

Fine-grained anomaly detection via multi-task self-supervision

Loïc Jézéquel^{1,2}

Ngoc-Son Vu¹

Jean Beaudet²

Aymeric Histace¹

¹ ETIS UMR 8051 (CY Cergy Paris Université, ENSEA, CNRS) F-95000

² Idemia Identity & Security, 95520 Osny France

{loic.jezequel, son.vu, aymeric.histace}@ensea.fr

Abstract

Detecting anomalies using deep learning has become a major challenge over the last years, and is becoming increasingly promising in several fields. The introduction of self-supervised learning has greatly helped many methods including anomaly detection where simple geometric transformation recognition tasks are used. However these methods do not perform well on fine-grained problems since they lack finer features. By combining both high-scale shape features and low-scale fine features in a multi-task framework, our method greatly improves fine-grained anomaly detection. It outperforms state-of-the-art with up to 31% relative error reduction measured with AUROC on various anomaly detection problems including one-vs-all, out-of-distribution detection and face presentation attack detection.

Detecting anomalies straying apart from a well-defined normal situation has always been a major challenge in many fields such as video surveillance [33, 43], intrusion detection [15], fraud detection [42], medical imaging [20] and more recently adversarial attack detection [25]. Deep visual anomaly detection has been introduced to tackle this problem and has proven to be more robust and reliable than classical binary classification. Rather than directly try to discriminate anomalies from normal samples, we only learn the normal class boundary and deem as anomalous any observation outside.

Recently, the introduction of self-supervised learning has greatly improved many one-class anomaly detection learning methods. It enables to discriminate anomalies from normal samples by learning to solve simple tasks such as geometric transformation classification [10]. However, even if this approach has greatly improved anomaly detection performance, it still suffers from limitations on more challenging problems with local and fine-grained differences between anomalies and normal samples.

In this given context, our main contributions in this paper

are the following:

- We improve the detection of fine-grained anomalies by independently solving in a multi-task self-supervised fashion high-scale geometric task and low-scale jigsaw puzzle task.
- We validate the efficiency of the proposed method using an exhaustive protocol for anomaly detection on one-vs-all, out-of-distribution detection and anti-spoofing problems.
- The proposed method obtains better overall results with up to 31% AUROC relative improvement from state of the art methods.

1. Related work

1.1. Anomaly detection

The main goal in anomaly detection is to classify a sample as normal or anomalous. Formally, we predict $P(\mathbf{x} \in \mathcal{X}_{\text{norm}})$ for an observation \mathbf{x} and a normal (or positive) class $\mathcal{X}_{\text{norm}}$. In practice, a proxy anomaly score function $s_a(\mathbf{x})$ is usually estimated instead. The anomalous (or negative) class is then defined implicitly as the complementary of the normal class in image space. We can generally categorize anomalies into three families:

1. **Object anomaly:** any object which is not included in the positive class, e.g., a cat is an object anomaly in regards to dogs.
2. **Style anomaly:** observations representing the same object as the positive class but with a different style or support, e.g., a realistic mask or a printed face represent faces but with a visible different style.
3. **Local anomaly:** observations representing and sharing the same style as the positive class, however a localized part of the image is different. Most of the time,

these anomalies are the superposition of two generative processes, e.g., a fake nose on a real face is a local anomaly.

Usually, we assume in anomaly detection that only normal samples are available during training, meaning that most methods are part of one-class learning scheme. The first introduced methods simply used a pre-trained neural network to extract features, on which a classical algorithm such as One-Class SVM [31] (OCSVM) or Isolation Forest [17] (IF) were trained.

There have been also semi-supervised anomaly detection methods such as DeepSAD [27] or deviation networks [23] where we assume some of the anomalies representing a few modes are available. These methods can achieve better accuracy on borderline cases given enough diverse anomalies, which is often less manageable in practice. In particular, these two methods directly learn representations by minimizing the distance of normal sample features to an hypersphere center, while maximizing the distance to the anomalies. It follows the compactness principle, where we minimize the normal class representations variance and maximize the inter-class representations variance.

1.2. Self-supervised learning

Self Supervised Learning (SSL) is a part of representation learning, where we want to learn useful and general representations from an unlabeled dataset $\mathcal{X} = \{\mathbf{x}_i\}_1^N$. We can then use the learned features for a different task such as classification.

