 | 95.93 | 91.79 |
| Mistral (mntp-unsup-simcse) | EOS | 88.88 | 76.12 | 88.74 | 87.11 | 88.50 | 89.02 | 80.95 | 93.96 | 80.06 | 90.71 | 86.40 |
| | Mean | 87.74 | 76.16 | 94.93 | 79.56 | 93.72 | 85.18 | 71.12 | 97.63 | 90.91 | 96.82 | 87.38 |
| | Weighted Mean | 86.46 | 72.79 | 95.93 | 76.88 | 93.81 | 83.30 | 68.24 | 97.56 | 87.75 | 96.18 | 85.89 |
| Mistral (mntp-supervised) | EOS | 98.14 | 82.72 | 98.23 | 94.92 | 98.13 | 98.06 | 89.56 | 98.94 | 89.69 | 96.16 | 94.46 |
| | Mean | 97.08 | 78.75 | 98.90 | 89.69 | 98.48 | 96.22 | 79.09 | 99.44 | 89.94 | 97.08 | 92.47 |
| | Weighted Mean | 97.00 | 77.61 | 99.03 | 87.84 | 98.77 | 95.96 | 75.91 | 99.47 | 88.26 | 96.90 | 91.68 |
| OpenAI-3-small | - | 99.47 | 88.72 | 99.73 | **99.02** | 99.05 | **98.99** | 96.65 | 99.50 | 99.64 | 50.00 | 93.08 |
| OpenAI-ada-002 | - | 98.97 | 91.13 | 99.77 | 98.84 | 99.00 | 98.75 | 95.60 | 99.26 | 99.18 | 50.00 | 93.05 |
| OpenAI-3-large | - | **99.46** | 88.82 | **99.81** | 98.95 | 99.33 | 98.90 | **96.97** | 99.42 | **99.69** | 52.96 | 93.43 |

32

Table 19: Average AUROC(%) on DBpedia-4.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 86.64 | 82.12 | 92.22 | 83.09 | 85.92 | 84.82 | 70.75 | 86.25 | 61.66 | 91.17 | 82.46 |
| BERT | CLS | 68.96 | 63.59 | 93.78 | 67.67 | 82.68 | 64.06 | 57.37 | 86.37 | 74.36 | 83.02 | 74.19 |
|  | Mean | 76.02 | 72.27 | 93.65 | 72.75 | 93.23 | 71.38 | 63.34 | 94.07 | 88.91 | 94.38 | 82.00 |
| LLaMA-2 (mntp) | EOS | 71.83 | 69.48 | 72.80 | 71.09 | 72.77 | 69.98 | 60.90 | 70.17 | 66.53 | 71.23 | 69.68 |
|  | Mean | 85.18 | 82.91 | 89.40 | 84.03 | 79.83 | 82.68 | 74.62 | 75.26 | 72.24 | 75.87 | 80.20 |
|  | Weighted Mean | 81.99 | 80.23 | 87.43 | 80.91 | 78.56 | 79.13 | 71.28 | 77.33 | 69.73 | 75.90 | 78.25 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 97.43 | 88.44 | 97.72 | 96.30 | 98.78 | 96.16 | 88.80 | 98.91 | 87.72 | 96.91 | 94.72 |
|  | Mean | 92.43 | 87.34 | 96.54 | 92.33 | 97.64 | 92.34 | 84.28 | 98.58 | 78.40 | **98.61** | 91.85 |
|  | Weighted Mean | 92.05 | 85.47 | 95.88 | 91.97 | 97.71 | 92.03 | 83.33 | 98.34 | 74.56 | 98.32 | 90.97 |
| LLaMA-2 (mntp-supervised) | EOS | 97.04 | **89.15** | 98.54 | 96.56 | 97.95 | 96.61 | 93.54 | 98.43 | 88.73 | 95.64 | 95.22 |
|  | Mean | 95.95 | 83.52 | 98.73 | 95.58 | 97.78 | 95.72 | 87.47 | 98.71 | 85.29 | 96.22 | 93.50 |
|  | Weighted Mean | 95.17 | 84.63 | 98.37 | 94.70 | 97.92 | 94.87 | 85.52 | 98.45 | 83.35 | 96.15 | 92.91 |
| LLaMA-3 (mntp) | EOS | 92.05 | 89.42 | 90.54 | 91.60 | 95.28 | 91.77 | 84.06 | 93.35 | 86.54 | 91.08 | 90.57 |
|  | Mean | 88.88 | 87.18 | 95.54 | 88.77 | 96.54 | 87.22 | 81.80 | 97.10 | 82.94 | 93.81 | 89.98 |
|  | Weighted Mean | 89.43 | 87.75 | 95.12 | 89.13 | 96.14 | 88.25 | 81.75 | 96.56 | 82.35 | 94.02 | 90.05 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 73.74 | 66.52 | 85.02 | 69.60 | 76.21 | 67.38 | 65.19 | 80.80 | 69.16 | 74.95 | 72.86 |
|  | Mean | 91.37 | 86.63 | 95.26 | 90.81 | 96.02 | 91.03 | 84.40 | 97.43 | 80.17 | 97.12 | 91.02 |
|  | Weighted Mean | 91.30 | 86.73 | 94.78 | 90.79 | 96.24 | 91.16 | 83.58 | 97.30 | 76.66 | 97.04 | 90.56 |
| LLaMA-3 (mntp-supervised) | EOS | 95.96 | 84.73 | 98.22 | 94.82 | 94.75 | 94.71 | 85.64 | 96.31 | 85.91 | 95.08 | 92.61 |
|  | Mean | 98.28 | 88.03 | **99.22** | 97.69 | **99.24** | 97.75 | 92.08 | 98.93 | 91.95 | 97.26 | **96.04** |
|  | Weighted Mean | 98.05 | 86.61 | 98.92 | 97.57 | 99.12 | 97.65 | 91.01 | 98.71 | 89.72 | 96.89 | 95.42 |
| Mistral (mntp) | EOS | 93.33 | 81.38 | 93.52 | 84.66 | 93.63 | 92.76 | 77.53 | 93.86 | 80.03 | 90.32 | 88.10 |
|  | Mean | 88.83 | 83.05 | 95.74 | 84.23 | 96.29 | 84.09 | 77.68 | 96.60 | 83.66 | 92.92 | 88.31 |
|  | Weighted Mean | 89.74 | 81.42 | 95.47 | 84.30 | 95.87 | 85.79 | 77.06 | 95.88 | 83.76 | 93.99 | 88.33 |
| Mistral (mntp-unsup-simcse) | EOS | 86.20 | 72.20 | 85.42 | 82.71 | 85.82 | 86.43 | 73.96 | 89.74 | 72.96 | 85.27 | 82.07 |
|  | Mean | 84.02 | 70.93 | 91.94 | 75.91 | 90.49 | 82.12 | 68.83 | 95.10 | 77.17 | 84.42 | 82.09 |
|  | Weighted Mean | 82.48 | 70.17 | 92.89 | 74.03 | 90.45 | 80.30 | 67.09 | 94.69 | 74.71 | 91.93 | 81.87 |
| Mistral (mntp-supervised) | EOS | 96.34 | 79.86 | 96.59 | 92.16 | 96.30 | 96.25 | 85.23 | 96.83 | 82.43 | 93.21 | 91.52 |
|  | Mean | 96.57 | 75.56 | 97.72 | 88.28 | 97.39 | 95.73 | 78.00 | 98.32 | 80.15 | 93.94 | 90.17 |
|  | Weighted Mean | 96.28 | 74.55 | 97.64 | 86.88 | 97.51 | 95.30 | 75.78 | 98.25 | 77.31 | 84.70 | 88.42 |
| OpenAI-3-small | - | 98.60 | 84.79 | 98.33 | 97.68 | 98.72 | 97.63 | 93.60 | 98.39 | 98.64 | 50.00 | 91.64 |
| OpenAI-ada-002 | - | 98.46 | 89.91 | 98.69 | **98.42** | 98.04 | **98.40** | 95.86 | 97.40 | 98.57 | 51.67 | 92.54 |
| OpenAI-3-large | - | **98.89** | 84.20 | 98.87 | 98.17 | 98.74 | 98.12 | 95.67 | **99.15** | **99.26** | 44.73 | 91.58 |

