

---

# Towards a Unified Framework of Clustering-based Anomaly Detection

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Unsupervised Anomaly Detection (UAD) plays a crucial role in identifying abnormal  
2        patterns within data without labeled examples, holding significant practical  
3        implications across various domains. Although the individual contributions of  
4        representation learning and clustering to anomaly detection are well-established,  
5        their interdependencies remain under-explored due to the absence of a unified  
6        theoretical framework. Consequently, their collective potential to enhance anomaly  
7        detection performance remains largely untapped. To bridge this gap, in this paper,  
8        we propose a novel probabilistic mixture model for anomaly detection to establish  
9        a theoretical connection among representation learning, clustering, and anomaly  
10        detection. By maximizing a novel anomaly-aware data likelihood, representation  
11        learning and clustering can effectively reduce the adverse impact of anomalous  
12        data and collaboratively benefit anomaly detection. Meanwhile, a theoretically sub-  
13        stantiated anomaly score is naturally derived from this framework. Lastly, drawing  
14        inspiration from gravitational analysis in physics, we have devised an improved  
15        anomaly score that more effectively harnesses the combined power of representa-  
16        tion learning and clustering. Extensive experiments, involving 17 baseline methods  
17        across 30 diverse datasets, validate the effectiveness and generalization capability  
18        of the proposed method, surpassing state-of-the-art methods.

## 19    1 Introduction

20    Unsupervised Anomaly Detection (UAD) refers to the task dedicated to identifying abnormal patterns  
21    or instances within data in the absence of labeled examples [8]. It has long received extensive  
22    attention in the past decades for its wide-ranging applications in numerous practical scenarios,  
23    including financial auditing [3], healthcare monitoring [44] and e-commerce sector [23]. Due to the  
24    lack of explicit label guidance, the key to UAD is to uncover the dominant patterns that widely exist  
25    in the dataset so that samples do not conform to these patterns can be recognized as anomalies. To  
26    achieve this, early works [7] have heavily relied on powerful unsupervised *representation learning*  
27    methods to extract the normal patterns from high-dimensional and complex data such as images, text,  
28    and graphs. More recent works [45, 2] have utilized *clustering*, a widely observed natural pattern in  
29    real-world data, to provide critical global information for anomaly detection and achieved tremendous  
30    success.

31    While the individual contributions of representation learning and clustering to anomaly detection  
32    are well-established, their interrelationships remain largely unexplored. Intuitively, *discriminative*  
33    *representation learning* can leverage accurate clustering results to differentiate samples from distinct  
34    clusters in the embedding space (i.e., ①). Similarly, it can utilize accurate anomaly detection to  
35    avoid preserving abnormal patterns (i.e., ②). For *accurate clustering*, it can gain advantages from  
36    representation learning by operating in the discriminative embedding space (i.e., ③). Meanwhile, it

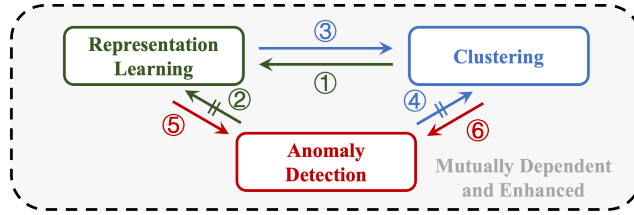


Figure 1: Interdependent relationships among representation learning, clustering, and anomaly detection.

37 can potentially benefit from accurate anomaly detection by excluding anomalies when formulating  
 38 clusters (i.e., ④). *Anomaly detection* can greatly benefit from both discriminative representation  
 39 learning and accurate clustering (i.e., ⑤ & ⑥). However, these benefits hinge on the successful  
 40 identification of anomalies and the reduction of their detrimental impact on the aforementioned  
 41 tasks. As depicted in Figure 1, the integration of these three elements exhibits a significant reciprocal  
 42 nature. In summary, representation learning, clustering, and anomaly detection are interdependent and  
 43 intricately intertwined. Therefore, it is crucial for anomaly detection to *fully leverage and mutually*  
 44 *enhance the relationships among these three components.*

45 Despite the intuitive significance of the interactions among representation learning, clustering, and  
 46 anomaly detection, existing methods have only made limited attempts to exploit them and fall short  
 47 of expectations. On one hand, some methods [58] have acknowledged the interplay among these  
 48 three factors, but their focus remains primarily on the interactions between two factors at a time,  
 49 making only targeted improvements. For instance, some strategies include explicitly removing outlier  
 50 samples during the clustering process [9] or designing robust representation learning methods [10] to  
 51 mitigate the influence of anomalies. On the other hand, recent methods [45] have begun to explore  
 52 the simultaneous optimization of these three factors within a single framework. However, these  
 53 attempts are still in the stage of merely superimposing the objectives of the three factors without a  
 54 unified theoretical framework. This lack of a guiding framework prevents the adequate modeling of  
 55 the interdependencies among these factors, thereby limiting their collective contribution to a unified  
 56 anomaly detection objective. Consequently, we aim to address the following question: *Is it possible*  
 57 *to employ a unified theoretical framework to jointly model these three interdependent objectives,*  
 58 *thereby leveraging their respective strengths to enhance anomaly detection?*

59 In this paper, we try to answer this question and propose a novel model named UniCAD for anomaly  
 60 detection. The proposed UniCAD integrates representation learning, clustering, and anomaly de-  
 61 tection into a unified framework, achieved through the theoretical guidance of maximizing the  
 62 anomaly-aware data likelihood. Specifically, we explicitly model the relationships between samples  
 63 and multiple clusters in the representation space using the probabilistic mixture models for the  
 64 likelihood estimation. Moreover, we creatively introduce a learnable indicator function into the  
 65 objective of maximum likelihood to explicitly attenuate the influence of anomalies on representation  
 66 learning and clustering. Under this framework, we can theoretically derive an anomaly score that  
 67 indicates the abnormality of samples, rather than heuristically designing it based on clustering results  
 68 as existing works do. Furthermore, building upon this theoretically supported anomaly score and  
 69 inspired by the theory of universal gravitation, we propose a more comprehensive anomaly metric that  
 70 considers the complex relationships between samples and multiple clusters. This allows us to better  
 71 utilize the learned representations and clustering results from this framework for anomaly detection.

72 To sum up, we underline our contributions as follows:

- 73 • We propose a unified theoretical framework to jointly optimize representation learning, clustering,  
 74 and anomaly detection, allowing their mutual enhancement and aid in anomaly detection.
- 75 • Based on the proposed framework, we derive a theoretically grounded anomaly score and further  
 76 introduce a more comprehensive score with the vector summation, which fully releases the power  
 77 of the framework for effective anomaly detection.
- 78 • Extensive experiments have been conducted on 30 datasets to validate the superior unsupervised  
 79 anomaly detection performance of our approach, which surpassed the state-of-the-art through  
 80 comparative evaluations with 17 baseline methods.

## 81 2 Related Work

82 Typical unsupervised anomaly detection (UAD) methods calculate a continuous score for each sample  
83 to measure its anomaly degree. Various UAD methods have been proposed based on different  
84 assumptions, making them suitable for detecting various types of anomaly patterns, including  
85 subspace-based models [24], statistical models [16], linear models [49, 32], density-based models [6,  
86 38], ensemble-based models [39, 29], probability-based models [40, 58, 28, 27], neural network-  
87 based models [42, 51], and cluster-based models [18, 9]. Considering the field of anomaly detection  
88 has progressed by integrating clustering information to enhance detection accuracy [26, 56], we  
89 primarily focus on and analyze anomaly patterns related to clustering, incorporating a global clustering  
90 perspective to assess the degree of anomaly. Notable methods in this context include CBLOF [18],  
91 which evaluates anomalies based on the size of the nearest cluster and the distance to the nearest large  
92 cluster. Similarly, DCFOD [45] introduces innovation by applying the self-training architecture of  
93 the deep clustering [50] to outlier detection. Meanwhile, DAGMM [58] combines deep autoencoders  
94 with Gaussian mixture models, utilizing sample energy as a metric to quantify the anomaly degree.  
95 In contrast, our approach introduces a unified theoretical framework that integrates representation  
96 learning, clustering, and anomaly detection, overcoming the limitations of heuristic designs and the  
97 overlooked anomaly influence in existing methods.

## 98 3 Methodology

99 In this section, we first define the problem we studied and the notations used in this paper. Then we  
100 elaborate on the proposed method UniCAD. More specifically, we first introduce a novel learning  
101 objective that optimizes representation learning, clustering, and anomaly detection within a unified  
102 theoretical framework by maximizing the data likelihood. A novel anomaly score with theoretical  
103 support is also naturally derived from this framework. Then, inspired by the concept of universal  
104 gravitation, we further propose an enhanced anomaly scoring approach that leverages the intricate  
105 relationship between samples and clustering to detect anomalies effectively. Finally, we present an  
106 efficient iterative optimization strategy to optimize this model and provide a complexity analysis for  
107 the proposed model.

108 **Definition 1** (Unsupervised Anomaly Detection). *Given a dataset  $\mathbf{X} \in \mathbb{R}^{N \times D}$  comprising  $N$*   
109 *instances with  $D$ -dimensional features, unsupervised anomaly detection aims to learn an anomaly*  
110 *score  $o_i$  for each instance  $\mathbf{x}_i$  in an unsupervised manner so that the abnormal ones have higher*  
111 *scores than the normal ones.*

### 112 3.1 Maximizing Anomaly-aware Likelihood

113 Previous research has demonstrated the importance of discriminative representation and accurate  
114 clustering in anomaly detection [45]. However, the presence of anomalous samples can significantly  
115 disrupt the effectiveness of both representation learning and clustering [12]. While some existing  
116 studies have attempted to integrate these three separate learning objectives, the lack of a unified  
117 theoretical framework has hindered their mutual enhancement, leading to suboptimal results.

118 To tackle this issue, in this paper, we propose a unified and coherent approach that considers  
119 representation learning, clustering, and anomaly detection by maximizing the likelihood of the  
120 observed data. Specifically, we denote the parameters of representation learning as  $\Theta$ , the clustering  
121 parameter as  $\Phi$ , and the dynamic indicator function for anomaly detection as  $\delta(\cdot)$ . These parameters  
122 are optimized simultaneously by maximizing the likelihood of the observed data  $\mathbf{X}$ :

$$\max \log p(\mathbf{X}|\Theta, \Phi) = \max \sum_{i=1}^N \delta(\mathbf{x}_i) \log p(\mathbf{x}_i|\Theta, \Phi) = \max \sum_{i=1}^N \delta(\mathbf{x}_i) \log \sum_{k=1}^K p(\mathbf{x}_i, c_i = k|\Theta, \Phi), \quad (1)$$

123 where  $c_i$  represents the latent cluster variable associated with  $\mathbf{x}_i$ , and  $c_i = k$  denotes the probabilistic  
124 event that  $\mathbf{x}_i$  belongs to the  $k$ -th cluster. The  $\delta(\mathbf{x}_i)$  is an indicator function that determines whether a  
125 sample  $\mathbf{x}_i$  is an anomaly of value 0 or a normal sample of value 1.

126 **3.1.1 Joint Representation Learning and Clustering with  $p(\mathbf{x}_i|\Theta, \Phi)$**

127 Based on the aforementioned advantages of MMs, we estimate the likelihood  $p(\mathbf{x}_i|\Theta, \Phi)$  with mixture  
128 models defined as:

$$\begin{aligned} p(\mathbf{x}_i|\Theta, \Phi) &= \sum_{k=1}^K p(\mathbf{x}_i, c_i = k|\Theta, \Phi) = \sum_{k=1}^K p(c_i = k) \cdot p(\mathbf{x}_i|c_i = k, \Theta, \boldsymbol{\mu}_k, \Sigma_k) \\ &= \sum_{k=1}^K \omega_k \cdot p(\mathbf{x}_i|c_i = k, \Theta, \boldsymbol{\mu}_k, \Sigma_k), \end{aligned} \quad (2)$$

129 where  $\Phi = \{\omega_k, \boldsymbol{\mu}_k, \Sigma_k\}$ . The mixture model is parameterized by the prototypes  $\boldsymbol{\mu}_k$ , covariance  
130 matrices  $\Sigma_k$ , and mixture weights  $\omega_k$  from all clusters.  $\sum_{k=1}^K \omega_k = 1$ , and  $k = 1, 2, \dots, K$ .

131 In practice, the samples are usually attributed to high-dimensional features and it is challenging to  
132 detect anomalies from the raw feature space [41]. Therefore, modern anomaly detection methods [42,  
133 58] often map raw data samples  $\mathbf{X} = \{\mathbf{x}_i\} \in \mathbb{R}^{N \times D}$  into a low-dimensional representation space  
134  $\mathbf{Z} = \{\mathbf{z}_i\} \in \mathbb{R}^{N \times d}$  with a representation learning function  $\mathbf{z}_i = f_{\Theta}(\mathbf{x}_i)$  and detect anomalies within  
135 this latent representation space.

136 Following this widely adopted practice, we model the distribution of samples in the latent represen-  
137 tation space with a multivariate Student's- $t$  distribution giving its cluster  $c_i = k$ . The Student's- $t$   
138 distribution is robust against outliers due to its heavy tails. Bayesian robustness theory leverages  
139 such distributions to dismiss outlier data, favoring reliable sources, making the Student's- $t$  process  
140 preferable over Gaussian processes for data with atypical information [1]. Thus the probability  
141 distribution of generating  $\mathbf{x}_i$  with latent representation  $\mathbf{z}_i$  given its cluster  $c_i = k$  can be expressed as:

$$p(\mathbf{x}_i|c_i = k, \Theta, \boldsymbol{\mu}_k, \Sigma_k) = \frac{\Gamma(\frac{\nu+1}{2})|\Sigma_k|^{-1/2}}{\Gamma(\frac{\nu}{2})\sqrt{\nu\pi}} \left(1 + \frac{1}{\nu}D_M(\mathbf{z}_i, \boldsymbol{\mu}_k)^2\right)^{-\frac{\nu+1}{2}}, \quad (3)$$

142 where  $\mathbf{z}_i = f_{\Theta}(\mathbf{x}_i)$  denotes the representation obtained from the data mapped through the neural  
143 network parameterized by  $\Theta$ .  $\Gamma$  denotes the gamma function while  $\nu$  is the degree of freedom.