We learn representations by solving from the data an auxiliary task \mathcal{T} , which is often unrelated to the final one. Therefore SSL consists of two steps:

1. Generating a labeled dataset $\mathcal{X}_{\mathcal{T}}$ aligned with \mathcal{T} , which for classification is usually done by applying c transformations T_j to our unlabeled samples

$$\mathcal{X}_{\mathcal{T}} = \{(T_j(\mathbf{x}_i), j)\}_{i,j} \quad (1)$$

2. Training a classification or regression network on this generated labeled set.

One of the final layers $\phi_{\mathcal{T}}$ can thus be used as a feature extractor. Some commonly used tasks are: 90° rotation prediction [9], jigsaw puzzle [22], distortions [7], colorization [41], image inpainting [24] or relative patches prediction [6].

1.3. SSL anomaly detection

Very recently, SSL has been adapted to the one-class anomaly detection framework. First we learn to solve an auxiliary task \mathcal{T} in a SSL fashion to obtain a pre-trained

network $\phi_{\mathcal{T}}$. Then, to classify at inference time an observation \mathbf{x} as anomalous or normal, we evaluate how well the network can solve the task. Indeed, the main assumption is that the network will perform relatively well on normal samples but will fail on anomalies. A task-independent metric L is computed on the generated labeled samples to compute the anomaly score function:

$$s_a(\mathbf{x}) = \{L(\phi_{\mathcal{T}}(T_i(\mathbf{x})), i) | i \in \llbracket 1, c \rrbracket\} \quad (2)$$

Unlike SSL, we are not directly interested in the intermediate features, but rather the final task outputs.

In **GeoTrans** [10], the auxiliary task is to classify which geometrical transformation has been applied to the input. A set of 72 transformations including identity is randomly sampled over all possible compositions of translations, rotations and symmetries. At the end of training, 72 Dirichlet distributions respectively parameterized by $\tilde{\alpha}_i$ are fitted over the normal class softmax outputs $\mathbf{y}(T_i(\mathbf{x}))$ for each transformation. The log-likelihood can then be used during inference as the task-independent metric L :

$$s_a(\mathbf{x}) = \sum_{i=1}^{72} (\tilde{\alpha}_i - 1) \cdot \log \mathbf{y}(T_i(\mathbf{x})) \quad (3)$$

In **MHRot** [12], the task is to simultaneously classify three types of transformations, each modeled by a softmax head: vertical translation, horizontal translations and 90° rotations. Accordingly, we are trying to predict the three following variables: vertical translations $(0, -t_y, +t_y)$, horizontal translations $(0, -t_x, +t_x)$ and 90° rotations $(0^\circ, 90^\circ, 180^\circ, 270^\circ)$.

During inference, we sum the three softmax of the known transformations for each transformation combination:

$$s_a(x) = \sum_{\substack{r \in \{0, 90, 180, 270\} \\ s \in \{0, -t_x, +t_x\} \\ t \in \{0, -t_y, +t_y\}}} \mathbf{y}(T_{r,s,t}(x))_{r,s,t} \quad (4)$$

2. Method overview

2.1. Anomaly detection pretext task

We present here a general rule of thumb regarding the choice of tasks for SSL anomaly detection. It is generally more restrictive than for simple representation learning [2].

Let \mathcal{T} be a task along its training loss $L_{\mathcal{T}}$. On the one hand, if the task is too hard on normal samples, meaning that the accuracy of our network remains close to random predictor throughout training (or that $\|\nabla_{\phi} L_{\mathcal{T}}\|$ is always small and that the minimum of L is high), then no meaningful representation will be reached at convergence. This will also result in poor accuracy on anomalies (Fig. 1.c) and yield unpredictable results during anomaly detection. On

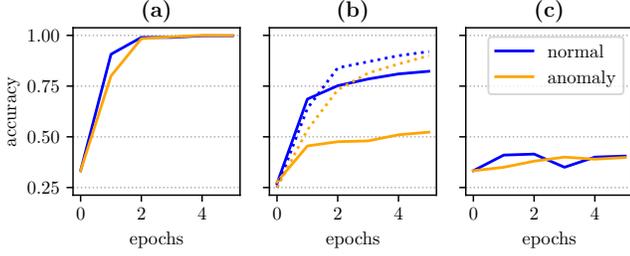


Figure 1. Tasks accuracy during training on coarse object AD (CIFAR-10 in plain line) and fine-grained AD (CaltechBirds in dotted line): (a) vertical translations, (b) 90° rotation, (c) unsolvable task

the other hand, if the task is too easy on normal samples, meaning that our model will converge to a perfect predictor in the first epochs, then the task loss will be minimized by many representations including trivial ones. Thus the network is more likely to learn such representations which will be unspecific to the normal class and encode very generic visual features. Since many anomalies will share these features, the task accuracy will be high on anomalies as well (Fig. 1.a).