33

Table 20: Average AUPRC(%) on DBpedia-0.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 75.71 | 70.37 | 81.18 | 71.46 | 78.23 | 74.33 | 57.78 | 76.00 | 49.73 | 79.74 | 71.45 |
| BERT | CLS | 80.55 | 80.16 | 94.03 | 81.80 | 88.41 | 78.31 | 74.76 | 91.09 | 70.10 | 88.12 | 82.73 |
|  | Mean | 75.26 | 73.05 | 89.84 | 73.91 | 84.82 | 73.11 | 64.81 | 88.09 | 78.47 | 88.49 | 78.99 |
| LLaMA-2 | EOS | 56.67 | 55.17 | 67.82 | 56.49 | 58.53 | 55.56 | 51.91 | 58.88 | 51.03 | 57.48 | 56.95 |
| (mntp) | Mean | 58.97 | 55.13 | 72.02 | 56.11 | 54.86 | 56.08 | 49.52 | 54.88 | 56.32 | 53.51 | 56.74 |
|  | Weighted Mean | 57.14 | 54.47 | 73.42 | 54.96 | 54.65 | 54.66 | 49.58 | 53.40 | 54.49 | 53.39 | 56.02 |
| LLaMA-2 | EOS | 94.68 | 81.07 | 96.03 | 94.25 | 96.33 | 94.43 | 81.18 | 96.07 | 81.45 | 92.40 | 90.79 |
| (mntp-unsup-simcse) | Mean | 63.24 | 59.69 | 83.47 | 62.01 | 84.27 | 62.60 | 48.93 | 87.45 | 68.42 | 90.69 | 71.08 |
|  | Weighted Mean | 63.40 | 59.51 | 81.40 | 62.32 | 82.92 | 62.71 | 50.39 | 86.75 | 66.20 | 88.81 | 70.44 |
| LLaMA-2 | EOS | 89.72 | 70.60 | 96.73 | 88.21 | 94.73 | 89.20 | 61.68 | 94.48 | 78.06 | 86.44 | 84.99 |
| (mntp-supervised) | Mean | 81.71 | 64.93 | 95.92 | 80.78 | 94.67 | 81.66 | 61.00 | 94.43 | 64.95 | 86.94 | 80.70 |
|  | Weighted Mean | 79.14 | 64.78 | 94.85 | 78.18 | 93.05 | 79.12 | 59.13 | 92.88 | 66.17 | 84.27 | 79.16 |
| LLaMA-3 | EOS | 68.72 | 64.56 | 79.29 | 67.88 | 75.63 | 65.35 | 50.81 | 72.98 | 64.43 | 66.37 | 67.60 |
| (mntp) | Mean | 54.94 | 51.43 | 69.99 | 52.30 | 61.18 | 52.12 | 41.88 | 66.09 | 52.83 | 66.13 | 56.89 |
|  | Weighted Mean | 54.16 | 52.46 | 72.65 | 52.58 | 61.92 | 51.10 | 43.51 | 67.94 | 53.23 | 65.73 | 57.53 |
| LLaMA-3 | EOS | 57.75 | 53.28 | 72.57 | 55.06 | 59.06 | 54.01 | 48.35 | 64.11 | 52.72 | 56.62 | 57.35 |
| (mntp-unsup-simcse) | Mean | 47.79 | 44.59 | 68.90 | 46.25 | 62.68 | 46.39 | 36.73 | 70.97 | 67.30 | 83.62 | 57.52 |
|  | Weighted Mean | 49.88 | 46.69 | 71.77 | 48.16 | 63.20 | 48.10 | 39.06 | 72.30 | 66.63 | 82.18 | 58.80 |
| LLaMA-3 | EOS | 90.61 | 78.77 | 96.21 | 88.64 | 92.29 | 88.65 | 75.71 | 93.15 | 76.48 | 87.20 | 86.77 |
| (mntp-supervised) | Mean | 87.01 | 69.91 | 96.93 | 85.78 | 96.33 | 86.77 | 66.12 | 95.22 | 69.91 | 88.33 | 84.23 |
|  | Weighted Mean | 87.94 | 69.85 | 96.62 | 86.84 | 95.82 | 87.79 | 69.86 | 94.77 | 69.17 | 87.42 | 84.61 |
| Mistral | EOS | 80.92 | 66.88 | 83.21 | 71.28 | 81.21 | 80.64 | 50.76 | 77.67 | 70.74 | 75.73 | 73.90 |
| (mntp) | Mean | 62.47 | 54.29 | 67.56 | 56.39 | 67.18 | 60.21 | 46.17 | 64.68 | 52.95 | 62.05 | 59.40 |
|  | Weighted Mean | 60.11 | 53.73 | 69.17 | 53.15 | 64.73 | 57.10 | 45.33 | 63.33 | 52.65 | 58.92 | 57.82 |
| Mistral | EOS | 79.77 | 63.05 | 83.92 | 79.71 | 78.99 | 80.73 | 60.25 | 84.27 | 63.50 | 73.31 | 74.75 |
| (mntp-unsup-simcse) | Mean | 52.76 | 49.51 | 66.80 | 47.96 | 63.55 | 51.16 | 42.89 | 69.32 | 59.91 | 67.52 | 57.14 |
|  | Weighted Mean | 52.11 | 48.24 | 66.61 | 48.00 | 61.74 | 50.76 | 43.77 | 68.77 | 59.21 | 66.33 | 56.55 |
| Mistral | EOS | 94.61 | 66.74 | 95.02 | 85.98 | 94.39 | 94.49 | 63.31 | 94.62 | 68.58 | 85.85 | 84.36 |
| (mntp-supervised) | Mean | 89.39 | 62.21 | 95.71 | 72.26 | 93.90 | 86.90 | 56.45 | 94.92 | 56.56 | 84.61 | 79.29 |
|  | Weighted Mean | 87.70 | 61.09 | 94.99 | 70.76 | 92.53 | 85.14 | 56.64 | 93.98 | 57.24 | 74.53 | 77.46 |
| OpenAI-3-small | - | 97.51 | 79.64 | 98.41 | 97.51 | 98.38 | 97.36 | 88.64 | 96.02 | 96.67 | 72.22 | 92.24 |
| OpenAI-ada-002 | - | 89.37 | 71.41 | 96.86 | 89.66 | 95.36 | 90.06 | 63.80 | 84.13 | 89.96 | 66.10 | 83.67 |
| OpenAI-3-large | - | 97.70 | 74.32 | 98.85 | 97.58 | 98.50 | 97.51 | 85.41 | 96.95 | 97.60 | 71.94 | 91.64 |

Table 21: Average AUPRC(%) on DBpedia-1.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 86.75 | 82.24 | 94.89 | 84.35 | 88.16 | 85.55 | 66.30 | 87.44 | 59.01 | 94.85 | 82.95 |
| BERT | CLS | 76.32 | 76.46 | 96.13 | 77.89 | 89.81 | 72.45 | 68.87 | 92.51 | 70.44 | 89.13 | 81.00 |
| | Mean | 72.73 | 68.32 | 94.03 | 72.09 | 89.50 | 69.72 | 61.34 | 91.53 | 90.35 | 93.37 | 80.30 |
| LLaMA-2 (mntp) | EOS | 64.35 | 62.71 | 70.95 | 64.42 | 68.97 | 63.06 | 59.05 | 70.60 | 56.56 | 68.32 | 64.90 |
| | Mean | 68.26 | 68.19 | 82.69 | 68.43 | 65.65 | 64.93 | 58.36 | 65.44 | 62.92 | 62.75 | 66.76 |
| | Weighted Mean | 67.72 | 66.05 | 82.57 | 67.03 | 65.76 | 64.30 | 59.28 | 64.95 | 59.89 | 63.56 | 66.11 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 99.19 | 93.01 | 97.80 | 99.10 | 99.42 | 98.93 | 94.22 | 99.53 | 94.74 | 98.74 | 97.47 |
| | Mean | 87.69 | 79.81 | 94.78 | 87.08 | 96.93 | 86.30 | 66.54 | 98.23 | 84.16 | 99.09 | 88.06 |
| | Weighted Mean | 88.51 | 79.98 | 92.00 | 87.65 | 96.64 | 87.10 | 69.58 | 97.97 | 83.38 | 98.75 | 88.16 |
| LLaMA-2 (mntp-supervised) | EOS | 98.64 | 91.44 | 98.82 | 98.35 | 99.33 | 98.42 | 90.69 | 99.56 | 93.23 | 97.86 | 96.63 |
| | Mean | 98.45 | 90.34 | 98.90 | 98.32 | 99.30 | 98.33 | 91.45 | 99.61 | 90.38 | 98.96 | 96.40 |
| | Weighted Mean | 97.99 | 88.19 | 98.26 | 97.83 | 99.19 | 97.89 | 89.79 | 99.52 | 88.64 | 98.51 | 95.58 |
| LLaMA-3 (mntp) | EOS | 83.34 | 78.38 | 86.21 | 82.66 | 92.76 | 80.23 | 66.02 | 92.08 | 80.44 | 85.20 | 82.73 |
| | Mean | 69.38 | 67.33 | 90.48 | 68.65 | 88.81 | 65.74 | 53.38 | 91.86 | 70.86 | 87.36 | 75.39 |
| | Weighted Mean | 71.73 | 69.22 | 90.65 | 70.99 | 88.68 | 67.88 | 56.64 | 91.65 | 69.84 | 86.87 | 76.42 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 72.76 | 64.13 | 84.03 | 67.20 | 77.27 | 65.14 | 57.85 | 84.48 | 69.24 | 73.59 | 71.57 |
| | Mean | 78.88 | 72.05 | 90.79 | 77.04 | 91.35 | 77.54 | 57.06 | 94.99 | 84.99 | 97.25 | 82.19 |
| | Weighted Mean | 81.46 | 72.85 | 88.87 | 79.52 | 91.18 | 79.94 | 61.44 | 94.85 | 78.24 | 96.30 | 82.46 |
| LLaMA-3 (mntp-supervised) | EOS | 98.63 | 91.99 | 98.27 | 97.98 | 98.19 | 97.75 | 92.91 | 99.03 | 94.73 | 98.10 | 96.76 |
| | Mean | 99.19 | 92.67 | 99.48 | 99.05 | 99.68 | 99.09 | 93.90 | 99.65 | 94.47 | 99.00 | 97.62 |
| | Weighted Mean | 99.15 | 92.12 | 99.14 | 99.02 | 99.60 | 99.06 | 94.29 | 99.57 | 92.98 | 98.74 | 97.37 |
| Mistral (mntp) | EOS | 94.78 | 74.90 | 96.18 | 82.58 | 94.81 | 94.28 | 60.89 | 95.39 | 83.53 | 91.91 | 86.92 |
| | Mean | 73.76 | 64.77 | 91.37 | 64.29 | 88.01 | 65.77 | 54.26 | 88.52 | 76.74 | 80.57 | 74.81 |
| | Weighted Mean | 75.00 | 62.49 | 90.42 | 64.71 | 87.25 | 67.33 | 55.76 | 88.09 | 76.99 | 80.74 | 74.88 |
| Mistral (mntp-unsup-simcse) | EOS | 95.05 | 85.81 | 95.04 | 92.77 | 94.69 | 95.24 | 82.87 | 96.35 | 86.50 | 93.43 | 91.78 |
| | Mean | 76.51 | 60.69 | 87.10 | 64.58 | 87.51 | 73.20 | 52.84 | 91.09 | 81.76 | 92.67 | 76.79 |
| | Weighted Mean | 76.03 | 60.73 | 85.11 | 65.51 | 86.34 | 73.13 | 55.61 | 90.69 | 77.89 | 91.61 | 76.26 |
| Mistral (mntp-supervised) | EOS | 98.72 | 82.23 | 96.14 | 97.58 | 98.64 | 98.76 | 85.79 | 99.07 | 90.12 | 97.92 | 94.50 |
| | Mean | 98.94 | 81.36 | 99.02 | 95.94 | 99.23 | 98.73 | 81.73 | 99.47 | 80.84 | 98.31 | 93.36 |
| | Weighted Mean | 98.85 | 82.36 | 98.17 | 95.33 | 99.09 | 98.62 | 81.85 | 99.38 | 80.78 | 98.10 | 93.25 |
| OpenAI-3-small | - | 99.72 | 94.78 | 99.78 | 99.67 | 99.79 | 99.65 | 98.53 | 98.85 | 99.73 | 72.22 | 96.27 |
| OpenAI-ada-002 | - | 98.73 | 90.29 | 99.55 | 98.78 | 99.50 | 98.72 | 88.16 | 82.64 | 99.08 | 64.63 | 92.01 |
| OpenAI-3-large | - | 99.50 | 86.86 | 99.69 | 99.38 | 99.68 | 99.39 | 98.10 | 97.07 | 99.54 | 51.02 | 93.02 |