144  $\Sigma_k$  is the scale parameter.  $D_M(\mathbf{z}_i, \boldsymbol{\mu}_k) = \sqrt{(\mathbf{z}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{z}_i - \boldsymbol{\mu}_k)}$  represents the Mahalanobis  
145 distance [33]. In the unsupervised setting, as cross-validating  $\nu$  on a validation set or learning it is  
146 unnecessary,  $\nu$  is set as 1 for all experiments [50, 48]. The overall marginal likelihood of the observed  
147 data  $\mathbf{x}_i$  can be simplified as:

$$p(\mathbf{x}_i|\Theta, \Phi) = \sum_{k=1}^K \omega_k \cdot \frac{\pi^{-1} \cdot |\Sigma_k|^{-1/2}}{1 + D_M(\mathbf{z}_i, \boldsymbol{\mu}_k)^2}. \quad (4)$$

148 **3.1.2 Anomaly Indicator  $\delta(\mathbf{x}_i)$  and Score  $o_i$**

149 As we have discussed, the indicator function  $\delta(\mathbf{x}_i)$  not only benefits both representation and clustering  
150 but also directly serves as the output of anomaly detection. Ideally, with the percentage of outliers  
151 denoted as  $l$ , an optimal solution for  $\delta(\mathbf{x}_i)$  that maximizes the objective function  $J(\Theta, \Phi)$  entails  
152 setting all  $\delta(\mathbf{x}_i) = 0$  for  $\mathbf{x}_i$  among the  $l$  percent of outliers with lowest generation possibility  
153  $p(\mathbf{x}_i|\Theta, \Phi)$ , and otherwise  $\delta(\mathbf{x}_i) = 1$  is set for the remaining normal samples. Therefore, the  
154 indicator function is determined as:

$$\delta(\mathbf{x}_i) = \begin{cases} 0, & \text{if } p(\mathbf{x}_i|\Theta, \Phi) \text{ is among the } l \text{ lowest,} \\ 1, & \text{otherwise.} \end{cases} \quad (5)$$

155 As this method involves sorting the samples based on the generation probability as being anomalous,  
156 the values of  $p(\mathbf{x}_i|\Theta, \Phi)$  can serve as a form of anomaly score, a classic approach within the mixture  
157 model framework [40, 58]. This suggests that the likelihood of a sample being anomalous is inversely  
158 related to its generative probability since a lower generative probability indicates a higher chance of  
159 the sample being an outlier. Thus the anomaly score of sample  $\mathbf{x}_i$  can be defined as:

$$o_i = \frac{1}{p(\mathbf{x}_i|\Theta, \Phi)} = \frac{1}{\sum_{k=1}^K \omega_k \cdot \frac{\pi^{-1} \cdot |\Sigma_k|^{-1/2}}{1 + D_M(\mathbf{z}_i, \boldsymbol{\mu}_k)^2}}. \quad (6)$$

## 160 3.2 Gravity-inspired Anomaly Scoring

161 In practical applications, it is proved that anomaly scores derived from generation probabilities often  
 162 yield suboptimal performance [17]. This observation prompts a reconsideration of *how to fully*  
 163 *leverage the complex relationships among samples or even across multiple clusters for anomaly*  
 164 *detection*. In this section, we first provide a brief introduction to the concept of Newton’s Law of  
 165 Universal Gravitation [35] and then demonstrate how the anomaly score is intriguingly similar to this  
 166 cross-field principle. Finally, we discuss the advantages of introducing the vector sum operation into  
 167 the anomaly score inspired by the analogy.

### 168 3.2.1 Analog Anomaly Scoring and Force Analysis

169 To begin with, Newton’s Law of Universal Gravitation [35] stands as a fundamental framework for  
 170 describing the interactions among entities in the physical world. According to this law, every object  
 171 in the universe experiences an attractive force from another object. In classical mechanics, force  
 172 analysis involves calculating the vector sum of all forces acting on an object, known as the **resultant**  
 173 **force**, which is crucial in determining an object’s acceleration or change in motion:

$$\vec{\mathbf{F}}_{i,\text{total}} = \sum_{k=1}^K \vec{\mathbf{F}}_{ik}, \quad \text{with } \vec{\mathbf{F}}_{ik} = \frac{G \cdot m_i m_k}{r_{ik}^2} \cdot \vec{\mathbf{r}}_{ik}, \quad (7)$$

174 where  $\vec{\mathbf{F}}_{ik}$  represents the  $k$ -th force acting on the object  $i$ . This force is proportional to the product of  
 175 their masses, ( $m_i$  and  $m_k$ ), and inversely proportional to the square of the distance  $r_{ik}$  between them.  
 176  $G$  represents the gravitational constant, and  $\vec{\mathbf{r}}_{ij}$  is the unit direction vector.

177 Similarly, if denoting:  $\tilde{\mathbf{F}}_{ik} = p(\mathbf{x}_i, c_i = k|\Theta, \Phi) = \omega_k \cdot \frac{\pi^{-1} \cdot |\Sigma_k|^{-1/2}}{1 + D_M(\mathbf{z}_i, \boldsymbol{\mu}_k)^2}$ , the score of Equation (6)  
 178 bears analogies to the summation of the magnitudes of forces as:

$$o_i = \frac{1}{\sum_{k=1}^K \tilde{\mathbf{F}}_{ik}}, \quad \text{with } \tilde{\mathbf{F}}_{ik} = \frac{\tilde{G} \cdot \tilde{m}_i \tilde{m}_k}{\tilde{r}_{ik}^2}, \quad (8)$$

179 where  $\tilde{G} = \pi^{-1}$ ,  $\tilde{m}_k = \omega_k |\Sigma_k|^{-1/2}$ ,  $\tilde{m}_i = 1$ , and  $\tilde{r}_{ik} = \sqrt{1 + D_M(\mathbf{z}_i, \boldsymbol{\mu}_k)^2}$ . Here,  $\tilde{r}_{ik}$  is taken as  
 180 the measure of distance within the representation space, modified slightly by an additional term for  
 181 smoothness. The constant  $\tilde{G}$  serves a role akin to the gravitational constant in this analogy, whereas  
 182  $\tilde{m}_k$  resembles the concept of mass for the cluster. The notation  $\tilde{m}_i$  suggests a standardization where  
 183 the mass of each data point is considered uniform and not differentiated.

### 184 3.2.2 Anomaly Scoring with Vector Sum

185 Comparing Equation (7) with Equation (8), what still differs is that, unlike a simple sum of the  
 186 scalar value, the resultant force  $\vec{\mathbf{F}}_{i,\text{total}}$  employs the vector sum and incorporates both the magnitude  
 187 and direction  $\hat{\mathbf{r}}_{ik}$  of each force. This distinction is crucial because forces in different directions  
 188 can neutralize each other with a large angle between them or enhance each other’s effects with a  
 189 small angle. Inspired by this difference, we consider modeling the relationship between samples and  
 190 clusters as a vector, and aggregating them through vector summation. The vector-formed anomaly  
 191 score  $o_i^V$  is defined as:

$$o_i^V = \frac{1}{\left\| \sum_{k=1}^K \tilde{\mathbf{F}}_{ik} \cdot \hat{\mathbf{r}}_{ik} \right\|}, \quad (9)$$

192 where  $\hat{\mathbf{r}}_{ik}$  represents the unit direction vector in the representation space from the sample  $\mathbf{z}_i$  to the  
 193 cluster prototype  $\boldsymbol{\mu}_k$ , and  $\| \cdot \|$  represents the  $L_2$  norm.

### 194 3.3 Iterative Optimization

195 Given the challenge posed by the interdependence of the parameters of the network  $\Theta$  and those of the  
 196 mixture model  $\{\omega_k, \boldsymbol{\mu}_k, \Sigma_k\}$  in joint optimization, we propose an iterative optimization procedure.  
 197 The pseudocode for training the model is presented in Algorithm 1 in the appendix.

#### 198 3.3.1 Update $\Phi$

199 To update the parameters of the mixture model  $\Phi = \{\omega_k, \boldsymbol{\mu}_k, \Sigma_k\}$ , we use the Expectation-  
 200 Maximization (EM) algorithm to maximize equation (1) [36]. The detailed derivation is included in  
 201 Appendix B.

202 **E-step.** During the E-step of iteration  $(t + 1)$ , our goal is to compute the posterior probabilities of  
 203 each data point belonging to the  $k$ -th cluster within the mixture model. Given the observed sample  
 204  $\mathbf{x}_i$  and the current estimates of the parameters  $\Theta^{(t)}$  and  $\Phi^{(t)}$ , the expected value of the likelihood  
 205 function of latent variable  $c_k$ , or the posterior possibilities, can be expressed as:

$$\tau_{ik}^{(t+1)} = p(c_i = k | \mathbf{x}_i, \Theta, \Phi^{(t)}) = \frac{p(\mathbf{x}_i, c_i = k | \Theta, \Phi^{(t)})}{\sum_{j=1}^K p(\mathbf{x}_i, c_i = j | \Theta, \Phi^{(t)})} = \frac{\tilde{\mathbf{F}}_{ik}^{(t)}}{\sum_{j=1}^K \tilde{\mathbf{F}}_{ij}^{(t)}}. \quad (10)$$

206 The scale factor[36] serving as an intermediate result for subsequent updates in the M-step is :

$$\mathbf{u}_{ik}^{(t+1)} = \frac{2}{1 + D_M(\mathbf{z}_i^{(t)}, \boldsymbol{\mu}_k^{(t)})}. \quad (11)$$

207 **M-step.** In the M-step of iteration  $(t + 1)$ , given the gradients  $\frac{\partial J(\Theta, \Phi)}{\partial \omega_k} = 0$ ,  $\frac{\partial J(\Theta, \Phi)}{\partial \boldsymbol{\mu}_k} = 0$ , and  
 208  $\frac{\partial J(\Theta, \Phi)}{\partial \Sigma_k} = 0$ , we derive the analytical solutions for the mixture model parameters  $\omega_k$ ,  $\boldsymbol{\mu}_k$ , and  $\Sigma_k$ .  
 209 Assume the anomalous ratio is  $l \in [0, 1]$ , the number of the normal samples is  $n = \text{int}(l * N)$ . The  
 210 updating process for  $\{\omega_k^{(t+1)}, \boldsymbol{\mu}_k^{(t+1)}, \Sigma_k^{(t+1)}\}$  is as follows:

- 211 • The mixture weights  $\omega_k$  are updated by averaging the posterior probabilities over all data points  
 212 with the number of samples , reflecting the relative presence of each component in the mixture:

$$\omega_k^{(t+1)} = \sum_{i=1}^n \tau_{ik}^{(t+1)} / n. \quad (12)$$

- 213 • The prototypes  $\boldsymbol{\mu}_k$  are updated to be the weighted average of the data points, where weights are the  
 214 posterior probabilities:

$$\boldsymbol{\mu}_k^{(t+1)} = \sum_{i=1}^n \left( \tau_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} \mathbf{z}_i \right) / \sum_{i=1}^n \left( \tau_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} \right). \quad (13)$$

- 215 • The covariance matrices  $\Sigma_k$  are updated by considering the dispersion of the data around the newly  
 216 computed prototypes:

$$\Sigma_k^{(t+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} (\mathbf{z}_i - \boldsymbol{\mu}_k^{(t+1)}) (\mathbf{z}_i - \boldsymbol{\mu}_k^{(t+1)})^\top}{\sum_{j=1}^K \tau_{ij}^{(t+1)}}. \quad (14)$$

#### 217 3.3.2 Update $\Theta$

218 We focus on anomaly-aware representation learning and use stochastic gradient descent to optimize  
 219 the network parameters  $\Theta$ , by minimizing the following joint loss:

$$\mathcal{L} = -J(\Theta, \Phi) + g(\Theta), \quad (15)$$

220 where  $J(\Theta, \Phi) = \log p(\mathbf{X} | \Theta, \Phi)$ . An additional constraint term  $g(\Theta)$  is introduced to prevent short-  
 221 cut solution [15]. In practice, an autoencoder architecture is implemented, utilizing a reconstruction  
 222 loss  $g(\Theta) = \|x - \hat{x}\|^2$  as the constraint.

223 These updates are iteratively performed until convergence, resulting in optimized model parameters  
 224 that best fit the given data according to the mixture model framework.

## 225 4 Experiments

### 226 4.1 Datasets & Baselines

227 We evaluated UniCAD on an extensive collection of datasets, comprising 30 tabular datasets that  
228 span 16 diverse fields. We specifically focused on naturally occurring anomaly patterns, rather  
229 than synthetically generated or injected anomalies, as this aligns more closely with real-world  
230 scenarios. The detailed descriptions are provided in Table 4 of Appendix D.1. Following the setup  
231 in ADBench [17], we adopt an inductive setting to predict newly emerging data, a highly beneficial  
232 approach for practical applications.

233 To assess the effectiveness of UniCAD, we compared it with 17 advanced unsupervised anomaly  
234 detection methods, including: (1) *traditional methods*: SOD [24] and HBOS [16]; (2) *linear methods*:  
235 PCA [49] and OCSVM [32]; (3) *density-based methods*: LOF [6] and KNN [38]; (4) *ensemble-based*  
236 *methods*: LODA [39] and IForest [29]; (5) *probability-based methods*: DAGMM [58], ECOD [28],  
237 and COPOD [27]; (6) *cluster-based methods*: DBSCAN [13], CBLOF [18], DCOD [45] and KMeans-  
238 - [9]; and (7) *neural network-based methods*: DeepSVDD [42] and DIF [51]. These baselines  
239 encompass the majority of the latest methods, providing a comprehensive overview of the state-of-  
240 the-art. For a detailed description, please refer to Appendix D.2.

### 241 4.2 Experiment Settings

242 In the unsupervised setting, we employ the default hyperparameters from the original papers for all  
243 comparison methods. Similarly, the UniCAD also utilizes a fixed set of parameters to ensure a fair  
244 comparison. For all datasets, we employ a two-layer MLP with a hidden dimension of  $d = 128$  and  
245 ReLU activation function as both encoder and decoder. We utilize the Adam optimizer [21] with a  
246 learning rate of  $1e^{-4}$  for 100 epochs. For the EM process, we set the maximum iteration number  
247 to 100 and a tolerance of  $1e^{-3}$  for stopping training when the objectives converge. The number of  
248 components in the mixture model is set as  $k = 10$ , and the proportion of the outlier is set as  $l = 1\%$ .  
249 We evaluate the methods using Area Under the Receiver Operating Characteristic (AUC-ROC) and  
250 Area Under the Precision-Recall Curve (AUC-PR) metrics [17], reporting the average ranking (Avg.  
251 Rank) across all datasets. All experiments are run 3 times with different seeds, and the mean results  
252 are reported.