To observe these effects, we train a network on several isolated tasks as described in Section 1.3. By monitoring its task classification accuracy on evaluation normal data and anomalous data during the first epochs, we empirically measure how well-suited a pretext task is for anomaly detection on a given dataset. We show that even though 90° rotation is more adapted than translations on coarse anomaly detection, it ultimately fails on fine-grained anomaly detection (Fig. 1.b). This confirms that basic geometric transformation recognition tasks, such as 90° rotations or translations, are only suited to simple object anomaly detection. Indeed, since these tasks are solvable accurately by learning high scale and shape features, it is unlikely the network will use finer characteristics that allow discriminating normal samples from more subtle anomalies.

2.2. Method overview

Finding a single task satisfying all the previous conditions is difficult, and must be highly dataset dependent. Therefore we resort to ensemble methods [5] by allowing the network to learn N tasks and merge their decision at inference. We learn richer features via multi-task learning [3], by sharing a common representation across all tasks. Our model is accordingly composed of a main feature extractor network ϕ and N dense layers $f_{\mathcal{T}_1}, \dots, f_{\mathcal{T}_N}$, where $f_{\mathcal{T}_i}(\phi(\mathbf{x}))$ is the output for the i^{th} task.

During inference, we aggregate the anomaly scores of all tasks into the final anomaly detection score. The whole

training and inference scheme is summarized in Figure 2. For each classification task \mathcal{T}_i , the task-independent metric $L_{\mathcal{T}_i}$ chosen is the softmax score corresponding to the true known class and we sum up these scores using the mean:

$$s_a(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N s_a^{(\mathcal{T}_i)}(\mathbf{x}) \quad (5)$$

where $s_a^{(\mathcal{T}_i)}$ is the anomaly score of the i^{th} task:

$$s_a^{(\mathcal{T}_i)}(\mathbf{x}) = \sum_j \text{softmax}(\phi \circ f_{\mathcal{T}_i}(T_j^{(i)}(\mathbf{x})))_j \quad (6)$$

We note that there is a caveat using the mean as anomaly score: adding new tasks can have a negative impact on the model performance. In practice if the task is not well suited to the normal class, it will add significant noise to the anomaly score and ultimately harm the anomaly detection accuracy.

To prevent our multi-task from being too easy on fine-grained problems, we introduce more challenging tasks. We choose here a simplified version of the jigsaw puzzle task. The jigsaw puzzle task consists in splitting an image into a grid of $n_h \times n_w$ patches, then randomly shuffling the different patches. The task is then to predict the original order of each patch. This task has proven in representation learning to provide a great challenge for extracting more local and finer features [22]. To avoid trivial solutions and force our model to understand pieces neighborhood, we are careful to add a margin between each patch with a random small offset.

Since we chose the softmax truth as the task-independent metric, we need to re-frame it into a classification problem by considering each permutation as a single class. This would greatly increase our model complexity, effectively adding $(n_w \cdot n_h)!$ classes. Therefore, we only consider $k < (n_w \cdot n_h)!$ randomly chosen permutations including the identity permutation. This quantity k becomes an additional parameter controlling the task difficulty.

The complete training loss for a single sample \mathbf{x} becomes

$$L(\mathbf{x}) = \sum_{i=1}^3 L_{CE}(\phi \circ f_v(T_i^{(v)}(\mathbf{x})), i) + \sum_{i=1}^3 L_{CE}(\phi \circ f_h(T_i^{(h)}(\mathbf{x})), i) + \sum_{i=1}^4 L_{CE}(\phi \circ f_{\text{rot}}(T_i^{(\text{rot})}(\mathbf{x})), i) + \sum_{i=1}^k L_{CE}(\phi \circ f_{\text{puzz}}(T_i^{(\text{puzz})}(\mathbf{x})), i) \quad (7)$$

where L_{CE} is the cross-entropy and $f_v, f_h, f_{\text{rot}}, f_{\text{puzz}}$ are respectively the dense layers for the vertical translations,