Table 22: Average AUPRC(%) on DBpedia-2.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 80.80 | 77.57 | 80.46 | 77.45 | 81.23 | 78.55 | 66.27 | 81.66 | 49.44 | 84.95 | 75.84 |
| BERT | CLS | 65.89 | 65.82 | 92.59 | 66.89 | 80.41 | 62.54 | 60.92 | 85.11 | 64.15 | 79.17 | 72.35 |
| | Mean | 72.93 | 69.45 | 92.30 | 71.37 | 88.32 | 68.77 | 64.33 | 91.91 | 83.29 | 89.92 | 79.26 |
| LLaMA-2 (mntp) | EOS | 59.23 | 56.31 | 70.00 | 59.26 | 61.23 | 57.93 | 52.93 | 59.69 | 55.93 | 56.55 | 58.91 |
| | Mean | 68.74 | 67.64 | 83.10 | 67.97 | 65.10 | 64.80 | 58.11 | 57.88 | 59.78 | 60.39 | 65.35 |
| | Weighted Mean | 66.79 | 66.50 | 81.56 | 65.82 | 63.97 | 63.06 | 56.63 | 57.28 | 58.34 | 59.57 | 63.95 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 90.13 | 72.78 | 94.16 | 88.24 | 93.55 | 87.99 | 75.20 | 95.82 | 79.57 | 91.20 | 86.86 |
| | Mean | 72.71 | 66.34 | 88.87 | 71.83 | 87.22 | 71.49 | 60.70 | 93.31 | 67.99 | 92.80 | 77.33 |
| | Weighted Mean | 69.35 | 62.07 | 88.28 | 68.09 | 86.00 | 67.73 | 56.58 | 92.01 | 65.00 | 91.25 | 74.64 |
| LLaMA-2 (mntp-supervised) | EOS | 84.89 | 72.15 | 92.71 | 83.99 | 90.96 | 84.24 | 74.35 | 91.33 | 74.64 | 84.93 | 83.42 |
| | Mean | 87.62 | 70.22 | 97.21 | 86.44 | 94.91 | 86.89 | 73.23 | 95.49 | 72.32 | 89.83 | 85.43 |
| | Weighted Mean | 86.51 | 68.50 | 96.70 | 85.26 | 94.65 | 85.78 | 72.05 | 94.59 | 68.76 | 89.43 | 84.22 |
| LLaMA-3 (mntp) | EOS | 83.54 | 79.07 | 87.28 | 82.84 | 90.40 | 82.43 | 71.68 | 87.28 | 75.51 | 82.18 | 82.22 |
| | Mean | 74.25 | 72.49 | 88.23 | 75.22 | 86.38 | 71.73 | 65.70 | 90.20 | 68.81 | 86.93 | 77.99 |
| | Weighted Mean | 73.81 | 72.61 | 88.56 | 73.77 | 85.63 | 71.31 | 62.93 | 88.54 | 68.65 | 82.80 | 76.86 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 58.41 | 51.81 | 73.50 | 53.66 | 59.23 | 52.69 | 50.26 | 68.65 | 54.12 | 60.34 | 58.27 |
| | Mean | 78.13 | 71.77 | 86.20 | 76.79 | 84.10 | 76.83 | 68.32 | 89.46 | 66.50 | 89.24 | 78.73 |
| | Weighted Mean | 75.70 | 67.31 | 85.99 | 73.56 | 83.57 | 73.77 | 63.61 | 88.95 | 66.50 | 89.42 | 76.84 |
| LLaMA-3 (mntp-supervised) | EOS | 81.23 | 61.06 | 92.78 | 75.94 | 85.27 | 76.25 | 57.49 | 89.62 | 74.66 | 82.59 | 77.69 |
| | Mean | 91.92 | 77.23 | 97.59 | 90.51 | 97.22 | 91.09 | 78.16 | 96.27 | 84.12 | 89.91 | 89.40 |
| | Weighted Mean | 91.22 | 72.61 | 97.23 | 89.66 | 96.75 | 90.40 | 75.79 | 95.84 | 82.13 | 89.37 | 88.10 |
| Mistral (mntp) | EOS | 89.95 | 74.33 | 91.18 | 80.78 | 90.30 | 89.42 | 73.13 | 92.12 | 77.06 | 86.92 | 84.52 |
| | Mean | 74.95 | 66.29 | 88.66 | 65.98 | 85.55 | 67.25 | 58.35 | 87.74 | 66.96 | 80.37 | 74.21 |
| | Weighted Mean | 75.04 | 63.23 | 88.60 | 64.12 | 84.03 | 67.21 | 55.95 | 85.84 | 65.40 | 78.00 | 72.74 |
| Mistral (mntp-unsup-simcse) | EOS | 77.22 | 59.67 | 75.84 | 72.47 | 77.09 | 77.36 | 62.71 | 80.77 | 62.16 | 75.05 | 72.03 |
| | Mean | 65.70 | 54.88 | 75.10 | 57.32 | 70.50 | 63.59 | 52.47 | 82.80 | 65.93 | 80.06 | 66.83 |
| | Weighted Mean | 60.17 | 52.27 | 74.74 | 52.62 | 66.31 | 58.04 | 48.66 | 80.45 | 63.03 | 76.56 | 63.28 |
| Mistral (mntp-supervised) | EOS | 93.35 | 70.83 | 95.09 | 86.16 | 93.55 | 93.09 | 77.55 | 94.28 | 74.09 | 88.26 | 86.62 |
| | Mean | 92.47 | 68.73 | 96.33 | 81.21 | 94.65 | 91.12 | 71.04 | 96.61 | 64.66 | 88.69 | 84.55 |
| | Weighted Mean | 91.63 | 66.34 | 96.25 | 78.16 | 94.74 | 89.97 | 67.99 | 96.41 | 64.20 | 88.48 | 83.42 |
| OpenAI-3-small | - | 95.13 | 73.48 | 96.57 | 93.18 | 95.37 | 93.04 | 84.27 | 95.12 | 95.12 | 72.22 | 89.35 |
| OpenAI-ada-002 | - | 93.23 | 79.13 | 95.80 | 93.09 | 92.55 | 92.83 | 85.63 | 92.03 | 93.68 | 66.76 | 88.47 |
| OpenAI-3-large | - | 91.52 | 60.66 | 95.83 | 88.18 | 94.46 | 88.46 | 80.80 | 91.25 | 91.93 | 52.44 | 83.55 |

Table 23: Average AUPRC(%) on DBpedia-3.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 84.38 | 76.84 | 95.78 | 79.43 | 91.92 | 81.01 | 62.90 | 90.64 | 59.40 | 94.51 | 81.68 |
| BERT | CLS | 56.39 | 53.00 | 93.20 | 53.46 | 82.94 | 52.16 | 43.97 | 85.30 | 77.76 | 82.14 | 68.03 |
| | Mean | 66.65 | 64.37 | 95.57 | 62.22 | 94.78 | 59.65 | 49.25 | 95.26 | 93.68 | 96.31 | 77.77 |
| LLaMA-2 (mntp) | EOS | 63.53 | 62.92 | 75.98 | 64.98 | 76.56 | 62.93 | 56.63 | 71.37 | 64.98 | 74.34 | 67.42 |
| | Mean | 77.93 | 78.58 | 89.65 | 78.81 | 79.86 | 76.45 | 67.97 | 76.15 | 71.52 | 73.70 | 77.06 |
| | Weighted Mean | 75.44 | 76.98 | 86.53 | 76.14 | 78.69 | 74.44 | 65.55 | 74.61 | 69.81 | 73.57 | 75.18 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 98.11 | 89.27 | 99.29 | 97.48 | 99.15 | 97.17 | 91.99 | 99.62 | 92.87 | 98.37 | 96.33 |
| | Mean | 93.58 | 85.59 | 99.12 | 93.39 | 98.82 | 93.21 | 83.83 | 99.48 | 91.94 | 99.45 | 93.84 |
| | Weighted Mean | 92.69 | 85.52 | 98.84 | 92.18 | 98.93 | 92.08 | 80.69 | 99.39 | 88.50 | 99.31 | 91.81 |
| LLaMA-2 (mntp-supervised) | EOS | 95.27 | 86.85 | 98.52 | 95.03 | 97.86 | 95.09 | 89.64 | 98.80 | 90.48 | 96.35 | 94.39 |
| | Mean | 93.62 | 78.90 | 99.39 | 92.83 | 98.66 | 92.98 | 79.76 | 99.49 | 89.48 | 97.46 | 92.26 |
| | Weighted Mean | 92.23 | 77.64 | 99.28 | 91.17 | 98.78 | 91.37 | 77.18 | 99.39 | 86.37 | 97.23 | 91.06 |
| LLaMA-3 (mntp) | EOS | 94.47 | 90.85 | 94.28 | 93.65 | 97.67 | 93.77 | 81.94 | 96.89 | 91.36 | 94.49 | 92.94 |
| | Mean | 90.97 | 89.47 | 98.51 | 91.29 | 98.50 | 88.50 | 80.95 | 98.98 | 87.38 | 97.67 | 92.22 |
| | Weighted Mean | 91.48 | 86.49 | 97.85 | 90.54 | 98.46 | 89.30 | 77.54 | 98.69 | 86.65 | 97.06 | 91.41 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 67.67 | 59.24 | 87.12 | 61.72 | 74.28 | 60.58 | 56.15 | 82.43 | 71.07 | 71.19 | 68.15 |
| | Mean | 93.80 | 88.44 | 98.35 | 93.12 | 98.00 | 93.14 | 86.48 | 98.81 | 90.62 | 98.53 | 93.93 |
| | Weighted Mean | 93.53 | 84.96 | 97.40 | 92.32 | 98.26 | 92.55 | 83.77 | 98.81 | 89.52 | 98.68 | 92.98 |
| LLaMA-3 (mntp-supervised) | EOS | 91.23 | 74.75 | 97.44 | 87.79 | 95.00 | 87.92 | 72.76 | 97.57 | 90.91 | 93.81 | 88.92 |
| | Mean | 96.08 | 82.14 | 99.72 | 95.00 | 99.57 | 95.06 | 82.75 | 99.65 | 94.02 | 97.98 | 94.20 |
| | Weighted Mean | 95.59 | 81.61 | 99.66 | 94.03 | 99.57 | 94.12 | 78.80 | 99.61 | 94.69 | 96.66 | 93.43 |
| Mistral (mntp) | EOS | 94.11 | 83.63 | 95.48 | 88.25 | 94.15 | 93.81 | 78.45 | 95.78 | 88.18 | 94.24 | 90.61 |
| | Mean | 91.69 | 83.24 | 97.71 | 86.06 | 97.86 | 84.67 | 73.51 | 98.34 | 86.42 | 95.97 | 89.55 |
| | Weighted Mean | 92.24 | 82.24 | 97.17 | 84.75 | 97.37 | 86.68 | 70.45 | 97.63 | 85.01 | 94.59 | 88.81 |
| Mistral (mntp-unsup-simcse) | EOS | 88.34 | 71.99 | 88.17 | 85.23 | 88.06 | 88.41 | 76.69 | 93.27 | 78.86 | 89.28 | 84.83 |
| | Mean | 85.15 | 71.33 | 93.98 | 75.44 | 91.81 | 81.86 | 66.00 | 97.08 | 90.69 | 96.29 | 84.69 |
| | Weighted Mean | 82.84 | 66.60 | 94.76 | 71.57 | 91.05 | 78.81 | 62.22 | 96.88 | 87.29 | 95.53 | 82.76 |
| Mistral (mntp-supervised) | EOS | 97.69 | 80.41 | 97.72 | 93.90 | 97.72 | 97.54 | 87.66 | 98.65 | 88.67 | 95.71 | 93.57 |
| | Mean | 96.16 | 74.42 | 98.62 | 87.00 | 98.14 | 95.06 | 74.35 | 99.26 | 89.07 | 96.70 | 90.88 |
| | Weighted Mean | 95.92 | 72.71 | 98.75 | 84.00 | 98.38 | 94.56 | 69.62 | 99.28 | 87.06 | 96.46 | 89.67 |
| OpenAI-3-small | - | 99.32 | 86.90 | 99.64 | 98.76 | 98.97 | 98.73 | 95.93 | 99.36 | 99.56 | 72.27 | 94.94 |
| OpenAI-ada-002 | - | 98.74 | 90.35 | 99.70 | 98.57 | 98.93 | 98.47 | 95.12 | 98.80 | 99.00 | 72.26 | 94.99 |
| OpenAI-3-large | - | 99.31 | 86.90 | 99.75 | 98.69 | 99.25 | 98.63 | 96.48 | 98.89 | 99.60 | 61.25 | 93.88 |