### 253 4.3 Performance and Analysis

254 **Performance Comparison.** Table 1 presents a comparison of UniCAD with 10 unsupervised  
255 baseline methods across 30 tabular datasets using the AUC-ROC metric. The experimental results,  
256 which encompass 17 baselines, are included in Tables 5 and 6 of Appendix D.3, with additional  
257 experiments on other data domains presented in Appendix E. Our proposed UniCAD achieves the  
258 top average ranking, exhibiting the best or near-best performance on a larger number of datasets  
259 and confirming advanced capabilities. It is noteworthy that there is no one-size-fits-all unsupervised  
260 anomaly detection method suitable for every type of dataset, as demonstrated by the observation that  
261 other methods have also achieved some of the best results on certain datasets. However, our model  
262 showcased a remarkable ability to generalize across most datasets featuring natural anomalies, as  
263 evidenced by statistical average ranking. As for clustering-based methods such as KMeans-, DCOD,  
264 and CBLOF, they mostly rank in the top tier among all baseline methods, supporting the advantage of  
265 combining deep clustering with anomaly detection. However, our method significantly outperformed  
266 these methods by mitigating their limitations and further providing a unified framework for joint  
267 representation learning, clustering, and anomaly detection.

268 **Effectiveness of Vector Sum in Anomaly Scoring.** As demonstrated in Table 1, we compare the  
269 anomaly score  $\mathbf{o}_i$  derived directly from the generation possibility with its vector summation form  $\mathbf{o}_i^V$ .  
270 According to our statistical findings, we observe that vector scores  $\mathbf{o}_i^V$  consistently outperform scalar  
271 scores  $\mathbf{o}_i$ . This indicates that the introduction of the vector summation, analogous to the concept  
272 of resultant force, makes a substantial difference in anomaly detection scenarios involving multiple  
273 clusters. The performance gains of the vector sum scores strongly demonstrate the effectiveness  
274 of the UniCAD in capturing the subtle differences in the distinctions among multiple clusters and  
275 underscore the utility of this factor in the context of anomaly detection based on clustering.

Table 1: AUCROC of 10 unsupervised algorithms on 30 tabular benchmark datasets. In each dataset, the algorithm with the highest AUCROC is marked in **red**, the second highest in **blue**, and the third highest in **green**.

Dataset	OC SVM	LOF	IForest	DA GMM	ECOD	DB SCAN	CBLOF	DCOD	KMeans--	DIF	UniCAD (Scalar)	UniCAD (Vector)
annthyroid	57.23	70.20	<b>82.01</b>	56.53	<b>78.66</b>	50.08	62.28	55.01	64.99	66.76	<b>75.27</b>	72.72
backdoor	85.04	85.79	72.15	55.98	86.08	76.55	81.91	79.57	<b>89.11</b>	<b>92.87</b>	87.28	<b>89.24</b>
breastw	80.30	40.61	98.32	N/A	<b>99.17</b>	85.20	96.86	<b>99.02</b>	97.05	77.45	98.15	<b>98.56</b>
campaign	65.70	59.04	71.71	56.03	<b>76.10</b>	50.60	64.34	63.16	63.51	67.53	<b>73.52</b>	<b>73.64</b>
celeba	70.70	38.95	70.41	44.74	76.48	50.36	73.99	<b>91.41</b>	56.76	65.29	<b>81.38</b>	<b>82.00</b>
census	54.90	47.46	59.52	59.65	67.63	58.50	60.17	<b>72.84</b>	63.33	59.66	<b>67.90</b>	<b>67.84</b>
glass	35.36	69.20	77.13	76.09	65.83	54.55	78.30	78.07	77.30	<b>84.57</b>	<b>79.52</b>	<b>82.17</b>
Hepatitis	67.75	38.06	69.75	54.80	<b>75.22</b>	68.12	73.05	48.38	64.64	74.24	<b>75.53</b>	<b>80.62</b>
htp	<b>99.59</b>	27.46	<b>99.96</b>	N/A	98.10	49.97	<b>99.60</b>	99.53	99.55	99.49	99.53	99.52
Ionosphere	75.92	90.59	84.50	73.41	73.15	81.12	<b>90.79</b>	57.78	<b>91.36</b>	89.74	<b>92.04</b>	90.37
landsat	36.15	53.90	47.64	43.92	36.10	50.17	<b>63.69</b>	33.40	<b>55.31</b>	54.84	49.60	<b>57.37</b>
Lymphography	99.54	89.86	<b>99.81</b>	72.11	99.52	74.16	<b>99.81</b>	81.19	<b>100.00</b>	83.67	99.29	<b>99.73</b>
mnist	82.95	67.13	80.98	67.23	74.61	50.00	79.96	65.23	82.45	<b>88.16</b>	<b>86.00</b>	<b>86.64</b>
musk	80.58	41.18	<b>99.99</b>	76.85	95.40	50.00	<b>100.00</b>	42.19	72.16	98.22	<b>99.92</b>	<b>100.00</b>
pendigits	93.75	47.99	94.76	64.22	93.01	55.33	<b>96.93</b>	94.33	94.37	93.79	<b>95.12</b>	<b>95.52</b>
Pima	66.92	65.71	<b>72.87</b>	55.93	63.05	51.39	71.49	72.16	70.44	67.28	<b>75.16</b>	<b>74.87</b>
satellite	59.02	55.88	70.43	62.33	58.09	55.52	71.32	55.97	67.71	<b>74.52</b>	<b>72.46</b>	<b>77.65</b>
satimage-2	97.35	47.36	99.16	96.29	96.28	75.74	<b>99.84</b>	86.01	<b>99.88</b>	99.63	<b>99.87</b>	<b>99.88</b>
shuttle	97.40	57.11	<b>99.56</b>	97.92	<b>99.13</b>	50.40	93.07	97.20	69.97	97.00	<b>99.15</b>	98.75
skin	49.45	46.47	<b>68.21</b>	N/A	49.08	50.00	68.03	64.34	65.47	66.36	<b>72.26</b>	<b>69.69</b>
Stamps	83.86	51.26	91.21	88.89	87.87	52.08	69.89	<b>93.41</b>	79.78	87.95	<b>91.37</b>	<b>94.18</b>
thyroid	87.92	86.86	<b>98.30</b>	79.75	<b>97.94</b>	53.57	94.74	78.55	92.26	96.26	<b>97.66</b>	97.48
vertebral	37.99	<b>49.29</b>	36.66	<b>53.20</b>	40.66	<b>49.74</b>	41.01	38.13	38.14	47.20	33.11	47.37
vowels	61.59	<b>93.12</b>	73.94	60.58	62.24	57.50	<b>92.12</b>	51.56	<b>93.45</b>	81.02	88.38	92.09
Waveform	56.29	73.32	71.47	49.35	62.36	66.41	71.27	63.47	<b>74.35</b>	<b>75.33</b>	71.81	<b>74.29</b>
WBC	<b>99.03</b>	54.17	<b>99.01</b>	N/A	<b>99.11</b>	87.43	96.88	94.92	97.45	81.27	97.68	98.93
Wilt	31.28	<b>50.65</b>	41.94	37.29	36.30	<b>49.96</b>	34.50	44.71	34.91	39.46	48.95	<b>52.56</b>
wine	73.07	37.74	80.37	61.70	77.22	40.33	27.14	<b>82.18</b>	27.36	41.69	<b>82.72</b>	<b>95.25</b>
WPBC	45.35	41.41	46.63	47.80	46.65	<b>52.22</b>	45.32	<b>49.67</b>	45.01	44.69	48.02	<b>49.90</b>
<b>Avg. Rank</b>	<b>7.8</b>	<b>8.9</b>	<b>5.1</b>	<b>8.7</b>	<b>6.4</b>	<b>9.3</b>	<b>5.7</b>	<b>7.4</b>	<b>6.0</b>	<b>5.8</b>	<b>3.7</b>	<b>2.6</b>

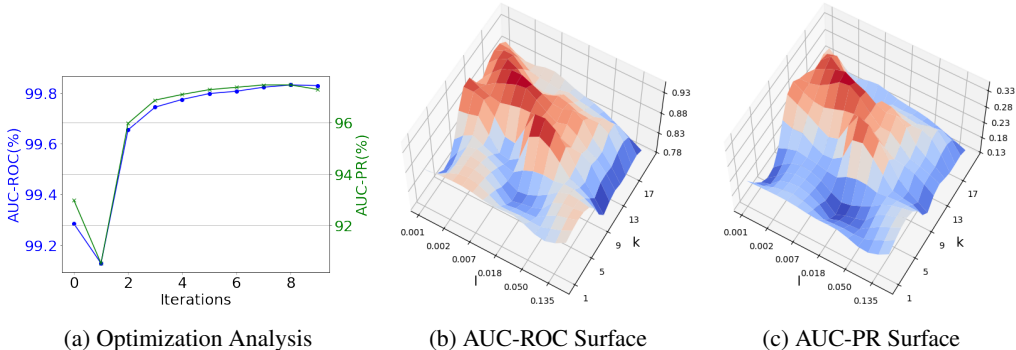


Figure 2: (a) demonstrates the performance variations during the optimization process on the satimage-2 dataset. (b) & (c) Analysis of cluster count  $k$ , anomaly ratio  $l$ .

276 **Analysis of EM Iterative Optimization.** To comprehend the iterative training within our model,  
 277 we have illustrated the performance variations accompanying the increase in iteration counts in  
 278 Figure 2a. Specifically, we monitored the iteration number  $t$  for the satimage-2 dataset, ranging  
 279 from 0 to 10, while maintaining other default parameters constant. Both AUC-ROC and AUC-PR  
 280 performance curves displayed consistent trends, with minor fluctuations only during the initial phase.  
 281 The performance remained relatively stable throughout the last steps, illustrating the effectiveness  
 282 and convergence of iterative EM optimization.

283 **Runtime Comparison.** We present a analysis of the runtime performance of various methods,  
 284 including our proposed approach, as detailed in Table 2. Our experiments, conducted on the backdoor  
 285 dataset, reveal that while non-deep learning methods exhibit lower runtime, they often simplify the  
 286 problem space excessively, failing to capture the complex non-linear relationships present in the  
 287 data. In contrast, our method, when compared to existing deep learning techniques, demonstrates  
 288 a significant reduction in computational time. This indicates that our approach not only manages



Table 2: Runtime Comparison. The runtime is reported in seconds (s).

Phase	IForest	KMeans--	DAGMM	DCOD	UniCAD
Fit	0.256	103.697	795.004	4548.634	246.113
Infer	0.0186	0.059	4.190	16.190	0.079

Table 3: Ablation study on AUC-ROC scores, calculated across 30 datasets.

Metric	w/ Gauss.	w/o $J(\Theta, \Phi)$	w/o $\delta(\mathbf{x}_i)$	Full Model
Avg. Rank (w/ baselines & variants)	6.2	6.6	5.0	<b>4.2</b>

289 to efficiently model complex patterns but also achieves an optimal balance between computational  
 290 efficiency and modeling capability.

#### 291 4.4 Ablation Studies

292 In this section, we examine the contributions of different components in UniCAD. Tables 3 reports the  
 293 results. We make three major observations. **Firstly**, the anomaly detection performance experiences a  
 294 significant drop when replacing the Student’s t distribution with a Gaussian distribution for the Mixture  
 295 Model, highlighting the robustness of the Student’s t distribution in unsupervised anomaly detection.  
 296 **Secondly**, omitting the likelihood maximization loss (w/o  $J(\Theta, \Phi)$ ) also results in a considerable  
 297 decrease in overall performance. This observation underscores the importance of deriving both  
 298 the optimization objectives and anomaly scores from the likelihood generation probability through  
 299 a theoretical framework, which allows for unified joint optimization of anomaly detection and  
 300 clustering in the representation space. **Furthermore**, the indicator function  $\delta(\mathbf{x}_i)$  also contributes to a  
 301 performance increase. These results further confirm the effectiveness of our UniCAD in mitigating the  
 302 negative influence of anomalies in the clustering process, as the existence of outliers may significantly  
 303 degrade the performance of clustering. In summary, all these ablation studies clearly demonstrate  
 304 the effectiveness of our theoretical framework in simultaneously considering representation learning,  
 305 clustering, and anomaly detection.

#### 306 4.5 Sensitivity of Hyperparameters

307 In this section, we conducted a sensitivity analysis on key hyperparameters of the model applied  
 308 to the donors dataset, focusing on the number of clusters  $k$  and the proportion of the outlier set  $l$ .  
 309 The results of this analysis are illustrated in Figure 2. Notably, the optimal range for  $l$  tends to be  
 310 lower than the actual proportion of anomalies in the dataset. Furthermore, a pattern was observed  
 311 with the number of clusters  $k$ , where the model performance initially improved with an increase in  $k$ ,  
 312 followed by a subsequent decline. This suggests the existence of an optimal range for the number of  
 313 clusters, which should be carefully selected based on the specific application context.

### 314 5 Conclusion

315 This paper presents UniCAD, a novel model for Unsupervised Anomaly Detection (UAD) that  
 316 seamlessly integrates representation learning, clustering, and anomaly detection within a unified  
 317 theoretical framework. Specifically, UniCAD introduces an anomaly-aware data likelihood based on  
 318 the mixture model with the Student-t distribution to guide the joint optimization process, effectively  
 319 mitigating the impact of anomalies on representation learning and clustering. This framework  
 320 enables a theoretically grounded anomaly score inspired by universal gravitation, which considers  
 321 complex relationships between samples and multiple clusters. Extensive experiments on 30 datasets  
 322 across various domains demonstrate the effectiveness and generalization capability of UniCAD,  
 323 surpassing 15 baseline methods and establishing it as a state-of-the-art solution in unsupervised  
 324 anomaly detection. Despite its potential, the proposed method’s applicability to broader fields like  
 325 time series and multimodal anomaly detection requires further exploration and validation, highlighting  
 326 a significant area for future work.

