

Anomaly Detection Requires Better Representations

Tal Reiss, Niv Cohen, Eliahu Horwitz, Ron Abutbul, and Yedid Hoshen

School of Computer Science and Engineering
The Hebrew University of Jerusalem, Israel
http://www.vision.huji.ac.il/ssrl_ad/

Abstract. Anomaly detection seeks to identify unusual phenomena, a central task in science and industry. The task is inherently unsupervised as anomalies are unexpected and unknown during training. Recent advances in self-supervised representation learning have directly driven improvements in anomaly detection. In this position paper, we first explain how self-supervised representations can be easily used to achieve state-of-the-art performance in commonly reported anomaly detection benchmarks. We then argue that tackling the next generation of anomaly detection tasks requires new technical and conceptual improvements in representation learning.

Keywords: Anomaly Detection, Self-Supervised Learning, Representation Learning

1 Introduction

Discovery commences with the awareness of anomaly, i.e., with the recognition that nature has somehow violated the paradigm-induced expectations that govern normal science.

—Kuhn, *The Structure of Scientific Revolutions* (1970)

I do not know what I may appear to the world, but to myself I seem to have been only like a boy playing on the seashore, and diverting myself in now and then finding a smoother pebble or a prettier shell than ordinary, whilst the great ocean of truth lay all undiscovered before me.

—Isaac Newton

Anomaly detection, discovering unusual patterns in data, is a core task for human and machine intelligence. The importance of the task stems from the centrality of discovering unique or unusual phenomena in science and industry. For example, the fields of particle physics and cosmology have, to large extent, been driven by the discovery of new fundamental particles and stellar objects. Similarly, the discovery of new, unknown, biological organisms and systems is

a driving force behind biology. The task is also of significant economic potential. Anomaly detection methods are used to detect credit card fraud, faults on production lines, and unusual patterns in network communications.

Detecting anomalies is essentially unsupervised as only "normal" data, but no anomalies, are seen during training. While the field has been intensely researched for decades, the most successful recent approaches use a very simple two-stage paradigm: (i) each data point is transformed to a representation, often learned in a self-supervised manner. (ii) a density estimation model, often as simple as a k nearest neighbor estimator, is fitted to the normal data provided in a training set. To classify a new sample as normal or anomalous, its estimated probability density is computed - low likelihood samples are denoted as anomalies.

In this position paper, we first explain that advances in representation learning are the main explanatory factor for the performance of recent anomaly detection (AD) algorithms. We show that this paradigm essentially "solves" the most commonly reported image anomaly detection benchmark (Sec. 4). While this is encouraging, we argue that existing self-supervised representations are unable to solve the next generation of AD tasks (Sec. 5). In particular, we highlight the following issues: (i) masked-autoencoders are much worse for AD than earlier self-supervised representation learning (SSRL) methods (ii) current approaches perform poorly in datasets with multiple objects per-image, complex background, fine-grained anomalies. (iii) in some cases SSRL performs worse than handcrafted representations (iv) for "tabular" datasets, no representation performed better than the original representation of the data (i.e. that data itself) (v) in the presence of nuisance factors of variation, it is unclear whether SSRL can *in-principle* identify the optimal representation for effective AD.

Anomaly detection presents both rich rewards as well as significant challenges for representation learning. Overcoming these issues will require significant progress, both technical and conceptual. We expect that increasing the involvement of the self-supervised representation learning community in anomaly detection will mutually benefit both fields.

2 Related Work

Classical AD approaches were typically based on either density estimation [9,20] or reconstruction [15]. With the advent of deep learning, classical methods were augmented by deep representations [23,38,19,24]. A prevalent way to learn these representations was to use self-supervised methods, e.g. autoencoder [30], rotation classification [10,13], and contrastive methods [36,35]. An alternative approach is to combine pretrained representations with anomaly scoring functions [25,32,27,28]. The best performing methods [27,28] combine pretraining on auxiliary datasets and a second finetuning stage on the provided normal samples in the training set. It was recently established [27] that given sufficiently powerful representations (e.g. ImageNet classification), a simple criterion based on the k NN distance to the normal training data achieves strong performance. We therefore limit the discussion of AD in this paper to this simple technique.