Table 24: Average AUPRC(%) on DBpedia-4.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 81.18 | 76.63 | 88.21 | 77.48 | 84.59 | 78.97 | 63.42 | 85.14 | 53.68 | 89.41 | 77.87 |
| BERT | CLS | 54.95 | 53.48 | 91.63 | 55.20 | 80.75 | 50.06 | 46.91 | 83.62 | 69.06 | 77.50 | 66.32 |
| | Mean | 63.21 | 60.98 | 91.18 | 60.88 | 90.85 | 58.34 | 52.17 | 91.38 | 87.34 | 92.40 | 74.87 |
| LLaMA-2 (mntp) | EOS | 64.64 | 61.95 | 68.04 | 63.36 | 66.47 | 62.36 | 53.91 | 64.18 | 62.12 | 64.99 | 63.20 |
| | Mean | 77.58 | 75.28 | 87.56 | 76.42 | 74.47 | 73.87 | 64.98 | 70.43 | 65.76 | 70.10 | 73.65 |
| | Weighted Mean | 73.56 | 71.47 | 85.04 | 72.50 | 72.78 | 69.67 | 61.31 | 72.37 | 62.59 | 70.23 | 71.15 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 96.38 | 84.68 | 96.88 | 94.76 | 98.37 | 94.60 | 84.25 | 98.50 | 84.37 | 96.00 | 92.88 |
| | Mean | 88.61 | 82.70 | 94.43 | 88.29 | 96.25 | 88.19 | 78.58 | 97.80 | 77.16 | 97.88 | 88.99 |
| | Weighted Mean | 88.12 | 80.58 | 92.67 | 87.80 | 96.22 | 87.78 | 77.58 | 97.45 | 72.26 | 97.50 | 87.80 |
| LLaMA-2 (mntp-supervised) | EOS | 95.53 | 86.69 | 97.93 | 94.91 | 97.56 | 94.95 | 90.97 | 98.07 | 86.18 | 94.80 | 93.76 |
| | Mean | 94.81 | 80.36 | 98.06 | 94.31 | 97.60 | 94.45 | 84.81 | 98.40 | 82.43 | 95.63 | 92.09 |
| | Weighted Mean | 93.79 | 81.21 | 97.57 | 93.17 | 97.58 | 93.36 | 82.50 | 98.09 | 80.14 | 95.49 | 91.29 |
| LLaMA-3 (mntp) | EOS | 87.73 | 85.19 | 84.18 | 87.17 | 93.16 | 87.17 | 78.02 | 91.04 | 80.95 | 87.86 | 86.25 |
| | Mean | 83.17 | 81.68 | 93.55 | 83.42 | 94.67 | 80.78 | 74.83 | 95.87 | 77.27 | 91.52 | 85.68 |
| | Weighted Mean | 84.41 | 83.17 | 91.88 | 84.26 | 94.21 | 82.77 | 75.13 | 95.26 | 76.27 | 91.93 | 85.93 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 63.32 | 57.09 | 82.21 | 58.58 | 70.83 | 57.31 | 54.35 | 77.00 | 63.63 | 66.61 | 65.09 |
| | Mean | 89.01 | 83.72 | 93.67 | 88.25 | 95.02 | 88.41 | 81.21 | 96.62 | 78.72 | 96.16 | 89.08 |
| | Weighted Mean | 88.63 | 82.89 | 92.28 | 87.86 | 95.08 | 88.21 | 79.80 | 96.38 | 75.06 | 95.99 | 88.22 |
| LLaMA-3 (mntp-supervised) | EOS | 93.65 | 80.20 | 97.59 | 91.88 | 93.84 | 91.93 | 79.57 | 95.63 | 82.51 | 93.33 | 90.01 |
| | Mean | 97.83 | 85.28 | 98.87 | 97.06 | 99.07 | 97.17 | 90.11 | 98.73 | 90.28 | 96.74 | 95.11 |
| | Weighted Mean | 97.60 | 83.74 | 98.40 | 96.94 | 98.93 | 97.08 | 88.95 | 98.49 | 87.83 | 96.31 | 94.43 |
| Mistral (mntp) | EOS | 92.25 | 76.55 | 93.23 | 81.31 | 92.57 | 91.63 | 72.25 | 93.09 | 77.12 | 89.02 | 85.90 |
| | Mean | 80.48 | 73.23 | 93.55 | 74.16 | 93.20 | 73.03 | 67.28 | 94.33 | 77.10 | 89.34 | 81.57 |
| | Weighted Mean | 82.42 | 71.93 | 92.55 | 75.00 | 92.51 | 76.52 | 67.23 | 93.26 | 77.29 | 90.45 | 81.92 |
| Mistral (mntp-unsup-simcse) | EOS | 85.05 | 66.72 | 83.64 | 78.77 | 84.71 | 85.15 | 67.26 | 88.29 | 70.74 | 82.82 | 79.31 |
| | Mean | 79.73 | 64.80 | 89.41 | 69.63 | 87.45 | 77.02 | 62.07 | 93.68 | 75.25 | 86.61 | 78.56 |
| | Weighted Mean | 77.22 | 62.70 | 88.99 | 66.63 | 86.46 | 74.15 | 59.51 | 92.89 | 72.63 | 90.07 | 77.12 |
| Mistral (mntp-supervised) | EOS | 95.71 | 76.27 | 95.25 | 90.71 | 95.71 | 95.58 | 82.85 | 96.35 | 79.82 | 92.45 | 90.07 |
| | Mean | 95.55 | 71.23 | 97.15 | 85.45 | 96.96 | 94.52 | 73.18 | 97.99 | 76.50 | 93.18 | 88.17 |
| | Weighted Mean | 95.11 | 69.45 | 96.96 | 82.96 | 97.03 | 93.85 | 69.47 | 97.92 | 73.20 | 87.36 | 86.33 |
| OpenAI-3-small | - | 98.09 | 81.37 | 96.59 | 96.90 | 98.40 | 96.83 | 91.47 | 97.37 | 98.33 | 72.22 | 92.76 |
| OpenAI-ada-002 | - | 97.94 | 88.21 | 97.99 | 97.85 | 97.79 | 97.86 | 94.83 | 92.09 | 98.14 | 68.95 | 93.16 |
| OpenAI-3-large | - | 98.48 | 81.20 | 98.10 | 97.54 | 98.49 | 97.51 | 94.24 | 98.76 | 98.99 | 48.41 | 91.17 |

