

Environmentally Robust Defect Classification with Domain Augmentation Framework

Sungho Lee^a, Jaewoong Shim^{a,*}

^a*Department of Data Science, Seoul National University of Science & Technology, 232 Gongneung-ro, Nowon-gu, Seoul 01811, Republic of Korea*

Abstract

Visual defect classification is a critical process in manufacturing systems, aiming to achieve high-quality production and reduce costs. Recently, deep learning-based defect classification models have achieved significant success. However, the performance of these models can be significantly reduced due to variations in manufacturing environments across multiple production lines. These variations, not present in the training data, result in a domain gap between training and test data. To address this challenge, we propose a domain augmentation framework for constructing a robust defect classification model. This model can deliver high performance across various manufacturing environments using a training dataset from only a single production line. The proposed framework first creates multiple augmented domains using image transformation functions. Then, a defect classification model is trained using a multi-source domain generalization (DG) method with these augmented domains. This approach mitigates the single-source DG problem to a multi-source DG problem, enabling the adoption of multi-source DG methods, which leads to performance improvements. The effectiveness of the proposed framework is demonstrated through experiments using a dataset provided by a Korean manufacturing company.

Keywords: domain generalization, image augmentation, visual defect classification, manufacturing system

1. Introduction

Defect classification is a vital process in manufacturing process that involves not only the detection of defects but also the classification of defect types. This process prevents defective products from being shipped and provides clues to the cause of defects through the analysis of defect types, significantly contributing to product quality assurance, reduction of product recall costs, and enhancement of customer satisfaction. Traditional defect classification involved direct human visual inspection of the produced goods, which incurs considerable time expenditure, decreased inspection accuracy due to the possibility of human error, and a finite volume of goods that can be inspected due to the daily labor time constraints of humans. To overcome these limitations, automatic defect classification was required (Bertolini et al., 2021).

In recent years, considerable research has been conducted to apply machine learning to automate defect classification (Jha & Babiceanu, 2023; Wagner, 2017). A classification model is constructed by learning from a training dataset composed of previously collected images labeled with their defect type. The trained model is then used to

*Corresponding author. Tel.: +82 2 970 6485

Email addresses: sean0310@seoultech.ac.kr (Sungho Lee), jaewoong@seoultech.ac.kr (Jaewoong Shim)

classify new product images into defect types. To utilize traditional machine learning algorithms on defect classification, feature extraction from images using various methods must precede (Park et al., 2016). However, with recent advancements in deep learning, the need for separate feature extraction has been eliminated, and this has demonstrated excellent performance in visual defect classification.

In manufacturing systems, the same product is typically produced through multiple production lines. Each production line has slightly different manufacturing environments, leading to variations in elements such as lighting intensity and type, camera positioning and angles, and so on. As a result, the images of the products captured are subtly different from each other. Moreover, even within the same production line, the imaging environment can change slightly over time. These variations, not present in the training data, can significantly reduce the performance of deep learning-based defect classification models. Therefore, as illustrated in Fig. 1, a defect classification model built on a single production line generally exhibits lower performance when deployed on multiple production lines. Ideally, a separate defect classification model should be built and deployed for each production line, and the model should be frequently retrained. However, due to high costs, it is practically challenging to acquire sufficient labeled training data for each production line to build an individual model. Especially, it becomes even more difficult to acquire labeled training data and build a model when a new production line is established. In real-world manufacturing systems, it is common to obtain a training dataset from a single environment due to various cost constraints. Therefore, there is a need for a robust defect classification model that can deliver high performance in various manufacturing environments, even if it is built only from labeled datasets from a single production line.

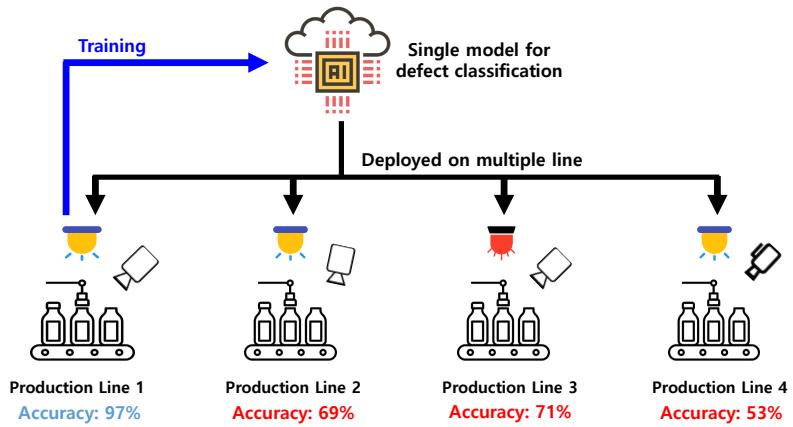


Fig. 1: Performance degradation of defect classification model in different manufacturing environments

To create a robust defect classification model, domain generalization (DG) methods can be utilized, treating each different manufacturing environment as a separate domain (Wang et al., 2022; Azamfar et al., 2020). DG is a research field that aims to overcome the domain gap between training and test data by constructing a domain-invariant model (Zhou et al., 2022; Zheng et al., 2019). One of the most common and basic methods is to augment training images during the training procedure. Although various DG methods have been developed, most of them assume a multi-source DG situation where data from multiple domains are available in the training data. Training a model using data from only a single production line corresponds to a single-source DG situation where data from only a single domain is available. Unfortunately, the single-source DG problem is much more challenging than the

multi-source DG problem, and therefore, the methodologies that can be attempted are limited (Wang et al., 2021; Zhou et al., 2022; Peng et al., 2022).

In this paper, we propose a domain augmentation framework for the construction of a robust defect classification model that can achieve high performance in various manufacturing environments using a training dataset from only a single production line. The proposed framework first creates multiple augmented domains using image transformation functions. Then, a defect classification model is trained using a multi-source DG method with these augmented domains. This simple framework mitigates the single-source situation to a multi-source situation, enabling the adoption of multi-source DG methods, which leads to performance improvements. This is a general framework that can be utilized with any image transformation function and multi-source DG method. The effectiveness of the proposed framework is investigated through experiments using a dataset provided by a Korean manufacturing company.

The main contributions of the proposed framework are summarized as follows.