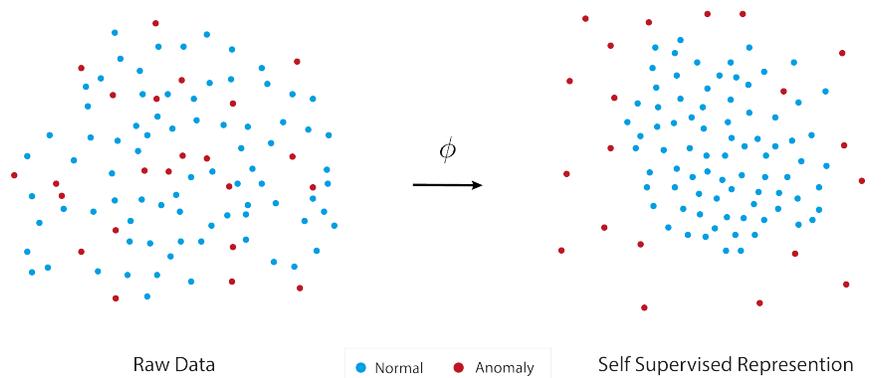


Fig. 1. Normal and Anomalous Representations: The self-supervised representations transform the raw data into a space in which normal and anomalous data can be easily separated using density estimation methods

3 Anomaly Detection as a Downstream Task for Representation Learning

In this section we describe the computational task, method, and evaluation setting for anomaly detection.

Task definition. We assume access to N random samples, denoted by $\mathcal{X}_{train} = \{x_1, x_2 \dots x_N\}$, from the distribution of the normal data $p_{norm}(x)$. At test time, the algorithm observes a sample \tilde{x} from the real-world distribution $p_{real}(x)$, which consists of a combination of the normal and anomalous data distributions: $p_{norm}(x)$ and $p_{anom}(x)$. The task is to classify the sample \tilde{x} as normal or anomalous.

Representations for anomaly detection. In AD, it is typically assumed that anomalies $a \sim p_{anom}$ have a low likelihood under the normal data distribution, i.e. that $p_{norm}(a)$ is small. Under this assumption, the PDF of normal data p_{norm} acts as an effective anomaly classifier. In practice, however, training an estimator q for scoring anomalies using p_{norm} is a challenging statistical task. The challenge is greater when: (i) the data are high-dimensional (e.g. images) (ii) p_{norm} is sparse or irregular (iii) normal and anomalous data are not separable using simple functions. Representation learning may overcome these issues by transforming the sample x into a representation $\phi(x)$, which is of lower dimension, where p_{norm} is relatively smooth and where normal and anomalous data are more separable. As no anomaly labels are provided, self-supervised representation learning is needed.

A two-stage anomaly detection paradigm. Given a self-supervised representation ϕ , we follow a simple two stage anomaly detection paradigm: (i) *Representation encoder*: each sample during training or test is mapped to a feature

descriptor using the mapping function ϕ . (ii) *Density estimation*: a probability estimator $q_{norm}(x)$ is fitted to the distribution of the normal sample features $\mathcal{X}_{train} = \{\phi(x_1), \phi(x_2) \dots \phi(x_N)\}$. A sample is scored at test time by first mapping it to the representation space $\phi(\tilde{x})$ and scoring it according to the density estimator. Given an estimator q_{norm} of the normal probability density p_{norm} , the anomaly score s is given by $s(\tilde{x}) = -q_{norm}(\phi(\tilde{x}))$. Normal data will typically obtain lower scores than anomalous samples. A user can then set a threshold for the prediction of anomalies based on an appropriate false positive rate.

4 Successful Representation Learning Enables Anomaly Detection

Detecting anomalies in images is probably the most researched task by the deep learning anomaly detection community. In this section, we show that the simple paradigm presented in Sec. 3 achieves state-of-the-art results. As most density estimators achieve very similar results, the anomaly detection performance is mostly determined by the quality of learned representation. This makes anomaly detection an excellent testing ground for representations. Furthermore, we discuss different approaches to finetune a representation on the normal train data and show significant gains.