Table 25: Average AUPRC(%) on WOS.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 31.51 | 25.02 | 45.53 | 25.09 | 23.64 | 28.64 | 24.04 | 26.10 | 36.96 | 29.00 | 29.55 |
| BERT | CLS | 29.69 | 29 99 | 38.01 | 29.06 | 30.06 | 28.99 | 27.90 | 30.19 | 31.78 | 31.39 | 30.71 |
| | Mean | 53.09 | 30.56 | 38.03 | 29.62 | 25.45 | 29.96 | 27.22 | 26.73 | 31.84 | 30.13 | 32.26 |
| LLaMA-2 (mntp) | EOS | 28.92 | 28.53 | 30.28 | 28.83 | 30.58 | 28.37 | 28.79 | 31.11 | 29.96 | 31.24 | 29.66 |
| | Mean | 30.87 | 29.41 | 31.08 | 28.98 | 27.53 | 29.89 | 28.17 | 26.33 | 29.00 | 28.61 | 28.99 |
| | Weighted Mean | 30.45 | 28.33 | 30.55 | 28.85 | 27.25 | 29.47 | 28.22 | 25.71 | 29.28 | 28.38 | 28.65 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 38.17 | 35.56 | 39.46 | 37.97 | 42.25 | 38.12 | 33.04 | 40.98 | 35.59 | 39.13 | 38.03 |
| | Mean | 26.58 | 25.76 | 40.63 | 24.32 | 26.71 | 24.04 | 22.77 | 28.52 | 31.04 | 32.62 | 28.30 |
| | Weighted Mean | 26.71 | 25.64 | 39.34 | 24.72 | 26.42 | 24.39 | 23.14 | 28.31 | 31.46 | 32.54 | 28.27 |
| LLaMA-2 (mntp-supervised) | EOS | 41.67 | 41.00 | 39.60 | 44.35 | 61.23 | 45.26 | 40.51 | 54.19 | 37.75 | 45.02 | 45.06 |
| | Mean | 26.32 | 27.50 | 49.36 | 24.62 | 37.38 | 24.27 | 22.88 | 38.23 | 32.98 | 36.86 | 32.04 |
| | Weighted Mean | 26.87 | 28.16 | 49.32 | 25.41 | 38.66 | 25.08 | 23.48 | 38.94 | 33.10 | 37.28 | 32.63 |
| LLaMA-3 (mntp) | EOS | 35.73 | 34.80 | 38.07 | 34.96 | 36.77 | 33.94 | 31.60 | 36.94 | 33.71 | 35.02 | 35.15 |
| | Mean | 27.58 | 25.43 | 41.89 | 24.66 | 26.26 | 24.33 | 22.65 | 31.01 | 33.41 | 32.93 | 29.02 |
| | Weighted Mean | 27.57 | 26.46 | 39.01 | 25.12 | 26.23 | 24.69 | 23.16 | 30.95 | 33.66 | 32.00 | 28.88 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 36.61 | 35.50 | 36.96 | 36.20 | 39.69 | 36.90 | 34.47 | 39.30 | 34.21 | 36.85 | 36.67 |
| | Mean | 25.05 | 23.63 | 45.34 | 22.59 | 25.76 | 22.23 | 21.58 | 28.87 | 30.93 | 31.77 | 27.78 |
| | Weighted Mean | 25.54 | 24.79 | 43.07 | 23.28 | 25.80 | 22.85 | 22.16 | 28.89 | 31.12 | 32.12 | 27.96 |
| LLaMA-3 (mntp-supervised) | EOS | 32.39 | 31.37 | 40.34 | 30.69 | 45.62 | 31.29 | 28.75 | 43.89 | 33.00 | 35.97 | 35.55 |
| | Mean | 36.77 | 36.70 | 43.71 | 36.90 | 62.51 | 38.19 | 32.49 | 56.01 | 37.74 | 44.33 | 42.54 |
| | Weighted Mean | 36.09 | 33.82 | 43.77 | 36.08 | 61.52 | 37.31 | 31.92 | 55.05 | 36.91 | 43.67 | 41.61 |
| Mistral (mntp) | EOS | 41.87 | 36.28 | 35.50 | 38.09 | 42.10 | 41.41 | 35.74 | 42.33 | 35.45 | 38.44 | 38.72 |
| | Mean | 28.00 | 27.37 | 36.12 | 28.36 | 26.94 | 25.03 | 25.68 | 30.83 | 31.08 | 32.63 | 29.20 |
| | Weighted Mean | 28.10 | 29.51 | 34.82 | 28.25 | 26.57 | 25.44 | 25.76 | 30.39 | 30.50 | 33.11 | 29.25 |
| Mistral (mntp-unsup-simcse) | EOS | 46.73 | 40.62 | 32.22 | 41.56 | 47.56 | 45.86 | 39.90 | 45.92 | 36.91 | 41.23 | 41.85 |
| | Mean | 23.52 | 22.77 | 44.12 | 21.91 | 24.90 | 20.96 | 21.44 | 27.94 | 29.60 | 28.84 | 26.60 |
| | Weighted Mean | 23.73 | 23.80 | 41.61 | 22.32 | 24.95 | 21.25 | 21.81 | 27.98 | 30.15 | 28.96 | 26.66 |
| Mistral (mntp-supervised) | EOS | 44.47 | 32.89 | 40.36 | 35.54 | 44.93 | 43.46 | 31.88 | 44.02 | 34.15 | 38.06 | 38.98 |
| | Mean | 31.23 | 28.32 | 43.12 | 26.89 | 38.44 | 29.15 | 24.84 | 39.50 | 33.63 | 36.56 | 33.17 |
| | Weighted Mean | 31.28 | 29.00 | 43.11 | 26.88 | 37.91 | 29.26 | 24.87 | 38.74 | 33.31 | 36.13 | 33.05 |
| OpenAI-3-small | - | 36.87 | 34.33 | 40.96 | 38.33 | 57.87 | 36.87 | 30.76 | 34.67 | 33.36 | 56.90 | 40.09 |
| OpenAI-ada-002 | - | 37.03 | 36.27 | 37.29 | 41.80 | 55.00 | 41.14 | 32.47 | 41.39 | 35.58 | 65.73 | 42.37 |
| OpenAI-3-large | - | 41.83 | 37.42 | 50.02 | 43.63 | 61.09 | 42.63 | 34.44 | 51.49 | 37.26 | 57.73 | 45.75 |

Table 26: Average AUPRC(%) on SST2.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 48.76 | 47.88 | 60.76 | 47.94 | 51.37 | 48.01 | 47.34 | 51.96 | 57.41 | 56.31 | 51.77 |
| BERT | CLS | 53.33 | 53.00 | 57.23 | 52.81 | 56.73 | 52.85 | 51.49 | 55.72 | 54.52 | 59.24 | 54.69 |
| | Mean | 52.39 | 51.76 | 57.69 | 52.02 | 57.90 | 51.91 | 50.49 | 56.52 | 53.76 | 60.52 | 54.50 |
| LLaMA-2 (mntp) | EOS | 51.25 | 51.52 | 54.59 | 51.14 | 56.92 | 50.89 | 50.55 | 58.29 | 52.29 | 62.98 | 54.04 |
| | Mean | 50.93 | 50.65 | 55.74 | 50.18 | 55.74 | 50.29 | 49.84 | 57.48 | 52.58 | 61.24 | 53.47 |
| | Weighted Mean | 50.41 | 50.46 | 55.68 | 49.87 | 55.82 | 49.85 | 49.48 | 56.98 | 51.97 | 61.97 | 53.25 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 58.39 | 55.96 | 67.03 | 58.43 | 71.81 | 58.41 | 53.16 | 76.37 | 60.25 | 69.90 | 62.97 |
| | Mean | 57.56 | 55.53 | 68.62 | 57.48 | 72.00 | 57.56 | 53.50 | 77.38 | 59.68 | 70.48 | 62.98 |
| | Weighted Mean | 57.33 | 56.28 | 67.83 | 57.47 | 71.56 | 57.56 | 53.40 | 77.05 | 60.53 | 69.60 | 62.86 |
| LLaMA-2 (mntp-supervised) | EOS | 65.79 | 63.93 | 65.37 | 64.92 | 69.68 | 64.38 | 58.96 | 72.44 | 64.67 | 71.36 | 66.15 |
| | Mean | 68.73 | 66.20 | 69.40 | 68.76 | 76.86 | 69.02 | 59.38 | 77.93 | 67.34 | 75.24 | 69.89 |
| | Weighted Mean | 66.82 | 64.44 | 67.20 | 66.45 | 74.09 | 66.54 | 58.12 | 75.92 | 65.33 | 73.64 | 67.85 |
| LLaMA-3 (mntp) | EOS | 55.20 | 54.64 | 58.17 | 54.67 | 61.13 | 54.91 | 53.12 | 64.89 | 55.91 | 62.91 | 57.55 |
| | Mean | 55.11 | 54.24 | 61.87 | 54.74 | 64.08 | 54.86 | 52.41 | 69.33 | 58.29 | 66.11 | 59.10 |
| | Weighted Mean | 55.32 | 54.66 | 59.58 | 55.27 | 65.01 | 55.59 | 52.70 | 69.98 | 57.56 | 67.31 | 59.30 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 55.73 | 54.91 | 59.83 | 55.29 | 60.09 | 55.43 | 53.99 | 64.23 | 57.20 | 62.55 | 57.92 |
| | Mean | 55.87 | 54.02 | 64.04 | 54.94 | 68.59 | 55.47 | 51.46 | 75.04 | 59.66 | 69.14 | 60.82 |
| | Weighted Mean | 55.54 | 54.65 | 63.09 | 54.84 | 66.18 | 55.18 | 52.16 | 73.52 | 59.28 | 67.66 | 60.21 |
| LLaMA-3 (mntp-supervised) | EOS | 67.52 | 64.57 | 65.43 | 65.72 | 70.51 | 65.54 | 58.28 | 73.81 | 63.76 | 72.35 | 66.75 |
| | Mean | 67.99 | 64.34 | 67.52 | 67.44 | 75.48 | 67.82 | 57.17 | 76.70 | 64.79 | 74.60 | 68.39 |
| | Weighted Mean | 67.70 | 63.21 | 66.81 | 66.93 | 74.18 | 67.21 | 57.19 | 76.19 | 64.63 | 74.59 | 67.86 |
| Mistral (mntp) | EOS | 64.76 | 51.30 | 58.21 | 51.78 | 59.38 | 63.94 | 49.48 | 61.70 | 55.86 | 60.67 | 57.71 |
| | Mean | 57.15 | 51.24 | 61.53 | 50.75 | 58.20 | 52.03 | 49.44 | 62.21 | 55.55 | 60.71 | 55.88 |
| | Weighted Mean | 57.83 | 50.50 | 60.07 | 50.02 | 57.72 | 52.73 | 48.56 | 61.43 | 55.81 | 60.48 | 55.52 |
| Mistral (mntp-unsup-simcse) | EOS | 74.50 | 55.19 | 61.85 | 54.09 | 65.91 | 73.43 | 49.87 | 68.59 | 58.93 | 65.14 | 62.75 |
| | Mean | 72.69 | 53.82 | 65.79 | 54.33 | 69.94 | 70.00 | 51.22 | 74.57 | 58.62 | 68.93 | 63.99 |
| | Weighted Mean | 72.67 | 53.18 | 65.05 | 54.62 | 69.59 | 70.11 | 51.09 | 73.64 | 59.47 | 68.20 | 63.76 |
| Mistral (mntp-supervised) | EOS | 79.81 | 61.59 | 64.52 | 62.99 | 71.83 | 79.50 | 57.27 | 72.87 | 65.25 | 72.12 | 68.78 |
| | Mean | 82.74 | 60.07 | 68.49 | 62.28 | 76.65 | 82.03 | 54.84 | 76.45 | 65.06 | 75.30 | 70.39 |
| | Weighted Mean | 81.97 | 58.75 | 68.08 | 61.70 | 74.97 | 81.29 | 54.90 | 75.37 | 64.93 | 73.48 | 69.54 |
| OpenAI-3-small | - | 67.77 | 61.04 | 71.04 | 67.22 | 77.75 | 67.32 | 55.83 | 71.69 | 63.72 | 75.95 | 67.93 |
| OpenAI-ada-002 | - | 64.27 | 61.13 | 69.52 | 64.78 | 76.34 | 64.64 | 56.86 | 65.86 | 63.08 | 75.95 | 66.24 |
| OpenAI-3-large | - | 68.75 | 61.39 | 73.62 | 68.44 | 79.29 | 68.27 | 57.46 | 72.77 | 65.21 | 75.95 | 69.12 |