- The framework can achieve high defect classification performance on new production lines using a dataset from a single existing production line, eliminating the need for additional data collection or model training costs.
- It alleviates the challenging single-source DG problem by converting it into a multi-source DG problem. This allows for the utilization of a wider range of methods, potentially leading to performance improvement.
- The proposed framework is a straightforward yet effective framework that can be utilized with any image transformation function and any multi-source DG method, enhancing its applicability.

The rest of this paper is structured as follows. In section 2, we offer a thorough review of the relevant literature. In section 3, we present our proposed framework. In section 4 and section 5, we outline the experimental setup and discuss the results, respectively. Finally, the conclusions of our research are summarized in section 6.

2. Related work

2.1. Image augmentation

Deep learning technologies have recently achieved remarkable results across a wide range of applications. Particularly, they have seen significant success in the field of computer vision, which has been made possible by the construction of large labeled datasets like ImageNet (Deng et al., 2009). However, creating such large labeled datasets in real-world situations, such as in manufacturing, is very challenging due to various practical constraints. Image augmentation is one of the primary methods to overcome these challenges. By applying various image transformations that can preserve labels to a given dataset, image augmentation effectively increases the amount of training data (Sousa et al., 2020). This alleviates the problem of having a small training dataset and improves the generalization performance of the trained model. In fact, image augmentation has been extensively studied for defect detection and classification in manufacturing systems (Yun et al., 2020; Jain et al., 2022).

Image augmentation methods can be divided into two main types: generative model-based methods and explicit transformation-based methods (Shorten & Khoshgoftaar, 2019). Generative model-based methods involve using another deep learning architecture to learn the distribution of training data and then generate new images based on

the trained model. Representative architectures for generative models in this area include variational autoencoders
75 (VAE) (Kingma & Welling, 2014), generative adversarial networks (GAN) (Goodfellow et al., 2014), and diffusion
models (Ho et al., 2020). These methods have the advantage of providing appropriate transformations automatically
without an explicit function. However, they require high computational cost and a substantial amount of dataset
needed for training the generative model. Moreover, they generate images only within the distribution range of
80 the original dataset, making it difficult to achieve diversity. On the other hand, explicit transformation-based
methods change given training images by applying various existing transformations. This includes 1) geometric
transformations such as rotation and flipping, 2) photometric transformations such as adjusting brightness, contrast,
and saturation, 3) noise injection, and 4) applying filters such as the Gaussian blur (Gedraite & Hadad, 2011) and
Laplacian filters (Paris et al., 2011). These augmentation methods are affordable and can be effective when selective
transformation methods are used based on domain knowledge.

85 Depending on when the transformation is applied to the training images, image augmentation can be categorized
into offline augmentation and online augmentation (Tang et al., 2020). Offline augmentation is a method where
augmentation is applied to the entire training set before model training begins. This method requires additional
storage space and has the disadvantage of reduced diversity compared to online augmentation. Therefore, in most
90 applications, online augmentation is primarily used, where image transformations are applied during model training
dynamically every epoch to greatly enhance the diversity of the data. However, in our proposed domain augmentation
framework, we have adopted the offline augmentation method to create augmented domains with explicit domain
labels.

2.2. Domain generalization

Deep learning models typically assume that the training and test datasets are drawn from the same distribution.
95 However, in real-world applications, mismatches between the distribution of the training dataset and the test dataset
are common. In such situations, the performance of deep learning models trained using conventional methods tends
to decline (Guo et al., 2024; Asutkar et al., 2023). To overcome this, DG method has been studied. The goal of DG
is to build a model that can perform well on test sets from various distributions (target domains), using only the
training data from some specific distributions (source domains) (Zhou et al., 2022; Gulrajani & Lopez-Paz, 2020).

100 DG can be categorized into single-source DG and multi-source DG, based on the number of domains in the
training dataset. Unlike multi-source DG, where the training dataset is composed of multiple domains, single-
source DG obtains the training dataset from only one domain. Naturally, single-source DG is more challenging than
multi-source DG, and therefore, the methodologies studied are relatively scarce (Wang et al., 2021; Zhou et al.,
2022). Nevertheless, in many real-world situations, obtaining a labeled training set from multiple domains is difficult
105 due to data acquisition issues, making single-source DG more practical.

DG methods can be divided into three strategies (Zhou et al., 2022):

- **Augmentation based strategy:** This strategy involves transforming the original image data to simulate
domain shifts (Zhou et al., 2020b,a). For instance, MixStyle (Zhou et al., 2020b) enhances generalization
performance by mixing training images that have the same label but belong to different domains.

- 110 • **Regularization strategy:** This strategy aims to minimize discrepancies across domains by constraining the model's weights and complexity during training (Pan et al., 2018; Sagawa et al., 2019; Krueger et al., 2021). VReX (Krueger et al., 2021) introduces a penalty to the loss function to reduce differences in risk across training domains, thereby reducing the model's sensitivity. GroupDRO (Sagawa et al., 2019) assigns individual weights to the loss for each domain to minimize the worst-case loss over a set of domains. IBN-Net
115 (Pan et al., 2018) employs instance normalization and batch normalization to limit the learning of domain-specific information. Due to their inherent characteristics, both VReX and IBN-Net can be seamlessly applied in a single-source DG situation.
- 120 • **Domain alignment strategy:** The essence of this strategy is to reduce the discrepancy between source domains, aiming to learn domain-invariant information (Sun & Saenko, 2016; Ganin et al., 2016; Long et al., 2015). DANN (Ganin et al., 2016) fosters the training of domain-invariant features by using adversarial training between a domain classifier and a class classifier. Deep CORAL (Sun & Saenko, 2016) trains a nonlinear function that can align correlations of layer activations between domains. DAN (Long et al., 2015) leverages an adaptation layer to extract domain-invariant features.