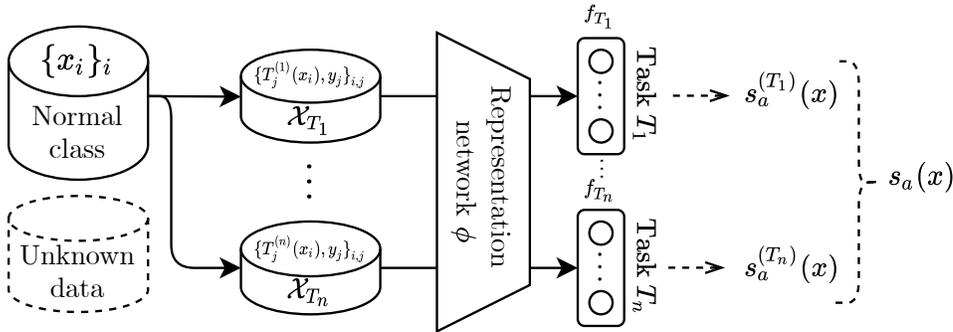


Figure 2. Multi-task self-supervised anomaly detection. In dotted line are additional steps during inference

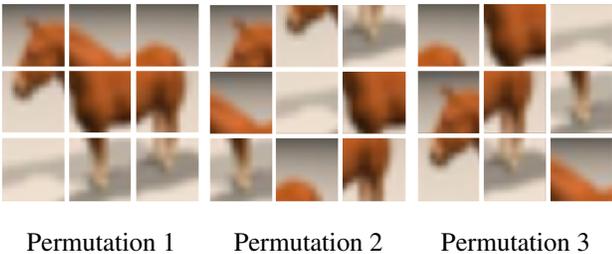


Figure 3. Example of simplified jigsaw puzzle task for $k = 3$

horizontal translations, rotations and puzzle tasks. Compared to the MHRot anomaly score in Equation 4 we evaluate each task *independently*, which allows us to greatly reduce the required amount of network forward pass to compute the anomaly score during inference. As a consequence, our method inference step is roughly 10 times faster.

By combining the base geometrical transformation recognition task with the jigsaw task, we allow the model to learn high-scale shape features more suited toward object anomaly detection as well as low-scale fine features more suited toward style anomaly and local anomaly detection.

3. Implementation details

The geometrical transformation task is composed as in [12] of horizontal translations, vertical translations and 90° rotations. As for the jigsaw puzzle task, we found best results with $n_w = n_h = 3$ and $k = 3$.

Regarding network architecture, we use a 16-4 WideResNet [40] ($\approx 10M$ parameters with a depth of 16) for the feature extractor network ϕ , along with two dense softmax layers respectively of size 10 for the geometrical transformation task and size 3 for the jigsaw puzzle task. Each of these dense layers have a dropout rate of 0.3 [32]. Training is performed under SGD optimizer with nesterov momentum [34], using a batch size of 32.

4. Results

4.1. Evaluation protocol

Until now, most of the anomaly detection literature have adopted the one-vs-all protocol to evaluate their method. In the one-vs-all protocol, we consider one class of a multi-class dataset, originally created for object recognition, as the normal class. All the other classes are then considered as anomalous, and we can in a leave-one-out cross-validation fashion evaluate the model on each possible normal class. The final reported result is the mean of each run.

Even though such datasets are easier to acquire and result in a highly multi-modal anomaly class, these might not be enough to fully evaluate anomaly detection methods. Indeed, these only cover coarse object anomalies which are now becoming too easy for state-of-the-art methods, and do not reflect realistic anomaly detection challenges.

	Dataset	Anomaly type		
		Object	Style	Local
Obj.classif	MNIST	✓	✗	✗
	F-MNIST	✓	✗	✗
	CIFAR-100	✓	✗	✗
Fine-grained	Caltech-Birds	✓	✓	✗
	FounderType	✗	✓	✗
Anti-spoofing	SiW-M	✓	✓	✓

Table 1. Summary of evaluation datasets.