Learning representations from the normal data. Perhaps the most common approach taken by recent AD methods is to learn the representations in a self-supervised manner using solely the normal samples (i.e. the training dataset). Examples of such methods are RotNet [13], CSI [36] and others. The main disadvantage of such methods is that most of the datasets are of small size and hence do not suffice for learning powerful representations.

Extracting representations from a pretrained model. A very simple alternative is to use an off-the-shelf pretrained model and extract features for the normal (i.e. training) data from it. The pretraining may be either supervised (e.g. using ImageNet labels [8,12]) or self-supervised (e.g. DINO), in both cases pretraining may be performed on ImageNet. These representations tend to perform much better than those extracted from models trained only on the normal data.

A hybrid approach. A natural extension to the above approaches is to combine the two. This is done by using the pretrained model as an initialization for a self-supervised finetuning phase (on the normal data). In this way, the powerful representation of the pretrained model can be used and refined within the context of the anomaly detection dataset and task. Multiple approaches [27,28] have been used for the self-supervised finetuning stage. However, in this paper we present what is possibly the simplest approach, using DINO’s objective for the finetuning stage. In this approach, a pretrained DINO model is used as an initialization. During the finetuning phase, the model is trained on the target anomaly detection training dataset (i.e. only normal data) in a self-supervised manner by simply using the original DINO objective.



Fig. 2. MAE vs. DINO nearest neighbors: For each image, the top 5 nearest neighbors are shown according to their order. Note how MAE neighbors are chosen mostly based on the colors and not their semantic contents, in contrast, DINO neighbors are semantically accurate.

In Fig. 1 the above process is demonstrated with a toy example. Tab. 1 presents anomaly detection results on the CIFAR-10 [18] dataset, which is the most commonly used dataset for evaluation. As can be seen, using representations extracted from a recent self-supervised method (i.e. DINO) following the hybrid approach and coupled with a trivial k NN estimator for the density estimation phase nearly solves this dataset. Although a possible conclusion could have been that the anomaly detection task has been solved, in the next section we show this is not the case.

5 Gaps in Anomaly Detection Point to Bottlenecks in Representations Learning

While Sec. 4 presented a very optimistic view of the ability of representation learning to solve anomaly detection, in this section we paint a more complex picture. We use this to highlight several limitations of current self-supervised representations.

Table 1. Image anomaly detection results: Mean ROC-AUC %. Bold denotes the best results, FT stands for finetuned

Approach	Self-supervised		Pretrained		Hybrid		
Method	RotNet [13]	CSI [36]	ResNet	DINO	PANDA [27]	MSAD [28]	DINO-FT
CIFAR-10	90.1	94.3	92.5	97.1	96.2	97.2	98.4

5.1 Masked-Autoencoder: Advances in self-supervised learning do not always imply better anomaly detection

Recently, masked-autoencoder (MAE) based methods achieved significant improvements on several self-supervised representation learning benchmarks [11]. Yet, the representations learnt by MAE underperform contrastive self-supervised methods on unsupervised tasks such as anomaly detection. A comparison between MAE to contrastive self-supervised method (DINO) is presented in Tab. 2 demonstrating the much better performance of DINO for AD. Finetuning on the normal training data improves both methods, however a large gap still remains. Implementation details for the experiments can be found in the App. A. In many papers, self-supervised methods are evaluated using supervised benchmarks, such as classification accuracy with finetuning. The key difference between anomaly detection and ordinary benchmarks where MAE excel is that anomaly detection is an unsupervised task. This is also suggested by MAE’s worse performance with linear probing (as reported by the original paper), where the supervised labels cannot be used to improve the backbone representations.

MAE’s optimization objective may explain why its strong representation does not translate into better anomaly detection capabilities. As MAE’s objective is to reconstruct patches, it may learn a representation that encodes local information needed for reconstructing the image, overlooking semantic object properties. Consequently, the nearest neighbors may pay more attention to local similarity than to global semantic properties (See Fig. 2). In contrast, the goal of contrastive-based objectives is to map semantically similar images to nearby representations, disregarding some of the local properties.

Conclusion. Better performance on supervised downstream tasks does not necessarily imply better representations across the board. In some cases, while the representation may excel in a supervised downstream task, it may underperform in an unsupervised counterpart. Looking forward, we suggest that new self-supervised representation learning methods present evaluations on unsupervised anomaly detection tasks alongside the common supervised benchmarks.

