

Image-based Novel Fault Detection with Deep Learning Classifiers using Hierarchical Labels

Nurettin Sergin^a, Jiayu Huang^a, Tzyy-Shuh Chang^b, Hao Yan^a

^aSchool of Computing and Augmented Intelligence, Arizona State University, Tempe, USA

^bOG Technology, Ann Arbor, MI, USA

Abstract

One important characteristic of modern fault classification systems is the ability to flag the system when faced with previously unseen fault types. This work considers the unknown fault detection capabilities of deep neural network-based fault classifiers. Specifically, we propose a methodology on how, when available, labels regarding the fault taxonomy can be used to increase unknown fault detection performance without sacrificing model performance. To achieve this, we propose to utilize soft label techniques to improve the state-of-the-art deep novel fault detection techniques during the training process and novel hierarchically consistent detection statistics for online novel fault detection. Finally, we demonstrated increased detection performance on novel fault detection in inspection images from the hot steel rolling process, with results well replicated across multiple scenarios and baseline detection methods.

Keywords: Novel fault detection, Hierarchical structure, Deep learning, Fault classification.

1 Introduction

Many manufacturing systems are instrumented with image-sensing systems to monitor process performance and product quality. The low cost and rich information of the image-based sensing systems have led to high-dimensional data streams that provide distinctive opportunities for performance improvement. Among these, accurate process monitoring and fault classification are among the benefits gained from the rich information these image sensors can provide. In literature, process monitoring often refers to the step of detecting and isolating abnormal samples in a certain process. Normally, after process monitoring, fault classification is performed, and the isolated fault is classified into one or more known types of fault. Fault classification is an essential step within the process monitoring loop, at which point the type of detected and identified faults are determined (Chiang et al., 2001). Accurate fault classification can provide engineers with favorable information to isolate and diagnose system faults and anomalies to improve quality and maximize system efficiency.

However, fault classification in manufacturing systems typically assumes a fixed set of fault modes. In this case, the existing fault classification model may make overconfident decisions or fail silently and, at certain times, dangerously for new unseen fault types. An additional policy should be developed to flag the emergence of unseen faults to alert the system and log the related instances for human evaluation and subsequent updates of

the model. Such practice has been implemented in some monitoring systems in modern machine learning models (Breck et al., 2017). We name this problem image-based novel fault detection, which will be the major focus of this paper.

The recent resurgence in machine learning research heavily influenced anomaly detection, novel fault detection, and fault classification literature (Liu et al., 2018). For example, feature-based and tensor-based anomaly detection and fault classification methods have been proposed (Yan et al., 2014). However, the traditional methods typically rely on the accurate definition of the handcrafted features, which often require much labor work and also lead to unsatisfying performance. Deep learning, a subclass of machine learning methods, has attracted a peculiar interest. What makes deep learning methods so attractive for anomaly detection (Kwon et al., 2019; Sergin and Yan, 2021) or fault classification (Pan et al., 2017) is their ability to learn useful feature representation for data representation or class discrimination automatically from data as opposed to requiring handcrafted ones, which is often challenging in fault classification. Despite the impressive performance increases, deep learning has brought the the problem of fault classification and anomaly detection. These models either assume that no labels are available (e.g., anomaly detection) or assume a closed set of fault types by using a fixed number of neurons in the last layers (e.g., classification), which cannot be used in novel fault detection.

To address the novel fault detection problem, some literature has recently proposed out-of-distribution methods to detect new emerging classes. In these works, a set of fault classes is assumed to be given, and the goal is to detect novel or unseen faults in the system. For more details about the novel fault classification, we will provide a complete literature review in Section 2.1.

However, in many existing complex systems, the classes typically can be represented in a hierarchical tree structure rather than a flat structure. In modern manufacturing systems, many fault classes exist and the fault classes can often be represented in a hierarchical structure. For example, for the defect inspection problem in the rolling manufacturing system, 14 different types of anomalies have been identified. These 14 anomalies can be classified into 8 subcategories, given the shapes of the defect areas. Here, many fault classes may have a very small sample size (Sahoo et al., 2003). Such a hierarchical fault structure is common in manufacturing systems due to the hierarchical root causes contributing to product quality often represented by the fishbone diagram. The hierarchical set of factors organized by the fishbone diagram may also lead to a similar set of hierarchical fault patterns. An example of such a hierarchical fault pattern in a steel rolling process inspection is shown in Figure 1. The first level determines which major category of defect this inspection image patch belongs to. Some examples of these defect categories are chip marks, cracks, and overfill. The second level defines the subcategory within each major category. For example, the overfilled major category contains black (A12), white (A11), and lite (A10) subcategories. The major shape of these categories are similar (showing verticle patterns), but the color is different for each subcategory.

In the literature, there are many existing works on hierarchical fault classification methods considering hierarchical fault structures. In these works, hierarchical classification methods have demonstrated improved classification accuracy over common flat classifiers. The improvement in performance is mostly due to the following two major reasons: First, it can guide the learning process due to the injected inductive bias and result in a more accurate estimation. Second, the knowledge learned from the other fault nodes can be transferred to

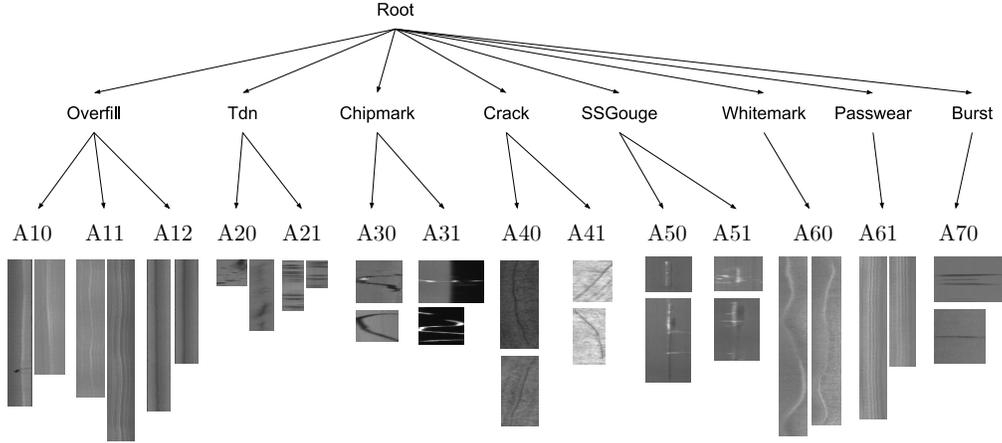


Figure 1: Graphical illustration of the hierarchy of defects in the hot steel rolling dataset.

the rare fault modes by utilizing the relationship between the fault modes within the same ancestor. Finally, the incorporation of the hierarchical information can help understand the severity of the mistakes in the training process, therefore improving the overall classification accuracy. The motivation behind this is that if the classifier fails to identify the right class, it may still predict it as a known class that is closer to the ground truth in the hierarchy (i.e., share the same parent node).

However, even though a large number of works for hierarchical classification exist, to the best of the authors' knowledge, there is currently no deep learning-based novel fault detection methods considering this hierarchical fault taxonomy.

In this work, we propose a novel approach such that the knowledge of the fault hierarchical taxonomy can be incorporated into both the training procedure and the online novel fault detection procedure by combining soft labeling and hierarchically consistent score in deep learning. We demonstrate how a model trained in consideration of the hierarchical relationship between the labels will be more effective at detecting previously unseen fault classes. We further demonstrate that the proposed method can improve upon state-of-the-art deep learning-based fault detection methods. While many fault sets are inherently hierarchical (e.g., fault trees), oftentimes this information is neglected in favor of the simplicity of using only the label information at leaf nodes. With this work, we aim to convince the practitioners to inject this information into the training procedure and the detection statistics for anomaly detection.

The main contribution of this paper can be summarized as follows:

- We propose a novel methodology by considering the hierarchical structure of taxonomy in real-world novel fault detection by incorporating it into a deep-learning model classifier. By injecting hierarchical information into the loss function in the deep learning model, we consider both classification accuracy and hierarchical structure consistency during model training.
- We propose a hierarchically consistent score to detect unknown faults online. The primary purpose of this score is to identify samples that are inconsistent with the known fault hierarchical taxonomy, improving the overall performance of the classifier-based out-of-distribution detector.

- We provide some methodology insights into why including such hierarchical information in offline training and introducing hierarchical consistency scores for online detection can improve the novel fault detection accuracy. Additionally, we propose two propositions to illustrate the insights of our proposed score and visualization of why hierarchical training would be beneficial. Finally, we discuss the practical implementation of the proposed methods for a rolling process with hierarchical structures on the anomaly taxonomy and demonstrate the improved accuracy for novel fault detection.

The rest of the article is organized as follows. In Section 2, we provide a background of related works in the out-of-distribution detection and hierarchical classification literature. In Section 3, we provide the details of the interventions that we propose at training and test time. In Section 4, we introduce our hot steel rolling image defect dataset, formalize our experiment design and define the performance evaluation criteria we consider. The results of the experiments and related discussions are presented in Section 4.3. Finally, we conclude the work in Section 5, and a brief data availability statement will also be included in Section 6.

2 Related Work

This section briefly reviews the literature related to out-of-distribution detection in Section 2.1 and the use of hierarchical classification to address the hierarchical labeling relationship in Section 2.2, as these two methods are the fundamental pillars of the proposed fault detection algorithm that is discussed in the next section.

2.1 Out-of-distribution Detection

As mentioned in (Bendale and Boult, 2016), deep learning classifiers lack the ability to detect unknown faults due to the implicit closed-set assumption, which states that all classes are known a priori (Bendale and Boult, 2016). This closed-set assumption prevents the deep-learning methodology from applying to safety-related applications. The research on out-of-distribution (OOD) detection focuses on detecting anything other than what lies in the training set.

The definition of what can be considered out-of-distribution is broad. Samples from previously unseen classes, adversarial examples (Goodfellow et al., 2015) or domain shift (Ben-David et al., 2010) can all be considered OOD.

2.1.1 Classifier-based models

The literature on OOD detection reveals two primary approaches: 1) classifier-based models, which differentiate between normal and anomaly based on a specific optimization function. 2) generative models, which involve modeling the data’s distribution and detecting anomalies based on deviations from learned distribution or reconstruction, such as the Variational Auto-encoder (VAE) and Generative Adversarial Networks (GAN).

In the out-of-distribution detection combined with the classification problem, the uncertainty of the labels, called aleatoric uncertainty (Kendall and Gal, 2017)(Sallak et al., 2013), is deserved to be considered since the samples we obtained could be a novel and unseen class

and the neural network outputs are often probabilities. This type of uncertainty can be used as an indicator of anomaly samples. The rest of the section will demonstrate the aleatoric uncertainty usage for out-of-distribution detection in some methods.

Two main types of classifier-based methods have been developed. The first class focuses on utilizing Maximum Softmax Probability (MSP) (Hendrycks and Gimpel, 2017). MSP is the largest probability value produced by the deep learning classifiers. Their major observation in (Hendrycks and Gimpel, 2017) was that probability assigned to the most likely class (i.e., MSP) differed considerably between in-distribution and out-of-distribution samples. For example, the probability distribution of the out-of-distribution samples tend to be closer to the uniform distribution compared to the in-distribution samples. In the MSP method, the idea of the aleatoric uncertainty difference between in-distribution and out-of-distribution samples is applied. Later, many efforts have been contributed to improving the use of MSP for out-of-distribution detection. Among those methods, the Out-of-Distribution detector for Neural networks (ODIN) (Liang et al., 2018) has improved the performance of MSP by modifying the training regime with temperature scaling and adversarial perturbation, which significantly improves the separation between in-distribution and out-of-distribution samples. According to the review paper (Shafaei et al., 2019), ODIN has achieved state-of-art performance in all MSP-based methods.