3. Proposed Framework

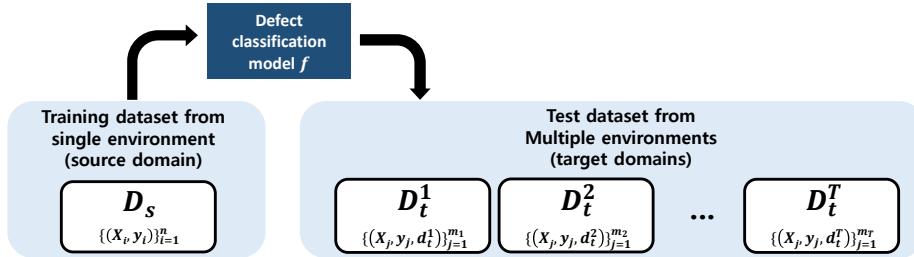


Fig. 2: Problem situation assumed in this study

Notation	Description
X_i	Image of the i -th product
y_i	Label (defect type) of the i -th product
D_s	Training dataset
D_s^k	Dataset of k -th augmented domain
D_t	Test dataset of target domain
d_s^k	Domain label of k -th augmented domain
d_t	Domain label of target domain
n	Number of instances in the source domain
m	Number of instances in the target domain
K	Number of augmented domains
h_k	Image transformation function for k -th augmented domain
f	Defect classification model

Table 1: Notations used in this paper.

125 The objective of this study is to construct a defect classification model that can robustly deliver high performance
 in various manufacturing environments. Each different manufacturing environment can be treated as a domain, and
 DG methods can be utilized to maintain high performance across multiple domains. In real-world manufacturing
 systems, it is common to obtain training dataset from a single environment due to various cost constraints. Therefore,
 we assume that the training dataset consists of only one single domain, but the test dataset consists of multiple
 130 domains. This single-source DG problem situation is depicted in Fig. 2 with a list of notations provided in Table 1. As
 depicted in the figure, a defect classification model is constructed using a training dataset from a single environment,
 and applied to multiple manufacturing environments (multiple domains).

135 The conventional approach in this situation is depicted in Fig. 3(a). The simplest and most widely used method
 is image augmentation, which involves applying various transformations to the training image X during the model
 training process. In addition to this, several single-source DG methods can be used to train the model f . How-
 ever, single-source DG methods are relatively limited in variety compared to multi-source DG methods, making it
 challenging to achieve high performance.

(a) Conventional method



(b) Proposed method

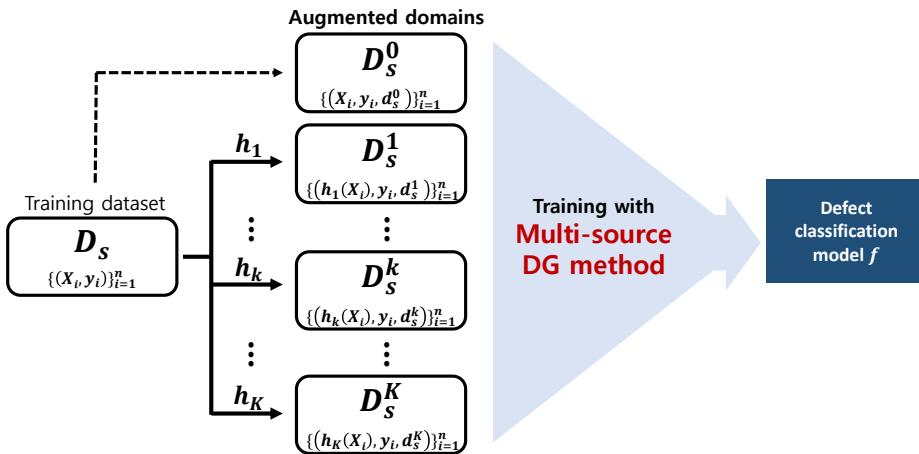


Fig. 3: Proposed domain augmentation framework

We propose a domain augmentation framework for the construction of a robust defect classification model.
 Fig. 3(b) illustrates the procedure of the proposed framework. This framework involves a two-step procedure, which
 140 will be elaborated in the next subsections: 1) Application of image transformation functions to generate multiple
 augmented domains. 2) Training a defect classification model via multi-source DG methods by using the various
 augmented domains. By augmenting the domain with image transformation functions, we can mitigate the single-
 source DG situation to a multi-source DG situation. This allows us to utilize multi-source DG methods, potentially
 leading to performance improvements. This strategy effectively broadens the variety of methodologies available for

145 model training, overcoming the limitations of single-source DG methods. This is a general framework that can be utilized with any image transformation function and multi-source DG method. The proposed domain augmentation framework is described in Algorithm 1.

Algorithm 1 Domain augmentation framework

input: training dataset D_s , Image transformation functions h_1, h_2, \dots, h_K

output: defect classification model f

```

1: procedure
2:    $D_{tr} \leftarrow \{(X_i, y_i, d_s^0) | (X_i, y_i) \in D_s\}$ 
3:   for  $k = 1, \dots, K$  do
4:      $D_{tr} \leftarrow D_{tr} \cup \{(h_k(X_i), y_i, d_s^k) | (X_i, y_i) \in D_s\}$ 
5:   end for
6:   Train  $f$  with  $D_{tr}$  using multi-source DG method
7: end procedure

```

3.1. Domain Augmentation

Given a single domain training dataset $D_s = \{(X_i, y_i)\}_{i=1}^n$, the objective of this step is to create multiple 150 augmented domains for training. The difference from conventional image augmentation methods is that different domain labels d_s^k are assigned for different types of image transformation functions h_k .

Firstly, the original dataset D_s is considered as dataset D_s^0 by assigning all instances with domain label d_s^0 , i.e., $D_s^0 = \{(X_i, y_i, d_s^0)\}_{i=1}^n$. Next, to create datasets for multiple domains, for all instances in D_s , we apply the image transformation function h_k to the image X , preserve the label y , and assign the domain label d_s^k to generate 155 $D_s^k = \{(X_i, y_i, d_s^k)\}_{i=1}^n$. If a total of K image transformation functions are used, K datasets are created. As a result, a total of $K + 1$ datasets with different domain labels, $D_s^0, D_s^1, D_s^2, \dots, D_s^K$, are combined to form a new training dataset D_{tr} , which is used to train the model.

For this domain augmentation step, any general image transformation can be used as h . If there is prior knowledge about various manufacturing environments, image transformation functions related to that can be included and 160 utilized on this step.