327 **References**

- 328 [1] J Ailton A Andrade. On the robustness to outliers of the student-t process. *Scandinavian*  
329 *Journal of Statistics*, 50(2):725–749, 2023.
- 330 [2] Caglar Aytakin, Xingyang Ni, Francesco Cricri, and Emre Aksu. Clustering and unsupervised  
331 anomaly detection with  $l_2$  normalized deep auto-encoder representations. In *2018 International*  
332 *Joint Conference on Neural Networks (IJCNN)*, pages 1–6. IEEE, 2018.
- 333 [3] Alexander Bakumenko and Ahmed Elragal. Detecting anomalies in financial data using machine  
334 learning algorithms. *Systems*, 10(5):130, 2022.
- 335 [4] Sambaran Bandyopadhyay, Saley Vishal Vivek, and MN Murty. Outlier resistant unsupervised  
336 deep architectures for attributed network embedding. In *Proceedings of the 13th international*  
337 *conference on web search and data mining*, pages 25–33, 2020.
- 338 [5] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko.  
339 Translating embeddings for modeling multi-relational data. *Advances in neural information*  
340 *processing systems*, 26, 2013.
- 341 [6] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. Lof: identifying  
342 density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference*  
343 *on Management of data*, pages 93–104, 2000.
- 344 [7] Raghavendra Chalapathy and Sanjay Chawla. Deep learning for anomaly detection: A survey.  
345 *arXiv preprint arXiv:1901.03407*, 2019.
- 346 [8] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM*  
347 *computing surveys (CSUR)*, 41(3):1–58, 2009.
- 348 [9] Sanjay Chawla and Aristides Gionis. k-means–: A unified approach to clustering and outlier  
349 detection. In *Proceedings of the 2013 SIAM international conference on data mining*, pages  
350 189–197. SIAM, 2013.
- 351 [10] Hyunsoo Cho, Jinseok Seol, and Sang-goo Lee. Masked contrastive learning for anomaly  
352 detection. *arXiv preprint arXiv:2105.08793*, 2021.
- 353 [11] Kaize Ding, Jundong Li, Rohit Bhanushali, and Huan Liu. Deep anomaly detection on attributed  
354 networks. In *Proceedings of the 2019 SIAM International Conference on Data Mining*, pages  
355 594–602. SIAM, 2019.
- 356 [12] Lian Duan, Lida Xu, Ying Liu, and Jun Lee. Cluster-based outlier detection. *Annals of*  
357 *Operations Research*, 168:151–168, 2009.
- 358 [13] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm  
359 for discovering clusters in large spatial databases with noise. In *kdd*, volume 96, pages 226–231,  
360 1996.
- 361 [14] Haoyi Fan, Fengbin Zhang, and Zuoyong Li. Anomalydae: Dual autoencoder for anomaly  
362 detection on attributed networks. In *ICASSP 2020-2020 IEEE International Conference on*  
363 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 5685–5689. IEEE, 2020.
- 364 [15] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,  
365 Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature*  
366 *Machine Intelligence*, 2(11):665–673, 2020.
- 367 [16] Markus Goldstein and Andreas Dengel. Histogram-based outlier score (hbos): A fast unsuper-  
368 vised anomaly detection algorithm. *KI-2012: poster and demo track*, 1:59–63, 2012.
- 369 [17] Songqiao Han, Xiyang Hu, Hailiang Huang, Minqi Jiang, and Yue Zhao. Adbench: Anomaly  
370 detection benchmark. *Advances in Neural Information Processing Systems*, 35:32142–32159,  
371 2022.
- 372 [18] Zengyou He, Xiaofei Xu, and Shengchun Deng. Discovering cluster-based local outliers.  
373 *Pattern recognition letters*, 24(9-10):1641–1650, 2003.

- 374 [19] Meng Jiang. Catching social media advertisers with strategy analysis. In *Proceedings of the*  
375 *First International Workshop on Computational Methods for CyberSafety*, pages 5–10, 2016.
- 376 [20] Ming Jin, Yixin Liu, Yu Zheng, Lianhua Chi, Yuan-Fang Li, and Shirui Pan. Anemone: Graph  
377 anomaly detection with multi-scale contrastive learning. In *Proceedings of the 30th ACM*  
378 *International Conference on Information & Knowledge Management*, pages 3122–3126, 2021.
- 379 [21] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
380 *arXiv:1412.6980*, 2014.
- 381 [22] Thomas N Kipf and Max Welling. Variational graph auto-encoders. *arXiv preprint*  
382 *arXiv:1611.07308*, 2016.
- 383 [23] Yufeng Kou, Chang-Tien Lu, Sirirat Sirwongwattana, and Yo-Ping Huang. Survey of fraud  
384 detection techniques. In *IEEE International Conference on Networking, Sensing and Control,*  
385 *2004*, volume 2, pages 749–754. IEEE, 2004.
- 386 [24] Hans-Peter Kriegel, Peer Kröger, Erich Schubert, and Arthur Zimek. Outlier detection in  
387 axis-parallel subspaces of high dimensional data. In *Advances in Knowledge Discovery and*  
388 *Data Mining: 13th Pacific-Asia Conference, PAKDD 2009 Bangkok, Thailand, April 27-30,*  
389 *2009 Proceedings 13*, pages 831–838. Springer, 2009.
- 390 [25] Srijan Kumar, Xikun Zhang, and Jure Leskovec. Predicting dynamic embedding trajectory  
391 in temporal interaction networks. In *Proceedings of the 25th ACM SIGKDD international*  
392 *conference on knowledge discovery & data mining*, pages 1269–1278, 2019.
- 393 [26] Jinbo Li, Hesam Izakian, Witold Pedrycz, and Iqbal Jamal. Clustering-based anomaly detection  
394 in multivariate time series data. *Applied Soft Computing*, 100:106919, 2021.
- 395 [27] Zheng Li, Yue Zhao, Nicola Botta, Cezar Ionescu, and Xiyang Hu. Copod: copula-based outlier  
396 detection. In *2020 IEEE international conference on data mining (ICDM)*, pages 1118–1123.  
397 IEEE, 2020.
- 398 [28] Zheng Li, Yue Zhao, Xiyang Hu, Nicola Botta, Cezar Ionescu, and George Chen. Ecod: Unsu-  
399 pervised outlier detection using empirical cumulative distribution functions. *IEEE Transactions*  
400 *on Knowledge and Data Engineering*, 2022.
- 401 [29] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In *2008 eighth ieee*  
402 *international conference on data mining*, pages 413–422. IEEE, 2008.
- 403 [30] Yixin Liu, Zhao Li, Shirui Pan, Chen Gong, Chuan Zhou, and George Karypis. Anomaly  
404 detection on attributed networks via contrastive self-supervised learning. *IEEE transactions on*  
405 *neural networks and learning systems*, 33(6):2378–2392, 2021.
- 406 [31] Xuexiong Luo, Jia Wu, Amin Beheshti, Jian Yang, Xiankun Zhang, Yuan Wang, and Shan Xue.  
407 Comga: Community-aware attributed graph anomaly detection. In *Proceedings of the Fifteenth*  
408 *ACM International Conference on Web Search and Data Mining*, pages 657–665, 2022.
- 409 [32] Larry M Manevitz and Malik Yousef. One-class svms for document classification. *Journal of*  
410 *machine Learning research*, 2(Dec):139–154, 2001.
- 411 [33] Geoffrey J McLachlan. Mahalanobis distance. *Resonance*, 4(6):20–26, 1999.
- 412 [34] Emmanuel Müller, Patricia Iglesias Sánchez, Yvonne Mülle, and Klemens Böhm. Ranking  
413 outlier nodes in subspaces of attributed graphs. In *2013 IEEE 29th international conference on*  
414 *data engineering workshops (ICDEW)*, pages 216–222. IEEE, 2013.
- 415 [35] Isaac Newton. *Philosophiae naturalis principia mathematica*, volume 1. G. Brookman, 1833.
- 416 [36] David Peel and Geoffrey J McLachlan. Robust mixture modelling using the t distribution.  
417 *Statistics and computing*, 10:339–348, 2000.
- 418 [37] Zhen Peng, Minnan Luo, Jundong Li, Luguo Xue, and Qinghua Zheng. A deep multi-view  
419 framework for anomaly detection on attributed networks. *IEEE Transactions on Knowledge*  
420 *and Data Engineering*, 34(6):2539–2552, 2020.

- 421 [38] Leif E Peterson. K-nearest neighbor. *Scholarpedia*, 4(2):1883, 2009.
- 422 [39] Tomáš Pevný. Loda: Lightweight on-line detector of anomalies. *Machine Learning*, 102:275–  
423 304, 2016.
- 424 [40] Douglas A Reynolds et al. Gaussian mixture models. *Encyclopedia of biometrics*, 741(659-663),  
425 2009.
- 426 [41] Lukas Ruff, Jacob R Kauffmann, Robert A Vandermeulen, Grégoire Montavon, Wojciech  
427 Samek, Marius Kloft, Thomas G Dietterich, and Klaus-Robert Müller. A unifying review of  
428 deep and shallow anomaly detection. *Proceedings of the IEEE*, 109(5):756–795, 2021.
- 429 [42] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui,  
430 Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In  
431 *International conference on machine learning*, pages 4393–4402. PMLR, 2018.
- 432 [43] Mayu Sakurada and Takehisa Yairi. Anomaly detection using autoencoders with nonlinear  
433 dimensionality reduction. In *Proceedings of the MLSDA 2014 2nd workshop on machine  
434 learning for sensory data analysis*, pages 4–11, 2014.
- 435 [44] Osman Salem, Yaning Liu, Ahmed Mehaoua, and Raouf Boutaba. Online anomaly detection in  
436 wireless body area networks for reliable healthcare monitoring. *IEEE journal of biomedical  
437 and health informatics*, 18(5):1541–1551, 2014.
- 438 [45] Hanyu Song, Peizhao Li, and Hongfu Liu. Deep clustering based fair outlier detection. In  
439 *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*,  
440 pages 1481–1489, 2021.
- 441 [46] Jianheng Tang, Jiabin Li, Ziqi Gao, and Jia Li. Rethinking graph neural networks for anomaly  
442 detection. In *International Conference on Machine Learning*, pages 21076–21089. PMLR,  
443 2022.
- 444 [47] Shantanu Thakoor, Corentin Tallec, Mohammad Gheshlaghi Azar, Mehdi Azabou, Eva L Dyer,  
445 Remi Munos, Petar Veličković, and Michal Valko. Large-scale representation learning on graphs  
446 via bootstrapping. *arXiv preprint arXiv:2102.06514*, 2021.
- 447 [48] Laurens Van Der Maaten. Learning a parametric embedding by preserving local structure. In  
448 *Artificial intelligence and statistics*, pages 384–391. PMLR, 2009.
- 449 [49] Svante Wold, Kim Esbensen, and Paul Geladi. Principal component analysis. *Chemometrics  
450 and intelligent laboratory systems*, 2(1-3):37–52, 1987.
- 451 [50] Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering  
452 analysis. In *International conference on machine learning*, pages 478–487. PMLR, 2016.
- 453 [51] Hongzuo Xu, Guansong Pang, Yijie Wang, and Yongjun Wang. Deep isolation forest for  
454 anomaly detection. *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- 455 [52] Xiaowei Xu, Nurcan Yuruk, Zhidan Feng, and Thomas AJ Schweiger. Scan: a structural  
456 clustering algorithm for networks. In *Proceedings of the 13th ACM SIGKDD international  
457 conference on Knowledge discovery and data mining*, pages 824–833, 2007.
- 458 [53] Zhiming Xu, Xiao Huang, Yue Zhao, Yushun Dong, and Jundong Li. Contrastive attributed  
459 network anomaly detection with data augmentation. In *Advances in Knowledge Discovery and  
460 Data Mining: 26th Pacific-Asia Conference, PAKDD 2022, Chengdu, China, May 16–19, 2022,  
461 Proceedings, Part II*, pages 444–457. Springer, 2022.
- 462 [54] Xu Yuan, Na Zhou, Shuo Yu, Huafei Huang, Zhikui Chen, and Feng Xia. Higher-order structure  
463 based anomaly detection on attributed networks. In *2021 IEEE International Conference on  
464 Big Data (Big Data)*, pages 2691–2700. IEEE, 2021.
- 465 [55] Yu Zheng, Ming Jin, Yixin Liu, Lianhua Chi, Khoa T Phan, and Yi-Ping Phoebe Chen. Generative  
466 and contrastive self-supervised learning for graph anomaly detection. *IEEE Transactions  
467 on Knowledge and Data Engineering*, 2021.

- 468 [56] Shuang Zhou, Xiao Huang, Ninghao Liu, Qiaoyu Tan, and Fu-Lai Chung. Unseen anomaly  
469 detection on networks via multi-hypersphere learning. In *Proceedings of the 2022 SIAM*  
470 *International Conference on Data Mining (SDM)*, pages 262–270. SIAM, 2022.
- 471 [57] Shuang Zhou, Qiaoyu Tan, Zhiming Xu, Xiao Huang, and Fu-lai Chung. Subtractive aggregation  
472 for attributed network anomaly detection. In *Proceedings of the 30th ACM International*  
473 *Conference on Information & Knowledge Management*, pages 3672–3676, 2021.
- 474 [58] Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Daeki Cho, and  
475 Haifeng Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection.  
476 In *International conference on learning representations*, 2018.

---

**Algorithm 1** Model training for UniCAD

---

**Input:** data points  $\mathbf{X}$ , cluster number  $K$ , outlier ratio  $l$ , tolerance  $\lambda$ , iterations  $t$

**Output:** network parameters  $\Theta$ , mixture parameters  $\{\omega_k, \boldsymbol{\mu}_k, \Sigma_k\}$

```
1: Initialize  $\Theta$  and  $\{\boldsymbol{\mu}_k, \omega_k, \Sigma_k\}$ ;
2: for  $i = 1$  to  $t$  do
3:   if  $i = 1$  then
4:      $\mathbf{X}_i \leftarrow \mathbf{X}$ ;
5:   else
6:     Re-order the point in  $\mathbf{X}$  such that  $o_1 \geq \dots \geq o_n$ ;
7:      $L_i \leftarrow \{x_1, \dots, x_{\lfloor N \cdot l \rfloor}\}$ ;
8:      $\mathbf{X}_i \leftarrow \mathbf{X} \setminus L_i$ ;
9:   end if
10:  Update  $\Theta$  with Equation (15);
11:  while  $|J(\Theta, \Phi) - J^{old}(\Theta, \Phi)| > \lambda$  do
12:     $J^{old}(\Theta, \Phi) = J(\Theta, \Phi)$ ;
13:    Calculate  $\boldsymbol{\tau}$  with Equation (10);
14:    Update  $\{\omega_k, \boldsymbol{\mu}_k, \Sigma_k\}$  with Equation (12), (13) and (14);
15:  end while
16:  Calculate  $o_i$  with Equation (9);
17: end for
18: return  $\Theta$  and  $\{\omega_k, \boldsymbol{\mu}_k, \Sigma_k\}$ 
```

---

## 477 A Iterative Training Algorithm

478 The pseudocode for training the model is presented in Algorithm 1. Initially, all parameters undergo  
479 random initialization. In subsequent iterations, following the initial round, the outlier set  $L$  undergoes  
480 updates based on the anomaly score  $o_i$ . This is succeeded by the adjustment of the network parameters  
481  $\Theta$  based on  $\mathbf{x}_i$ , further optimizing the performance of  $\Theta$  through the utilization of the estimated  
482 parameters  $\boldsymbol{\mu}_k, \omega_k, \Sigma_k$ . The essence of the algorithm is embedded in its alternating optimization  
483 strategy, iteratively refining the accuracy of representation learning and mixed model parameter  
484 estimation, thereby augmenting the overall training effectiveness of the model.