The second line of work focuses on the neuron response values at different levels to detect samples from novel classes. This method is originally proposed by Lee et al. (2018), where the authors observe that neuron response outputs at various levels of a neural network differ significantly from one class to the other, as well as the in-distribution class to the out-of-distribution class. Therefore, the neuron response of a new sample can be compared with the distribution of the response from the existing known classes, where the sample-to-distribution distance can be used as an abnormality score. For example, multivariate Gaussian distributions with tied covariances at each level are used for modeling the neuron response for each known class, and the Mahalanobis distance is used to define a potential new class in the samples. While experiments suggest further improvements in accuracy, this method relies on additional out-of-distribution datasets for the tuning of the hyper-parameters.

We also would like to mention that there are some other frameworks besides these two lines of research by utilizing the out-of-distribution samples to improve the out-of-distribution detection performance. Among these methods, Hendrycks et al. (2019) proposed a complementary method to the existing scoring-based out-of-distribution detection methods. The method, named Outlier Exposure, simply fine-tunes an already trained model, again with the use of the auxiliary outlying dataset, to encourage even more separation between in-distribution and out-of-distribution, regardless of the scoring method used. However, these methods assume that the out-of-distribution samples are available, which is not applicable in the application of novel fault detection.

2.1.2 Generative models

Generative models such as VAE and GAN have been widely used for anomaly detection via reconstruction error. However, their performance in detecting out-of-distribution (OOD) data has been questioned, particularly for image data, given their focus on detecting global patterns rather than local anomalies (Nalisnick et al., 2018; Havtorn et al., 2021). Additionally, VAE methods have been criticized for generating blurry images with little detail,

limiting their effectiveness for OOD detection. While VAE-based OOD detection techniques have been developed, they have primarily been evaluated on natural image datasets with very clear intra- and inter-class variations, such as MNIST, FashionMNIST, etc (Guo et al., 2021). GAN methods have been criticized for being prone to mode collapse (Zhang, 2021; Murray and Rawat, 2021), where they only generate a subset of the possible variations of the target distribution, leading to a limited representation of the data.

Some initial works have been applied to OOD for industrial image inspection (Sergin and Yan, 2021; Yan et al., 2019). However, in many classification-based OOD problems, where the class information is available, we are interested in detecting the unseen anomaly, given that multiple classes exist. Traditional VAE or GAN lacks the ability to encode the class information.

Even though in literature, there are some existing works on encoding the class information into generative models, such as conditional VAE (Mirza and Osindero, 2014), Knowledge-oriented VAE (Shen et al., 2019) and conditional GAN (Sohn et al., 2015). However, in literature, the class information is often represented as the flat structures, instead of the hierarchical structure used in this paper.

Several comparison studies in the literature have shown that generative models generally have poorer performance than most other classifier-based methods, such as ODIN, MSP, and DMD, especially when class information is available (Li et al., 2022). Furthermore, training a generative model such as VAE or GAN is typically more expensive than training a classifier, making them less practical for practitioners.

To the best of our knowledge, there has been no exploration of incorporating hierarchical label structure into deep learning models for improving the performance of existing classifiers or other generative methods in OOD detection.

In this paper, we propose a method that leverages the statistical properties of deep learning classifiers by injecting hierarchical knowledge into their optimization functions. The proposed hierarchical models have the potential to perform better at detecting unseen anomalies.

2.2 Hierarchical Classification

In this subsection, we will first make a clear definition of hierarchical classification. We will then introduce the existing literature as well as the advantages of a hierarchical classifier over a flat classifier.

First, a formal definition of the hierarchical classification model, originally proposed by Jr. and Freitas (2011), can be given with the following three assumptions:

1. A hierarchy can be modeled either by a tree or a directed acyclic graph, depicting subsumption relationships among classes. Each class is subsumed under its parent(s) in this structure.
2. The structure is pre-defined and provided by the modeler. It is not meant to be learned from the data.
3. Each data point belongs to one and only one of the leaf nodes in the hierarchy.

Given this structure, a simple approach to classification would be flattening out all the leaf nodes and treating the classification problem as an N -way flat classification task where

there are, in total, N leaf nodes. However, this structure treats each class as equally different from the other and completely ignores the rich information provided by the hierarchy.

In literature, many existing methods utilize hierarchical classification to improve the classification accuracy of different tasks, such as bioinformatics (Freitas and Carvalho, 2007), community data (Gauch Jr and Whittaker, 1981), web content (Dumais and Chen, 2000), accident report (Zhao et al., 2021), e-commerce (Shen et al., 2012), and visual recognition (Yan et al., 2015). Overall, Silla and Freitas (2011) divided the existing approaches into the top-down local classifier approach and the global-classifier approach.

The top-down local classifier is originally proposed by Koller and Sahami (1997). More specifically, these approaches can be further divided into a local classifier per node (Fagni and Sebastiani, 2007), a local classifier per parent node (Brown and Mues, 2012), and a local classifier per level (Clare and King, 2003), depending on how the local information is considered. However, one of the major disadvantages of the top-down approach is that any error at a certain class level will be propagated down the hierarchy. Please see (Silla and Freitas, 2011) for a more detailed review of the local classifier approach.

To overcome the drawback of the local classifier approach, global classifiers are proposed by Xiao et al. (2007), which aims to consider a single classification model trained from the training set, taking into consideration of the class hierarchy as a whole during a single run of the classification algorithm. Typically, learning a single model for all classes has the advantage that the size of the global classification model is considerably smaller, which is especially important for deep-learning-based models, given that deep-learning models are already large. Moreover, the dependencies between the classes can be taken into account in a more natural and straightforward way. Currently, many approaches exist, such as multi-label classification methods (Kiritchenko et al., 2005) and distance-based methods (Rocchio, 1971). The multi-label classification approach represents each label and its structure using multiple attributes, and the algorithm aims to predict all attributes simultaneously. Distance-based methods aim to assign the new sample to the nearest class by computing the distance between the new test example and each class, considering the hierarchical structure in the distance measure. Recently, a new global hierarchical classification model has been proposed (Bertinetto et al., 2020) by defining the soft labels to represent the hierarchical structure of the class. Such a method is simple to implement while achieving state-of-the-art performance in the ImageNet classification. Therefore, this paper focuses on incorporating the “soft-label” approach toward the OOD performance with the hierarchical label structure.

3 Proposed Methodology

We present the mathematical framework related to the key elements of the neural network-based novel fault detection systems that will be used throughout this section.

Assume we have a dataset $\mathcal{D} = \{(\mathbf{x}^n, y^n) \mid n \in \{1 \dots N\}\}$ where \mathbf{x}^n are inputs and $y^n \in \{1 \dots K\}$ denote which of the K labels that \mathbf{x}^n belong to. In the context of fault classification, \mathbf{x} can be various inputs, such as images collected from automated optical inspection cameras or readings from sensors. Each input is associated with one and only one fault label, y , which is assumed to be drawn from a categorical distribution with K distinct fault types. We can assume that this dataset is generated by a random process governed by the joint distribution $p(\mathbf{x}, y)$.

The aim of a neural network-based fault classifier is to approximate the conditional

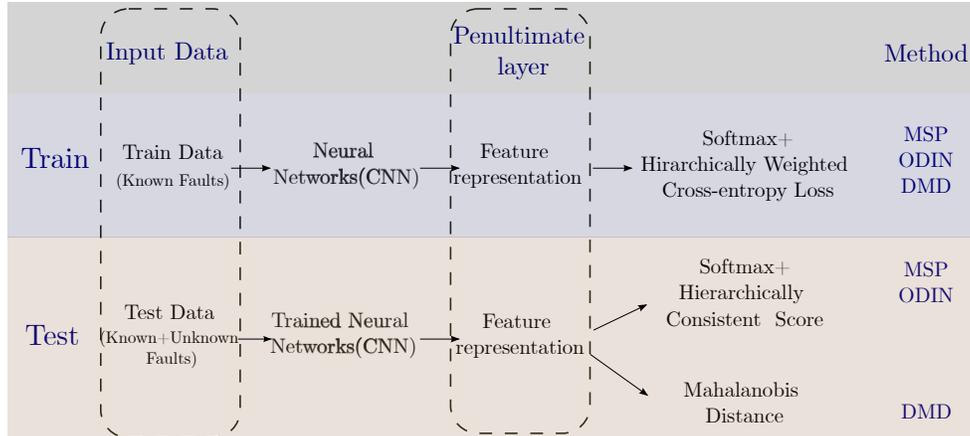


Figure 2: Methodology Flowchart

distribution $p(y|\mathbf{x})$ with a function $g(\mathbf{x}, \boldsymbol{\theta})$ parameterized by neural network weights $\boldsymbol{\theta}$ such that the class with the highest predicted probability \hat{y} will be assigned to the input \mathbf{x} . The function $g(\mathbf{x}, \boldsymbol{\theta})$ represents the penultimate layer of the neural network. Here, the penultimate layer often refers to the feature layer before the Softmax layer. A common practice is to project this representation into the probability simplex by using a Softmax layer. We shall denote this function whose codomain is the probability simplex as $f(\mathbf{x}, \boldsymbol{\theta})$. To simplify notation, $f(\mathbf{x})$ and $g(\mathbf{x})$ will be used instead of $f(\mathbf{x}, \boldsymbol{\theta})$ and $g(\mathbf{x}, \boldsymbol{\theta})$ respectively, which implies that the model has been trained on the training data and the parameters are fixed. The output of the penultimate layer for samples in class k is denoted by the subscript, $g_k(\mathbf{x})$, and the respective probability prediction is $f_k(\mathbf{x})$. In another word, we have $f_k(\mathbf{x}) = \text{Softmax}(g_k(\mathbf{x}))$.

This section will start with the review of three baseline methods in Section 3.1, including the Maximum Softmax Probability (MSP), Out-of-Distribution detector (ODIN), and Deep Mahalanobis Detection (DMD). In Section 3.2, we demonstrate that in the training stage, hierarchical regularization via soft labeling modification in the loss function will be incorporated. In Section 3.3, a novel hierarchically consistent score function for novel fault detection will be illustrated for three baseline methods. Figure 2 illustrates a summary of the methodology in this section.

3.1 Formulation of baseline OOD detection scoring methods used in this study

In this subsection, we will review three state-of-the-art OOD classifiers in three subsections, namely the Maximum Softmax Probability (MSP), Out-of-Distribution detector (ODIN), and Deep Mahalanobis Detection (DMD).

3.1.1 Maximum Softmax Probability

Maximum Softmax probability (MSP) was introduced in (Hendrycks and Gimpel, 2017). The idea of MSP is intuitive, where the in-distribution samples turned out to have a smaller uncertainty or a larger prediction probability for a certain class. On the other hand, out-of-

distribution samples may have a much more uniform distribution for all the known classes. Here, the maximum of the Softmax probability (MSP) can be used as the monitoring statistics, which represent how certain the neural network is about the specific sample \mathbf{x} is $\max_k f_k(\mathbf{x})$. In other words, a sample is detected to be out-of-distribution if its negative MSP score is larger than a certain threshold c as $-\max_k f_k(\mathbf{x}) > c$.

3.1.2 Out-of-Distribution detector for Neural networks (ODIN)

The ODIN method was introduced in (Liang et al., 2018) as an improvement over the MSP method. ODIN is composed of two additional improvements to further improve its OOD Detection accuracy. The first improvement is the temperature scaling, which was originally proposed by Hinton et al. (2015) for knowledge distillation of large neural networks and later used to calibrate the classification uncertainty of the modern neural network in (Guo et al., 2017) to solve the overconfident issues in large deep neural networks. The intuition is that by calibrating the classification uncertainties, the MSP score could be more accurate. The formulation of temperature scaling is defined as follows: for a given input, temperature scaled output of the Softmax for a class k is formulated as below, with the single hyperparameter being the temperature T . The adjusted softmax function for class k can be formulated as follows:

$$f_k(\mathbf{x}; T) = \frac{\exp(g_k(\mathbf{x})/T)}{\sum_{j=1}^K \exp(g_j(\mathbf{x})/T)}. \quad (1)$$

The second improvement to the original MSP method is the adversarial perturbation. Recent research shows that most deep neural networks are not robust to adversarial perturbation (Goodfellow et al., 2015) and highlights the failure of some anomaly detection algorithms caused by adversarial attack (Biehler et al., 2024). ODIN method takes advantage of this property and claims that adversarial perturbation trained for in-distribution samples will be far more effective for the in-distribution samples compared to the out-of-distribution samples, which further increases the MSP score difference for these two groups. The formulation of adversarial perturbation is given as follows: we first add adversarial perturbing into \mathbf{x} , and the perturbed version of input is denoted by $\tilde{\mathbf{x}}$, which is computed as follows:

$$\tilde{\mathbf{x}} = \mathbf{x} - \epsilon \cdot \text{sign}(-\nabla_{\mathbf{x}} \log f_{\hat{y}}(\mathbf{x}; T)), \quad (2)$$

where $\hat{y} = \arg \max_k f_k(\mathbf{x}; T)$ and ϵ is the perturbation magnitude.