5.2 Complex datasets: Current representations struggle on scenes, finegrained classes, multiple objects

Current representations are very effective for anomaly detection on datasets with a single object occupying a large portion of the image. Furthermore, these methods typically perform well when the number of object categories in the

Table 2. Anomaly detection comparison of MAE and DINO: Mean ROC-AUC %. Bold denotes the best results

Method	CIFAR-10	CUB-200	INet-S
MAE	78.1	73.1	83.2
DINO	97.1	93.9	99.3

normal train set is small and have coarse differences (e.g. "cat" and "ship"). A prime example is CIFAR-10, which is virtually solved. On the other hand, anomaly detection accuracy is much lower on more complex datasets containing multiple small objects, complex backgrounds; and when anomalies consist of related object categories (e.g. "sofa" and "armchair"). We modified the MS-COCO [21] dataset by using all images from a single super-category ('vehicles') as normal data, apart from a single category ('bicycle') which are used as anomalies. We experiment both with cropping just the object bounding boxes or using the entire image (including the background and other objects). Similarly, we report results for a multi-modal CUB-200 [37] anomaly detection benchmark. The results are presented in Tab. 3 (implementation details can be found in the Appendix). It is clear that these datasets are far from solved and that better representations are needed to achieve acceptable performance.

Conclusion. While current representations are effective for relatively easy datasets, more realistic cases with small objects, backgrounds and many object categories call for the development of new SSRL methods.

Table 3. Multi-modal datasets: Mean ROC-AUC %. "MS-COCO-I" / "MS-COCO-O" indicates MS-COCO image / object level benchmarks (respectively).

Method	MS-COCO-I	MS-COCO-O	CUB-200
PANDA [27]	61.5	77.0	78.4
MSAD [28]	61.7	76.9	80.1
DINO	61.5	73.4	74.5

5.3 Unidentifiability: Representations for anomaly detection may be ambiguous without further guidance

In some settings, we would like our representation to focus on specific attributes (which we denote as *relevant*) while ignoring nuisance attributes that might bias the model. Consider two different companies interested in anomaly detection in cars. The first company may be interested in detecting novel car models, while the second is interested in unusual driving behaviors. Although both may wish to apply density estimation using a state-of-the-art self-supervised representation,

Table 4. Summary of the findings of from Horwitz and Hoshen [14]: Average metrics across all MVTec3D-AD classes, "INet" indicates ImageNet [8] pretrained features, PC indicates point cloud. I-ROC indicates image level ROC-AUC % [4], P-ROC indicates pixel level ROC-AUC %. Higher score indicates better the results

Modality	RGB	Depth	Depth	Depth	Depth	Depth	PC	RGB+PC
Method	INet	INet	NSA [33]	Raw	HoG [7]	SIFT [22]	FPFH [31]	RGB+FPFH
PRO [2]	87.6	58.6	57.2	19.1	61.4	86.6	92.4	96.4
I-ROC	78.5	63.7	69.6	52.8	56.0	71.4	75.3	86.5
P-ROC	96.6	82.1	81.7	54.8	84.5	95.4	98.0	99.3

each will view the ground truth anomalies of the other company as a false-positive case. As each company is interested in different anomalies, they may require different representations. One company would require the representation to contain only the driving patterns and be agnostic to the car model, at the same time, the other company would strive for the opposite. As these preferences are not present at the time of pretraining the self-supervised backbone, the correct solution is often unidentifiable.

One initial effort is RedPANDA [6] that proposed providing labels for nuisance attributes. We note that only the attributes to be ignored are labeled, while the other attributes (the ones in which characterize anomalies) are not provided. Representation learning is then performed using domain-supervised disentanglement [16], resulting in a representation only describing the unlabelled attributes. Yet, the field of domain-supervised disentanglement is still in its infancy, and the assumption of nuisance attribute labels is often not applicable.

Conclusion. Self-supervised representation learning methods are designed to focus on semantic attributes of images, but choosing the most relevant ones is unidentifiable without further guidance. Incorporating guidance may be achieved by a careful choice of inductive bias [16] (e.g. augmentations) or using concept-based representation techniques [17].