## 485 B Derivation of EM Algorithm

486 This appendix provides the detailed derivation of the Expectation-Maximization (EM) algorithm  
487 for optimizing the parameters of a mixture model based on Student's t-distribution. The focus is  
488 on deriving analytical solutions for the maximization of the parameters  $\Phi = \{\boldsymbol{\mu}_k, \Sigma_k, \omega_k\}$  of the  
489 mixture components. The EM algorithm alternates between two steps:

490 **In the E-step**, we calculate the posterior probabilities  $\tau_{ik}$ , representing the probability of data point  
491  $i$  belonging to cluster  $k$ , given the current parameters. The posterior probabilities for a Student's  
492 t-distribution mixture model are formulated as:

$$\tau_{ik} = \frac{\omega_k \cdot p(\mathbf{z}_i | \boldsymbol{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \omega_j \cdot p(\mathbf{z}_i | \boldsymbol{\mu}_j, \Sigma_j)}, \quad (16)$$

493 where  $\tau(\mathbf{z}_i | \boldsymbol{\mu}_k, \Sigma_k)$  denotes the Student's t-distribution for data point  $i$  with respect to cluster  $k$ , and  
494  $K$  is the number of mixture components.

495 The Student's t-distribution is depicted as a hierarchical conditional probability, resembling a Gaussian  
496 distribution with an accuracy scale factor  $\mathbf{u}$ , where its latent variable follows a gamma distribution.  
497 Adopting a degree of freedom  $\nu = 1$ , the value of  $\mathbf{u}_{ik}$  is given by:

$$\mathbf{u}_{ik} = \frac{\nu + 1}{\nu + D_M(z_i, \boldsymbol{\mu}_k)} = \frac{2}{1 + D_M(z_i, \boldsymbol{\mu}_k)} \quad (17)$$

498 **In the M-step**, we update the parameters  $\Phi = \{\omega_k, \boldsymbol{\mu}_k, \Sigma_k\}$  using the derivatives obtained in  
499 the previous steps. In our model, the likelihood function for a Student's-t Distribution Mixture Model

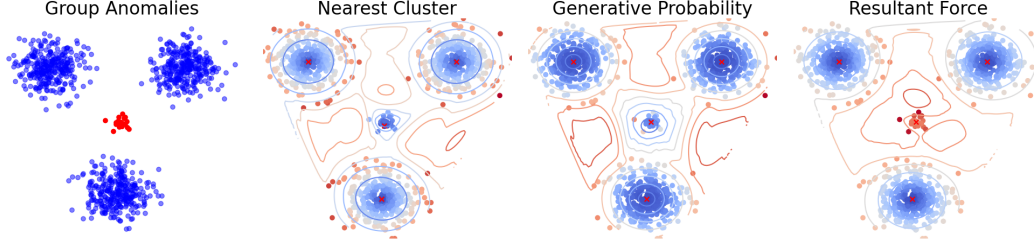


Figure 3: Score comparison with other methods.

500 (SMM) is represented as:

$$L(\omega, \boldsymbol{\mu}, \Sigma) = \sum_{i=1}^N \sum_{k=1}^K \omega_k \cdot \frac{\pi^{-1} \cdot |\Sigma_k|^{-\frac{1}{2}}}{1 + (\mathbf{z}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{z}_i - \boldsymbol{\mu}_k)}, \quad (18)$$

501 where  $\omega_k$  are the mixture weights,  $\Sigma_k$  the covariance matrices,  $\boldsymbol{\mu}_k$  the means, and  $\mathbf{z}_i$  the data points.

502 The derivative with respect to  $\omega_k$  must consider the constraint that the sum of the mixture weights  
 503 equals 1, i.e.,  $\sum_k \omega_k = 1$ . Hence, we introduce a Lagrange multiplier  $\lambda$  to address this constraint  
 504 and construct the Lagrangian  $L'$ :

$$L'(\omega, \boldsymbol{\mu}, \Sigma, \lambda) = L(\omega, \boldsymbol{\mu}, \Sigma) + \lambda \left( 1 - \sum_{k=1}^K \omega_k \right), \quad (19)$$

505 The derivative with respect to  $\omega_k$  is:

$$\frac{\partial L'}{\partial \omega_k} = \frac{\partial L}{\partial \omega_k} - \lambda, \quad (20)$$

506 Substituting the definition of  $L(\omega, \boldsymbol{\mu}, \Sigma)$ , we obtain:

$$\frac{\partial L}{\partial \omega_k} = \sum_i \frac{p(\mathbf{z}_i | \boldsymbol{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \omega_j \cdot p(\mathbf{z}_i | \boldsymbol{\mu}_j, \Sigma_j)} = \sum_i \frac{\tau_{ik}}{\omega_k}, \quad (21)$$

507 To solve for  $\omega_k$ , we first multiply both sides of the equation by  $\omega_k$  and apply the constraint condition:

$$\sum_k \omega_k \left( \sum_i \frac{\tau_{ik}}{\omega_k} - \lambda \right) = 0, \quad (22)$$

508 Upon further organization, we find that the Lagrange multiplier  $\lambda$  actually equals the total number of  
 509 data points  $N$  (since  $\sum_i \tau_{ik} = N_k$ , where  $N_k$  is the expected total number of data points belonging  
 510 to the  $k$ th component, and the sum of all  $N_k$  equals the total number of data points  $N$ ).

511 Finally, we can solve for  $\omega_k$ :

$$\omega_k = \frac{\sum_i \tau_{ik}}{N}, \quad (23)$$

512 This result indicates that the weight  $\omega_k$  of each mixture component equals the proportion of the  
 513 posterior probabilities of the data points it contains relative to all data points.

514 To update  $\boldsymbol{\mu}_k$  and  $\Sigma_k$ , we consider the conditional expectation of the data log-likelihood function:

$$Q(\boldsymbol{\mu}_k, \Sigma_k) = \sum_{i=1}^N \tau_{ik} \left( -\log(\pi) - \frac{1}{2} \log |\Sigma_k| + \frac{1}{2} \log u_{ik} - \frac{1}{2} \mathbf{u}_{ik} (\mathbf{z}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{z}_i - \boldsymbol{\mu}_k) \right) \quad (24)$$

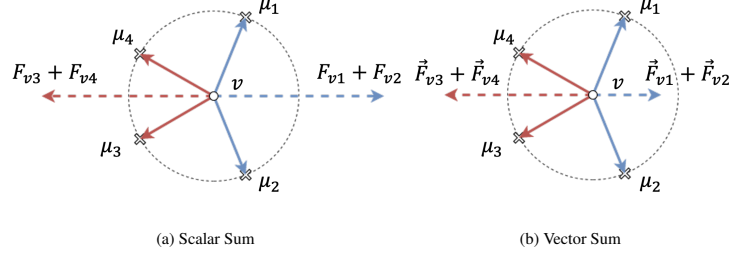


Figure 4: Analysis of gravitational force.

515 Maximizing  $Q(\boldsymbol{\mu}_k, \Sigma_k)$  with respect to  $\boldsymbol{\mu}_k$  leads to:

$$\frac{\partial Q}{\partial \boldsymbol{\mu}_k} = \frac{1}{2} \sum_{i=1}^N \tau_{ik} \mathbf{u}_{ik} (2\Sigma_k^{-1} \boldsymbol{\mu}_k - 2\Sigma_k^{-1} \mathbf{z}_{ik}) \quad (25)$$

516 Setting  $\frac{\partial Q}{\partial \boldsymbol{\mu}_k} = 0$  results in the updated mean  $\boldsymbol{\mu}_k^{(t+1)}$ :

$$\boldsymbol{\mu}_k^{(t+1)} = \sum_{i=1}^n \left( \tau_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} \mathbf{z}_i \right) / \sum_{i=1}^n \left( \tau_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} \right). \quad (26)$$

517 Considering the derivative of  $Q(\boldsymbol{\mu}_k, \Sigma_k)$  with respect to  $\Sigma_k^{-1}$ :

$$\frac{\partial Q}{\partial \Sigma_k^{-1}} = \frac{1}{2} \sum_{i=1}^N \tau_{ik} (\Sigma_k - \mathbf{u}_{ik} (\mathbf{z}_i - \boldsymbol{\mu}_k) \times (\mathbf{z}_i - \boldsymbol{\mu}_k)^T). \quad (27)$$

518 Setting  $\frac{\partial Q}{\partial \Sigma_k} = 0$  yields the updated covariance matrix  $\Sigma_k^{(t+1)}$ :

$$\Sigma_k^{(t+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(t+1)} \mathbf{u}_{ik}^{(t+1)} (\mathbf{z}_i - \boldsymbol{\mu}_k^{(t+1)}) (\mathbf{z}_i - \boldsymbol{\mu}_k^{(t+1)})^T}{\sum_{j=1}^K \tau_{ij}^{(t+1)}}. \quad (28)$$

## 519 C Anomaly Score with Vector Sum

### 520 C.1 Advantages

521 Here we discuss the advantages of employing vector sum in anomaly score with a toy example.

522 The application of the vector sum principle extends beyond physical mechanics and finds relevance  
 523 in various domains. In relational embedding [5], for example, relationships can be represented as  
 524 vectors. Aggregating these vectors allows for capturing complexities like transitivity, symmetry, and  
 525 antisymmetry.

526 Similarly, in our context, the vector sum can help capture more complex relationships along clusters.  
 527 Consider Figure 4 as an example, where a sample  $v$  is attracted by two groups of cluster  
 528 prototypes ( $\{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2\}$ ,  $\{\boldsymbol{\mu}_3, \boldsymbol{\mu}_4\}$ ) with the same mass and sample-prototype distances ( $\tilde{m}_1 = \tilde{m}_2 =$   
 529  $\tilde{m}_3 = \tilde{m}_4, \tilde{r}_{v1} = \tilde{r}_{v2} = \tilde{r}_{v3} = \tilde{r}_{v4}$ ). Without considering the direction of the forces, the two groups  
 530 of prototypes would attract the sample with equal forces. However, we argue that the two groups of  
 531 prototypes should exert different influences. A sample close to two clusters with a large difference  
 532 ( $\{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2\}$ ) is more likely to be an anomaly compared to a sample that is close to two clusters with  
 533 a smaller difference ( $\{\boldsymbol{\mu}_3, \boldsymbol{\mu}_4\}$ ). For example, in a social network, a user who equally likes two  
 534 extremely different communities, like money-saving tips and luxury items, is more anomalous than  
 535 a user who equally likes two similar communities, like private jets and luxury items. Applying  
 536 the vector sum, the total force of  $\{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2\}$  is much smaller than that of  $\{\boldsymbol{\mu}_3, \boldsymbol{\mu}_4\}$ . As the anomaly  
 537 score is inversely related to the total force, it is more anomalous when equally attracted by  $\{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2\}$   
 538 with large difference. This indicates that *the vector sum successfully captures subtle differences*  
 539 *in the distinctions among multiple clusters, thereby assisting in the identification of more accurate*  
 540 *anomalies.*



## 541 C.2 Toy Example

542 In the appendix, as illustrated in Figure 3, we investigated a toy example. We discussed a specific  
543 pattern of anomalies termed *group anomalies*, where a small number of anomalous samples cluster  
544 together. It is crucial to note that we do not claim this anomaly pattern is common in real-world data;  
545 our goal is merely to point out a specific anomaly pattern that is challenging for traditional cluster-  
546 based anomaly detection methods to detect. Specifically, we utilize three Gaussian distributions with  
547 high variance (each generating 300 data samples) and one with lower variance (generating 30 data  
548 samples). Because the samples from the smaller Gaussian follow a different generative mechanism  
549 and represent a minority in the dataset, we consider them anomalies.

550 We set the cluster number for KMeans-- and GMM at four, indicating that the Gaussian distribution  
551 comprising anomalous samples was also recognized as a cluster. KMeans-- employs a cluster-based  
552 approach, using the distance to the nearest cluster center as the anomaly score, while GMM uses  
553 a probability-based approach, considering the samples' likelihood in the mixture model as the  
554 anomaly score. However, both approaches are ineffective in this scenario. Rather than identifying the  
555 small cluster as anomalous, they tend to misidentify samples on the peripheries of larger clusters as  
556 anomalies.

557 By contrast, our scoring method views the entire small cluster as more likely anomalous, followed by  
558 outlier samples on the margins of the larger clusters. This visualization provides a perspective that  
559 distinguishes our method from previous efforts.

## 560 D Experimental Supplementary

### 561 D.1 Benchmark Datasets Details

562 Due to space constraints in the main text, we utilized 30 public datasets from ADBench [17], covering  
563 all different types of data. The details of the 30 datasets are presented in Table 4.

### 564 D.2 Baselines Details

565 A comprehensive overview of the unsupervised anomaly detection methods is presented below.

#### 566 D.2.1 Traditional Models

- 567 • **Subspace Outlier Detection (SOD) [24]:** Identifies outliers in varying subspaces of a high-  
568 dimensional feature space, targeting anomalies that emerge in lower-dimensional projections.
- 569 • **Histogram-based Outlier Detection (HBOS) [16]:** Assumes feature independence and calculates  
570 outlyingness via histograms, offering scalability and efficiency.

#### 571 D.2.2 Linear Models

- 572 • **Principal Component Analysis (PCA) [49]:** Utilizes singular value decomposition for dimension-  
573 ality reduction, with anomalies indicated by reconstruction errors.
- 574 • **One-class SVM (OCSVM) [32]:** Defines a decision boundary to separate normal samples from  
575 outliers, maximizing the margin from the data origin.

#### 576 D.2.3 Density-based Models

- 577 • **Local Outlier Factor (LOF) [6]:** Measures local density deviation, marking samples as outliers if  
578 they lie in less dense regions compared to their neighbors.
- 579 • **K-Nearest Neighbors (KNN) [38]:** Anomaly scores are assigned based on the distance to the k-th  
580 nearest neighbor, embodying a simple yet effective approach.

#### 581 D.2.4 Ensemble-based Models

- 582 • **Lightweight On-line Detector of Anomalies (LODA) [39]:** An ensemble method suitable for  
583 real-time processing and adaptable to concept drift through random projections and histograms.
- 584 • **Isolation Forest (IForest) [29]:** Isolates anomalies by randomly selecting features and split values,  
585 leveraging the ease of isolating anomalies to identify them efficiently.

Table 4: Statistics of tabular benchmark datasets.