Similar to the MSP method, an OOD sample is detected if the negative MSP score is larger than a threshold c , as shown below

$$-\max f_k(\tilde{\mathbf{x}}; T) > c. \quad (3)$$

Throughout the experiments in this study, T and ϵ are fixed to 1000 and 0.0012, respectively, as these are the values are suggested by (Liang et al., 2018).

3.1.3 Deep Mahalanobis Detection

Deep Mahalanobis Detection (DMD) was proposed in (Lee et al., 2018) based on the scaled distance metric (Kim et al., 2011) and it has a different formulation than the aforementioned baseline methods, given it is not based on the Softmax probability. It infers a class-conditional Gaussian distribution over the outputs of the penultimate layer $g(\mathbf{x})$. For

each class k , its respective mean $\hat{\mu}_k$ and the global covariance $\hat{\Sigma}$ for each class is calculated as below:

$$\begin{aligned}\hat{\mu}_k &= \frac{1}{N_k} \sum_{\{n:y_n=k\}} g(\mathbf{x}_n) \\ \hat{\Sigma} &= \frac{1}{N} \sum_k \sum_{\{n:y_n=k\}} (g(\mathbf{x}_n) - \hat{\mu}_k)(g(\mathbf{x}_n) - \hat{\mu}_k)^T\end{aligned}$$

where N_k is the number of samples predicted as class k -based on the MSP function, y_n is the predicted label for sample \mathbf{x}_n .

The squared sample’s Mahalanobis distance to class k can be formulated as:

$$D_k(\mathbf{x}) = (g(\mathbf{x}_n) - \hat{\mu}_k)^T \hat{\Sigma}^{-1} (g(\mathbf{x}_n) - \hat{\mu}_k). \quad (4)$$

Finally, for a given output and threshold c , the Minimum Mahalanobis distance of the feature embedding $g(\mathbf{x}^n)$ and the closest class center $\hat{\mu}_k$ can be used as the OOD statistics and an OOD sample is detected if and only if

$$\min_k D_k(\mathbf{x}) > c. \quad (5)$$

3.2 Hierarchical Regularization in Hierarchical Classification

Typically, the classification objective is to minimize the total cross-entropy loss over the training dataset by optimizing the parameters of the neural network. To construct the cross-entropy loss, we first have to define the one-hot embedding function. Assume a label y_n from the dataset. The one-hot embedding $l_k : k \in \{1 \dots K\} \rightarrow \{0, 1\}^K$ maps a label to a K -length binary vector, where it attains only one ‘1’ on the dimension of the label k and K is the total number of classes. In other words, $l_k(y^n) = 1$ if and only if $k = y_n$ and is zero otherwise.

The cross-entropy objective function over the training dataset \mathcal{D} can be formalized as follows:

$$\min_{\theta} E(\theta) = \min_{\theta} - \sum_{n=1}^N \sum_{k=1}^K l_k(y_n) \log f_k(\mathbf{x}_n, \theta) \quad (6)$$

Note that in this formulation, the model is only penalized for the prediction made for the actual label. The problem is that when it makes a mistake in predicting certain classes, it is completely indifferent to the prediction of similar classes. This is the core challenge in the hierarchical classification literature. In this paper, we will use one of the state-of-the-art techniques from that literature for hierarchical classification, namely the soft label formulation from (Bertinetto et al., 2020) to model the relationship of labels. The soft label is a label representation trick to transform the loss function in a way that is sensitive to the predictions of the other classes, too. Specifically, the new loss formulation forces the prediction mistakes of the model to be hierarchically consistent. A closer look at how the soft label embeddings are constructed is illustrated as follows:

$$l_k^{\text{soft}}(i) = \frac{\exp(-\beta d(k, i))}{\sum_{j=1}^K \exp(-\beta d(j, i))}. \quad (7)$$

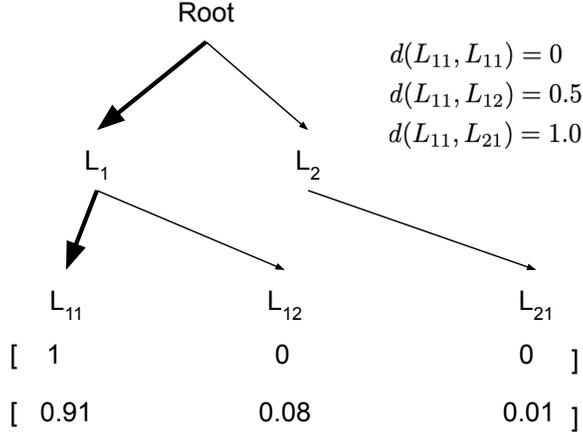


Figure 3: Graphical illustration of the soft labeling logic. A hypothetical two-level fault taxonomy is given. The distances show how the least common ancestor-based distances manifest themselves for this structure. Using these distances and Equation (7), and taking $\beta = 5$, we obtain the soft labels in the second row under the leaf labels, given the real label is L_{11} . For comparison, one-hot labeling is also shown for the same case.

Here β is a constant parameter, and $d(i, j)$ can be any proper distance function that determines how close two labels are in the hierarchy. Here, we choose to use the normalized lowest common ancestor distance, which is a common measure of the distance between two nodes i and j , on a taxonomy tree \mathcal{T} (Bertinetto et al., 2020).

$$d(i, j) = \frac{LCA(i, j)}{h_{\mathcal{T}}}. \quad (8)$$

Here, $LCA(i, j)$ is the lowest common ancestor, which is the lowest ancestor that has both i and j as descendants. Furthermore, to normalize this distance, we can divide the distance by the height of the tree $h_{\mathcal{T}}$, so that the distance is normalized to 1. Note the label probabilities still sum up to one, but the value assigned to the real label is amortized by a fraction, which is redistributed to other classes in a hierarchically consistent manner. See Figure 3 for an illustration of the soft labeling mechanism. For example, in Figure 3, the depth of the tree $h_{\mathcal{T}} = 2$ and $LCA(L_{11}, L_{12}) = 1$. Therefore, $d(L_{11}, L_{12}) = \frac{LCA(L_{11}, L_{12})}{h_{\mathcal{T}}} = 0.5$. Similarly, $LCA(L_{11}, L_{21}) = 2$, therefore, $d(L_{11}, L_{21}) = \frac{LCA(L_{11}, L_{21})}{h_{\mathcal{T}}} = 1$.

Replacing one-hot embedding with soft labels, we obtain the updated loss function as follows:

$$E(\boldsymbol{\theta}) = - \sum_{n=1}^N \sum_{k=1}^K l_k^{\text{soft}}(y_n) \log f_k(\mathbf{x}_n, \boldsymbol{\theta}) \quad (9)$$

Note that with the updated version of the loss function, the prediction function $f(\mathbf{x}, \boldsymbol{\theta})$ is not only encouraged to make the correct prediction but also to do so in a way that reflects the hierarchical relationships between the classes. We call this mechanism hierarchical regularization. The strength of this regularization is determined by the parameter β in Equation (7) where lower values of β assign greater importance to consistency.

3.3 Proposed Methods

In this section, we discuss the proposed methods in detail. We will first demonstrate the hierarchically consistent score function in Section 3.3.1 and demonstrate how it can be combined with three state-of-the-art OOD methods: MSP, ODIN, and DMD. Finally, we would like to discuss the intuition behind the proposed hierarchically consistent score.

3.3.1 Proposed Hierarchically Consistent Score Function

In this section, we will show how we can extend three established methods in the literature using hierarchical regularization and also propose a novel hierarchically consistent score function for anomaly detection. The three methods are Maximum Softmax Probability (MSP) (Hendrycks and Gimpel, 2017), Out-of-Distribution detector for neural networks (ODIN) (Liang et al., 2018) and Deep Mahalanobis Detector (DMD) from (Lee et al., 2018), introduced in Section 3.1.

MSP and ODIN use the negative maximum softmax probability as an out-of-distribution (OOD) score, and detect a sample \mathbf{x} is out-of-distribution if and only if

$$-\max\{f_k(\mathbf{x}, \boldsymbol{\theta}^*) | k \in \{1 \dots K\}\} > c, \quad (10)$$

where $\boldsymbol{\theta}^*$ are the parameters optimized over the training set, and thus, they are fixed during test time, and c is the threshold for OOD detection.

Under the flat classification assumption, the formulation in Equation (10) is reasonable, given that the model assumes that other classes are equally unimportant given the predicted class, and there is no apparent need to incorporate information outside the primarily predicted class. However, we argue that if hierarchical regularization is used for training, the OOD score function should be updated accordingly to reflect hierarchical consistency. The most straightforward formulation to obtain a weighted score would be using the soft labels. Let $\hat{y} = \arg \max_{k \in \{1 \dots K\}} \{f_k(\mathbf{x}, \boldsymbol{\theta}^*)\}$ be the predicted label with the highest prediction score. Here, we detect a sample \mathbf{x} as an unknown fault if and only if

$$-\sum_{k=1}^K l_k^{\text{soft}}(\hat{y}) \log f_k(\mathbf{x}, \boldsymbol{\theta}^*) > c, \quad (11)$$

where $l_k^{\text{soft}}(\hat{y})$ is the soft label defined in Equation (7) and c is the threshold for the OOD detection.

The effectiveness of this formulation depends upon a fundamental assumption, which is the samples generated from the in-distribution process yield hierarchically consistent predictions because they are trained with hierarchical regularization, whereas no such guarantee exists for out-of-distribution samples. Therefore, the proposed OOD score is able to detect OOD samples not only based on the maximum Softmax probability but also on whether the prediction Softmax probability $f_k(\mathbf{x}, \boldsymbol{\theta}^*)$ is hierarchically consistent. In our experiments, we will replace Equation (10) with Equation (11), whenever hierarchical regularization is used during training. Furthermore, it is easy to show that when the predicted label distribution follows the hierarchical structured defined by the soft labels, $f_k(\mathbf{x}, \boldsymbol{\theta}^*) = l_k^{\text{soft}}(\hat{y})$, the OOD score $-\sum_{k=1}^K l_k^{\text{soft}}(\hat{y}) \log f_k(\mathbf{x}, \boldsymbol{\theta}^*)$ can be minimized. Therefore, this detector aims to detect the outlier samples that are hierarchically inconsistent.

3.3.2 Insights to improve upon the state-of-the-art OOD Methods

In this section, we delve into the insights explaining why hierarchical treatment can enhance the performance of state-of-the-art Out-of-Distribution (OOD) methods, including MSP, ODIN, and DMD. A summary of extensions to these base methods, both in the training and monitoring stages, is provided in Table 1. We particularly focus on the improvements realized at the monitoring stage.

Improvement over the flat MSP Method

Maximum Softmax Probability (MSP) employs the maximum of the Softmax Probability, denoted as $\max_k f_k(x)$ for anomaly detection. In contrast, the hierarchical version of MSP utilizes Equation (11) for anomaly detection. As previously mentioned, the hierarchical consistent score aims to leverage the hierarchical structure to enhance the detection score. We elucidate the intuition behind this improvement through Proposition 1.