3.2. Multi-Source Domain Generalization

The new training dataset D_{tr} , created in subsection 3.1, contains $K + 1$ different domains, with each instance having its corresponding domain label. As the original single-source DG situation has been circumvented to a multi-source DG situation, we now employ a multi-source DG method for training the defect classification model f . For 165 this step, any general multi-source DG method can be utilized. This allows for a wide range of methodologies to be employed, enhancing the potential for achieving a high-performance model.

4. Experimental Design

4.1. Dataset Description

To demonstrate the effectiveness of the proposed framework, we used a dataset provided by a Korean manufacturing company. This dataset comprises images of D-Sub connectors, each with a size of (512, 288), and each image 170

corresponds to one of six classes: OK, Scratch, FM, Pin, Dent, and Glue. A detailed description of each class is provided in Table 2, and representative samples for each class are illustrated in Fig. 4.

class	description
OK	normal product with no defect
Scratch	scraped or dug into by a sharp object
FM	foreign materials in pin hole
Pin	pin bent or deformed
Dent	hollow on black surface
Glue	some sticky on black surface

Table 2: Description of each defect class.

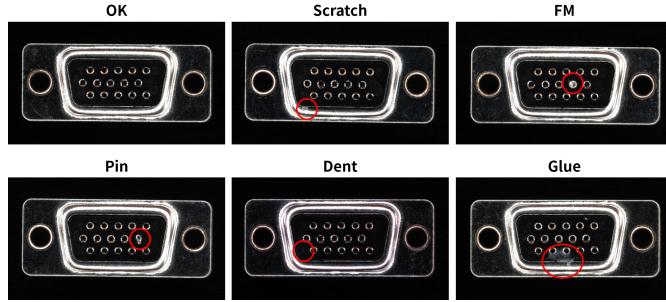


Fig. 4: Representative examples of each defect class. For clarity, the defect area has been indicated with a red circle.

This D-Sub connector dataset contains images taken in a ‘Default’ environment, and it also contains images from various other shooting environments. Specifically, there are three environmental factors that can affect the images:
 175 ‘Color’ (changes due to light scattering), ‘Brightness’ (changes due to illumination levels), ‘Focus’ (changes due to camera focus). Based on a ‘Default’ environment, each factor independently has four configurations depending on its intensity. Representative examples of each configuration are shown in Fig. 5. With the four configurations for each of the three factors, a total of 12 different environments are provided. Along with the ‘Default’ environment, there are a total of 13 environments. These environments are treated as domains in the experiments.

180 This dataset is structured into training, validation, and test sets to align with the problem situation we are considering. The distribution of each class is summarized in Table 3. The training and validation datasets contain the images from only the ‘Default’ environment. In contrast, the test dataset includes the images from all 13 environments. Specifically, 40% of the test set corresponds to ‘Default’, and the remaining 60% is evenly distributed across the 12 environments, for all classes. Therefore, if the model achieves high performance on this test dataset,
 185 it demonstrates its robustness against various environments in real-world manufacturing situations.

4.2. Experimental settings

To demonstrate the effectiveness of the domain augmentation framework, we set four different baseline methods to be compared.

- **Vanilla:** The model was trained based on the cross-entropy loss of the given training dataset, as in a conventional situation that does not consider domain shifts. It does not involve any image augmentation or DG

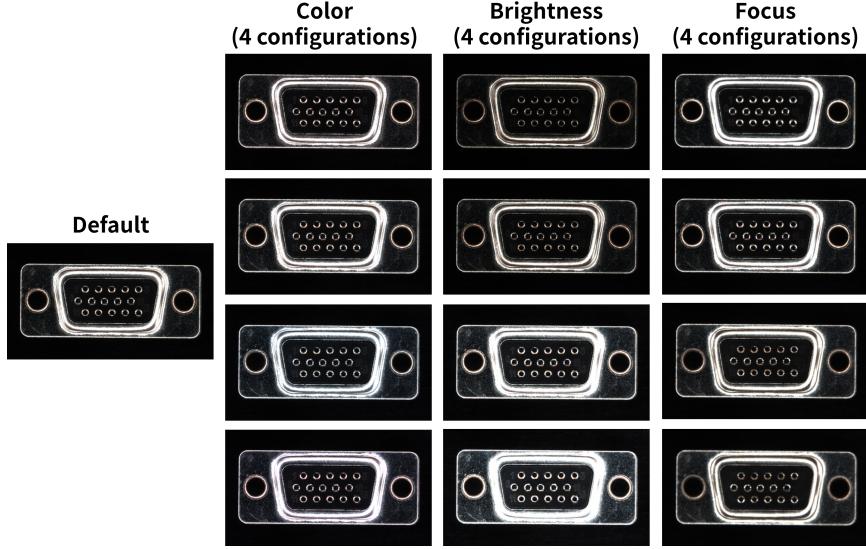


Fig. 5: Representative examples of D-Sub connectors in 13 different environments (domains). Based on the ‘Default’, there are three environmental factors, each with four configurations.

Dataset	Class						Number of Domains
	Ok	Scratch	FM	Pin	Dent	Glue	
Training	160	20	20	20	20	20	260
Validation	40	4	4	4	4	4	60
Test	1000	120	120	120	120	120	1600
							13

Table 3: Distribution of the D-Sub connector dataset.

methods.

- **Offline aug:** The model was trained using only offline image augmentation without any additional DG method. This can be seen as an ablation study from our proposed framework, where only the first step is applied.
- **Online aug:** The model was trained using only online image augmentation without any additional DG method.
- **Single DG:** The model was trained by using a single-source DG method.

195

196

197

198

199

200

For multi-source DG methods to be used for our proposed framework, we adopted ‘Deep CORAL (Sun & Saenko, 2016)’, ‘ERM (Vapnik, 1999)’, ‘GroupDRO (Sagawa et al., 2019)’, IRM (Arjovsky et al., 2019), ‘MixStyle (Zhou et al., 2020b)’, ‘IBN-net (Pan et al., 2018)’, ‘VReX (Krueger et al., 2021)’, and ‘DANN (Ganin et al., 2016)’. For the baseline ‘Single DG’, ‘IBN-net’ and ‘VReX’ methods were utilized since they are also applicable on single-source situation.

As an image transformation function h , we adopted the following eight types of transformation from the open-

source library, ImgAug¹. Each transformation will be used for creating augmented domains.