5.4 3D Point Clouds: Self-supervised representations do not always improve over handcrafted ones

In an empirical investigation [14], we evaluated representative methods designed for different modalities on the MVTec3D-AD dataset [3]. The paper showed that currently, handcrafted features for 3D surface matching outperform learning-based methods designed either for images or for 3D point clouds. A key insight was that rotation invariance is very beneficial in this modality, and is often overlooked. A summary of the findings, taken from the original paper, is found in Tab. 4.

Conclusion. When dealing with 3D point-cloud, self-supervised representations are yet to outperform handcrafted features for anomaly detection. For

Table 5. Tabular results: Mean F1 & ROC-AUC % from the ODDS benchmarks results. Bold denotes the best results

Method	GOAD [1]		ICL [34]		Raw
Scoring	Auxiliary	k NN	Auxiliary	k NN	k NN
F1	54.4	63.2	68.1	69.8	69.9
ROC-AUC	78.2	87.6	88.9	89.4	90.2

modalities less mature than images, domain specific priors may still need to be integrated into the architecture or objective. This stresses the need for better 3D point-cloud representations.

5.5 Tabular Data: When representations do not improve over the original data

The tabular setting is probably the most general anomaly detection setting, where each sample in the dataset consists of a set of numerical and categorical variables. This is strictly harder than any other setting as no regularity in the data can be assumed. Such data are frequently encountered, as unstructured databases are very common. In recent years, self-supervised methods have been proposed for tabular anomaly detection [39,1,34,26]. These methods differ by the auxiliary task that they use for representation learning (and potentially also for anomaly scoring). Two representative deep learning approaches are GOAD [1] which predicts geometric transformations, and use the prediction errors to detect anomalies, and ICL [34] which adopts the contrastive learning task for training and for anomaly scoring by differentiating between in-window and out-window features. As part of our evaluation, we used both their standard pipeline (i.e. their auxiliary tasks for anomaly scoring) and our AD density estimation paradigm (see Appendix). These results were then compared with k NN on the original raw features without any modifications. The results are presented in Tab. 5. Self-supervised representation learning did not improve performance in comparison with the original raw features.

Conclusion. Representation learning for general datasets is an open research question. Some prior knowledge of the dataset must be used in order to learn non-trivial data representations, at least in the context of anomaly detection.

6 Final Remarks

In this position paper, we advocated the study of self-supervised representations for the task of anomaly detection. We explained that advances in representation learning have been the main driving force behind progress in anomaly detection. On the other hand, we demonstrated that current self-supervised representation learning methods often fall short in challenging anomaly detection settings. Our hope is that interplay between the self-supervised representation learning and anomaly detection fields will result in mutual benefits for both communities.

7 Acknowledgements

This work was partially supported by the Malvina and Solomon Pollack Scholarship, a Facebook award, the Israeli Cyber Directorate, the Israeli Higher Council and the Israeli Science Foundation. We also acknowledge support of Oracle Cloud credits and related resources provided by the Oracle for Research program.

A Appendix

In this paper we report anomaly detection results using the standard uni-modal protocol, which is widely used in the anomaly detection community. In the uni-modal protocol, multi-class datasets are converted to anomaly detection by setting a class as normal and all other classes as anomalies. The process is repeated for all classes, converting a dataset with C classes into C datasets. Finally, we report the mean ROC-AUC % over all C datasets as the anomaly detection results.

A.1 Anomaly detection comparison of MAE and DINO

We compare between DINO [5] and MAE [11] as a representation for a k NN based anomaly detection algorithm. For MAE, we experimented both with k NN and reconstruction error for anomaly scoring and found that the latter works badly, therefore we report just the k NN results. We evaluate using a variety of datasets, in the uni-modal setting described above. We used the following datasets:

INet-S [29]: The dataset is subset of 10 animal classes taken from ImageNet21k (e.g "petrel", "tyrannosaur", "rat snake", "duck", "bee fly", "sheep", "beer cub", "red deer", "silverback", "opossum rat") that do not appear in ImageNet1K dataset. The dataset is coarse-grained and contains images relatively close to ImageNet1K dataset. It intended to convey that even for easy tasks the MAE doesn't achieve as good results as DINO.