Proposition 1. *If we define $s(f_k) = -\sum_{k=1}^K l_k^{\text{soft}}(\hat{y}) \log f_k(\mathbf{x}, \boldsymbol{\theta}^*)$ as a function of f_k , where f_k is the output of the Softmax layer, which satisfies $\sum_k f_k = 1$. The hierarchically consistent score $s(f_k)$ will be minimized if and only if $f_k = l_k^{\text{soft}}(\hat{y})$.*

The proof is straightforward and can be seen in Appendix A.1. Proposition 1 demonstrates that the proposed hierarchically consistent score will be minimized if and only if the Softmax output f_k aligns with the soft labels. This is inherently hierarchically consistent, as the soft labels are defined based on the hierarchical tree structure $l_k^{\text{soft}}(\hat{y})$. Therefore, test samples that violate hierarchical consistency will yield in higher scores and will be flagged as unknown faults. The validation of this core assumption for the improvement of the MSP will be discussed in Section 4.4.1.

Improvement over the flat ODIN Method

The ODIN method enhances the original MSP in two significant ways: temperature scaling and adversarial perturbation.

Temperature scaling framework corrects the overconfident Softmax score by incorporating a temperature term in Equation (1). Using the soft label approach, the predicted probability aims to align with the soft label probability $l_k^{\text{soft}}(\hat{y})$.

Proposition 2. *When the temperature T is large, $f_k(\mathbf{x}; T) = \frac{1}{K - \frac{1}{T} \sum_{j \neq k} (g_k(\mathbf{x}) - g_j(\mathbf{x}))}$*

The proof of Proposition 2 is elaborated in Appendix A.2. For in-distribution data, the $f_k(\mathbf{x}; T)$ is more hierarchically consistent, where $g_{\hat{y}}(\mathbf{x})$ is largest and other $g_k(\mathbf{x})$ is smaller. For outliers, $g_{\hat{y}}(\mathbf{x})$ and $g_k(\mathbf{x})$ becomes much closer.

Adversarial Perturbation is added to data \mathbf{x} as shown in Equation (2), and the perturbed version of input is denoted by $\tilde{\mathbf{x}}$. The label is determined by the original data \mathbf{x} as $\hat{y} = \arg \max_k f_k(\mathbf{x}; T)$. However, the hierarchically consistent score is computed through the perturbed data $s(f_k) = -\sum_{k=1}^K l_k^{\text{soft}}(\hat{y}) \log f_k(\tilde{\mathbf{x}})$. We aim to provide intuition on why this perturbation can enhance the OOD performance through the Proposition 3 and the subsequent remarks.

Proposition 3. *The perturbation by $\tilde{\mathbf{x}} = \mathbf{x} - \epsilon \text{sign}(-\nabla_{\mathbf{x}} \log f_{\hat{y}}(\mathbf{x}; T))$ leads to the following Taylor expansion of the hierarchically consistent score $-\sum_k l_k^{\text{soft}}(\hat{y}) \log f_k(\tilde{\mathbf{x}}) = -\sum_k l_k^{\text{soft}}(\hat{y}) \log f_k(\mathbf{x}) + \epsilon U_1 + \epsilon U_2 + O(\epsilon^2)$, where $U_1 = -l_{\hat{y}}^{\text{soft}}(\hat{y}) \|\nabla_{\mathbf{x}} \log f_{\hat{y}}(\mathbf{x})\|_1$ and $U_2 = -\sum_{k \neq \hat{y}} l_k^{\text{soft}}(\hat{y}) \text{sign}(-\nabla_{\mathbf{x}} \log f_{\hat{y}}(\mathbf{x})) \cdot \nabla_{\mathbf{x}} \log f_k(\mathbf{x})$ and $U_2 \geq -\sum_{k \neq \hat{y}} l_k^{\text{soft}}(k) \|\nabla_{\mathbf{x}} \log f_k(\mathbf{x})\|_1$.*

The proof of Proposition 3 is shown in Appendix A.3. Proposition 3 demonstrated that the first-order effect of the perturbation on the proposed hierarchically consistent score can be decomposed into two parts U_1 and U_2 . U_1 illustrates the effect of the label \hat{y} itself, while U_2 represents the effects of other labels k on the label \hat{y} . In general, $U_2 = 0$ if there is no hierarchical structure on the labels. From Proposition 3, we would like to illustrate the following two important observations.

Remark 1. U_1 is typically smaller for in-distribution data compared to outliers.

Remark 1 is demonstrated in the original paper (Liang et al., 2018). For more validation of how U_1 and U_2 differ the abnormal samples from the normal samples, please refer to Appendix A.4 for the validation using the real case study in Section 4. This behavior leads to a better detection performance of ODIN compared to the baseline MSP methods.

Remark 2. U_2 is typically smaller for in-distribution data with hierarchical structures compared to outliers.

Appendix A.4 provides an illustrative example for Remark 2 using our dataset. We will briefly illustrate the insight of Remark 2 here. The reason is that for in-distribution data with hierarchical structures, there exists a set of classes $k \in \mathcal{K}$ that is close to \hat{y} where $l_k^{\text{soft}}(k)$ is large. For these $k \in \mathcal{K}$, $\text{sign}(-\nabla_x \log f_{\hat{y}}(x)) \approx \text{sign}(-\nabla_x \log f_k(x))$, given these two classes are similar (or hierarchically consistent). Therefore, for in-distribution data, U_2 can reach the lower bound $\sum_{k \neq \hat{y}} l_k^{\text{soft}}(k) \|\nabla_x \log f_k(x)\|_1$ much easier. On the other hand, for outliers, $\text{sign}(-\nabla_x \log f_{\hat{y}}(x))$ will be random and can have different signs compared to $\nabla_x \log f_k(x)$. Therefore the inner product $\text{sign}(-\nabla_x \log f_{\hat{y}}(x)) \cdot \nabla_x \log f_k(x)$ will be much smaller or even becomes 0 in U_2 .

In conclusion, Remark 1 and Remark 2 guarantee that the adversarial perturbation will increase the hierarchical consistent score much larger for in-distribution data compared to outliers. This core assumption of the improvement over the ODIN methods will be validated in Section 4.4.2.

Improvement over the flat DMD method

The DMD method does not necessarily require additional treatment for its OOD score formulation. By definition, it fits a class-conditional Gaussian distribution over the entire prediction output of a neural network, not just the maximum prediction score. However, hierarchical consistency will tighten the spread of the conditional distributions of each class, especially if these two classes belong to the same parent. This should result in more conservative anomaly thresholds, thereby reducing Type-I errors. In our experiments, we will employ a simpler variant of this method, fitting distributions only on the outputs of the penultimate layer. For further details, readers are referred to Equation 1 and Equation 2 in (Lee et al., 2018).

The numerical and real-case study validation of these core assumptions for DMD improvement will be discussed in Section 4.4.3.

3.4 Practical Guidelines and Tuning Parameter Selections

3.4.1 Parameter β Selection

The parameter β determines the importance of the hierarchical information in training and testing, and therefore, we want to explore its impact on our methodology. When the value of β is small, the methodology places less emphasis on the knowledge of hierarchical structures.

Table 1: Summary of proposed extensions over base methods, at training and at the monitoring stage.

Base Method	At Training	At Monitoring
MSP	Equation (9)	Equation (11)
ODIN	Equation (9)	Equation (11)
DMD	Equation (9)	—

As $\beta \rightarrow 0$, the soft-label embedding in Equation (7) will be a uniform distribution, which makes the label useless. If β is large, the proposed hierarchically consistent OOD score will ignore any prior assumptions or knowledge about the hierarchical information and becomes the standard Maximum Softmax Probability as illustrated in Proposition 2.

Proposition 4. *When $\beta \rightarrow +\infty$, for \mathbf{x} , the hierarchical consistent score can be used to detect anomalies as:*

$$-\sum_{k=1}^K l_k^{soft}(\hat{y}) \log f_k(\mathbf{x}, \boldsymbol{\theta}^*) > c$$

is equivalent to $\max_k \{f_k(\mathbf{x}, \boldsymbol{\theta}^)\} < c'$, where c' is a positive constant threshold. This implies that the proposed hierarchically consistent OOD score will become the traditional Maximum Softmax Probability (MSP) criterion used in MSP and ODIN methods when $\beta \rightarrow +\infty$.*

The proof is given in Appendix A.5. Proposition 4 shows the connection between the proposed hierarchically consistent score and the MSP. In summary, it reveals that the proposed hierarchically consistency score will become equivalent to the detection score in the traditional MSP method when $\beta \rightarrow +\infty$. In practice, we may need to select a proper value of β to incorporate the hierarchical structure information.

3.4.2 OOD Detection Procedure and Threshold c

In the context of industrial application, our proposed method can be effectively applied by the following steps:

1. In data preparation, obtain in control dataset \mathcal{D} and partition it to two sets: training and validation sets \mathcal{D}_{trn} and \mathcal{D}_{val} .
2. Train the hierarchical classifiers using the training data alone \mathcal{D}_{trn} .
3. Computing the testing statistic for all validation samples and takes its $1 - \alpha$ percentile as c . Remove the out-of-control process and update the percentile until convergence.
4. Compute the test statistics for new samples and start to identify samples as out-of-control if they are above the threshold c

In practice, some evaluation metrics, such as the type I error of the classification results in the validation set can be used to obtain a proper threshold c depending on different industrial applications. For example, we can try to use validation data to get the critical values for different methods via the process shown in Figure 5.

4 Case Study

In this section, we implement the proposed methodology in real word data as described in Section 4.1 and introduce the details of the experimental design and the evaluation criteria in Section 4.2. Finally, the result comparison is presented in Section 4.3.

4.1 Hot Steel Rolling Dataset

Our proposed extensions will be evaluated on a real-life case study image dataset collected from a hot steel rolling process. There are, in total, 3732 image patches that have been detected as potential anomalies obtained by an industrial vision monitoring system. These images have been carefully labeled by domain engineers. Each defect is assigned a two-level label, where there are 8 major categories and 13 subcategories. Please refer to Figure 1 for the full tree of defects hierarchy for this dataset, as well as example images for each leaf category.

The following details of the dataset are provided: 1) Image dimension: The original image sizes for different images could have different dimensions. To make them as the uniform input size 224×224 , we first center crop the figure into an 80×80 square and resize them to the required size; 3) They are all grayscale images; 4) Sample sizes: 18 to 135 samples are available for different defect types. The detailed sample size table is listed in Appendix A.6.

For this study, we randomly partition the data set into train, validation, and test partitions with sizes 60% – 20% – 20% of the original dataset, respectively. The proportions of class prevalence remain the same across all partitions. In other words, we employ a stratified split based on class proportions.

4.2 Experimental Design and Evaluation Criteria

We design an experiment in which we control all possible aspects of a neural network training process except for the hierarchical regularization in Section 3.2 and unknown OOD score functions discussed in Section 3.3. With this design, we aim to answer the question of whether detection performance improves considering hierarchical information of labels.

To simulate the emergence of a new class, we leave out a class from the tree for each experiment. Singling out of each class can be considered a new scenario. There will also be a set of testing samples from known classes to evaluate how the scoring differs between known and unknown classes.

An overview of the steps of each experiment can be listed as follows:

1. Leave out one of the classes as the emerging anomaly class. Only collect the anomaly sample in a single test partition.
2. Train a deep convolutional neural network (CNN) for the classification task with samples in the training set over a set of predetermined hyperparameters.
3. Choose the CNN model with the optimal hyperparameter that produces the best loss in the validation set, freeze its parameters, and proceed with it to the testing stage.
4. Compute the anomaly score function for the validation and testing samples.