Thus we propose to use fine-grained classification datasets in the same one-vs-all protocol. Since discrimination between these classes is mostly based on local and fine patterns, we can have a good coverage of style anomalies and local anomalies. Also we note that because of the increased shift in object recognition toward fine-grained classification, such datasets have become readily available. For one-vs-all datasets, we used MNIST [16], Fashion MNIST

Model		CIFAR-100	MNIST	F-MNIST	Caltech-Birds 200	Fonts	SiW-M
Semi-Supervised	Deep-SAD (75%) [27]	88.7	<u>99.9</u>	<u>98.1</u>	73.6	<u>99.8</u>	85.4
	Deep-SAD (25%)	87.9	98.5	95.4	70.9	99.4	76.0
	Deep-SAD (10%)	<u>89.1</u>	96.5	88.2	66.1	98.0	80.6
One-class	ADGAN [4]	54.7	94.7	88.4	-	-	-
	GANomaly [1]	56.5	92.8	80.9	-	-	-
	ARNet [8]	78.8	98.3	93.9	-	-	-
	OCSVM [31]	-	84.7	74.2	76.3	-	-
	IF [17]	-	87.1	84.0	74.2	-	-
	PIAD [36]	78.8	98.1	94.3	63.5	-	81.2
	GeoTrans [10]	84.7	96.9	92.6	66.6	92.3	81.1
	MHRot [12]	83.6	95.2	92.5	77.6	96.7	83.1
	Ours	85.8	96.0	92.8	83.2	96.9	88.4

Table 2. Comparison with the state-of-the-art AUROC over several datasets, underline indicates best result, bold indicates best one-class learning result. We re-implemented all the methods except the three one-class methods in the first block (results are from original papers [4, 8]).

Model	Airplane	Automobile	Bird	Cat	Deer	Dog	Frog	Horse	Ship	Truck	Avg
VAE [13]	70.0	38.6	67.9	53.5	74.8	52.3	68.7	49.3	69.6	38.6	58.3
OCSVM [31]	63.0	44.0	64.9	48.7	73.5	50.0	72.5	53.3	64.9	50.8	58.5
AnoGAN [30]	67.1	54.7	52.9	54.5	65.1	60.3	58.5	62.5	75.8	66.5	61.8
PixelCNN [37]	53.1	99.5	47.6	51.7	73.9	54.2	59.2	78.9	34.0	66.2	61.8
Deep-SVDD [28]	61.7	65.9	50.8	59.1	60.9	65.7	67.7	67.3	75.9	73.1	64.8
OCGAN [26]	75.7	53.1	64.0	62.0	72.3	62.0	72.3	57.5	82.0	55.4	65.6
Puzzle-AE [29]	78.9	78.0	69.9	54.8	75.4	66.0	74.7	73.3	83.3	69.9	72.4
DROCC [11]	81.7	76.7	66.7	67.1	73.6	74.4	74.4	71.4	80.0	76.2	74.2
GeoTrans [10]	74.7	95.7	78.1	72.4	87.8	87.8	83.4	95.5	93.3	91.3	86.0
Ours	75.1	96.3	84.8	74.2	91.1	89.9	88.7	95.5	94.7	91.9	88.2

Table 3. Detailed comparison with one-class state-of-the-art AUROC on CIFAR-10 dataset.

[39], CIFAR-100 [14]. For the fine-grained dataset, we chose the Caltech-Birds 200 database [38].

We also put forward datasets from real anomaly detection problems over different fields. First, we use font recognition challenges as they provide shape-focused style anomaly detection. Indeed two different fonts represent the same characters albeit with a distinctive style. Even though these images lie on a low dimensional manifold compared to natural images, they still provide insight into how well the model can capture small shape hints. In particular, we use **FounderType-200** [18] introduced for novelty detection and containing 6700 images per font. Furthermore, we choose a dataset from face anti-spoofing, where the goal is to discriminate real faces from fake representations of someone’s face. Due to the richness and high variability of such frauds, this problem effectively encompasses all three types of anomalies. We use here the **Spoof in the Wild Multiple** (SiW-M) [19] database which contains more than

1600 short videos of real faces and presentation attacks. There are 493 real identities along with several types of attacks: paper print, screen replay, masks and partial attacks where only a localized area of the face is fake. The masks are composed of half-masks, paper masks, silicone mask and transparent masks. All evaluation datasets are summarized in Table 1.