CIFAR-10 [18]: Consists of low-resolution 32×32 images from 10 different classes.

CUB-200 [37]: Bird species image dataset which contains 11,788 images of 200 subcategories. In the experiment we calculated mean ROC-AUC % over the 20 first categories.

A.2 Multi-modal datasets

In these experiment we specify a single class as anomalous, and treat all images which does not contain it as normal.

MS-COCO-I [21]: We build a multi-modal anomaly detection dataset comprised of scenes benchmarks, where each image is evaluated against other images featuring similar scenes. We choose 10 object categories ("bicycle", "traffic light", "bird", "backpack", "frisbee", "bottle", "banana", "chair", "tv", "microwave",

”book”) from different MS-COCO super-categories. To construct a multi-modal anomaly detection benchmark, we designate an object category from the list as the anomalous class, and training images of a similar super-category that do not contain it as our normal train set. Our test set contains all the test images from that super-category, where images containing the anomalous object are labelled as anomalies. This process is repeated for the 10 object categories resulting in 10 different evaluations. We report their average ROC-AUC %.

MS-COCO-O: We introduce a similar benchmark to MS-COCO-I, focusing on single objects rather than scenes. We crop all objects from our 10 super-categories (described above) according to the MS-COCO supplied bounding boxes. We repeat a similar process, using a similar object category as normal and the rest as anomalies.

CUB-200 [37]: We create a multi-modal anomaly detection benchmark based on the CUB-200 dataset. We focus on the 20 first categories, designating only one as an anomaly each time.

A.3 Tabular domain

Various datasets used for tabular data anomaly detection were used for the experiments. A total of 31 datasets from Outlier Detection DataSets (ODDS)¹ are employed. For the evaluation of GOAD and ICL we used the official repositories and made an effort to select the best configuration available. For all density estimation evaluations we used k NN with $k = 5$ nearest neighbors. To convert GOAD and ICL into the standard paradigm of representation learning followed by density estimation: i) we use the original approaches to train a feature encoder (followed by a classifier which we discard) ii) we use the feature encoder to represent each sample iii) density estimation is performed on the representations using k NN exactly as in Sec. 3.

¹ <http://odds.cs.stonybrook.edu/>

References

1. Bergman, L., Hoshen, Y.: Classification-based anomaly detection for general data. In: ICLR (2020) [9](#)
2. Bergmann, P., Fauser, M., Sattlegger, D., Steger, C.: Mvtec ad—a comprehensive real-world dataset for unsupervised anomaly detection. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 9592–9600 (2019) [8](#)
3. Bergmann, P., Jin, X., Sattlegger, D., Steger, C.: The mvtec 3d-ad dataset for unsupervised 3d anomaly detection and localization. arXiv preprint arXiv:2112.09045 (2021) [8](#)
4. Bradley, A.P.: The use of the area under the roc curve in the evaluation of machine learning algorithms. *Pattern recognition* **30**(7), 1145–1159 (1997) [8](#)
5. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9650–9660 (2021) [10](#)
6. Cohen, N., Kahana, J., Hoshen, Y.: Red panda: Disambiguating anomaly detection by removing nuisance factors. arXiv preprint arXiv:2207.03478 (2022) [8](#)
7. Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection. In: 2005 IEEE computer society conference on computer vision and pattern recognition (CVPR’05). vol. 1, pp. 886–893. Ieee (2005) [8](#)
8. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009) [4](#), [8](#)
9. Eskin, E., Arnold, A., Prerau, M., Portnoy, L., Stolfo, S.: A geometric framework for unsupervised anomaly detection. In: Applications of data mining in computer security, pp. 77–101. Springer (2002) [2](#)
10. Golan, I., El-Yaniv, R.: Deep anomaly detection using geometric transformations. In: NeurIPS (2018) [2](#)
11. He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16000–16009 (2022) [6](#), [10](#)
12. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016) [4](#)
13. Hendrycks, D., Mazeika, M., Kadavath, S., Song, D.: Using self-supervised learning can improve model robustness and uncertainty. In: NeurIPS (2019) [2](#), [4](#), [6](#)
14. Horwitz, E., Hoshen, Y.: An empirical investigation of 3d anomaly detection and segmentation (2022) [8](#)
15. Jolliffe, I.: Principal component analysis. Springer (2011) [2](#)
16. Kahana, J., Hoshen, Y.: A contrastive objective for learning disentangled representations. arXiv preprint arXiv:2203.11284 (2022) [8](#)
17. Koh, P.W., Nguyen, T., Tang, Y.S., Musmann, S., Pierson, E., Kim, B., Liang, P.: Concept bottleneck models. In: International Conference on Machine Learning. pp. 5338–5348. PMLR (2020) [8](#)
18. Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images (2009) [5](#), [10](#)
19. Larsson, G., Maire, M., Shakhnarovich, G.: Learning representations for automatic colorization. In: ECCV (2016) [2](#)