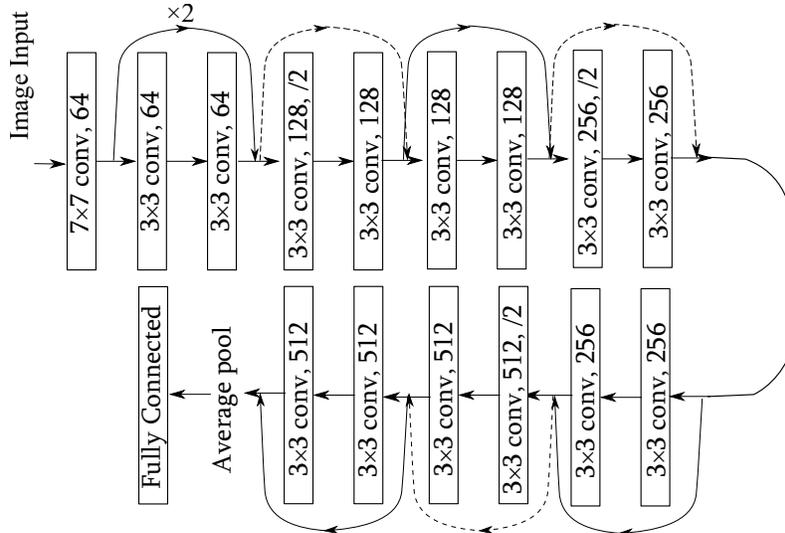


Figure 4: ResNet18 Architecture (He et al., 2016).

5. Report summary evaluation statistics over all test samples.

The experiments are factorized in steps (2) and (4). In step (2), we will compare the training with or without hierarchical regularization. In step (4), we use MSP, ODIN, or DMD as explained in Section 3.2. For MSP and ODIN, the hierarchically consistent OOD scoring will be used in step (4).

In step (1), we replicate the same experiment for four different scenarios. In all scenarios, we remove all instances from one of the classes in the training samples from A12, A31, A61 and A40. From the hierarchical structure tree as shown in Figure 1, there is 1) only one parent node having three child fault types, 2) four parent nodes having two child fault types, 3) three parent nodes having only one child fault type. In the experiment setting, we aim to cover all these three scenarios in picking up our OOD class for testing.

In the first scenario, we select class A12 as the novel fault to observe the behavior of the detectors. This selection can represent the model’s performance on a class with two sibling leave nodes, and we assume that hierarchical algorithms can help improve the detection of this novel fault class (see Figure 1). For the next two scenarios, we select classes A31 and A61 as the novel fault classes. Both classes have only one sibling leave node. In the last scenario, we select class A40 as the novel fault class, as it has no sibling leave node in the structure. Compared to other non-sibling fault classes, the dissimilarity of this fault from all others is relatively distinguishable. To evaluate the efficacy of our proposed method, we have deliberately chosen this diverse range of novel fault classes. Our intention is to ensure that our approach is robust and capable of detecting novel faults in various scenarios. In all of the experiment runs, we utilize the 18-layer variant of the ResNet architecture for the classifier (He et al., 2016). The detailed architecture is given in Figure 4.

As advised by Hendrycks and Gimpel (2017) and employed by most other papers in the literature, we use Area Under the Receiving Operating Characteristics curve (AUROC) given that it doesn’t need to involve the careful selection of the threshold c and type-I error. Here, in literature, the AUROC can be interpreted as the likelihood of a detector to score a random sample from an unknown fault higher than that is randomly picked from known

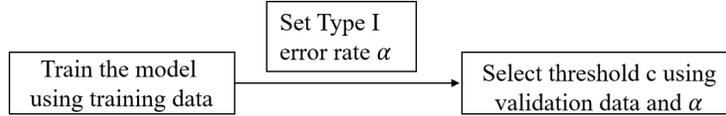


Figure 5: Example flowchart for critical value

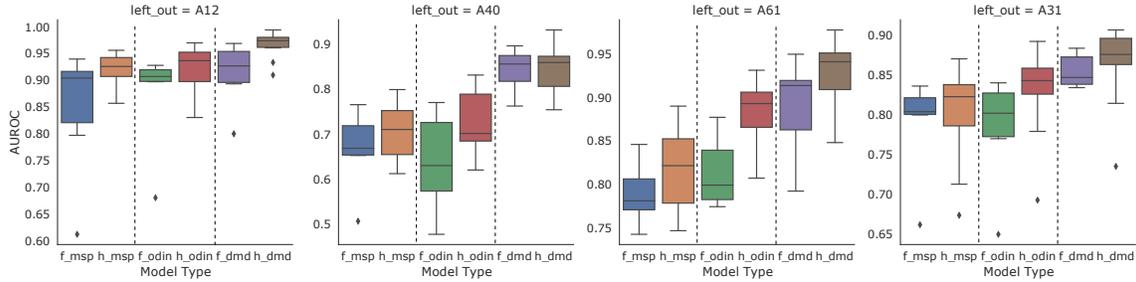


Figure 6: AUROC Comparison: This figure includes all experiments with $\beta = 10$ or $\beta = 100$. Each subplot shows six values on the x-axis, representing the model type and the baseline method used. Specifically, the prefix `f_` denotes `flat model` while `h_` denotes `hierarchical model`. The baseline methods `misp`, `odin`, and `dmd`, which correspond to the MSP, ODIN, and DMD methods reviewed in Section 3, are also represented on the x-axis. To aid in visual clarity, dotted lines are used to separate the results of different baseline methods for each scenario in each subplot.

fault types.

4.3 Results & Discussion

In this section, we will discuss the comparison of the detection performance of the hierarchical method with the flat method for all three aforementioned benchmark methods, including MSP, ODIN, and DMD in all four different scenarios.

4.3.1 Detection Performance Comparisons

In order to compare the performance via Area Under ROC Curve (AUROC) between flat model and hierarchical model clearly, box-plots are shown in Figure 6 given proper β values. Further sensitivity analysis regarding the hyper-parameter β will be investigated in the subsequent subsection. The evaluation dataset of our models encompasses the complete test set, including the known and unknown fault classes, across multiple replications with different learning rates and train seed configurations.

Regarding the model based on MSP, we can refer to the `f_misp` and `h_misp` to find the performance comparison between flat and hierarchical models. We first observe a relatively high novel fault class detection AUC score for ‘A12’ with AUROC larger than 0.9, given that this class is fairly different from the existing classes. For example, ‘A12’ has a black line within the defect region, and ‘A11’ and ‘A10’ have white defective regions. ‘A31’ and ‘A61’ have detection AUROC around 0.8, while ‘A40’ has AUROC about 0.7. Compared to other classes, ‘A40’ has a very small sample size as shown in Table 2 and may affect

the test accuracy robustness. We then observe an increase in both median and highest detection performance in all four left-out scenarios when the hierarchical regularization and consistency are employed at training and testing time, respectively.

Table 2: Sample size summary for left-out classes

Label	Sample Size
A12	75
A31	115
A40	18
A61	75

Regarding ODIN and hierarchical ODIN performance, we can refer to f_{odin} and h_{odin} in Figure 6 as well. First, we have observed an improvement in the baseline detection performance of ODIN over the MSP detector in most classes except class ‘A40’ with 18 samples in testing, which has been demonstrated in (Liang et al., 2018). Second, we have observed a similar performance increase when the hierarchical model structure is considered. This is not a surprising result, given that both ODIN and MSP use the maximum softmax probability. More specifically, the performance improvement of ‘A61’ is the largest when ODIN is used as the base detector compared to MSP, given that the temperature scaling helps reduce the over-confident issue in the neural network and helps achieve a more reasonable hierarchical consistent score.

Finally, we will show the novel fault detection result when DMD is employed as the baseline detection. We also observe results strongly in favor of employing a hierarchical treatment to increase the novel fault detection performance. This reinforces the confidence in the results as DMD employs a rather different strategy at detection time as opposed to MSP and ODIN. As mentioned in Section 3, it is based on the feature representation without a softmax layer. Based on the above observation and overall review of Figure 6, we can find the promising improvement of the hierarchical classifiers in novel fault detection problems even without the re-design of the deep learning architecture.

4.3.2 Sensitivity Analysis

To investigate the sensitivity of the model performance concerning alternations in the soft-label embedding hyper-parameter β , we have conducted a series of experiments in which we varied the values of β and recorded their impact on the model’s performance. For each experiment, β is selected from the set $\beta \in \{0.1, 1, 10, 100\}$. Figure 7, Figure 8, and Figure 9 reflect the hierarchical model performance built on different baseline methods.

First, through an examination of Figure 7, the results generated from the model with $\beta = 0.1$ are the poorest among all scenarios, while models with $\beta \in \{10, 100\}$ perform fairly well in the MSP hierarchical model, except for the A12 scenario. Thus, based on this observation, we can conclude that a very low value of β could negatively impact the hierarchical model performance.

Figure 8 illustrates the sensitivity of ODIN hierarchical models with respect to β . The outcomes of this sensitivity analysis exhibit a closely similar pattern to that observed in the previously discussed results of MSP hierarchical models. Specifically, except for A12, models with $\beta \in \{10, 100\}$ perform fairly well in the ODIN hierarchical model. Additionally, models

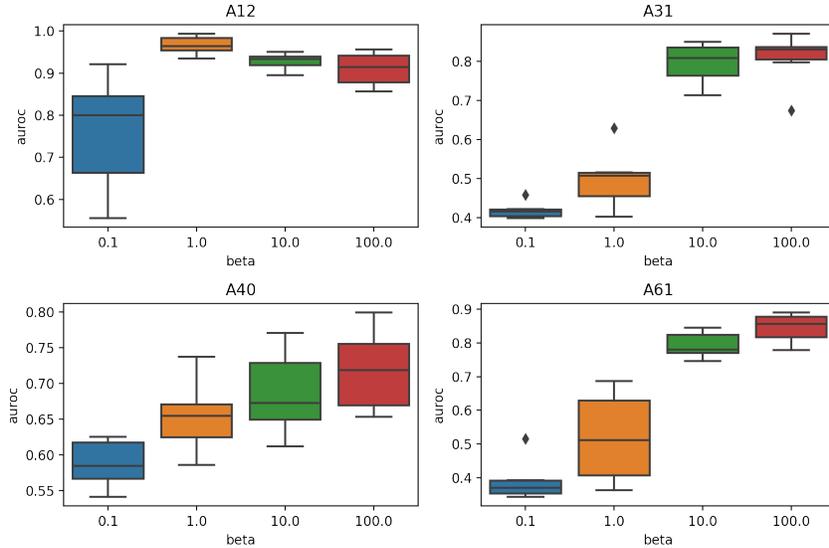


Figure 7: Sensitivity analysis for hierarchical MSP model on different β

with $\beta = 10$ perform very consistently over all classes. Thus, based on this observation, we also can conclude that a very low value of β could negatively impact the hierarchical performance, and β around 10 could be a good option in parameter selection.

Upon examination of Figure 9, we can also observe the poor performance of the hierarchical model with $\beta = 0.1$. Conversely, hierarchical DMD models with $\beta \in \{1, 10, 100\}$ exhibit levels of performance that are comparable to one another, demonstrating low sensitivity across different values of β between $[1, 100]$.

After the examination of these visualizations of the sensitivity analysis over hyperparameter β , we can make a brief conclusion that extremely low values of β are sub-optimal for all the three methods. Instead, selecting a value of β from the set 10, 100 yields relatively stable and satisfactory model performance. Furthermore, among these 4 different selections for β , our experimental results suggest that $\beta = 10$ is associated with the most stable and optimal performance of the proposed hierarchical models.

4.4 Validation of Core Assumptions

In this section, we would like to check the core assumptions that hierarchically consistent scores can improve upon many state-of-the-art OOD methods, including MSP, ODIN, and DMD. The insight into why hierarchically consistent score works have been discussed in Section 3.3.2, which we would like to validate the major assumption that we made.