Data	# Samples	# Features	# Anomaly	% Anomaly	Category
anthyroid	7200	6	534	7.42	Healthcare
backdoor	95329	196	2329	2.44	Network
breastw	683	9	239	34.99	Healthcare
campaign	41188	62	4640	11.27	Finance
celeba	202599	39	4547	2.24	Image
census	299285	500	18568	6.20	Sociology
glass	214	7	9	4.21	Forensic
Hepatitis	80	19	13	16.25	Healthcare
http	567498	3	2211	0.39	Web
Ionosphere	351	33	126	35.90	Oryctognosy
landsat	6435	36	1333	20.71	Astronautics
Lymphography	148	18	6	4.05	Healthcare
magic.gamma	19020	10	6688	35.16	Physical
mnist	7603	100	700	9.21	Image
musk	3062	166	97	3.17	Chemistry
pendigits	6870	16	156	2.27	Image
Pima	768	8	268	34.90	Healthcare
satellite	6435	36	2036	31.64	Astronautics
satimage-2	5803	36	71	1.22	Astronautics
shuttle	49097	9	3511	7.15	Astronautics
skin	245057	3	50859	20.75	Image
Stamps	340	9	31	9.12	Document
thyroid	3772	6	93	2.47	Healthcare
vertebral	240	6	30	12.50	Biology
vowels	1456	12	50	3.43	Linguistics
Waveform	3443	21	100	2.90	Physics
WBC	223	9	10	4.48	Healthcare
Wilt	4819	5	257	5.33	Botany
wine	129	13	10	7.75	Chemistry
WPBC	198	33	47	23.74	Healthcare

### 586 D.2.5 Probability-based Models

- 587 • **Deep Autoencoding Gaussian Mixture Model (DAGMM) [58]:** Combines a deep autoencoder  
588 with a GMM for anomaly scoring, utilizing both low-dimensional representation and reconstruction  
589 error.
- 590 • **Empirical-Cumulative-distribution-based Outlier Detection (ECOD) [28]:** Uses ECDFs to  
591 estimate feature densities independently, targeting outliers in distribution tails.
- 592 • **Copula Based Outlier Detector (COPD) [27]:** A hyperparameter-free method leveraging  
593 empirical copula models for interpretable and efficient outlier detection.

### 594 D.2.6 Cluster-based Models

- 595 • **DBSCAN [13]:** A density-based clustering algorithm that identifies clusters based on the density  
596 of data points, effectively separating high-density clusters from low-density noise, and is widely  
597 used for anomaly detection in spatial data.
- 598 • **Clustering Based Local Outlier Factor (CBLOF) [18]:** Calculates anomaly scores based on  
599 cluster distances, using global data distribution.
- 600 • **KMeans-- [45]:** Extends k-means to include outlier detection in the clustering process, offering an  
601 integrated approach to anomaly detection.
- 602 • **Deep Clustering-based Fair Outlier Detection (DCFOD) [9]:** Enhances outlier detection with a  
603 focus on fairness, combining deep clustering and adversarial training for representation learning.

Table 5: AUCROC of 17 unsupervised algorithms on 30 tabular benchmark datasets. In each dataset, the algorithm with the highest AUCROC is marked in **red**, the second highest in **blue**, and the third highest in **green**.

Dataset	SOD	HBOS	PCA	OC SVM	LOF	KNN	LODA	IForest	DA GMM	ECOD	COPOD	DB SCAN	CBLOF	DCOD	KMeans--	Deep SVDD	DIF	UniCAD (Scalar)	UniCAD (Vector)
anthyroid	77.38	60.15	66.24	57.23	70.20	71.69	41.02	<b>82.01</b>	56.53	<b>78.66</b>	76.80	50.08	62.28	55.01	64.99	76.09	66.76	75.27	72.72
backdoor	68.77	71.56	80.16	85.04	85.79	80.58	66.38	72.15	53.98	86.08	80.97	76.55	81.91	79.57	89.11	78.83	<b>92.87</b>	87.28	<b>89.24</b>
breastw	93.97	98.94	95.13	80.30	40.61	97.01	98.49	98.32	N/A	<b>99.17</b>	<b>99.68</b>	85.20	96.86	99.02	97.05	63.36	77.45	98.15	98.56
campaign	69.16	<b>78.55</b>	72.78	65.70	59.04	72.27	51.67	71.71	56.03	76.10	<b>77.69</b>	50.60	64.34	63.16	63.51	54.42	67.53	73.52	73.64
celeba	48.44	76.18	79.38	70.70	38.95	59.63	60.17	70.41	44.74	76.48	75.68	50.36	73.99	<b>91.41</b>	56.76	45.17	65.29	81.38	<b>82.00</b>
census	62.12	64.89	<b>68.74</b>	54.90	47.46	66.88	37.14	59.52	59.65	<b>67.63</b>	<b>69.07</b>	58.50	60.17	<b>72.84</b>	63.33	54.16	59.66	67.90	67.84
glass	73.36	77.23	66.29	35.36	69.20	<b>82.29</b>	73.13	77.13	76.09	65.83	72.43	54.55	78.30	78.07	77.30	55.71	<b>84.57</b>	79.52	82.17
Hepatitis	67.83	79.85	75.95	67.75	38.06	52.76	64.87	69.75	54.80	75.22	<b>82.05</b>	68.12	73.05	48.38	64.64	57.45	74.24	75.53	<b>80.62</b>
hnp	78.04	99.53	<b>99.72</b>	99.59	27.46	3.37	12.48	<b>99.96</b>	N/A	98.10	99.29	49.97	99.60	99.53	99.55	60.38	99.49	99.53	99.52
Ionosphere	86.37	62.49	79.19	75.92	90.59	88.26	78.42	84.50	73.41	73.15	79.34	81.12	90.79	57.78	<b>91.36</b>	53.94	89.74	<b>92.04</b>	90.37
landsat	59.54	55.14	35.76	36.15	53.90	57.95	38.17	47.64	43.92	36.10	41.55	50.17	<b>63.69</b>	33.40	55.31	<b>62.48</b>	54.84	49.60	57.37
Lymphography	71.22	99.49	<b>99.82</b>	99.54	89.86	55.91	85.55	<b>99.81</b>	72.11	99.52	99.48	74.16	99.81	81.19	<b>100.00</b>	71.91	83.67	99.29	99.73
mnist	60.10	60.42	85.29	82.95	67.13	80.58	72.27	80.98	67.23	74.61	77.74	50.00	79.96	65.23	82.45	50.98	<b>88.16</b>	<b>86.00</b>	<b>86.64</b>
musk	74.09	<b>100.00</b>	<b>100.00</b>	80.58	41.18	69.89	95.11	<b>99.99</b>	76.85	95.40	94.20	50.00	<b>100.00</b>	42.19	72.16	66.02	98.22	99.92	<b>100.00</b>
pendigits	66.29	93.04	93.73	93.75	47.99	72.95	89.10	94.76	64.22	93.01	90.68	55.33	<b>96.93</b>	94.33	94.37	27.32	93.79	95.12	<b>95.52</b>
Pima	61.25	71.07	70.77	66.92	65.71	73.43	65.93	72.87	55.93	63.05	69.10	51.39	71.49	72.16	70.44	49.49	67.28	75.16	74.87
satellite	63.96	<b>74.80</b>	59.62	59.02	55.88	65.18	61.98	70.43	62.33	58.09	63.20	55.52	71.32	55.97	67.71	57.40	74.52	72.46	<b>77.65</b>
satimage-2	83.08	97.65	97.62	97.35	47.36	92.60	97.56	99.16	96.29	96.28	97.21	75.74	<b>99.84</b>	86.01	<b>99.88</b>	55.68	<b>99.63</b>	<b>99.87</b>	<b>99.88</b>
shuttle	69.51	98.63	98.62	97.40	57.11	69.64	60.95	<b>99.56</b>	97.92	99.13	<b>99.35</b>	50.40	93.07	97.20	69.97	51.81	97.00	99.15	98.75
skin	60.35	60.15	45.26	49.45	46.47	<b>71.46</b>	45.75	68.21	N/A	49.08	47.55	50.00	68.03	64.34	65.47	45.69	66.36	<b>72.26</b>	<b>69.69</b>
Stamps	73.26	90.73	91.47	83.86	51.26	68.61	87.18	91.21	88.89	87.87	93.40	52.08	69.89	<b>93.41</b>	79.78	59.48	87.95	91.37	<b>94.18</b>
thyroid	92.81	95.62	96.34	87.92	86.86	95.93	74.30	<b>98.30</b>	79.75	<b>97.94</b>	94.30	53.57	94.74	78.55	92.26	52.14	96.26	97.66	97.48
vertebral	40.32	28.56	37.06	37.99	49.29	33.79	30.57	36.66	<b>53.20</b>	40.66	25.64	<b>49.74</b>	41.01	38.13	38.14	37.81	47.20	33.11	47.37
vowels	92.65	72.21	65.29	61.59	93.12	<b>97.26</b>	70.36	73.94	60.58	62.24	53.15	57.50	92.12	51.56	<b>93.45</b>	49.87	<b>81.02</b>	88.38	92.09
Waveform	68.57	68.77	65.48	56.29	73.32	73.78	60.13	71.47	49.35	62.36	<b>75.03</b>	66.41	71.27	63.47	74.35	53.94	<b>75.33</b>	71.81	74.29
WBC	94.60	98.72	98.20	<b>99.03</b>	54.17	90.56	96.91	<b>99.01</b>	N/A	<b>99.11</b>	<b>99.11</b>	87.43	96.88	94.92	97.45	62.46	81.27	97.68	98.93
Wilt	<b>53.25</b>	32.49	20.39	31.28	50.65	48.42	26.42	41.94	37.29	36.30	33.40	49.96	34.50	44.71	34.91	45.90	39.46	48.95	<b>52.56</b>
wine	46.11	<b>91.36</b>	84.37	73.07	37.74	44.98	90.12	80.37	61.70	77.22	88.65	40.33	27.14	82.18	27.36	64.26	41.69	82.72	<b>95.25</b>
WPPC	<b>51.28</b>	51.24	46.01	45.35	41.41	46.59	49.31	46.63	47.80	46.65	49.34	<b>52.22</b>	45.32	49.67	45.01	44.01	44.69	48.02	49.90
Avg. Rank	<b>11.00</b>	<b>8.26</b>	<b>8.98</b>	<b>11.59</b>	<b>13.59</b>	<b>10.00</b>	<b>13.24</b>	<b>7.09</b>	<b>13.24</b>	<b>9.19</b>	<b>8.29</b>	<b>14.21</b>	<b>8.07</b>	<b>10.90</b>	<b>8.71</b>	<b>15.48</b>	<b>8.38</b>	<b>5.41</b>	<b>3.59</b>

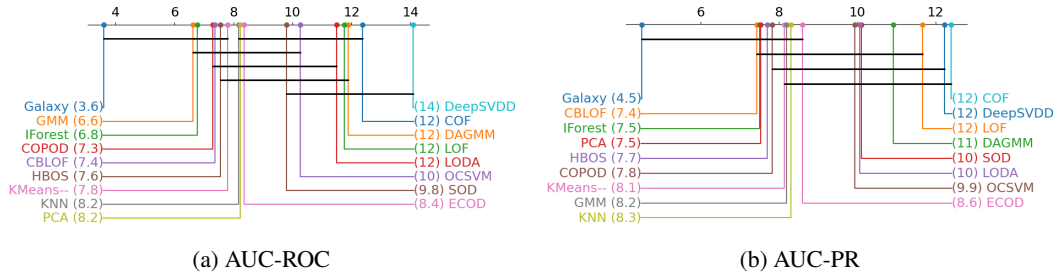


Figure 5: Critical difference diagrams for AUC-ROC and AUC-PR.

## 604 D.2.7 Neural Network-based Models

- 605 • **Deep Support Vector Data Description (DeepSVDD) [42]:** Minimizes the volume of a hyper-
- 606 sphere enclosing network data representations, isolating anomalies outside this sphere.
- 607 • **Deep Isolation Forest for Anomaly Detection (DIF) [51]:** Utilizes deep learning to enhance
- 608 traditional isolation forest techniques, offering improved anomaly detection in complex datasets
- 609 with minimal parameter tuning.

610 Each method’s unique mechanism and application context provide a rich landscape of techniques  
 611 for unsupervised anomaly detection, illustrating the field’s diverse methodologies and the breadth of  
 612 approaches to tackling anomaly detection challenges.

## 613 D.3 Supplementary Experimental Results

614 In the appendix, we detail the statistical analysis conducted to compare the performance of various  
 615 anomaly detectors. We obtained this diagram by conducting a Friedman test (p-value: 4.657e-19),  
 616 indicating significant differences among different detectors. We utilized average ranks and the  
 617 Nemenyi test to generate the critical difference diagram, as shown in Figure 5. It is noteworthy that  
 618 the vector version exhibits significantly superior performance compared to the scalar version across  
 619 more methods. The detailed outcomes for the AUCROC and AUCPR metrics, spanning 30 datasets  
 620 and against 17 baseline approaches, are showcased in Table 5 and Table 6.

## 621 D.4 Complexity Analysis

622 The complexity of each iteration in UniCAD involves three parts: constructing the outlier set,  
 623 updating the network parameters  $\Theta$ , and optimizing the mixture model using the EM algorithm.

Table 6: AUCPR of 17 unsupervised algorithms on 30 tabular benchmark datasets. In each dataset, the algorithm with the highest AUCPR is marked in **red**, the second highest in **blue**, and the third highest in **green**.