20. Latecki, L.J., Lazarevic, A., Pokrajac, D.: Outlier detection with kernel density functions. In: International Workshop on Machine Learning and Data Mining in Pattern Recognition. pp. 61–75. Springer (2007) [2](#)
21. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014) [7](#), [10](#)
22. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. International journal of computer vision **60**(2), 91–110 (2004) [8](#)
23. Mathieu, M., Couprie, C., LeCun, Y.: Deep multi-scale video prediction beyond mean square error. ICLR (2016) [2](#)
24. Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: ECCV (2016) [2](#)
25. Perera, P., Patel, V.M.: Learning deep features for one-class classification. IEEE Transactions on Image Processing **28**(11), 5450–5463 (2019) [2](#)
26. Qiu, C., Pfrommer, T., Kloft, M., Mandt, S., Rudolph, M.: Neural transformation learning for deep anomaly detection beyond images. In: International Conference on Machine Learning. pp. 8703–8714. PMLR (2021) [9](#)
27. Reiss, T., Cohen, N., Bergman, L., Hoshen, Y.: Panda: Adapting pretrained features for anomaly detection and segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2806–2814 (2021) [2](#), [4](#), [6](#), [7](#)
28. Reiss, T., Hoshen, Y.: Mean-shifted contrastive loss for anomaly detection. arXiv preprint arXiv:2106.03844 (2021) [2](#), [4](#), [6](#), [7](#)
29. Ridnik, T., Ben-Baruch, E., Noy, A., Zelnik-Manor, L.: Imagenet-21k pretraining for the masses (2021) [10](#)
30. Ruff, L., Gornitz, N., Deecke, L., Siddiqui, S.A., Vandermeulen, R., Binder, A., Müller, E., Kloft, M.: Deep one-class classification. In: ICML (2018) [2](#)
31. Rusu, R.B., Blodow, N., Beetz, M.: Fast point feature histograms (fpfh) for 3d registration. In: 2009 IEEE International Conference on Robotics and Automation. pp. 3212–3217 (2009) [8](#)
32. Salehi, M., Sadjadi, N., Baselizadeh, S., Rohban, M.H., Rabiee, H.R.: Multiresolution knowledge distillation for anomaly detection. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 14902–14912 (2021) [2](#)
33. Schlüter, H.M., Tan, J., Hou, B., Kainz, B.: Self-supervised out-of-distribution detection and localization with natural synthetic anomalies (nsa). arXiv preprint arXiv:2109.15222 (2021) [8](#)
34. Shenkar, T., Wolf, L.: Anomaly detection for tabular data with internal contrastive learning. In: International Conference on Learning Representations (2021) [9](#)
35. Sohn, K., Li, C.L., Yoon, J., Jin, M., Pfister, T.: Learning and evaluating representations for deep one-class classification. arXiv preprint arXiv:2011.02578 (2020) [2](#)
36. Tack, J., Mo, S., Jeong, J., Shin, J.: Csi: Novelty detection via contrastive learning on distributionally shifted instances. NeurIPS (2020) [2](#), [4](#), [6](#)
37. Welinder, P., Branson, S., Mita, T., Wah, C., Schroff, F., Belongie, S., Perona, P.: Caltech-ucsd birds 200 (2010) [7](#), [10](#), [11](#)
38. Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: ECCV (2016) [2](#)
39. Zong, B., Song, Q., Min, M.R., Cheng, W., Lumezanu, C., Cho, D., Chen, H.: Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In: International conference on learning representations (2018) [9](#)