4.4.1 Validation of Hierarchically consistent score improves upon MSP

To validate the core assumption outlined in Section 3.2, we present Figure 10, which illustrates the relationship between the “distance to prediction” and “the prediction rank”. The term “distribution to prediction” refers to the normalized lowest common ancestor distance

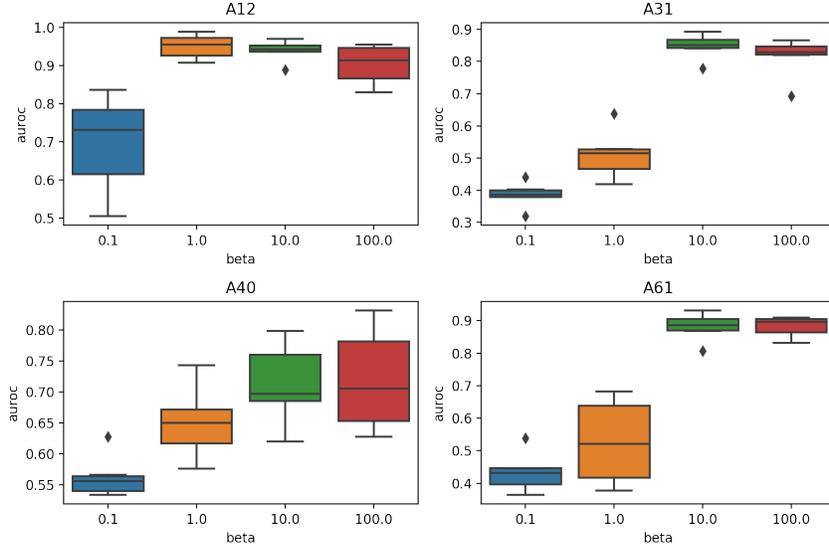


Figure 8: Sensitivity analysis for hierarchical ODIN model on different β

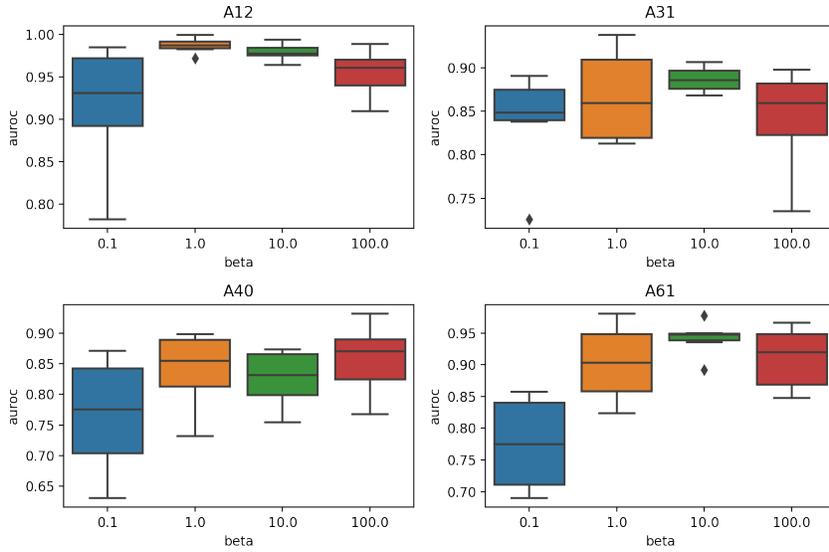


Figure 9: Sensitivity analysis for hierarchical DMD model on different β

of the predicted label, as defined in Equation (8). The “prediction rank” represents the ranked Softmax probability. We start counting from rank 2 on the x-axis, as the “distribution to prediction” is always 0 for the predicted label (i.e., the rank one probability of the Softmax probability); the “distance to prediction” signifies the distance from a given label in the prediction to the true label. The experiments were conducted over 300 iterations and the mean and 95% confidence interval are presented for all four scenarios both for anomaly/novel fault samples (i.e., shown in purple) and normal/known samples (i.e., shown in green).

The objective is to investigate whether a monotonic relationship exists between the

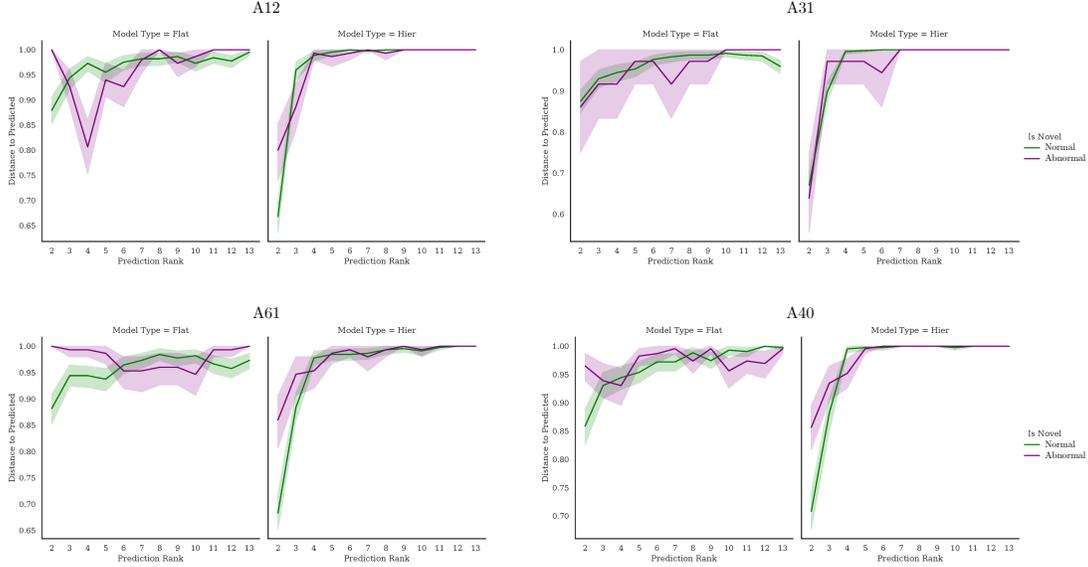


Figure 10: Evolution of distance to predicted by prediction rank. Each subfigure is created by an exemplar run from a scenario where one of the leaf classes is left out. For each scenario, each element in a test is assigned a prediction. The predictions are ranked and their lowest common ancestor-based distance is recorded. The curves depict the evolution of distances with increasing rank with a 95% confidence interval band. The line color denotes whether samples from known(normal) or unknown(abnormal) faults were used to create summary statistics. Figures on the left denoted when hierarchical regularization was not applied while for the figures on the right, it was.

“distance to prediction” and “the prediction rank”. If such a relationship is observed, it would imply that the mean distance to the predicted label (i.e., Rank 1 prediction label) in the hierarchy increases monotonically as the prediction rank decreases. This would further indicate that the prediction is “hierarchically consistent”.

From this study, we can draw several interesting conclusions: 1) In the model without the hierarchical treatment, the prediction is not “hierarchically consistent”. This suggests that a hierarchical structure, if not explicitly enforced, is not inherently present in traditional flat detection methods. For the model with the hierarchical treatment, only the normal samples are hierarchically consistent. However, for the anomaly samples, the monotonicity breaks on a few ranks. 2) For the flat MSP, both the normal and abnormal samples are not clearly separable. However, when hierarchical treatment is applied during model training, the separation between the curves becomes more distinct. This is particularly important for the proposed hierarchically consistent anomaly score function, as it relies on this hierarchical consistency to distinguish between normal and abnormal samples. Specifically, the curve representing anomaly samples diverges from the curve for normal samples at earlier ranks rather than later ones. Given that the hierarchically consistent score function places greater emphasis on earlier ranks, this behavior further ensures the separability of anomaly samples. This is crucial for the success of MSP and ODIN, as both methods depend on this score function for effective anomaly detection.

4.4.2 Validation of Hierarchically Consistent Score Improves upon ODIN

The ODIN model, an enhanced version of the Maximum Softmax Probability (MSP) method, incorporates temperature scaling and perturbation techniques to address issues of overconfidence and improve the separation between normal and abnormal samples. In this study, we extend the ODIN model by incorporating hierarchical structure information and employing a hierarchical consistency score for both training and prediction. Specifically, we examine whether the hierarchical training approach can make ODIN more effective in distinguishing between normal and abnormal instances. In Figure 11 we compare the performance of the newly proposed Hierarchically Trained ODIN model with that of the flat ODIN model, while maintaining consistent hyper-parameter settings. All the scores are standardized using the mean and standard deviation from the normal validation samples.

From Figure 11, the abnormal class scores from hierarchically trained models (hier) validate the enhancement achieved by incorporating hierarchical structure into ODIN. The 'hier' model's abnormal class scores exhibit greater deviation from the standardized normal class scores, centered at zero, when compared to 'flat' model abnormal class scores. This enhancement underscores the utility of the hierarchical approach in bolstering the precision of out-of-distribution detection and showcases its potential to advance the state-of-the-art in anomaly detection models.

4.4.3 Validation of Hierarchically consistent score improves upon DMD

Figure 12 and Figure 13 present the t-Distributed Stochastic Neighbor Embedding (t-SNE) plots (Maaten and Hinton, 2008) of the flat and hierarchical classifiers, respectively. A comparison of these two t-SNE embeddings reveals that the hierarchical model more effectively groups child classes under the same parent class. This suggests that the hierarchical model captures the underlying taxonomy structure more accurately than the baseline flat model does. Additionally, the classes in the hierarchical model are more distinctly separated compared to those in the flat baseline model.

The improvement in performance of the hierarchical model, as evidenced by the AUROC scores, further substantiates that the hierarchically consistent loss function is effective in learning a more hierarchically consistent feature representation. This feature representation is subsequently used to calculate the Mahalanobis distance in DMD.

5 Conclusion

We showed how, when available, hierarchical labels can be used to improve novel fault class detection when deep neural network-based fault classifiers are used. Integration of hierarchical information can be intimidating, especially in the context of neural network training. This is why we employ the most straightforward approach we could find in the literature, the soft labels, as a proof of concept. We extended the existing literature on novel class detection in deep learning to its hierarchical counterpart.

Our claims are experimentally demonstrated on a real-life hot steel rolling defect image dataset. Our experiment design emulates an environment where a fault class has never been

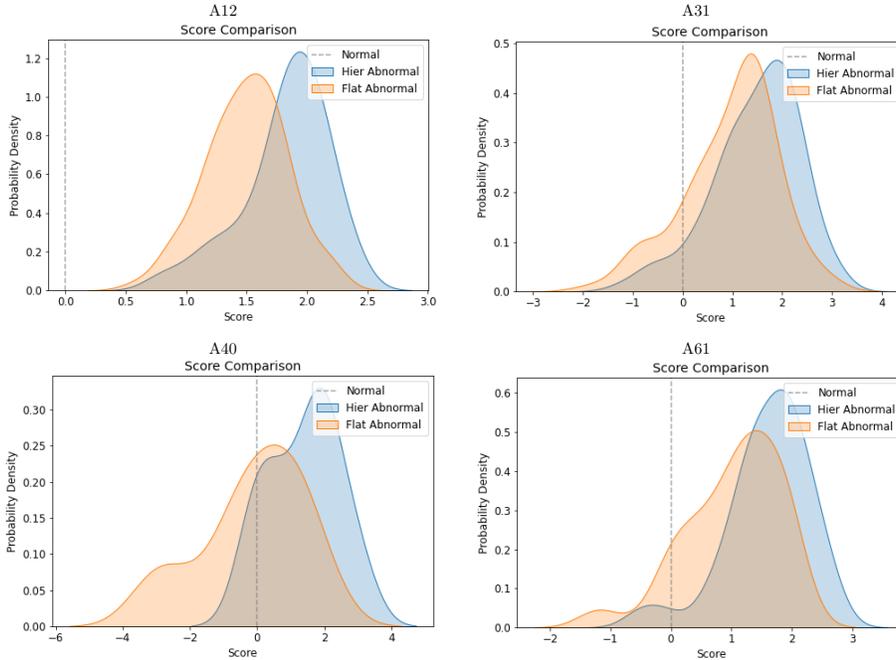


Figure 11: Abnormal class score comparison for flat and hierarchical ODIN: 1) The grey dotted lines in subplots indicate that after standardization, the normal class score is near 0. 2) Subplots with different colors indicate the standardized hierarchical consistency score comparison for 4 different left-out classes. Higher scores signify a greater likelihood of an instance belonging to the unknown/abnormal class.

observed during training time and is introduced at testing time. Accounting for variations in weight initializations and hyperparameter optimizations, we demonstrate that our approach increases the AUROC metric for all possible scenarios and methodology extensions.