Dataset	SOD	HBOS	PCA	OC SVM	LOF	KNN	LODA	IForest	DA GMM	ECOD	COPOD	DB SCAN	CBLOF	DCOD	KMeans--	Deep SVDD	DIF	UniCAD (Scalar)	UniCAD (Vector)
anthyroid	18.84	16.99	16.12	10.37	15.71	16.74	7.06	<b>30.47</b>	9.64	25.35	16.58	7.60	13.74	10.01	15.41	21.75	18.93	<b>26.37</b>	25.03
backdoor	37.07	4.96	31.29	8.79	26.14	<b>44.37</b>	13.84	4.75	5.47	10.72	7.69	21.04	7.03	6.77	15.47	<b>55.70</b>	41.46	37.77	36.36
breastw	84.88	97.71	95.11	82.70	28.55	92.19	97.04	96.04	N/A	<b>98.54</b>	<b>99.40</b>	78.42	91.94	96.83	92.25	48.60	50.65	94.47	95.90
campaign	19.14	<b>38.01</b>	27.90	29.25	14.59	27.18	14.11	32.26	14.54	36.65	<b>38.58</b>	11.43	20.88	19.61	18.86	16.75	26.52	27.66	27.12
celeba	2.36	13.82	<b>15.89</b>	10.73	1.73	3.14	4.04	8.96	1.95	13.96	13.69	2.32	11.22	<b>17.48</b>	3.19	2.73	5.44	15.12	14.66
census	8.54	8.68	<b>10.02</b>	6.82	5.48	9.04	5.03	7.78	9.03	9.46	9.92	7.52	7.52	<b>10.92</b>	8.13	8.42	7.42	9.70	9.75
glass	18.73	11.82	10.05	8.02	<b>20.11</b>	<b>20.26</b>	13.37	10.99	<b>24.58</b>	15.35	9.78	6.88	11.57	9.66	14.66	8.46	18.86	13.29	15.33
Hepatitis	24.73	37.73	36.65	29.44	13.67	21.95	30.90	26.25	22.93	32.80	<b>41.50</b>	22.31	36.54	19.53	25.14	30.04	34.93	36.08	<b>43.37</b>
http	8.32	44.79	<b>56.43</b>	46.86	3.82	0.70	0.67	<b>90.83</b>	N/A	16.61	35.19	0.37	47.53	44.03	45.09	13.39	41.72	45.53	45.52
Ionosphere	85.88	41.78	73.92	74.54	88.07	<b>90.41</b>	73.04	80.41	64.97	64.69	69.89	63.04	89.77	47.63	<b>91.36</b>	43.24	87.45	89.55	87.61
landsat	26.38	22.03	16.18	16.21	24.69	24.65	18.86	19.81	24.48	16.24	17.48	20.80	<b>31.05</b>	15.57	22.40	<b>36.92</b>	24.35	20.84	23.27
Lymphography	22.00	91.83	<b>97.02</b>	93.59	23.08	38.69	44.54	<b>97.31</b>	19.52	90.87	88.68	7.66	<b>97.31</b>	12.34	<b>100.00</b>	34.58	32.84	91.69	96.66
mnist	19.15	12.51	39.93	33.20	20.90	35.53	25.86	27.71	23.75	17.45	21.35	9.21	30.60	23.59	37.12	20.18	<b>44.55</b>	<b>41.19</b>	<b>41.94</b>
musk	7.59	<b>100.00</b>	<b>99.89</b>	10.61	2.82	9.65	47.60	99.61	32.76	50.13	34.79	3.16	<b>100.00</b>	2.87	37.55	8.78	70.70	97.65	99.96
pendigits	4.46	29.27	23.65	23.52	3.78	6.50	18.71	26.05	4.67	30.65	21.22	2.94	<b>32.87</b>	22.21	<b>32.67</b>	1.53	23.75	24.86	21.68
Phoneme	48.24	<b>36.61</b>	54.03	50.00	47.18	55.14	44.09	<b>55.82</b>	41.55	50.45	55.19	36.65	52.99	50.24	53.50	35.02	46.34	54.66	54.23
satellite	47.23	67.25	59.64	57.61	37.68	50.01	61.94	65.92	58.33	32.22	56.58	37.56	61.43	43.31	54.68	41.77	68.92	<b>71.68</b>	<b>75.13</b>
satimage-2	26.11	78.04	85.69	82.71	4.30	39.14	80.52	93.45	22.07	64.49	76.55	12.08	97.09	8.12	97.13	2.58	72.90	<b>97.33</b>	<b>97.31</b>
shuttle	20.27	<b>96.40</b>	92.35	85.29	13.76	20.38	48.75	<b>97.62</b>	93.20	90.45	<b>96.56</b>	7.68	79.89	81.82	32.66	12.41	67.23	92.05	<b>92.36</b>
skin	24.61	23.70	17.40	19.03	18.25	<b>28.72</b>	18.44	26.08	N/A	18.37	17.99	20.89	<b>28.34</b>	26.29	25.58	19.06	25.36	<b>28.87</b>	<b>28.72</b>
Stamps	20.28	35.24	41.09	31.39	21.29	23.53	34.60	39.49	43.73	33.21	43.10	11.03	24.46	<b>47.36</b>	35.63	12.07	34.68	42.39	<b>50.94</b>
thyroid	23.56	50.98	44.34	21.23	20.81	34.98	14.68	<b>63.11</b>	16.06	51.06	19.64	9.44	29.88	10.56	31.69	2.70	50.36	<b>60.99</b>	60.06
vertebral	11.79	9.23	10.49	10.94	14.24	10.57	9.68	10.46	<b>15.24</b>	11.84	8.89	13.11	11.43	11.58	10.54	10.62	<b>14.31</b>	9.78	12.96
vowels	<b>38.88</b>	13.41	8.92	8.24	34.42	<b>63.41</b>	13.82	15.12	12.22	10.56	4.14	13.27	35.14	3.58	<b>49.10</b>	4.58	14.97	26.52	32.42
Waveform	9.66	5.86	5.79	4.37	11.33	13.04	4.71	6.24	3.11	4.76	6.90	5.33	<b>17.93</b>	4.26	<b>19.74</b>	4.41	11.28	6.49	7.83
WBC	54.00	73.56	82.29	<b>89.87</b>	5.57	66.55	78.67	<b>90.49</b>	N/A	86.19	<b>86.19</b>	30.25	67.31	33.43	71.88	8.99	13.32	68.69	83.14
Wilt	<b>5.53</b>	3.84	3.13	3.62	5.05	4.73	3.36	4.23	4.00	3.93	3.69	<b>5.33</b>	3.74	4.62	3.76	4.65	4.05	4.80	<b>5.19</b>
wine	7.95	43.08	30.87	21.56	7.77	8.43	<b>48.82</b>	25.96	17.51	23.54	45.71	8.11	5.98	24.44	6.27	18.78	8.38	21.40	<b>49.59</b>
WPBC	<b>25.62</b>	23.04	23.01	22.93	20.29	21.49	<b>25.39</b>	22.42	22.49	21.24	22.81	23.86	21.08	22.86	20.58	<b>25.00</b>	20.73	22.71	24.90
Avg. Rank	<b>10.83</b>	<b>8.19</b>	<b>8.31</b>	<b>11.14</b>	<b>13.24</b>	<b>9.36</b>	<b>11.79</b>	<b>7.29</b>	<b>11.96</b>	<b>9.36</b>	<b>9.53</b>	<b>14.91</b>	<b>8.53</b>	<b>11.97</b>	<b>9.03</b>	<b>13.41</b>	<b>9.10</b>	<b>6.31</b>	<b>4.74</b>

624 Constructing the outlier set requires a sorting operation, for which we use Numpy’s built-in quantile  
625 calculation with a time complexity of  $\mathcal{O}(N \log N)$ . Considering the number of network parameters  
626 along with the computation of the loss function, the computational complexity for optimizing  $\Theta$  is  
627 approximately  $\mathcal{O}(TNDd + TNKd)$ . The EM algorithm for the Student’s t mixture model includes  
628 two main steps: the E-step, where the complexity for computing the probability (or responsibility)  
629 of each data point belonging to each component is approximately  $\mathcal{O}(NKd)$ , and the M-step, where  
630 the full computational complexity of updating the parameters (mean, covariance matrix) of each  
631 component is  $\mathcal{O}(NKd^2)$ . In practice, we use diagonal covariance matrices, which reduces the  
632 update complexity to roughly  $\mathcal{O}(NKd)$ . If the EM algorithm requires  $T$  round to converge, its  
633 time complexity is approximately  $\mathcal{O}(TNKd)$ . Therefore, the time complexity for  $t$ -iterations is  
634  $\mathcal{O}(tN(\log N + Td(D + K)))$ .

## 635 E Additional Experiments on Graph

### 636 E.1 Baselines

637 Our proposed method was compared with 16 graph domain baseline methods grouped into three  
638 categories as follows:

- 639 • **Contrastive Learning-based Methods:** This group includes CoLA [30], SLGAD [55],  
640 CONAD [53], and ANEMONE [20]. These methods primarily assume that the contrastive loss  
641 between anomalous nodes and their neighborhoods is more significant.
- 642 • **Autoencoder-based Methods:** This category consists of MLPAE [43], GCNAE [22], DOMI-  
643 NANT [11], GUIDE [54], ComGA [31], AnomalyDAE [14], ALARM [37], DONE/AdONE [4]  
644 and AAGNN [57]. These methods focus on the reconstruction errors of anomalous nodes during  
645 the process of reconstructing the graph structure or features.
- 646 • **Clustering-based Methods:** This category of methods encompasses SCAN [52], CBLOF [18],  
647 and DCFOF [45]. These methods generally identify anomalies by detecting if a sample deviates  
648 from the clustering.

### 649 E.2 Datasets

650 We assess the performance of our model using four graph benchmark datasets containing organic  
651 anomalies. Table 7 presents the statistical summary for each dataset. These datasets contain naturally  
652 occurring real-world anomalies and are valuable for assessing the performance of anomaly detection  
653 algorithms in real-world scenarios. The sources and compositions of these datasets are as follows:

Table 7: Statistics of graph benchmark datasets.

Dataset	# Nodes	# Edges	# Features	# Anomaly	Category
Disney	124	670	28	6	co-purchase network
Weibo	8,405	407,963	400	868	social media network
Reddit	10,984	168,016	64	366	user-subreddit network
T-Finance	39,357	42,445,086	10	1,803	trading network

- 654 • **Weibo**[19] is a labeled graph comprising user posts extracted from the social media platform  
655 Tencent Weibo. The user-user graph establishes connections between users who exhibit similar  
656 topic labels. A user is considered anomalous if they have engaged in a minimum of five suspicious  
657 events, whereas normal nodes represent users who have not.
- 658 • **Reddit**[25] consists of a user-subreddit graph extracted from the popular social media platform  
659 Reddit. This publicly accessible dataset encompasses user posts within various subreddits over  
660 a month. Each user is assigned a binary label indicating whether they have been banned on the  
661 platform. Our assumption is that banned users exhibit anomalous behavior compared to regular  
662 Reddit users.
- 663 • **Disney**[34] is a co-purchase network of movies that includes attributes such as price, rating, and the  
664 number of reviews. The ground truth labels, indicating whether a movie is considered anomalous  
665 or not, were assigned by high school students through majority voting.
- 666 • **T-Finance**[46] aims to identify anomalous accounts within a trading network. The nodes in  
667 this network represent unique anonymous accounts, each characterized by ten features related to  
668 registration duration, recorded activity, and interaction frequency. Graph edges denote transaction  
669 records between accounts. If a node is associated with activities such as fraud, money laundering,  
670 or online gambling, human experts will designate it as an anomaly.

### 671 E.3 Experiment Settings

Table 8: AUC-ROC and AUC-PR of 16 unsupervised algorithms on 4 graph benchmark datasets.

Group	Method	Weibo		Reddit		Disney		T-Finance	
		AUC-ROC	AUC-PR	AUC-ROC	AUC-PR	AUC-ROC	AUC-PR	AUC-ROC	AUC-PR
CL-Based	CoLA	0.382	0.087	0.527	0.036	0.455	0.060	0.243	0.031
	SL-GAD	0.421	0.109	<b>0.594</b>	0.040	0.494	0.061	0.442	0.041
	ANEMONE	0.320	0.082	0.536	0.036	0.454	0.068	0.226	0.030
	CONAD	0.806	0.432	0.551	0.037	0.600	0.138	N/A	N/A
AE-Based	MLPAE	0.880	0.629	0.501	0.035	0.563	0.064	0.299	0.030
	GCNAE	0.847	0.567	0.526	0.033	0.517	0.059	0.295	0.030
	GUIDE	0.897	0.692	0.566	0.040	0.521	0.060	N/A	N/A
	DOMINANT	0.927	0.797	0.561	0.037	0.590	0.077	N/A	N/A
	ComGA	0.925	0.809	0.568	0.037	0.494	0.058	N/A	N/A
	AnomalyDAE	0.892	0.694	0.560	0.037	0.520	0.070	N/A	N/A
	ALARM	0.952	0.843	0.559	0.037	0.595	0.123	N/A	N/A
	DONE	0.856	0.579	0.551	0.037	0.517	0.061	0.550	0.046
AAGNN	0.804	0.530	0.564	<b>0.045</b>	0.479	0.059	N/A	N/A	
Cluster-Based	SCAN	0.701	0.186	0.496	0.033	0.548	0.053	N/A	N/A
	CBLOF*	0.972	0.875	0.503	0.035	0.574	<b>0.146</b>	0.524	0.046
	DCFOD*	0.684	0.196	0.552	0.038	0.675	0.119	0.521	0.066
	UniCAD *	<b>0.985</b>	<b>0.927</b>	0.560	0.040	<b>0.701</b>	0.130	<b>0.876</b>	<b>0.422</b>

672 In this experiment, we compared graph-based methods on relational data. For methods originally  
673 designed around feature vectors, including CBLOF, DCFOD, and our approach, we uniformly  
674 employed the same graph representation learning technique as described in BGRL [47]. Specifically,  
675 we used a two-layer Graph Convolutional Network (GCN) for encoding, which produced output  
676 embeddings with a dimensionality of 128. The training epochs were set to 3000, including a warm-up  
677 period of 300 epochs. The hidden size of the predictor was set to 512, and the momentum was fixed  
678 at 0.99.

#### 679 **E.4 Performance Analysis**

680 The performance of UniCAD compared to 16 baseline methods on the four datasets are summarized  
681 in Table 8. From the results, we have the following observations: Our model consistently outperforms  
682 the baseline methods on most datasets, underlining its effectiveness in anomaly detection even within  
683 graph data contexts. This highlights the superiority of UniCAD in detecting anomalies in real-world  
684 graph data.

685 When comparing UniCAD with the four contrastive learning-based methods, it exhibits a distinct  
686 advantage, outperforming them by a substantial margin across all metrics. Unlike contrastive learning  
687 methods that rely on the local neighborhood for anomaly detection, UniCAD leverages the global  
688 clustering distribution. This key difference contributes to its consistently superior performance.  
689 Although CONAD incorporates human prior knowledge about anomalies, enabling it to outperform  
690 other similar methods on the Weibo and Disney datasets, it still falls short compared to our proposed  
691 UniCAD.

692 Compared to the autoencoder-based methods, UniCAD offers the advantage of lower memory  
693 requirements along with better performance. Graph autoencoders typically reconstruct the entire  
694 adjacency matrix during full graph training, resulting in memory usage of at least  $\mathcal{O}(N^2)$ . In contrast,  
695 UniCAD, as a clustering-based method, only requires  $\mathcal{O}(N \times K)$ . Among the autoencoder-based  
696 methods, GCNAE, DONE, and AdONE can be extended to the T-Finance dataset as they only  
697 reconstruct the sampled subgraphs rather than the entire adjacency matrix. However, UniCAD still  
698 showcases superior performance while being more memory-efficient.

699 UniCAD also demonstrates superior performance compared to various other clustering-based methods,  
700 including traditional structural clustering (SCAN) methods that treat the embedding from BGRL as  
701 tabular data (CBLOF, DCFOD).

## 702 **NeurIPS Paper Checklist**

703 The checklist is designed to encourage best practices for responsible machine learning research,  
704 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
705 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
706 follow the references and precede the (optional) supplemental material. The checklist does NOT  
707 count towards the page limit.