- 205 • **CoarseDropout:** This transformation sets a certain fraction of pixels in images to zero, with a probability ranging from 0.005 to 0.025%.
- 210 • **AdditiveGaussianNoise:** This transformation adds Gaussian noise $N(0, s)$ to each pixel of the image, where s is varied between 0 and 2.55.
- 215 • **RandomFlip:** This transformation randomly applies a flip to input images either horizontally or vertically.
- 220 • **GaussianBlur:** This transformation applies a Gaussian kernel-based blur to the image, with the blur intensity, defined by σ , varying randomly between 0.2 and 2.
- 225 • **Affine:** This transformation first applies random rotations between -25 and 25 degrees to make an rotated image. Then, it creates final image with alpha-blending, $0.25 \cdot [\text{rotated image}] + 0.75 \cdot [\text{original image}]$.
- 230 • **MultiplyHue:** This transformation operates in the HSV color space, multiplying the H channel pixel values by a random factor between 0.25 and 4. Then, it also adds a random value to each pixel with ‘EnhanceBrightness’.
- 235 • **GammaContrast:** This transformation modifies the contrast of images using a formula $255 \times \left(\frac{v}{255}\right)^\gamma$, where v represents the pixel value and γ is randomly sampled from [0.4, 2].
- 240 • **HomomorphicFilter:** This transformation, called Homomorphic Filtering (Seow & Asari, 2006), works by balancing light and dark areas in an image, making it clearer and more detailed by adjusting light levels and enhancing contrasts.

Representative examples of these transformations are shown in Fig. 6. For the main results presented in subsection 5.1, three methods - ‘MultiplyHue’, ‘Affine’, and ‘RandomFlip’ - were utilized to create augmented domains. These methods were chosen because they demonstrated the highest performance for the proposed framework. A detailed discussion regarding the types and number of transformation functions will be presented in subsection 5.2.

245 In all experiments, ResNet50 was adopted as a backbone network architecture. The SGD optimizer was set with a learning rate of 0.01, and cosine annealing was applied. The batch size was set to 32. Training was conducted for a total of 100 epochs, and early stopping was applied if there was no improvement in validation loss for 4 consecutive epochs. All other settings were consistent with the default settings as specified in the original papers. All experiments are repeated 5 times and the average was reported. We used PyTorch 1.71, running on a GPU of
250 RTX 3090.

5. Results and Discussion

5.1. Main results

Table 4 presents the results of comparing the proposed framework with four other baselines. As shown in the figure, the IBN-net using the proposed domain augmentation framework demonstrated the highest overall accuracy

¹<https://github.com/aleju/imgaug>



Fig. 6: Representative examples of each type of image transformation to be used for creating augmented domains.

of 0.8941 on the test dataset. Moreover, all other methods using the domain augmentation framework also showed relatively high performance, with the average accuracy of the seven methods being 0.8705, which was more effective compared to all other baseline methods. When examining the accuracy for each environment, all methodologies tend to decrease in accuracy when there are changes in ‘Color’, ‘Brightness’, and ‘Focus’ compared to the ‘Default’ environment directly learned from the training set. However, the decrease was smallest when using the domain augmentation framework. This demonstrates that the domain augmentation framework is effective in constructing a robust defect classification model that can handle environmental changes.

Environmental factors	Vanilla	Offline aug	Online aug	Single DG			Domain augmentation framework (proposed)							
				IBN-net	VReX	Avg	Deep Coral	ERM	GroupDRO	IBN-net	MixStyle	VReX	DANN	Avg
Default	87.0±3.14	77.01±2.24	79.01±1.02	88.46±1.96	86.21±2.64	87.34	88.34±0.75	91.23±1.20	91.46±3.67	93.43±2.42	92.11±2.42	84.68±1.94	91.81±2.77	90.44
Color	66.84±4.49	73.14±2.50	72.16±1.47	86.81±2.99	81.46±2.24	84.14	88.16±0.98	89.13±0.63	89.46±3.26	93.19±2.73	93.07±3.49	83.98±1.26	88.17±2.38	89.31
Brightness	63.22±3.12	72.19±3.03	73.44±1.36	76.03±4.15	72.24±7.41	74.14	86.91±1.02	87.5±3.37	87.66±2.87	90.17±2.10	88.18±2.10	81.67±0.89	90.39±2.95	87.5
Focus	52.49±16.75	69.71±2.06	70.67±2.33	85.64±3.19	79.07±3.16	82.16	84.31±1.17	81.64±1.72	81.69±3.19	83.88±3.52	85.41±2.15	80.61±3.72	80.67±19.2	82.6
Overall	73.29±3.66	78.75±1.67	82.1±1.36	84.56±1.41	80.2±1.47	82.38	87.61±1.10	87.36±0.98	87.31±2.79	89.41±2.42	86.31±1.90	84.54±2.04	86.84±2.29	87.05

Table 4: Comparison of accuracy for each environmental factor on test dataset (mean \pm standard deviation)

Fig. 7 presents the confusion matrices for Offline aug, Single DG (IBN-net), and the proposed framework (IBN-net). These matrices represent the best outcomes from multiple experiments conducted for each method. The top part shows the results on the Default environment while the bottom part shows the results on non-Default environments (total of 12 environments). Similar to the results in Table 4, the proposed framework demonstrates relatively robust performance compared to other baselines in the non-Default environment. Notably, while Single DG exhibits competitive performance in the Default environment comparable to the proposed framework, it experiences a more significant performance decrease in non-Default environments.