Table 27: Average AUPRC(%) on SMS-SPAM.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 26.52 | 22.75 | 25.21 | 23.05 | 23.86 | 24.57 | 22.44 | 22.00 | 49.75 | 23.83 | 26.40 |
| BERT | CLS | 32.75 | 21.45 | 28.19 | 22.44 | 19.23 | 27.88 | 20.63 | 21.01 | 34.94 | 22.18 | 25.07 |
| | Mean | 28.47 | 25.40 | 30.48 | 26.59 | 21.24 | 26.97 | 20.99 | 23.91 | 74.98 | 33.71 | 31.27 |
| LLaMA-2 (mntp) | EOS | 56.19 | 48.91 | 41.10 | 50.68 | 22.32 | 49.64 | 30.43 | 28.45 | 25.98 | 35.36 | 38.91 |
| | Mean | 63.63 | 56.71 | 30.73 | 54.21 | 19.45 | 57.92 | 30.47 | 19.16 | 27.26 | 20.75 | 38.03 |
| | Weighted Mean | 71.68 | 65.50 | 43.99 | 64.01 | 19.60 | 68.09 | 33.36 | 20.57 | 24.20 | 22.68 | 43.37 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 61.66 | 50.00 | 63.14 | 61.91 | 61.57 | 59.93 | 36.33 | 64.90 | 50.76 | 63.95 | 57.42 |
| | Mean | 62.99 | 46.33 | 48.39 | 61.06 | 64.61 | 64.74 | 29.74 | 65.77 | 69.10 | 79.97 | 59.27 |
| | Weighted Mean | 79.89 | 57.05 | 54.27 | 78.50 | 74.56 | 80.91 | 42.13 | 75.06 | 62.38 | 85.07 | 68.98 |
| LLaMA-2 (mntp-supervised) | EOS | 84.45 | 80.54 | 50.55 | 85.05 | 84.44 | 86.20 | 69.80 | 83.95 | 66.07 | 86.31 | 77.74 |
| | Mean | 60.93 | 53.46 | 41.79 | 64.26 | 66.77 | 64.70 | 36.08 | 69.68 | 49.51 | 69.85 | 57.70 |
| | Weighted Mean | 72.32 | 62.84 | 42.60 | 75.99 | 77.48 | 76.51 | 48.46 | 78.71 | 50.19 | 79.22 | 66.43 |
| LLaMA-3 (mntp) | EOS | 82.54 | 74.90 | 77.28 | 79.21 | 85.08 | 80.84 | 62.23 | 82.88 | 78.11 | 82.48 | 78.56 |
| | Mean | 48.64 | 42.12 | 40.25 | 41.28 | 46.30 | 53.41 | 27.09 | 39.11 | 80.10 | 48.88 | 46.72 |
| | Weighted Mean | 63.85 | 51.23 | 43.49 | 54.51 | 58.34 | 68.75 | 34.05 | 49.34 | 79.84 | 64.32 | 56.77 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 80.98 | 75.79 | 68.32 | 80.48 | 81.12 | 80.69 | 70.80 | 83.18 | 69.99 | 81.10 | 77.25 |
| | Mean | 57.89 | 46.62 | 54.39 | 56.49 | 59.68 | 59.59 | 32.27 | 55.07 | 85.59 | 78.69 | 58.63 |
| | Weighted Mean | 73.86 | 63.81 | 52.10 | 72.90 | 73.87 | 76.10 | 45.45 | 69.05 | 89.69 | 87.39 | 70.42 |
| LLaMA-3 (mntp-supervised) | EOS | 87.20 | 79.38 | 73.64 | 85.61 | 85.59 | 85.94 | 73.28 | 88.27 | 71.90 | 86.45 | 81.73 |
| | Mean | 80.35 | 69.59 | 66.37 | 79.66 | 82.28 | 80.11 | 59.53 | 82.63 | 72.60 | 82.60 | 75.57 |
| | Weighted Mean | 80.25 | 67.09 | 57.96 | 79.71 | 82.72 | 80.12 | 59.12 | 82.78 | 71.62 | 83.67 | 74.50 |
| Mistral (mntp) | EOS | 77.70 | 61.84 | 87.20 | 65.32 | 76.59 | 77.94 | 48.68 | 78.38 | 76.39 | 77.55 | 72.76 |
| | Mean | 68.80 | 51.29 | 71.74 | 54.98 | 48.62 | 64.86 | 32.43 | 38.91 | 63.48 | 54.56 | 54.97 |
| | Weighted Mean | 65.89 | 54.45 | 74.03 | 58.65 | 46.02 | 63.07 | 35.36 | 41.51 | 74.71 | 66.47 | 58.02 |
| Mistral (mntp-unsup-simcse) | EOS | 77.92 | 49.97 | 79.25 | 69.17 | 75.09 | 78.38 | 45.10 | 76.93 | 59.67 | 72.02 | 68.35 |
| | Mean | 67.61 | 39.54 | 51.01 | 46.39 | 54.71 | 66.82 | 26.13 | 56.59 | 86.00 | 82.81 | 57.76 |
| | Weighted Mean | 83.56 | 52.81 | 57.07 | 69.16 | 71.18 | 83.39 | 36.67 | 70.82 | 86.68 | 83.77 | 69.51 |
| Mistral (mntp-supervised) | EOS | 73.02 | 54.45 | 61.60 | 69.16 | 68.55 | 73.90 | 48.16 | 72.51 | 78.85 | 78.03 | 67.82 |
| | Mean | 38.66 | 33.45 | 44.74 | 38.33 | 37.60 | 37.32 | 23.82 | 46.16 | 65.17 | 59.67 | 42.49 |
| | Weighted Mean | 45.98 | 48.04 | 43.99 | 52.08 | 42.87 | 45.54 | 28.78 | 51.40 | 61.39 | 67.71 | 48.78 |
| OpenAI-3-small | - | 33.86 | 29.66 | 66.32 | 31.71 | 42.08 | 33.22 | 19.34 | 32.84 | 30.04 | 64.73 | 38.38 |
| OpenAI-ada-002 | - | 63.49 | 48.74 | 67.03 | 63.55 | 60.52 | 64.81 | 36.99 | 66.92 | 39.20 | 32.56 | 54.38 |
| OpenAI-3-large | - | 44.12 | 34.80 | 65.74 | 37.84 | 51.06 | 43.22 | 19.89 | 41.39 | 44.83 | 37.07 | 42.00 |

Table 28: Average AUPRC(%) on Reuters21578.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 84.40 | 84.32 | 86.11 | 84.53 | 87.87 | 83.97 | 75.58 | 87.37 | 68.23 | 87.91 | 83.03 |
| BERT | CLS | 67.87 | 64.83 | 83.76 | 64.11 | 65.62 | 64.13 | 58.90 | 69.48 | 60.42 | 69.37 | 66.85 |
|  | Mean | 82.12 | 82.44 | 92.20 | 81.75 | 87.02 | 79.84 | 71.82 | 87.01 | 82.39 | 90.75 | 83.73 |
| LLaMA-2 (mntp) | EOS | 76.00 | 76.47 | 69.21 | 77.28 | 73.92 | 76.74 | 73.87 | 75.79 | 76.00 | 77.65 | 75.29 |
|  | Mean | 72.41 | 75.13 | 73.37 | 73.73 | 72.38 | 71.49 | 68.02 | 78.75 | 68.66 | 72.81 | 72.68 |
|  | Weighted Mean | 72.38 | 75.38 | 72.96 | 73.83 | 72.24 | 71.96 | 68.32 | 73.98 | 69.47 | 73.22 | 72.37 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 95.64 | 92.78 | 90.29 | 95.62 | 96.33 | 95.27 | 88.33 | 97.07 | 91.50 | 95.81 | 93.86 |
|  | Mean | 88.99 | 87.31 | 86.94 | 87.30 | 89.34 | 87.28 | 74.16 | 92.17 | 96.04 | 94.97 | 88.45 |
|  | Weighted Mean | 88.05 | 84.49 | 85.66 | 86.46 | 88.77 | 86.49 | 73.00 | 91.39 | 95.40 | 94.88 | 87.46 |
| LLaMA-2 (mntp-supervised) | EOS | 96.33 | 90.37 | 85.10 | 96.11 | 97.01 | 96.05 | 91.61 | 97.39 | 88.14 | 95.32 | 93.34 |
|  | Mean | 94.93 | 89.93 | 87.57 | 95.18 | 95.97 | 95.19 | 86.36 | 97.49 | 92.37 | 95.51 | 93.05 |
|  | Weighted Mean | 94.39 | 88.56 | 86.19 | 94.76 | 95.44 | 94.89 | 85.19 | 96.94 | 91.52 | 94.88 | 92.28 |
| LLaMA-3 (mntp) | EOS | 87.96 | 85.53 | 86.73 | 87.88 | 92.44 | 87.31 | 74.79 | 91.52 | 74.10 | 89.32 | 85.76 |
|  | Mean | 89.44 | 88.82 | 90.67 | 89.72 | 90.85 | 87.49 | 75.96 | 92.53 | 79.37 | 92.04 | 87.69 |
|  | Weighted Mean | 86.10 | 86.41 | 89.34 | 86.16 | 89.62 | 83.18 | 72.19 | 91.39 | 76.59 | 90.77 | 85.18 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 85.72 | 80.32 | 81.40 | 83.25 | 86.27 | 83.96 | 76.75 | 86.52 | 80.93 | 83.34 | 82.85 |
|  | Mean | 90.13 | 87.92 | 88.05 | 88.33 | 88.72 | 88.24 | 77.47 | 91.17 | 94.85 | 93.31 | 88.82 |
|  | Weighted Mean | 88.75 | 86.77 | 87.05 | 86.77 | 88.07 | 87.06 | 75.00 | 90.27 | 92.99 | 92.20 | 87.49 |
| LLaMA-3 (mntp-supervised) | EOS | 94.97 | 89.83 | 94.94 | 93.98 | 92.89 | 93.41 | 89.10 | 95.18 | 92.99 | 93.54 | 93.08 |
|  | Mean | 96.88 | 92.35 | 91.75 | 96.86 | 96.61 | 96.78 | 89.91 | 97.75 | 94.54 | 95.88 | 94.93 |
|  | Weighted Mean | 96.26 | 90.88 | 90.64 | 96.30 | 95.79 | 96.23 | 88.97 | 97.32 | 93.13 | 95.40 | 94.09 |
| Mistral (mntp) | EOS | 85.42 | 87.24 | 81.89 | 90.52 | 84.50 | 85.56 | 77.16 | 86.44 | 79.52 | 89.43 | 84.77 |
|  | Mean | 83.79 | 77.66 | 90.04 | 75.48 | 82.27 | 77.54 | 61.81 | 84.78 | 77.38 | 83.89 | 79.46 |
|  | Weighted Mean | 81.61 | 72.97 | 89.65 | 72.16 | 81.48 | 74.20 | 61.98 | 84.14 | 75.92 | 82.37 | 77.65 |
| Mistral (mntp-unsup-simcse) | EOS | 93.80 | 82.86 | 84.58 | 91.27 | 93.71 | 93.81 | 82.21 | 93.63 | 86.30 | 91.08 | 89.32 |
|  | Mean | 86.03 | 76.53 | 85.31 | 78.52 | 85.82 | 85.26 | 66.56 | 88.21 | 93.69 | 91.60 | 83.75 |
|  | Weighted Mean | 84.72 | 75.24 | 84.77 | 76.05 | 85.02 | 83.98 | 64.90 | 87.55 | 90.15 | 90.52 | 82.29 |
| Mistral (mntp-supervised) | EOS | 95.64 | 88.88 | 83.58 | 95.17 | 95.46 | 95.78 | 87.37 | 96.11 | 92.06 | 94.20 | 92.42 |
|  | Mean | 97.70 | 89.84 | 89.21 | 96.60 | 96.62 | 97.80 | 87.06 | 97.74 | 92.70 | 96.26 | 94.15 |
|  | Weighted Mean | 96.74 | 85.92 | 87.43 | 95.63 | 95.47 | 96.91 | 84.73 | 96.97 | 93.47 | 95.17 | 92.84 |
| OpenAI-3-small | - | 98.79 | 95.85 | 93.30 | 98.75 | 97.20 | 98.58 | 93.02 | 98.20 | 98.61 | 81.16 | 95.35 |
| OpenAI-ada-002 | - | 98.17 | 92.22 | 90.86 | 97.62 | 95.85 | 97.36 | 86.44 | 93.20 | 97.95 | 81.16 | 93.08 |
| OpenAI-3-large | - | 99.52 | 93.80 | 96.38 | 99.43 | 98.18 | 99.38 | 95.33 | 98.36 | 99.45 | 81.16 | 96.10 |