708 Please read the checklist guidelines carefully for information on how to answer these questions. For  
709 each question in the checklist:

- 710 • You should answer [Yes] , [No] , or [NA] .
- 711 • [NA] means either that the question is Not Applicable for that particular paper or the relevant  
712 information is Not Available.
- 713 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

714 **The checklist answers are an integral part of your paper submission.** They are visible to the  
715 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it  
716 (after eventual revisions) with the final version of your paper, and its final version will be published  
717 with the paper.

718 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
719 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a  
720 proper justification is given (e.g., "error bars are not reported because it would be too computationally  
721 expensive" or "we were unable to find the license for the dataset we used"). In general, answering  
722 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we  
723 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
724 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
725 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
726 please point to the section(s) where related material for the question can be found.

727 IMPORTANT, please:

- 728 • **Delete this instruction block, but keep the section heading “NeurIPS paper checklist”,**
- 729 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 730 • **Do not modify the questions and only use the provided macros for your answers.**

### 731 1. **Claims**

732 Question: Do the main claims made in the abstract and introduction accurately reflect the  
733 paper’s contributions and scope?

734 Answer: [Yes]

735 Justification: The main claims presented in the abstract and introduction are consistent with  
736 the paper’s contributions and accurately outline the scope.

737 Guidelines:

- 738 • The answer NA means that the abstract and introduction do not include the claims made  
739 in the paper.
- 740 • The abstract and/or introduction should clearly state the claims made, including the  
741 contributions made in the paper and important assumptions and limitations. A No or NA  
742 answer to this question will not be perceived well by the reviewers.
- 743 • The claims made should match theoretical and experimental results, and reflect how much  
744 the results can be expected to generalize to other settings.
- 745 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
746 are not attained by the paper.

### 747 2. **Limitations**

748 Question: Does the paper discuss the limitations of the work performed by the authors?

749 Answer: [Yes]

750 Justification: The limitations of the method’s application scope are discussed in Section 5 of  
751 the paper, along with considerations for future work.

752 Guidelines:

- 753 • The answer NA means that the paper has no limitation while the answer No means that  
754 the paper has limitations, but those are not discussed in the paper.
- 755 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 756 • The paper should point out any strong assumptions and how robust the results are to  
757 violations of these assumptions (e.g., independence assumptions, noiseless settings, model  
758 well-specification, asymptotic approximations only holding locally). The authors should  
759 reflect on how these assumptions might be violated in practice and what the implications  
760 would be.
- 761 • The authors should reflect on the scope of the claims made, e.g., if the approach was only  
762 tested on a few datasets or with a few runs. In general, empirical results often depend on  
763 implicit assumptions, which should be articulated.
- 764 • The authors should reflect on the factors that influence the performance of the approach.  
765 For example, a facial recognition algorithm may perform poorly when image resolution is  
766 low or images are taken in low lighting. Or a speech-to-text system might not be used  
767 reliably to provide closed captions for online lectures because it fails to handle technical  
768 jargon.
- 769 • The authors should discuss the computational efficiency of the proposed algorithms and  
770 how they scale with dataset size.
- 771 • If applicable, the authors should discuss possible limitations of their approach to address  
772 problems of privacy and fairness.
- 773 • While the authors might fear that complete honesty about limitations might be used by  
774 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
775 limitations that aren’t acknowledged in the paper. The authors should use their best  
776 judgment and recognize that individual actions in favor of transparency play an important  
777 role in developing norms that preserve the integrity of the community. Reviewers will be  
778 specifically instructed to not penalize honesty concerning limitations.

### 779 3. Theory Assumptions and Proofs

780 Question: For each theoretical result, does the paper provide the full set of assumptions and  
781 a complete (and correct) proof?

782 Answer: [Yes]

783 Justification: The full set of assumptions and complete proofs for each theoretical result  
784 are provided and can be found in the appendix, specifically in Section B, ensuring that the  
785 theoretical framework is transparent and verifiable.

786 Guidelines:

- 787 • The answer NA means that the paper does not include theoretical results.
- 788 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
789 referenced.
- 790 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 791 • The proofs can either appear in the main paper or the supplemental material, but if they  
792 appear in the supplemental material, the authors are encouraged to provide a short proof  
793 sketch to provide intuition.
- 794 • Inversely, any informal proof provided in the core of the paper should be complemented  
795 by formal proofs provided in appendix or supplemental material.
- 796 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 797 4. Experimental Result Reproducibility

798 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
799 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
800 of the paper (regardless of whether the code and data are provided or not)?

801 Answer: [Yes]



802 Justification: All necessary information for reproducing the main experimental results,  
803 including dataset links, baseline comparisons, and the methodologies of our proposed  
804 approach, are comprehensively included within the submission files, facilitating transparency  
805 and reproducibility of the research findings.

806 Guidelines:

- 807 • The answer NA means that the paper does not include experiments.
- 808 • If the paper includes experiments, a No answer to this question will not be perceived well  
809 by the reviewers: Making the paper reproducible is important, regardless of whether the  
810 code and data are provided or not.
- 811 • If the contribution is a dataset and/or model, the authors should describe the steps taken to  
812 make their results reproducible or verifiable.
- 813 • Depending on the contribution, reproducibility can be accomplished in various ways.  
814 For example, if the contribution is a novel architecture, describing the architecture fully  
815 might suffice, or if the contribution is a specific model and empirical evaluation, it may be  
816 necessary to either make it possible for others to replicate the model with the same dataset,  
817 or provide access to the model. In general, releasing code and data is often one good  
818 way to accomplish this, but reproducibility can also be provided via detailed instructions  
819 for how to replicate the results, access to a hosted model (e.g., in the case of a large  
820 language model), releasing of a model checkpoint, or other means that are appropriate to  
821 the research performed.
- 822 • While NeurIPS does not require releasing code, the conference does require all submis-  
823 sions to provide some reasonable avenue for reproducibility, which may depend on the  
824 nature of the contribution. For example
  - 825 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to  
826 reproduce that algorithm.
  - 827 (b) If the contribution is primarily a new model architecture, the paper should describe  
828 the architecture clearly and fully.
  - 829 (c) If the contribution is a new model (e.g., a large language model), then there should  
830 either be a way to access this model for reproducing the results or a way to reproduce  
831 the model (e.g., with an open-source dataset or instructions for how to construct the  
832 dataset).
  - 833 (d) We recognize that reproducibility may be tricky in some cases, in which case authors  
834 are welcome to describe the particular way they provide for reproducibility. In the  
835 case of closed-source models, it may be that access to the model is limited in some  
836 way (e.g., to registered users), but it should be possible for other researchers to have  
837 some path to reproducing or verifying the results.

## 838 5. Open access to data and code

839 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
840 tions to faithfully reproduce the main experimental results, as described in supplemental  
841 material?

842 Answer: [Yes]

843 Justification: The paper ensures open access to both the data and code necessary for  
844 reproducing the main experimental results, complemented by detailed instructions in the  
845 supplemental material.

846 Guidelines:

- 847 • The answer NA means that paper does not include experiments requiring code.
- 848 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
849 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 850 • While we encourage the release of code and data, we understand that this might not be  
851 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
852 including code, unless this is central to the contribution (e.g., for a new open-source  
853 benchmark).
- 854 • The instructions should contain the exact command and environment needed to run to  
855 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
856 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.

- 857 • The authors should provide instructions on data access and preparation, including how to  
858 access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 859 • The authors should provide scripts to reproduce all experimental results for the new  
860 proposed method and baselines. If only a subset of experiments are reproducible, they  
861 should state which ones are omitted from the script and why.
- 862 • At submission time, to preserve anonymity, the authors should release anonymized ver-  
863 sions (if applicable).
- 864 • Providing as much information as possible in supplemental material (appended to the  
865 paper) is recommended, but including URLs to data and code is permitted.

## 866 6. Experimental Setting/Details

867 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
868 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
869 results?

870 Answer: [Yes]

871 Justification: Detailed information regarding the training and test setups, including data  
872 splits, hyperparameters and their selection process, the type of optimizer used, and other  
873 relevant details, are thoroughly documented in Section 4.2 of the paper.

874 Guidelines:

- 875 • The answer NA means that the paper does not include experiments.
- 876 • The experimental setting should be presented in the core of the paper to a level of detail  
877 that is necessary to appreciate the results and make sense of them.
- 878 • The full details can be provided either with the code, in appendix, or as supplemental  
879 material.

## 880 7. Experiment Statistical Significance

881 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
882 information about the statistical significance of the experiments?

883 Answer: [Yes]

884 Justification: The paper reports the statistical significance of the experiments by detailing  
885 the results of the Friedman test and the Nemenyi test in Appendix D.3.

886 Guidelines:

- 887 • The answer NA means that the paper does not include experiments.
- 888 • The authors should answer "Yes" if the results are accompanied by error bars, confidence  
889 intervals, or statistical significance tests, at least for the experiments that support the main  
890 claims of the paper.
- 891 • The factors of variability that the error bars are capturing should be clearly stated (for  
892 example, train/test split, initialization, random drawing of some parameter, or overall run  
893 with given experimental conditions).
- 894 • The method for calculating the error bars should be explained (closed form formula, call  
895 to a library function, bootstrap, etc.)
- 896 • The assumptions made should be given (e.g., Normally distributed errors).
- 897 • It should be clear whether the error bar is the standard deviation or the standard error of  
898 the mean.
- 899 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably  
900 report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality  
901 of errors is not verified.
- 902 • For asymmetric distributions, the authors should be careful not to show in tables or figures  
903 symmetric error bars that would yield results that are out of range (e.g. negative error  
904 rates).
- 905 • If error bars are reported in tables or plots, The authors should explain in the text how they  
906 were calculated and reference the corresponding figures or tables in the text.

## 907 8. Experiments Compute Resources

908 Question: For each experiment, does the paper provide sufficient information on the com-  
909 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
910 the experiments?

911 Answer: [Yes]

912 Justification: Detailed information on the compute resources, including the type of compute  
913 workers (CPU/GPU), memory, and execution time for each experiment, is provided in the  
914 supplementary materials, enabling accurate reproduction of the experiments.

915 Guidelines:

- 916 • The answer NA means that the paper does not include experiments.
- 917 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or  
918 cloud provider, including relevant memory and storage.
- 919 • The paper should provide the amount of compute required for each of the individual  
920 experimental runs as well as estimate the total compute.
- 921 • The paper should disclose whether the full research project required more compute than  
922 the experiments reported in the paper (e.g., preliminary or failed experiments that didn't  
923 make it into the paper).

## 924 9. Code Of Ethics

925 Question: Does the research conducted in the paper conform, in every respect, with the  
926 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

927 Answer: [Yes]

928 Justification: The research adheres to the NeurIPS Code of Ethics, including considerations  
929 for anonymity, fairness, and transparency, with no deviations reported or necessary under  
930 current laws or regulations.

931 Guidelines:

- 932 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 933 • If the authors answer No, they should explain the special circumstances that require a  
934 deviation from the Code of Ethics.
- 935 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration  
936 due to laws or regulations in their jurisdiction).

## 937 10. Broader Impacts

938 Question: Does the paper discuss both potential positive societal impacts and negative  
939 societal impacts of the work performed?

940 Answer: [Yes]

941 Justification: The paper thoroughly discusses both the potential positive impacts, such as  
942 enhancements in anomaly detection for critical applications.

943 Guidelines:

- 944 • The answer NA means that there is no societal impact of the work performed.
- 945 • If the authors answer NA or No, they should explain why their work has no societal impact  
946 or why the paper does not address societal impact.
- 947 • Examples of negative societal impacts include potential malicious or unintended uses  
948 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g.,  
949 deployment of technologies that could make decisions that unfairly impact specific groups),  
950 privacy considerations, and security considerations.
- 951 • The conference expects that many papers will be foundational research and not tied to  
952 particular applications, let alone deployments. However, if there is a direct path to any  
953 negative applications, the authors should point it out. For example, it is legitimate to point  
954 out that an improvement in the quality of generative models could be used to generate  
955 deepfakes for disinformation. On the other hand, it is not needed to point out that a  
956 generic algorithm for optimizing neural networks could enable people to train models that  
957 generate Deepfakes faster.

- 958
- 959
- 960
- 961
- 962
- 963
- 964
- 965
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

966

967

968

969

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

970

Answer: [NA]

971

Justification: The paper poses no such risks.

972

Guidelines:

- 973
- 974
- 975
- 976
- 977
- 978
- 979
- 980
- 981
- 982
- The answer NA means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

983

984

985

986

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

987

Answer: [Yes]

988

989

990

Justification: The paper appropriately credits the creators of the utilized assets, including code, data, and models, and explicitly mentions the licenses and terms of use, ensuring compliance with the original terms set by the asset owners.

991

Guidelines:

- 992
- 993
- 994
- 995
- 996
- 997
- 998
- 999
- 1000
- 1001
- 1002
- 1003
- 1004
- 1005
- The answer NA means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New Assets

1006

1007

1008

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

1009

Answer: [NA]

1010 Justification: The paper does not release new assets.

1011 Guidelines:

- 1012 • The answer NA means that the paper does not release new assets.
- 1013 • Researchers should communicate the details of the dataset/code/model as part of their sub-  
1014 missions via structured templates. This includes details about training, license, limitations,  
1015 etc.
- 1016 • The paper should discuss whether and how consent was obtained from people whose asset  
1017 is used.
- 1018 • At submission time, remember to anonymize your assets (if applicable). You can either  
1019 create an anonymized URL or include an anonymized zip file.

1020 **14. Crowdsourcing and Research with Human Subjects**

1021 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1022 include the full text of instructions given to participants and screenshots, if applicable, as  
1023 well as details about compensation (if any)?

1024 Answer: [NA]

1025 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1026 Guidelines:

- 1027 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1028 human subjects.
- 1029 • Including this information in the supplemental material is fine, but if the main contribution  
1030 of the paper involves human subjects, then as much detail as possible should be included  
1031 in the main paper.
- 1032 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or  
1033 other labor should be paid at least the minimum wage in the country of the data collector.

1034 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
1035 Subjects**

1036 Question: Does the paper describe potential risks incurred by study participants, whether  
1037 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1038 approvals (or an equivalent approval/review based on the requirements of your country or  
1039 institution) were obtained?

1040 Answer: [NA]

1041 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1042 Guidelines:

- 1043 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1044 human subjects.
- 1045 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1046 may be required for any human subjects research. If you obtained IRB approval, you  
1047 should clearly state this in the paper.
- 1048 • We recognize that the procedures for this may vary significantly between institutions  
1049 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1050 guidelines for their institution.
- 1051 • For initial submissions, do not include any information that would break anonymity (if  
1052 applicable), such as the institution conducting the review.