5.2. Analysis on the augmented domains

In this subsection, we examined how the performance of the domain augmentation framework varies depending on the type of image transformation function h and the number of the functions K . Firstly, we investigated the effect of the type of image transformation function h . To evaluate the effect of each of the eight image transformation

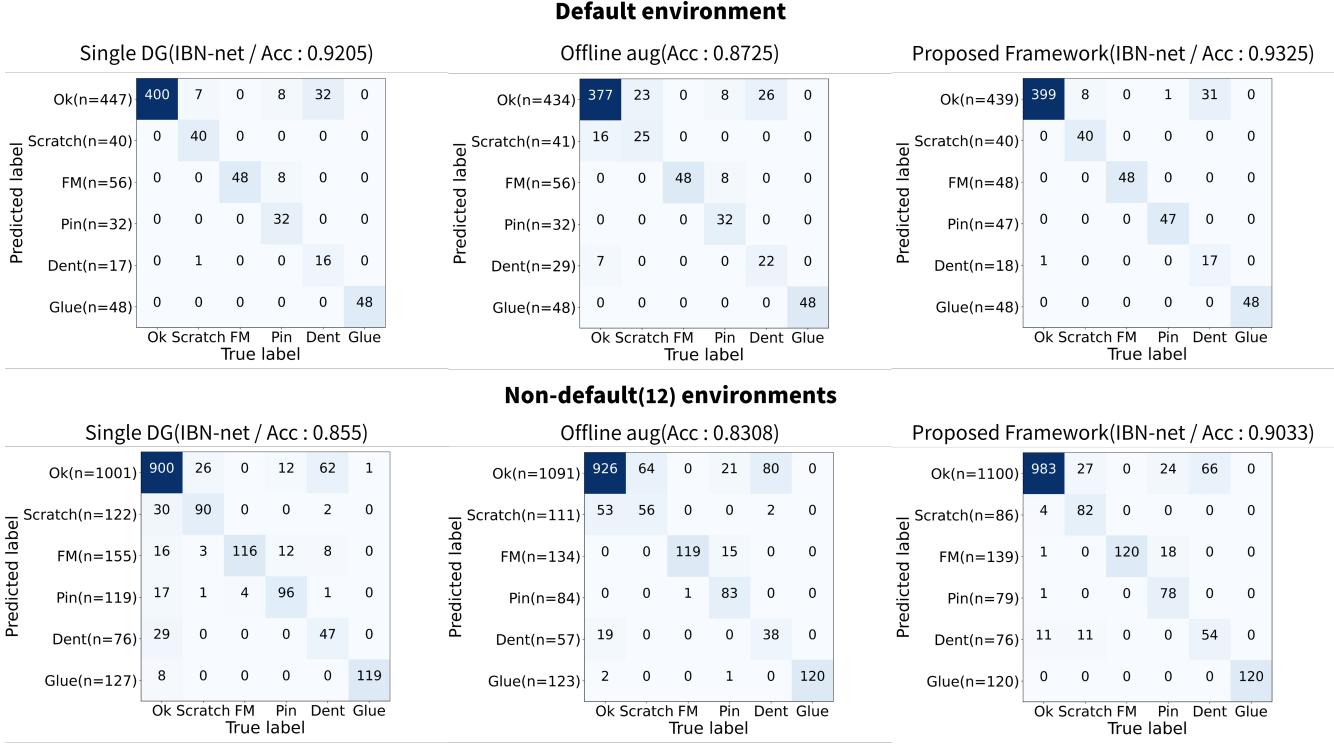


Fig. 7: Confusion matrices for each method on the test dataset: (Top) results on the ‘Default’ environment, which is the same environment as the training dataset. (Bottom) results on environments other than the ‘Default’ environment

functions introduced in subsection 4.2, we trained models with domain augmentation framework where $K = 1$, i.e., creating a single augmented domain with each of the eight image transformation functions.

DG method	Image transformation function							
	Coarse		Additive		Random		Gaussian	Affine
	Dropout	GaussianNoise	Flip	Blur	Hue	Multiply	Gamma	Homomorphic
Deep Coral	80.73±1.43	80.14±2.62	83.66±3.23	80.04±3.80	84.58±1.83	88.41±0.80	85.15±1.95	84.47±2.26
ERM	80.46±3.41	79.53±4.05	84.51±2.68	81.72±1.85	81.45±3.85	89.60±1.82	82.79±1.82	83.19±1.24
GroupDRO	81.14±3.51	74.13±4.43	85.93±1.71	80.90±2.92	85.90±2.10	87.54±3.39	63.77±1.04	83.64±2.01
IBN-Net	77.70±6.53	79.29±2.53	85.68±1.60	78.75±4.69	86.65±2.39	86.67±4.28	84.64±1.68	83.41±2.63
Mix-Style	83.56±1.96	80.49±1.36	85.13±2.5	84.87±3.09	86.43±1.25	87.66±1.97	83.96±1.37	86.29±2.04
VReX	80.65±3.09	71.00±12.19	77.12±5.15	83.08±2.15	79.48±4.14	81.84±4.28	54.06±17.36	76.72±4.19
DANN	81.08±3.32	78.47±5.31	84.11±3.53	81.17±4.90	83.95±2.64	86.87±3.14	83.88±2.31	85.67±1.96
Avg	80.76	77.58	83.73	81.51	84.06	86.94	76.89	83.34
Avg Rank	6.14	7.29	3.43	5.43	3.43	1.14	5.29	3.86

Table 5: Comparison of each image transformation function. In the domain augmentation framework, single augmented domain was added to the original source domain. (mean ± standard deviation)

Based on the results in Table 5, the effectiveness of each transformation function type was ranked as follows: MultiplyHue, Affine, Randomflip, HomomorphicFilter, GammaContrast, GaussianBlur, CoarseDropout, and AdditiveGaussianNoise. This order was established based on the average rank, and in the event of a tie, as was the case with Affine and Randomflip, the average accuracy was used. MultiplyHue showed the highest performance as

a single transformation function in the domain augmentation framework. MultiplyHue can create an augmented domain very similar to the change of the ‘Color’ environmental factor, which appears to have had a positive effect. Therefore, if there is prior knowledge about the actual manufacturing environment, it can be reflected in the image transformation function of the domain augmentation framework to effectively improve performance.

Subsequently, we analyzed the performance changes depending on the number of newly generated domains K , that is, the number of transformation function h . Following the effectiveness order of image transformation functions identified in the previous experiment, we adjusted K by incrementally adding augmented domains. For $K = 1$, only MultiplyHue was utilized. For $K = 2$, Affine was added. For $K = 3$, Randomflip was included. This pattern persisted until $K = 8$, where all eight transformation functions were employed.