We additionally evaluate our anomaly detection model on out-of-distribution (OOD) protocol. OOD detection, which is broader than anomaly detection, aims at discriminating the training dataset from other data distributions. The “normal” distribution in OOD is therefore usually more diverse and highly multi-modal. We also have a greater overlap in term of class between the in-distribution samples and out-of-distribution samples compared to anomaly detection. Nevertheless it gives us great insight into the multi-modality limits of our model. The most common evaluation setup is to discriminate one training multi-class dataset from

other datasets. Here we choose to learn on CIFAR-10 and discriminate CIFAR-100 and the easier Street View House Numbers (SVHN) dataset [21].

For all of the evaluations, the metric used is the area under the ROC curve (AUROC), averaged over all possible normal classes in the case of one-vs-all datasets.

4.2. Ablation study

We evaluate in Table 4 how combining the two tasks of geometric transformation recognition and jigsaw puzzle improves the anomaly detection. We drastically improve performances with a relative error reduction regarding AUROC of 13% on CIFAR-100, 25% on Caltech-Birds 200 and 31% on SiW-M. This validates our statement in Section 2: the finer the differences between anomaly and normal class, the greater the improvement is by adding the jigsaw task.

Auxiliary Task	CIFAR-100	Caltech-Birds	SiW-M
Geometric (G)	83.6	77.6	83.1
Jigsaw (J)	80.1	78.5	76.3
Ours (G+J)	85.8	83.2	88.4

Table 4. AUROC for different tasks, best result is in bold.

4.3. Comparison to the state-of-the-art

We compare our method with different one-class learning state-of-the-art approaches to anomaly detection: reconstruction error generative models with the PIAD model [36], self-supervised methods with GeoTrans [10] and MHRot [12]. As an addition, we include a semi-supervised learning anomaly detection method DeepSAD [27], which has access to a portion of the anomalies during training. As such, we train it with the same normal samples but three different ratio of the anomaly subclasses: 10%, 25% and 75%.

For the sake of fair comparison in the same conditions, we take the existing implementations or re-implement each method and evaluate each, except for the ADGAN, GANomaly and ARNet which we reference results from their original papers [4, 1, 8].

The results are gathered in Table 2 and 3. First of all, we can see our method generally maintains among the best accuracies on simple object anomaly detection, and even improves it on more challenging datasets such as CIFAR-100. Moreover, it greatly improves fine-grained anomaly detection and outperforms state-of-the-art methods which could not be realistically be used for this problem. We also show that our method, without further tuning, improves anti-spoofing detection performances on SiW-M. Finally, we notice our one-class learning model generally reduces the gap with semi-supervised method, and even outperforms

these on Caltech-Birds 200 and SiW-M, even though these take advantage of a significant amount of additional anomalous data.

Metrics	AUROC	EER	APCER (5%BPCER)
MHRot [12]	83.0	21.6	77.5
Ours	88.4	18.7	39.1

Table 5. AUROC, EER and APCER at 5% BPCER of MHRot and our model with jigsaw task on SiW-M dataset, best result is in bold.

We compare in Table 5 our method with the second best self-supervised method MHRot on SiW-M. We use metrics more adapted to face presentation attack detection with equal error rate (EER) and the false acceptance rate for the rate of false reject fixed at 5% (APCER@5%BPCER). Our comparison does not include other face anti-spoofing methods since we only use real faces training images while these all use a set of presentation attacks during training. Using our method, the APCER@5%BPCER drops from 77.5 to 39.1 thus also showing promising usage of anomaly detection methods in fraud detection.

OOD	SVHN	CIFAR-100	Avg
VAE [13]	2.4	52.8	27.6
Deep-SVDD [28]	14.5	52.1	33.3
PixelCNN [37]	15.8	52.4	34.1
RotNet [9]	97.9	81.2	89.5
Ours	98.8	83.4	91.1
CSI [35]	99.8	89.2	<u>94.5</u>

Table 6. Comparison with state-of-the-art on the Out-Of-Distribution detection protocol with CIFAR-10 as in-distribution, best result is underlined, best pretext task driven method is in bold.

Lastly, we compare our model with state-of-the-art on OOD detection in Table 6. Although not designed specifically for such complex normal class, we obtain better detection rates than other self-supervised anomaly detection methods with pretext tasks.

5. Conclusion and Future Work

In this paper, we investigate the power of multi-task self supervision for anomaly detection and show the limits of simple geometric tasks. In more details, we combine two complementary tasks of jigsaw puzzle and geometric transformation recognition. Through an ablation study, we show that this enables it to learn much complex and finer features and therefore better detect anomalies. Finally, we provide a more comprehensive evaluation protocol than previously

used datasets in the anomaly detection literature. It presents more challenging datasets and covers object, style and local anomalies. Our method outperforms state-of-the-art, including a semi-supervised method, on most of the fine-grained datasets.