Table 29: Average AUPRC(%) on IMDB.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 59.62 | 59.62 | 62.26 | 59.50 | 59.62 | 59.62 | 57.34 | 60.77 | 60.38 | 60.98 | 59.97 |
| BERT | CLS | 59.02 | 59.64 | 59.57 | 58.94 | 60.02 | 59.19 | 56.97 | 59.39 | 62.74 | 60.45 | 59.59 |
| | Mean | 58.52 | 58.70 | 61.10 | 58.24 | 60.29 | 58.07 | 55.39 | 59.87 | 59.81 | 59.88 | 58.99 |
| LLaMA-2 (mntp) | EOS | 61.30 | 61.22 | 62.38 | 61.30 | 61.69 | 61.36 | 61.09 | 61.60 | 61.64 | 62.11 | 61.57 |
| | Mean | 60.30 | 60.38 | 61.40 | 60.27 | 61.23 | 60.57 | 59.92 | 61.26 | 61.37 | 61.91 | 60.86 |
| | Weighted Mean | 60.53 | 60.52 | 61.84 | 60.47 | 61.33 | 60.78 | 60.14 | 61.32 | 61.36 | 61.96 | 61.02 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 59.52 | 59.14 | 65.87 | 58.66 | 64.56 | 58.63 | 54.35 | 65.64 | 61.06 | 62.37 | 60.98 |
| | Mean | 59.80 | 60.34 | 65.22 | 59.68 | 65.02 | 59.66 | 55.75 | 65.30 | 61.84 | 62.09 | 61.47 |
| | Weighted Mean | 60.25 | 59.90 | 65.85 | 60.17 | 64.75 | 60.16 | 56.39 | 65.37 | 62.14 | 62.57 | 61.76 |
| LLaMA-2 (mntp-supervised) | EOS | 63.15 | 61.64 | 71.39 | 60.90 | 67.06 | 60.61 | 54.33 | 65.46 | 63.53 | 63.65 | 63.17 |
| | Mean | 62.26 | 62.64 | 68.26 | 61.50 | 74.37 | 61.52 | 53.82 | 70.51 | 63.71 | 65.35 | 64.39 |
| | Weighted Mean | 62.25 | 61.97 | 68.42 | 61.47 | 74.17 | 61.49 | 53.85 | 70.06 | 63.52 | 65.42 | 64.26 |
| LLaMA-3 (mntp) | EOS | 58.50 | 58.08 | 61.20 | 58.08 | 59.74 | 58.14 | 56.36 | 61.36 | 60.43 | 60.75 | 59.26 |
| | Mean | 58.80 | 58.94 | 62.60 | 58.86 | 61.90 | 58.91 | 56.02 | 62.97 | 60.99 | 62.27 | 60.23 |
| | Weighted Mean | 59.05 | 58.87 | 62.93 | 59.12 | 61.66 | 59.14 | 56.54 | 62.92 | 61.25 | 62.23 | 60.37 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 59.57 | 59.75 | 63.77 | 59.23 | 60.89 | 59.10 | 57.13 | 61.49 | 61.16 | 60.61 | 60.27 |
| | Mean | 59.94 | 59.57 | 64.40 | 59.97 | 65.21 | 60.04 | 55.05 | 64.48 | 61.78 | 62.31 | 61.28 |
| | Weighted Mean | 60.00 | 60.26 | 64.29 | 60.05 | 64.60 | 60.12 | 55.30 | 64.11 | 61.78 | 62.07 | 61.26 |
| LLaMA-3 (mntp-supervised) | EOS | 63.81 | 63.68 | 69.31 | 63.06 | 68.75 | 62.75 | 55.51 | 66.63 | 64.99 | 64.49 | 64.30 |
| | Mean | 59.23 | 60.14 | 67.79 | 59.17 | 74.23 | 59.16 | 53.23 | 70.69 | 62.44 | 64.65 | 63.07 |
| | Weighted Mean | 59.36 | 59.79 | 67.60 | 59.30 | 74.12 | 59.29 | 53.11 | 70.46 | 62.38 | 64.98 | 63.04 |
| Mistral (mntp) | EOS | 57.59 | 57.09 | 59.22 | 57.35 | 57.66 | 57.52 | 55.71 | 58.28 | 59.90 | 58.61 | 57.89 |
| | Mean | 61.57 | 61.61 | 64.62 | 61.53 | 64.07 | 61.70 | 58.81 | 65.18 | 62.14 | 64.08 | 62.53 |
| | Weighted Mean | 61.79 | 61.35 | 64.80 | 61.71 | 63.81 | 61.85 | 59.25 | 64.99 | 62.26 | 63.71 | 62.55 |
| Mistral (mntp-unsup-simcse) | EOS | 67.81 | 61.45 | 67.25 | 60.69 | 67.14 | 66.53 | 54.89 | 65.41 | 63.34 | 63.92 | 63.84 |
| | Mean | 64.12 | 62.88 | 67.74 | 63.00 | 67.53 | 63.73 | 59.61 | 67.39 | 63.95 | 64.85 | 64.48 |
| | Weighted Mean | 64.25 | 63.60 | 68.19 | 63.19 | 67.14 | 63.86 | 59.83 | 67.19 | 63.80 | 64.90 | 64.60 |
| Mistral (mntp-supervised) | EOS | 71.63 | 65.36 | 71.56 | 66.39 | 69.71 | 70.77 | 55.88 | 66.90 | 65.43 | 67.58 | 67.12 |
| | Mean | 63.51 | 60.36 | 66.10 | 58.42 | 69.91 | 61.39 | 52.85 | 67.11 | 62.50 | 64.72 | 62.69 |
| | Weighted Mean | 63.88 | 61.53 | 66.13 | 58.81 | 69.75 | 61.82 | 53.02 | 66.82 | 62.81 | 64.77 | 62.93 |
| OpenAI-3-small | - | 62.69 | 61.83 | 72.47 | 61.90 | 70.79 | 61.77 | 53.42 | 64.98 | 61.77 | 81.25 | 65.29 |
| OpenAI-ada-002 | - | 61.87 | 62.27 | 69.91 | 61.37 | 73.20 | 61.21 | 54.76 | 62.70 | 60.73 | 76.26 | 64.43 |
| OpenAI-3-large | - | 64.63 | 62.85 | 73.06 | 63.92 | 77.58 | 63.75 | 51.36 | 68.31 | 62.91 | 81.25 | 66.96 |