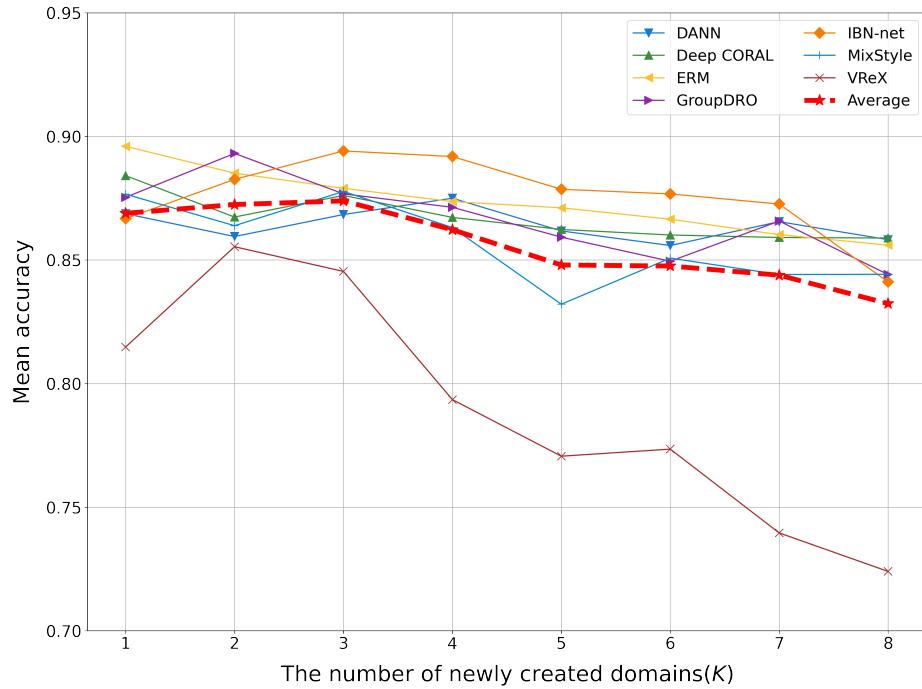


Fig. 8: Performance of the domain augmentation framework as the number of newly created domains K increased.

The result is depicted in Fig. 8. The performance changes for each multi-source DG method are shown as solid lines, and the average is depicted as a red dotted line. The IBN-net, which showed the highest performance, performed best when K was 2, 3, or 4. Observing the average performance trend, there is an initial increase in performance as K grows up to 3, followed by a gradual decline due to the inclusion of less effective transformation functions. However, the extent of this performance variation is relatively minor, indicating that the performance is not significantly impacted by the number of augmented domains employed. Notably, when employing the VReX, there was a considerable fluctuation in performance based on the type and number of transformation functions used. This is likely because the VReX assumption (Krueger et al., 2021)—that variation across training domains is representative of test set variation—does not align well with our situation, where variations between domains are subtle yet diverse.

6. Conclusion

This study primarily focuses on the performance degradation of the defect classification model in diverse manufacturing environments. Aiming for an environmentally-robust defect classification model across diverse production lines, we introduced the domain augmentation framework. In this framework, by augmenting the domain using image transformation functions, the single-source DG problem is mitigated to a multi-source DG problem. Then, more diverse and effective multi-source DG methods become available, overcoming the limitations of single-source DG methods. This is a general framework that can be utilized with any image transformation function and multi-source DG method, so it is capable of incorporating domain knowledge of manufacturing environments. Through experiments using real-world data, we demonstrated that the proposed framework can achieve high defect classification performance across diverse environments.

In future work, we plan to examine the introduction of conditional GAN (He et al., 2016) and diffusion model (Shorten & Khoshgoftaar, 2019) for creating augmented domains to enhance diversity. We anticipate that incorporating these advanced techniques into our domain augmentation framework could further improve performance. Moreover, we aim to develop a test time domain adaptation method to maximize the use of limited information from new manufacturing environments. Specifically, we plan to devise a methodology that assists the model in immediately adapting to a new environment at test time by providing limited information about the new environment in text form.

295 Acknowledgements

This paper was supported by Korea Institute for Advancement of Technology(KIAT) grant funded by the Korea Government (MOTIE) (P0017123, The Competency Development Program for Industry Specialist) and the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT; Ministry of Science and ICT) (No. RS-2022-00165783).

300 References

- Arjovsky, M., Bottou, L., Gulrajani, I., & Lopez-Paz, D. (2019). Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, . doi:<https://doi.org/10.48550/arXiv.1907.02893>.
- Asutkar, S., Chalke, C., Shivgan, K., & Tallur, S. (2023). Tinyml-enabled edge implementation of transfer learning framework for domain generalization in machine fault diagnosis. *Expert Systems with Applications*, 213, 119016.
- Azamfar, M., Li, X., & Lee, J. (2020). Deep learning-based domain adaptation method for fault diagnosis in semiconductor manufacturing. *IEEE transactions on semiconductor manufacturing*, 33, 445–453. doi:<https://doi.org/10.1109/tsm.2020.2995548>.
- Bertolini, M., Mezzogori, D., Neroni, M., & Zammori, F. (2021). Machine learning for industrial applications: A comprehensive literature review. *Expert Systems with Applications*, 175, 114820.

- 310 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 248–255). doi:<https://doi.org/10.1109/cvpr.2009.5206848>.
- Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., & Lempitsky, V. (2016). Domain-adversarial training of neural networks. *The journal of machine learning research*, 17, 2096–2030. doi:<https://doi.org/10.48550/arXiv.1505.07818>.
- 315 Gedraite, E. S., & Hadad, M. (2011). Investigation on the effect of a gaussian blur in image filtering and segmentation. In *Proceedings of the IEEE international symposium on electronics in marine* (pp. 393–396).
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. *Advances in neural information processing systems*, 27. doi:https://doi.org/10.1007/978-3-658-40442-0_9.
- 320 Gulrajani, I., & Lopez-Paz, D. (2020). In search of lost domain generalization. In *Proceedings of international conference on learning representations*. doi:<https://doi.org/10.48550/arXiv.2007.01434>.
- Guo, C., Zhao, Z., Ren, J., Wang, S., Liu, Y., & Chen, X. (2024). Causal explaining guided domain generalization for rotating machinery intelligent fault diagnosis. *Expert Systems with Applications*, 243, 122806.
- 325 He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770–778). doi:<https://doi.org/10.1109/cvpr.2016.90>.
- Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33, 6840–6851. doi:<https://doi.org/10.1109/powertech55446.2023.10202713>.
- 330 Jain, S., Seth, G., Paruthi, A., Soni, U., & Kumar, G. (2022). Synthetic data augmentation for surface defect detection and classification using deep learning. *Journal of intelligent manufacturing*, (pp. 1–14). doi:<https://doi.org/10.1007/s10845-020-01710-x>.
- Jha, S. B., & Babiceanu, R. F. (2023). Deep cnn-based visual defect detection: Survey of current literature. *Computers in Industry*, 148, 103911. doi:<https://doi.org/10.1016/j.compind.2023.103911>.
- 335 Kingma, D., & Welling, M. (2014). Auto-encoding variational bayes international. In *Proceedings of the international conference on learning representations*. doi:<https://doi.org/10.48550/arXiv.1312.6114>.
- Krueger, D., Caballero, E., Jacobsen, J.-H., Zhang, A., Binas, J., Zhang, D., Le Priol, R., & Courville, A. (2021). Out-of-distribution generalization via risk extrapolation (REx). In *Proceedings of the international conference on machine learning* (pp. 5815–5826). PMLR. doi:<https://proceedings.mlr.press/v139/krueger21a.html>.
- 340 Long, M., Cao, Y., Wang, J., & Jordan, M. (2015). Learning transferable features with deep adaptation networks. In *Proceedings of the international conference on machine learning* (pp. 97–105). PMLR. doi:<https://doi.org/10.48550/arXiv.1502.02791>.