In this work, we are only focusing on utilizing the soft-label approaches to encode hierarchical information as well as define the monitoring statistics. There are some other hierarchical models and advanced anomaly detection models, and how to propose proper test statistics for those would be another interesting future work.

Additionally, we are also considering a number of open questions in future work. Whether additional post-training calibration methods can drive the detection performance up higher or not has been answered in this work. Although we deliberately avoided the use of auxiliary datasets in the context of fault classification because of its impracticality in real-life scenarios, we suspect that simulated data may address this issue. Another question we are interested in is whether a partially unknown class (e.g., A12) can be detected along with the correct classification at a coarse level. This is an important question for the practitioner may choose to take different actions based on obliviousness to various levels of the hierarchy.

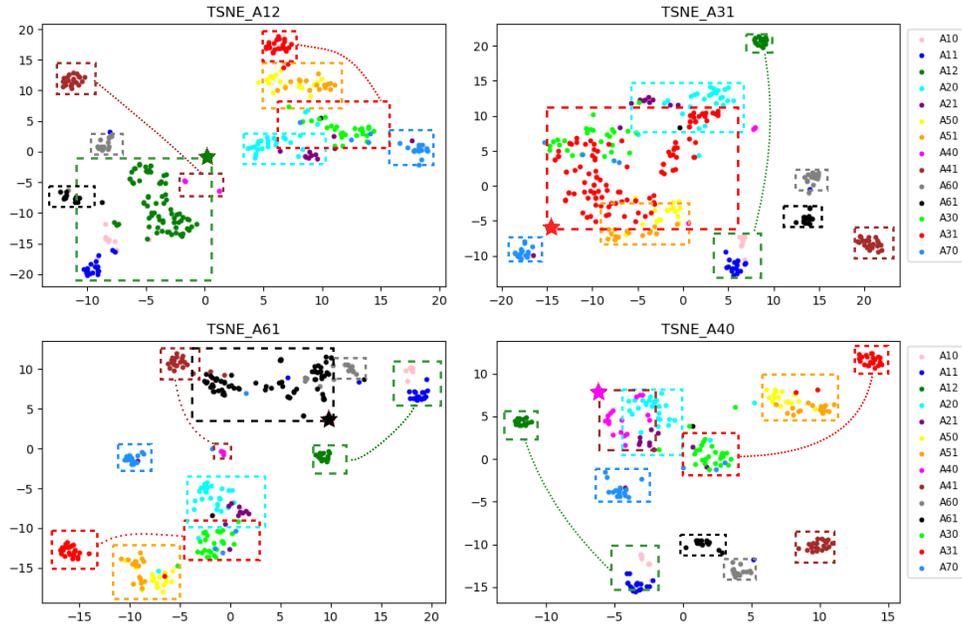


Figure 12: t-SNE plots generated by baseline model: In each subplot, a dashed box roughly squares the dots under the same first-level class (parent class), and the star with the corresponding color denotes the left-out class in that experiment.

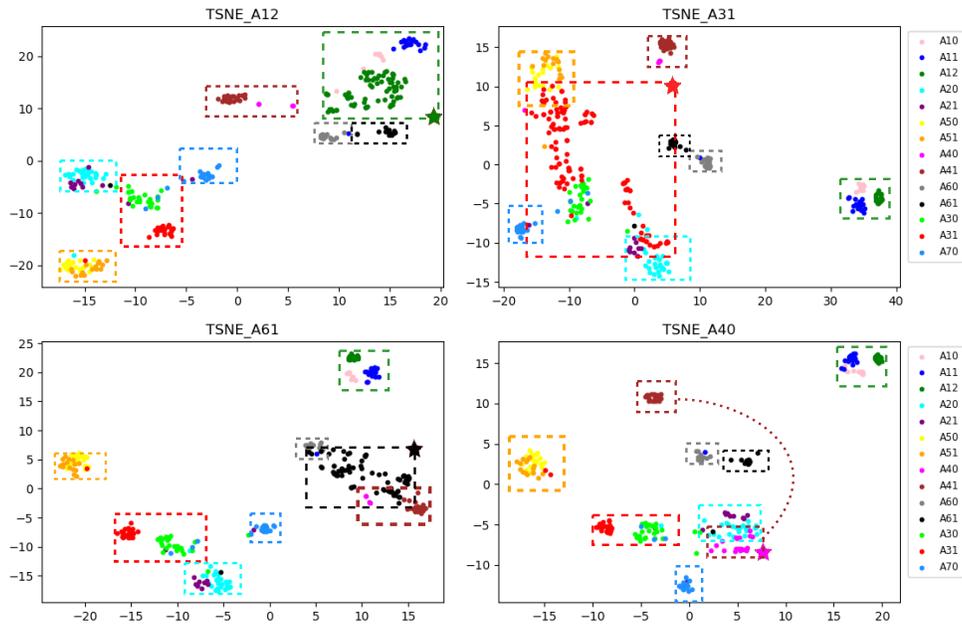


Figure 13: t-SNE plots generated by proposed hierarchical model: In each subplot, a dashed box roughly squares the dots under the same first-level class (parent class). The star with the corresponding color denotes the left-out class in that experiment.

6 Data Availability Statement

Due to the proprietary nature of the collected dataset from the industry, supporting data cannot be made openly available. However, the code will be made available on the Github repo. Further information about the codes will be made available at the Github repo after the paper acceptance.

References

- Ben-David, S., J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan (2010). A theory of learning from different domains. *Machine learning* 79(1), 151–175.
- Bendale, A. and T. E. Boulton (2016). Towards open set deep networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 1563–1572.
- Bertinetto, L., R. Mueller, K. Tertikas, S. Samangooei, and N. A. Lord (2020). Making better mistakes: Leveraging class hierarchies with deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12506–12515.
- Bertinetto, L., R. Müller, K. Tertikas, S. Samangooei, and N. A. Lord (2020). Making better mistakes: Leveraging class hierarchies with deep networks. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pp. 12503–12512.
- Biehler, M., Z. Zhong, and J. Shi (2024). Sage: S tealthy a ttack ge neration in cyber-physical systems. *IJSE Transactions* 56(1), 54–68.
- Breck, E., S. Cai, E. Nielsen, M. Salib, and D. Sculley (2017). The ML test score: A rubric for ML production readiness and technical debt reduction. In J. Nie, Z. Obradovic, T. Suzumura, R. Ghosh, R. Nambiar, C. Wang, H. Zang, R. Baeza-Yates, X. Hu, J. Kepner, A. Cuzzocrea, J. Tang, and M. Toyoda (Eds.), *2017 IEEE International Conference on Big Data, BigData 2017, Boston, MA, USA, December 11-14, 2017*, pp. 1123–1132.
- Brown, I. and C. Mues (2012). An experimental comparison of classification algorithms for imbalanced credit scoring data sets. *Expert Systems with Applications* 39(3), 3446–3453.
- Chiang, L. H., E. L. Russell, and R. D. Braatz (2001). *Fault Detection and Diagnosis in Industrial Systems* (1 ed.). Springer, London.
- Clare, A. and R. D. King (2003). Predicting gene function in *saccharomyces cerevisiae*. *Bioinformatics*, 19 no. suppl_2, ii42–ii49.
- Dumais, S. and H. Chen (2000). Hierarchical classification of web content. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 256–263.
- Fagni, T. and F. Sebastiani (2007). On the selection of negative examples for hierarchical text categorization. In *Proceedings of the 3rd language technology conference*, pp. 24–28.

- Freitas, A. and A. Carvalho (2007). A tutorial on hierarchical classification with applications in bioinformatics. *Research and trends in data mining technologies and applications*, 175–208.
- Gauch Jr, H. G. and R. H. Whittaker (1981). Hierarchical classification of community data. *The Journal of Ecology*, 537–557.
- Goodfellow, I. J., J. Shlens, and C. Szegedy (2015). Explaining and harnessing adversarial examples. In Y. Bengio and Y. LeCun (Eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Guo, C., G. Pleiss, Y. Sun, and K. Q. Weinberger (2017). On calibration of modern neural networks. In *International Conference on Machine Learning*, pp. 1321–1330.
- Guo, X., J. W. Gichoya, S. Purkayastha, and I. Banerjee (2021). Cvad: A generic medical anomaly detector based on cascade vae. *arXiv preprint arXiv:2110.15811*.
- Havtorn, J. D., J. Frellsen, S. Hauberg, and L. Maaløe (2021). Hierarchical vaes know what they don’t know. In *International Conference on Machine Learning*, pp. 4117–4128. PMLR.
- He, K., X. Zhang, S. Ren, and J. Sun (2016). Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 770–778.
- Hendrycks, D. and K. Gimpel (2017). A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*.
- Hendrycks, D., M. Mazeika, and T. G. Dietterich (2019). Deep anomaly detection with outlier exposure. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*.
- Hinton, G., O. Vinyals, and J. Dean (2015). Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Jr., C. N. S. and A. A. Freitas (2011). A survey of hierarchical classification across different application domains. *Data Min. Knowl. Discov.* 22(1-2), 31–72.
- Kendall, A. and Y. Gal (2017). What uncertainties do we need in bayesian deep learning for computer vision? *arXiv preprint arXiv:1703.04977*.
- Kim, S. B., T. Sukchotrat, and S.-K. Park (2011). A nonparametric fault isolation approach through one-class classification algorithms. *IIE Transactions* 43(7), 505–517.
- Kiritchenko, S., S. Matwin, A. F. Famili, et al. (2005). Functional annotation of genes using hierarchical text categorization. In *Proc. of the ACL Workshop on Linking Biological Literature, Ontologies and Databases: Mining Biological Semantics*.

- Koller, D. and M. Sahami (1997). Hierarchically classifying documents using very few words. Technical report, Stanford InfoLab.
- Kwon, D., H. Kim, J. Kim, S. C. Suh, I. Kim, and K. J. Kim (2019). A survey of deep learning-based network anomaly detection. *Cluster Computing* 22(1), 949–961.
- Lee, K., K. Lee, H. Lee, and J. Shin (2018). A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In S. Bengio, H. M. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 7167–7177.
- Li, X., C. Desrosiers, and X. Liu (2022). Deep neural forest for out-of-distribution detection of skin lesion images. *IEEE Journal of Biomedical and Health Informatics*.
- Liang, S., Y. Li, and R. Srikant (2018). Enhancing the reliability of out-of-distribution image detection in neural networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*.
- Liu, R., B. Yang, E. Zio, and X. Chen (2018). Artificial intelligence for fault diagnosis of rotating machinery: A review. *Mechanical Systems and Signal Processing* 108, 33–47.
- Maaten, L. v. d. and G. Hinton (2008). Visualizing Data using t-SNE. *Journal of Machine Learning Research* 9(86), 2579–2605.
- Mirza, M. and S. Osindero (2014). Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*.
- Murray, A. and D. B. Rawat (2021). On the performance of generative adversarial network by limiting mode collapse for malware detection systems. *Sensors* 22(1), 264.
- Nalisnick, E., A. Matsukawa, Y. W. Teh, D. Gorur, and B. Lakshminarayanan (2018). Do deep generative models know what they don’t know? *arXiv preprint arXiv:1810.09136*.
- Pan, J., Y. Zi, J. Chen, Z. Zhou, and B. Wang (2017). Liftingnet: A novel deep learning network with layerwise feature learning from noisy mechanical data for fault classification. *IEEE Transactions on Industrial Electronics* 65(6), 4973–4982.
- Rocchio, J. J. (1971). The smart retrieval system: Experiments in automatic document processing. *Relevance feedback in information retrieval*, 313–323.
- Sahoo, R. K., A. J. Oliner, I. Rish, M. Gupta, J. E. Moreira, S. Ma, R. Vilalta, and A. Sivasubramaniam (2003). Critical event prediction for proactive management in large-scale computer clusters. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 426–435.
- Sallak, M., W. Schön, and F. Aguirre (2013). Reliability assessment for multi-state systems under uncertainties based on the dempster–shafer theory. *IIE Transactions* 45(9), 995–1007.