Table 30: Average AUPRC(%) on Enron.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 37.93 | 36.93 | 40.82 | 37.03 | 44.10 | 38.10 | 34.42 | 41.41 | 34.71 | 49.17 | 39.46 |
| BERT | CLS | 32.16 | 32.18 | 57.63 | 32.31 | 41.17 | 31.81 | 31.74 | 37.95 | 27.95 | 57.25 | 38.22 |
|  | Mean | 31.62 | 31.79 | 57.16 | 31.83 | 35.86 | 31.86 | 30.67 | 34.47 | 34.99 | 52.55 | 37.28 |
| LLaMA-2 (mntp) | EOS | 34.63 | 34.84 | 44.46 | 34.71 | 36.72 | 35.47 | 33.59 | 39.80 | 32.79 | 65.89 | 39.29 |
|  | Mean | 32.36 | 33.56 | 57.62 | 32.55 | 37.46 | 32.05 | 31.93 | 36.79 | 32.99 | 47.38 | 37.47 |
|  | Weighted Mean | 32.56 | 33.13 | 56.84 | 32.49 | 35.75 | 32.16 | 31.98 | 36.02 | 33.88 | 47.08 | 37.19 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 54.05 | 44.27 | 83.79 | 46.85 | 72.03 | 46.89 | 42.63 | 95.97 | 66.67 | 91.42 | 64.46 |
|  | Mean | 38.49 | 37.37 | 49.46 | 37.91 | 49.34 | 38.35 | 34.83 | 89.67 | 63.77 | 83.36 | 52.26 |
|  | Weighted Mean | 36.62 | 36.34 | 51.16 | 35.72 | 47.99 | 35.92 | 32.97 | 89.19 | 62.83 | 84.83 | 51.36 |
| LLaMA-2 (mntp-supervised) | EOS | 55.39 | 47.70 | 78.19 | 47.78 | 73.45 | 45.84 | 44.82 | 96.12 | 65.16 | 93.65 | 64.81 |
|  | Mean | 44.31 | 44.35 | 71.67 | 43.44 | 59.99 | 42.76 | 41.61 | 94.61 | 61.62 | 90.76 | 59.51 |
|  | Weighted Mean | 44.89 | 43.73 | 72.18 | 43.52 | 64.10 | 42.57 | 41.42 | 95.27 | 59.66 | 91.11 | 59.85 |
| LLaMA-3 (mntp) | EOS | 43.86 | 42.10 | 54.53 | 42.32 | 66.56 | 42.59 | 39.57 | 89.34 | 56.05 | 87.60 | 56.45 |
|  | Mean | 36.76 | 36.60 | 54.71 | 36.99 | 46.04 | 36.78 | 34.72 | 71.43 | 44.84 | 72.35 | 47.12 |
|  | Weighted Mean | 34.17 | 34.73 | 53.76 | 34.26 | 43.67 | 33.80 | 32.40 | 70.51 | 43.06 | 70.64 | 45.10 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 67.41 | 46.72 | 74.17 | 46.90 | 74.98 | 52.98 | 43.72 | 94.54 | 69.06 | 90.13 | 66.06 |
|  | Mean | 38.74 | 36.92 | 51.74 | 37.09 | 47.95 | 37.94 | 33.32 | 86.17 | 62.29 | 82.33 | 51.45 |
|  | Weighted Mean | 37.01 | 34.93 | 53.45 | 34.61 | 46.58 | 35.06 | 31.28 | 86.20 | 62.23 | 82.01 | 50.34 |
| LLaMA-3 (mntp-supervised) | EOS | 58.44 | 48.15 | 75.28 | 49.34 | 61.85 | 48.50 | 46.40 | 91.00 | 65.31 | 89.76 | 63.40 |
|  | Mean | 58.44 | 48.20 | 75.28 | 49.34 | 61.85 | 48.50 | 46.40 | 91.00 | 62.76 | 90.18 | 63.20 |
|  | Weighted Mean | 58.44 | 48.50 | 75.28 | 49.34 | 61.85 | 48.50 | 46.40 | 91.00 | 63.96 | 89.96 | 63.32 |
| Mistral (mntp) | EOS | 94.16 | 40.55 | 60.58 | 41.59 | 64.22 | 93.79 | 37.99 | 80.67 | 55.98 | 85.80 | 65.53 |
|  | Mean | 42.34 | 36.85 | 55.35 | 36.87 | 45.61 | 39.65 | 34.17 | 62.80 | 44.85 | 65.97 | 46.45 |
|  | Weighted Mean | 40.89 | 34.49 | 57.19 | 34.88 | 43.04 | 36.54 | 33.24 | 60.47 | 42.85 | 64.79 | 44.84 |
| Mistral (mntp-unsup-simcse) | EOS | 97.30 | 53.27 | 84.83 | 56.86 | 80.74 | 96.84 | 51.76 | 93.32 | 73.20 | 92.87 | 78.10 |
|  | Mean | 89.95 | 33.26 | 50.40 | 33.26 | 47.74 | 61.91 | 30.34 | 84.12 | 58.23 | 85.15 | 57.44 |
|  | Weighted Mean | 90.04 | 31.22 | 53.07 | 31.26 | 46.20 | 62.02 | 29.05 | 83.05 | 58.46 | 84.15 | 56.85 |
| Mistral (mntp-supervised) | EOS | 94.47 | 40.33 | 69.56 | 40.91 | 57.48 | 94.00 | 38.50 | 82.82 | 67.20 | 88.87 | 67.41 |
|  | Mean | 92.61 | 40.60 | 70.03 | 41.95 | 55.81 | 88.20 | 39.76 | 90.56 | 60.06 | 89.10 | 66.87 |
|  | Weighted Mean | 92.87 | 41.47 | 69.11 | 41.19 | 56.01 | 90.10 | 38.92 | 90.44 | 59.61 | 88.68 | 66.84 |
| OpenAI-3-small | - | 46.08 | 47.25 | 67.22 | 45.77 | 62.80 | 45.32 | 42.15 | 61.19 | 42.70 | 65.93 | 52.64 |
| OpenAI-ada-002 | - | 52.83 | 51.17 | 84.73 | 52.33 | 74.67 | 50.96 | 48.17 | 61.25 | 52.28 | 65.92 | 59.43 |
| OpenAI-3-large | - | 51.20 | 48.38 | 76.08 | 50.12 | 67.18 | 49.62 | 45.96 | 62.01 | 52.27 | 65.94 | 56.88 |

Table 31: Average AUPRC(%) on 20Newsgroups.

| Embedding | Pooling | OCSVM | IForest | LOF | PCA | KNN | KDE | ECOD | AE | DSVDD | DPAD | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GloVe | Mean | 13.13 | 13.61 | 14.79 | 13.52 | 15.12 | 13.16 | 13.22 | 13.46 | 12.16 | 14.07 | 13.62 |
| BERT | CLS | 13.26 | 13.82 | 21.37 | 13.90 | 14.39 | 13.44 | 13.63 | 15.17 | 11.60 | 14.50 | 14.51 |
| | Mean | 13.09 | 13.74 | 16.96 | 13.47 | 13.84 | 13.32 | 12.98 | 13.98 | 14.34 | 14.80 | 14.05 |
| LLaMA-2 (mntp) | EOS | 11.46 | 11.13 | 10.56 | 10.83 | 11.70 | 10.55 | 10.89 | 11.07 | 10.57 | 11.10 | 10.99 |
| | Mean | 13.89 | 12.18 | 12.80 | 12.41 | 11.77 | 12.82 | 12.33 | 11.79 | 11.59 | 11.56 | 12.31 |
| | Weighted Mean | 13.82 | 12.03 | 12.74 | 12.30 | 11.67 | 12.58 | 12.24 | 11.41 | 11.29 | 11.43 | 12.15 |
| LLaMA-2 (mntp-unsup-simcse) | EOS | 15.61 | 16.01 | 13.98 | 15.43 | 14.46 | 15.19 | 13.84 | 14.16 | 14.10 | 14.11 | 14.69 |
| | Mean | 11.63 | 11.82 | 11.53 | 12.03 | 14.61 | 12.27 | 11.54 | 13.32 | 14.51 | 14.52 | 12.78 |
| | Weighted Mean | 11.84 | 13.01 | 11.14 | 12.26 | 14.55 | 12.55 | 11.73 | 13.31 | 14.82 | 14.88 | 13.01 |
| LLaMA-2 (mntp-supervised) | EOS | 16.31 | 13.70 | 19.20 | 14.99 | 16.62 | 14.98 | 12.94 | 21.36 | 16.41 | 17.91 | 16.44 |
| | Mean | 13.36 | 12.48 | 15.30 | 14.01 | 22.88 | 14.01 | 12.99 | 23.73 | 14.48 | 16.74 | 16.00 |
| | Weighted Mean | 13.98 | 13.41 | 15.99 | 14.50 | 25.15 | 14.45 | 13.36 | 24.81 | 14.30 | 16.78 | 16.67 |
| LLaMA-3 (mntp) | EOS | 19.05 | 17.70 | 18.89 | 18.93 | 18.55 | 18.92 | 17.24 | 18.49 | 16.70 | 19.87 | 18.43 |
| | Mean | 14.17 | 13.48 | 13.69 | 14.10 | 15.36 | 14.45 | 13.23 | 14.69 | 16.12 | 16.01 | 14.53 |
| | Weighted Mean | 13.78 | 13.43 | 13.16 | 13.67 | 14.65 | 14.08 | 12.90 | 13.83 | 15.23 | 14.97 | 13.97 |
| LLaMA-3 (mntp-unsup-simcse) | EOS | 13.33 | 13.15 | 13.84 | 12.94 | 13.92 | 12.72 | 12.62 | 13.96 | 13.46 | 12.86 | 13.28 |
| | Mean | 11.16 | 11.79 | 10.88 | 11.21 | 12.52 | 11.44 | 10.71 | 12.37 | 14.64 | 13.98 | 12.07 |
| | Weighted Mean | 10.72 | 9.984 | 10.74 | 10.71 | 11.78 | 10.94 | 10.26 | 11.76 | 14.54 | 13.88 | 11.53 |
| LLaMA-3 (mntp-supervised) | EOS | 18.02 | 16.92 | 22.38 | 17.26 | 20.32 | 17.47 | 14.78 | 24.17 | 13.87 | 18.58 | 18.38 |
| | Mean | 16.13 | 16.64 | 23.24 | 16.21 | 26.00 | 16.46 | 14.49 | 25.72 | 15.33 | 18.78 | 18.90 |
| | Weighted Mean | 15.11 | 13.41 | 20.47 | 15.05 | 23.70 | 15.26 | 13.57 | 23.88 | 14.02 | 17.94 | 17.24 |
| Mistral (mntp) | EOS | 24.22 | 23.74 | 22.74 | 23.39 | 23.67 | 24.48 | 21.89 | 24.32 | 19.67 | 25.73 | 23.38 |
| | Mean | 17.49 | 17.97 | 14.57 | 17.64 | 17.42 | 18.54 | 16.73 | 15.86 | 15.45 | 16.70 | 16.84 |
| | Weighted Mean | 16.88 | 16.61 | 15.27 | 16.81 | 16.07 | 17.50 | 16.12 | 15.02 | 15.18 | 15.35 | 16.08 |
| Mistral (mntp-unsup-simcse) | EOS | 13.34 | 16.78 | 13.02 | 23.25 | 12.48 | 13.71 | 19.31 | 13.73 | 14.42 | 16.50 | 15.65 |
| | Mean | 11.93 | 11.78 | 11.39 | 11.63 | 12.94 | 11.87 | 11.37 | 12.79 | 14.55 | 13.45 | 12.37 |
| | Weighted Mean | 11.99 | 12.03 | 11.38 | 11.69 | 12.66 | 11.93 | 11.42 | 12.67 | 13.94 | 13.54 | 12.32 |
| Mistral (mntp-supervised) | EOS | 13.75 | 13.78 | 14.55 | 13.80 | 13.26 | 13.73 | 11.39 | 17.79 | 13.56 | 16.40 | 14.20 |
| | Mean | 13.36 | 10.94 | 11.13 | 11.17 | 15.12 | 12.72 | 10.58 | 15.39 | 12.22 | 12.87 | 12.55 |
| | Weighted Mean | 13.25 | 11.66 | 10.78 | 11.18 | 14.63 | 12.64 | 10.59 | 15.30 | 12.62 | 12.64 | 12.53 |
| OpenAI-3-small | - | 14.41 | 14.25 | 11.85 | 15.16 | 21.53 | 15.25 | 12.88 | 16.37 | 12.60 | 56.08 | 19.04 |
| OpenAI-ada-002 | - | 16.11 | 15.47 | 14.53 | 18.66 | 29.65 | 18.30 | 16.12 | 19.60 | 13.78 | 56.08 | 21.83 |
| OpenAI-3-large | - | 12.88 | 12.09 | 14.03 | 11.30 | 19.36 | 11.96 | 9.939 | 15.59 | 12.59 | 56.08 | 17.58 |