For future work we could explore the combination of more tasks, including generative tasks (in contrast to discriminative tasks used here). Such tasks could range from re-colorization to image in-painting.

References

- [1] S. Akcay, A. Atapour-Abarghouei, and T. P. Breckon. Ganomaly: Semi-supervised anomaly detection via adversarial training. In *Asian Conference on Computer Vision*, pages 622–637. Springer, 2018. 5, 6
- [2] Y. M. Asano, C. Rupprecht, and A. Vedaldi. A critical analysis of self-supervision, or what we can learn from a single image. In *International Conference on Learning Representations*, 2020. 2
- [3] R. Caruana. Multitask learning. *Machine Learning*, 28(1):41–75, Jul 1997. 3
- [4] L. Deecke, R. Vandermeulen, L. Ruff, S. Mandt, and M. Kloft. Image Anomaly Detection with Generative Adversarial Networks. In M. Berlingerio, F. Bonchi, T. Gärtner, N. Hurley, and G. Ifrim, editors, *Machine Learning and Knowledge Discovery in Databases*, Lecture Notes in Computer Science, pages 3–17. Springer International Publishing, 2019. 5, 6
- [5] T. G. Dietterich. Ensemble methods in machine learning. In *Multiple Classifier Systems*, pages 1–15. Berlin, Heidelberg, 2000. Springer Berlin Heidelberg. 3
- [6] C. Doersch, H. Mulam, and A. Efros. Unsupervised visual representation learning by context prediction. *IEEE International Conference on Computer Vision (ICCV)*, pages 1422–1430, 05 2015. 2
- [7] A. Dosovitskiy, J. Springenberg, M. Riedmiller, and T. Brox. Discriminative unsupervised feature learning with exemplar convolutional neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1, 06 2014. 2
- [8] Y. Fei, C. Huang, C. Jinkun, M. Li, Y. Zhang, and C. Lu. Attribute Restoration Framework for Anomaly Detection. *IEEE Transactions on Multimedia*, pages 1–1, 2020. 5, 6
- [9] S. Gidaris, P. Singh, and N. Komodakis. Unsupervised representation learning by predicting image rotations. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018. 2, 6
- [10] I. Golan and R. El-Yaniv. Deep anomaly detection using geometric transformations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 9758–9769. Curran Associates, Inc., 2018. 1, 2, 5, 6
- [11] S. Goyal, A. Raghunathan, M. Jain, H. V. Simhadri, and P. Jain. Drocc: Deep robust one-class classification. In *International Conference on Machine Learning*, 2020. 5
- [12] D. Hendrycks, M. Mazeika, S. Kadavath, and D. Song. Using self-supervised learning can improve model robustness and uncertainty. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2, 4, 5, 6
- [13] D. Kingma and M. Welling. Auto-encoding variational bayes. In *Conference proceedings: papers accepted to the International Conference on Learning Representations (ICLR)*, 2014. 5, 6
- [14] A. Krizhevsky. Learning multiple layers of features from tiny images. 2009. 5
- [15] D. Kwon, H. Kim, J. Kim, S. Suh, I. Kim, and K. Kim. A survey of deep learning-based network anomaly detection. *Cluster Computing*, 22, 01 2019. 1
- [16] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 4
- [17] F. T. Liu, K. Ting, and Z.-H. Zhou. Isolation forest. In *2008 Eighth IEEE International Conference on Data Mining*, pages 413 – 422, 01 2009. 2, 5
- [18] J. Liu, Z. Lian, Y. Wang, and J. Xiao. Incremental kernel null space discriminant analysis for novelty detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4123–4131, 07 2017. 5
- [19] Y. Liu, J. Stehouwer, A. Jourabloo, and X. Liu. Deep tree learning for zero-shot face anti-spoofing. In *In Proceedings of IEEE Computer Vision and Pattern Recognition*, Long Beach, CA, June 2019. 5
- [20] A. S. Lundervold and A. Lundervold. An overview of deep learning in medical imaging focusing on mri. *Zeitschrift für Medizinische Physik*, 29(2):102 – 127, 2019. Special Issue: Deep Learning in Medical Physics. 1
- [21] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011. 6
- [22] M. Noroozi and P. Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In B. Leibe, J. Matas, N. Sebe, and M. Welling, editors, *Computer Vision – ECCV 2016*, pages 69–84, Cham, 2016. Springer International Publishing. 2, 3
- [23] G. Pang, C. Shen, and A. V. D. Hengel. Deep anomaly detection with deviation networks. *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019. 2
- [24] D. Pathak, P. Krähenbühl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2536–2544, 2016. 2
- [25] A. Paudice, L. Muñoz-González, A. György, and E. C. Lupu. Detection of adversarial training examples in poisoning attacks through anomaly detection. *CoRR*, abs/1802.03041, 2018. 1