- Sergin, N. D. and H. Yan (2021). Toward a better monitoring statistic for profile monitoring via variational autoencoders. *Journal of Quality Technology* 53(5), 454–473.
- Shafaei, A., M. Schmidt, and J. J. Little (2019). A less biased evaluation of out-of-distribution sample detectors. In *30th British Machine Vision Conference 2019, BMVC Cardiff, UK, September 9-12, 2019*.
- Shen, D., J. Ruvini, and B. Sarwar (2012). Large-scale item categorization for e-commerce. In X. Chen, G. Lebanon, H. Wang, and M. J. Zaki (Eds.), *21st ACM International Conference on Information and Knowledge Management, CIKM'12, Maui, HI, USA, October 29 - November 02, 2012*, pp. 595–604.
- Shen, Y., K. Yuan, M. Yang, B. Tang, Y. Li, N. Du, and K. Lei (2019). Kmr: knowledge-oriented medicine representation learning for drug–drug interaction and similarity computation. *Journal of cheminformatics* 11, 1–16.
- Silla, C. N. and A. A. Freitas (2011). A survey of hierarchical classification across different application domains. *Data Mining and Knowledge Discovery* 22(1), 31–72.
- Sohn, K., H. Lee, and X. Yan (2015). Learning structured output representation using deep conditional generative models. *Advances in neural information processing systems* 28.
- Xiao, Z., E. Dellandrea, W. Dou, and L. Chen (2007). Hierarchical classification of emotional speech. Technical report, RR-LIRIS-2007-006, LIRIS UMR 5205 CN.
- Yan, H., K. Paynabar, and J. Shi (2014). Image-based process monitoring using low-rank tensor decomposition. *IEEE Transactions on Automation Science and Engineering* 12(1), 216–227.
- Yan, H., H.-M. Yeh, and N. Sergin (2019). Image-based process monitoring via adversarial autoencoder with applications to rolling defect detection. In *2019 IEEE 15th International Conference on Automation Science and Engineering (CASE)*, pp. 311–316. IEEE.
- Yan, Z., H. Zhang, R. Piramuthu, V. Jagadeesh, D. DeCoste, W. Di, and Y. Yu (2015). HD-CNN: hierarchical deep convolutional neural networks for large scale visual recognition. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pp. 2740–2748.
- Zhang, K. (2021). On mode collapse in generative adversarial networks. In *Artificial Neural Networks and Machine Learning–ICANN 2021: 30th International Conference on Artificial Neural Networks, Bratislava, Slovakia, September 14–17, 2021, Proceedings, Part II 30*, pp. 563–574. Springer.
- Zhao, X., H. Yan, and Y. Liu (2021). Hierarchical tree-based sequential event prediction with application in the aviation accident report. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*, pp. 1925–1930.

A Appendix

A.1 Proof of Proposition 1

Proof. The hierarchical score is given as

$$\begin{aligned}
s(f_k) &= - \sum_{k=1}^K l_k^{\text{soft}}(\hat{y}) \log f_k(\mathbf{x}, \boldsymbol{\theta}^*) \\
&\text{Let } l(k) \text{ denote } l_k^{\text{soft}}(\hat{y}), \text{ and } f(k) \text{ denote } f_k(\mathbf{x}, \boldsymbol{\theta}^*) \\
&= - \sum_k l(k) \log f(k) \\
&= - \sum_k l(k) \log \frac{f(k)l(k)}{l(k)} \\
&= - \sum_k l(k) \log \frac{f(k)}{l(k)} - \sum_k l(k) \ln l(k) \\
&= D_{KL}(l \parallel f) + H(l) \geq H(L)
\end{aligned}$$

The equality condition is satisfied if and only if $f(k) = l(k)$.

where $D_{KL}(l \parallel f)$ denotes the KL divergence between l and f , $H(l)$ denotes the entropy of $l(k)$, which is a constant given the specified hierarchical structure and hyper-parameter β . \square

A.2 Proof of Proposition 2

The softmax function with temperature scaling is given by

$$\begin{aligned}
f_k(\mathbf{x}; T) &= \frac{\exp(g_k(\mathbf{x})/T)}{\sum_{j=1}^K \exp(g_j(\mathbf{x})/T)} \\
&= \frac{1}{1 + \sum_{j \neq k} \exp(\frac{g_j(\mathbf{x}) - g_k(\mathbf{x})}{T})} \\
&\approx \frac{1}{1 + (K-1) + \frac{1}{T} \sum_{j \neq k} (g_j(\mathbf{x}) - g_k(\mathbf{x})) + o(\frac{1}{T^2})} \\
&\approx \frac{1}{K + \frac{1}{T} \sum_{j \neq k} (g_j(\mathbf{x}) - g_k(\mathbf{x}))}
\end{aligned}$$

A.3 Proof of Proposition 3

The perturbed version of the hierarchically consistent score is given by

$$\begin{aligned}
-s(f_k) &= \sum_{k=1}^K l_k^{\text{soft}}(\hat{y}) \log f_k(\tilde{\mathbf{x}}) \\
&= l_{\hat{y}}^{\text{soft}}(\hat{y}) \log f_{\hat{y}}(\tilde{\mathbf{x}}) + \sum_{k \neq \hat{y}} l_k^{\text{soft}}(\hat{y}) \log f_k(\tilde{\mathbf{x}})
\end{aligned}$$

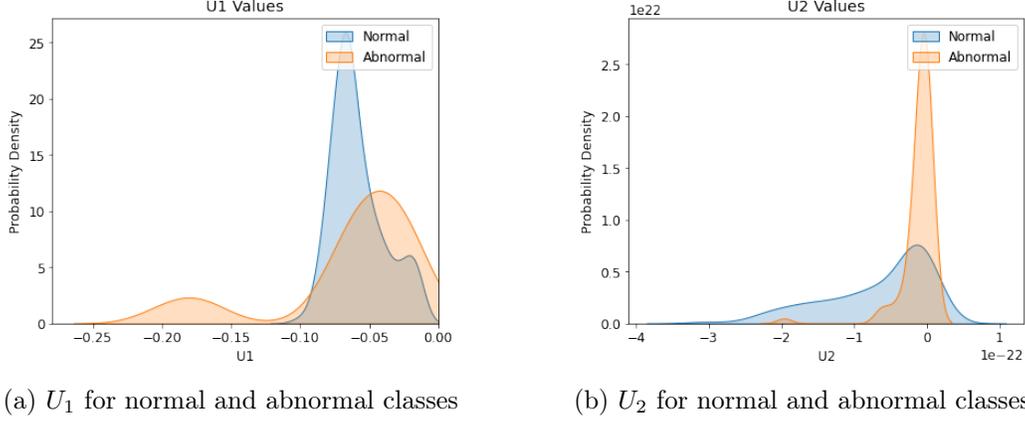


Figure 14: U_1 and U_2 values comparison

$$\begin{aligned}
& l_{\hat{y}}^{\text{soft}}(\hat{y}) \log f_{\hat{y}}(\tilde{\mathbf{x}}) \\
&= l_{\hat{y}}^{\text{soft}}(\hat{y}) \log f_{\hat{y}}(\mathbf{x}) - \epsilon \text{sign}(-\nabla_{\mathbf{x}} \log f_{\hat{y}}(\mathbf{x})) \cdot (\nabla_{\mathbf{x}} \log f_{\hat{y}}(\mathbf{x})) \\
&= l_{\hat{y}}^{\text{soft}}(\hat{y}) \log f_{\hat{y}}(\mathbf{x}) + \underbrace{\|\nabla_{\mathbf{x}} \log f_{\hat{y}}(\mathbf{x})\|_1}_{-U_1}
\end{aligned}$$

Similarly,

$$\begin{aligned}
& \sum_{k \neq \hat{y}} l_k^{\text{soft}}(\hat{y}) \log f_k(\tilde{\mathbf{x}}) \\
&= \sum_{k \neq \hat{y}} l_k^{\text{soft}}(\hat{y}) \log f_k(\mathbf{x}) - \sum_{k \neq \hat{y}} l_k^{\text{soft}}(\hat{y}) \epsilon \text{sign}(-\nabla_{\mathbf{x}} \log f_{\hat{y}}(\mathbf{x})) \cdot (\nabla_{\mathbf{x}} \log f_k(\mathbf{x})) \\
&= \sum_{k \neq \hat{y}} l_k^{\text{soft}}(\hat{y}) \log f_k(\mathbf{x}) - \underbrace{\sum_{k \neq \hat{y}} l_k^{\text{soft}}(\hat{y}) \text{sign}(-\nabla_{\mathbf{x}} \log f_{\hat{y}}(\mathbf{x})) \cdot \nabla_{\mathbf{x}} \log f_k(\mathbf{x})}_{-U_2}
\end{aligned}$$

Also, it is easy to get the lower bound of U_2 :

$$U_2 = \sum_{k \neq \hat{y}} l_k^{\text{soft}}(\hat{y}) \text{sign}(-\nabla_{\mathbf{x}} \log f_{\hat{y}}(\mathbf{x})) \cdot \nabla_{\mathbf{x}} \log f_k(\mathbf{x}) \geq - \sum_{k \neq \hat{y}} l_k^{\text{soft}}(\hat{y}) \|\nabla_{\mathbf{x}} \log f_k(\mathbf{x})\|_1.$$

A.4 U_1 and U_2 values

Figure 14a and Figure 14b present an illustrative example from the dataset in this study to elucidate Remark 1 and Remark 2. The normal class is more likely to exhibit smaller values than the abnormal classes for both U_1 and U_2 as discussed in Remark 1 and Remark 2. This validates the assumptions that combining the hierarchical treatment and the perturbation can indeed separate the abnormal class from the normal class better.

A.5 Proof of Proposition 4

Proof. \hat{y} is defined as $\hat{y} = \arg \max_{k \in \{1 \dots K\}} \{f_k(\mathbf{x}, \boldsymbol{\theta}^*)\}$. The soft label embedding can be considered as the weights $l_k^{\text{soft}}(\hat{y}) = \frac{\exp(-\beta d(k, \hat{y}))}{\sum_j \exp(-\beta d(j, \hat{y}))}$. When $\beta \rightarrow +\infty$, we can get:

$$\lim_{\beta \rightarrow \infty} l_k^{\text{soft}}(\hat{y}) = \lim_{\beta \rightarrow \infty} \frac{\exp(-\beta d(k, \hat{y}))}{\sum_j^K \exp(-\beta d(j, \hat{y}))} = \begin{cases} 1 & \text{if } k = \hat{y} \\ 0 & \text{if } k \neq \hat{y} \end{cases}$$

$$\lim_{\beta \rightarrow \infty} - \sum_{k=1}^K l_k^{\text{soft}}(\hat{y}) \log f_k(\mathbf{x}, \boldsymbol{\theta}^*)$$

$$= -\log f_{k=\hat{y}}(\mathbf{x}, \boldsymbol{\theta}^*) = -\log \max f_k(\mathbf{x}, \boldsymbol{\theta}^*) > c$$

Then, the hierarchical score for novel fault detection can be equivalently written as $\max f_k(\mathbf{x}, \boldsymbol{\theta}^*) < \exp(-c) = c'$. \square

A.6 Sample Size Summary

In our dataset used in the case study, the sample sizes for classes vary from 18 to 135. The detailed sample sizes for all classes are shown in Table 3.

Table 3: Sample size summary for all defects

Label	Sample Size
A10	44
A11	92
A12	75
A20	135
A21	54
A30	127
A31	115
A40	18
A41	105
A50	89
A51	78
A60	72
A61	75
A70	96