- Pan, X., Luo, P., Shi, J., & Tang, X. (2018). Two at once: Enhancing learning and generalization capacities via ibn-net. In *Proceedings of the european conference on computer vision* (pp. 464–479). doi:https://doi.org/10.1007/978-3-030-01225-0_29.
345
- Paris, S., Hasinoff, S. W., & Kautz, J. (2011). Local laplacian filters: Edge-aware image processing with a laplacian pyramid. *ACM transactions on graphics*, 30, 68. doi:<https://dx.doi.org/10.1145/2010324.1964963>.
- Park, J.-K., Kwon, B.-K., Park, J.-H., & Kang, D.-J. (2016). Machine learning-based imaging system for surface defect inspection. *International journal of precision engineering and manufacturing-green technology*, 3, 303–310.
350 doi:<https://doi.org/10.1007/s40684-016-0039-x>.
- Peng, X., Qiao, F., & Zhao, L. (2022). Out-of-domain generalization from a single source: An uncertainty quantification approach. *IEEE transactions on pattern analysis and machine intelligence*, . doi:<https://doi.org/10.1109/tpami.2022.3184598>.
- Sagawa, S., Koh, P. W., Hashimoto, T. B., & Liang, P. (2019). Distributionally robust neural networks. In *Proceedings of the international conference on learning representations*. doi:<https://doi.org/10.48550/arXiv.1911.08731>.
355
- Seow, M.-J., & Asari, V. K. (2006). Ratio rule and homomorphic filter for enhancement of digital colour image. *Neurocomputing*, 69, 954–958. doi:<https://doi.org/10.1016/j.neucom.2005.07.003>.
- Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on image data augmentation for deep learning. *Journal of big data*, 6, 1–48. doi:<https://doi.org/10.1186/s40537-019-0197-0>.
360
- Sousa, M. J., Moutinho, A., & Almeida, M. (2020). Wildfire detection using transfer learning on augmented datasets. *Expert Systems with Applications*, 142, 112975.
- Sun, B., & Saenko, K. (2016). Deep coral: Correlation alignment for deep domain adaptation. In *Proceedings of the european conference on computer vision* (pp. 443–450). Springer. doi:<https://doi.org/10.48550/arXiv.1607.01719>.
365
- Tang, Z., Gao, Y., Karlinsky, L., Sattigeri, P., Feris, R., & Metaxas, D. (2020). Onlineaugment: Online data augmentation with less domain knowledge. In *Proceedings of the european conference on computer vision* (pp. 313–329). Springer. doi:https://doi.org/10.1007/978-3-030-58571-6_19.
- Vapnik, V. (1999). *The nature of statistical learning theory*. Springer science & business media.
- Wagner, W. P. (2017). Trends in expert system development: A longitudinal content analysis of over thirty years of expert system case studies. *Expert systems with applications*, 76, 85–96.
370
- Wang, D., Liu, Q., Wu, D., & Wang, L. (2022). Meta domain generalization for smart manufacturing: Tool wear prediction with small data. *Journal of manufacturing systems*, 62, 441–449. doi:<https://doi.org/10.1016/j.jmsy.2021.12.009>.

- 375 Wang, Z., Luo, Y., Qiu, R., Huang, Z., & Baktashmotlagh, M. (2021). Learning to diversify for single domain
generalization. In *Proceedings of the IEEE/CVF international conference on computer vision* (pp. 834–843).
doi:<https://doi.org/10.1109/iccv48922.2021.00087>.
- Yun, J. P., Shin, W. C., Koo, G., Kim, M. S., Lee, C., & Lee, S. J. (2020). Automated defect inspection system for
metal surfaces based on deep learning and data augmentation. *Journal of manufacturing systems*, 55, 317–324.
380 doi:<https://doi.org/10.1016/j.jmsy.2020.03.009>.
- Zheng, H., Wang, R., Yang, Y., Li, Y., & Xu, M. (2019). Intelligent fault identification based on multisource
domain generalization towards actual diagnosis scenario. *IEEE transactions on industrial electronics*, 67, 1293–
1304. doi:<https://doi.org/10.1109/tie.2019.2898619>.
- Zhou, K., Liu, Z., Qiao, Y., Xiang, T., & Loy, C. C. (2022). Domain generalization: A survey. *IEEE transactions
385 on pattern analysis and machine intelligence*, . doi:<https://doi.org/10.1109/TPAMI.2022.3195549>.
- Zhou, K., Yang, Y., Hospedales, T., & Xiang, T. (2020a). Learning to generate novel domains for domain gen-
eralization. In *Proceedings of the european conference on computer vision* (pp. 561–578). Springer, Cham.
doi:https://doi.org/10.1007/978-3-030-58517-4_33.
- Zhou, K., Yang, Y., Qiao, Y., & Xiang, T. (2020b). Domain generalization with mixstyle. In *Proceedings of the
390 international conference on learning representations*. doi:<https://doi.org/10.48550/arXiv.2104.02008>.