- [26] P. Perera, R. Nallapati, and B. Xiang. Ocgan: One-class novelty detection using gans with constrained latent representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2898–2906, 2019. 5
- [27] L. Ruff, R. A. Vandermeulen, N. Görnitz, A. Binder, E. Müller, K.-R. Müller, and M. Kloft. Deep semi-supervised anomaly detection. In *International Conference on Learning Representations*, 2020. 2, 5, 6
- [28] L. Ruff, R. A. Vandermeulen, N. Görnitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft. Deep one-class classification. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 4393–4402, 2018. 5, 6
- [29] M. Salehi, A. Eftekhari, N. Sadjadi, M. H. Rohban, and H. R. Rabiee. Puzzle-ae: Novelty detection in images through solving puzzles. abs/2008.12959. 5
- [30] T. Schlegl, P. Seeböck, S. M. Waldstein, U. Schmidt-Erfurth, and G. Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *International conference on information processing in medical imaging*, pages 146–157. Springer, 2017. 5
- [31] B. Schölkopf, R. Williamson, A. Smola, J. Shawe-Taylor, and J. Platt. Support vector method for novelty detection. In *Proceedings of the 12th International Conference on Neural Information Processing Systems, NIPS’99*, page 582–588, Cambridge, MA, USA, 1999. MIT Press. 2, 5
- [32] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.*, 15:1929–1958, 2014. 4
- [33] W. Sultani, C. Chen, and M. Shah. Real-World Anomaly Detection in Surveillance Videos. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6479–6488, 06 2018. ISSN: 2575-7075. 1
- [34] I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and momentum in deep learning. In S. Dasgupta and D. McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1139–1147, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. 4
- [35] J. Tack, S. Mo, J. Jeong, and J. Shin. CSI: novelty detection via contrastive learning on distributionally shifted instances. *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, virtual*, 2020. 6
- [36] N. Tuluptceva, B. Bakker, I. Fedulova, and A. Konushin. Perceptual Image Anomaly Detection. In S. Palaiahnakote, G. Sanniti di Baja, L. Wang, and W. Q. Yan, editors, *Pattern Recognition*, Lecture Notes in Computer Science, pages 164–178. Springer International Publishing, 2020. 5, 6
- [37] A. Van den Oord, N. Kalchbrenner, L. Espeholt, O. Vinyals, A. Graves, et al. Conditional image generation with pixelcnn decoders. In *Advances in neural information processing systems*, pages 4790–4798, 2016. 5, 6
- [38] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-UCSD Birds 200. Technical Report CNS-TR-2010-001, California Institute of Technology, 2010. 5
- [39] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017. 5
- [40] S. Zagoruyko and N. Komodakis. Wide residual networks. In R. C. Wilson, E. R. Hancock, and W. A. P. Smith, editors, *Proceedings of the British Machine Vision Conference 2016, BMVC 2016, York, UK, September 19-22, 2016*. BMVA Press, 2016. 4
- [41] R. Zhang, P. Isola, and A. A. Efros. Colorful image colorization. In B. Leibe, J. Matas, N. Sebe, and M. Welling, editors, *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III*, volume 9907 of *Lecture Notes in Computer Science*, pages 649–666. Springer, 2016. 2
- [42] Z. Zhang, X. Zhou, X. Zhang, L. Wang, and P. Wang. A model based on convolutional neural network for online transaction fraud detection. *Security and Communication Networks*, 2018:1–9, 08 2018. 1
- [43] S. Zhu, C. Chen, and W. Sultani. Video anomaly detection for smart surveillance. *CoRR*, abs/2004.00222, 2020. 1