

Image Data Augmentation Approaches: A Comprehensive Survey and Future directions

*Corresponding author(s)

1st Teerath Kumar*
CRT-AI, ADAPT Research Centre
School of Computing,
Dublin City University, Ireland;
teerath.menghwar2@mail.dcu.ie

2nd Alessandra Mileo
INSIGHT & I-Form Research Centre
School of Computing,
Dublin City University, Ireland;
alessandra.mileo@dcu.ie

3rd Rob Brennan
ADAPT Research Centre
School of Computer Science,
University College Dublin, Ireland;
rob.brennan@adaptcentre.ie

4th Malika Bendeache
ADAPT & Lero Research Centres,
School of Computer Science,
University of Galway, Ireland;
malika.bendeache@universityofgalway.ie

Abstract—Deep learning algorithms have demonstrated remarkable performance in various computer vision tasks, however, limited labeled data can lead to overfitting problems, hindering the network’s performance on unseen data. To address this issue, various generalization techniques have been proposed, including dropout, normalization, and advanced data augmentation. Among these techniques, image data augmentation - which increases the dataset size by incorporating sample diversity - has received significant attention in recent times. In this survey, we focus on advanced image data augmentation techniques. We provide an overview of data augmentation, present a novel and comprehensive taxonomy of the reviewed data augmentation techniques, and discuss their strengths and limitations. Furthermore, we provide comprehensive results of the impact of data augmentation on three popular computer vision tasks: image classification, object detection, and semantic segmentation. For results reproducibility, the available codes of all data augmentation techniques have been compiled. Finally, we discuss the challenges and difficulties, as well as possible future directions for the research community. This survey provides several benefits: i) readers will gain a deeper understanding of how data augmentation can help address overfitting problems, ii) researchers will save time searching for comparison results, iii) the codes for the data augmentation techniques are available for result reproducibility, and iv) the discussion of future work will spark interest in the research community.

Index Terms—Computer vision, Data Augmentation, Deep learning, Image classification, Object detection, Semantic segmentation, Survey Data Augmentation

This research was supported by Science Foundation Ireland under grant numbers 18/CRT/6223 (SFI Centre for Research Training in Artificial intelligence), SFI/12/RC/2289/P_2 (Insight SFI Research Centre for Data Analytics), 13/RC/2094/P_2 (Lero SFI Centre for Software) and 13/RC/2106/P_2 (ADAPT SFI Research Centre for AI-Driven Digital Content Technology). For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission.

I. INTRODUCTION & MOTIVATION

Deep learning models have gained popularity and achieved tremendous progress in computer vision (CV) tasks such as image classification [11], [55], [67], [75], [78], [79], [118], [125], object detection [46], [54], , image segmentation [82], [87], [89], [95] and medical imaging [6], [113]–[116], [139]. This advancement has been propelled by various deep neural network architectures, powerful computation resources and extensive availability of data [123]. Convolutional Neural Networks (CNNs) have demonstrated remarkable performance in CV tasks among all deep learning models. CNNs learn different features of an image by applying the convolution operation with the input image and kernel. The initial layers of CNN learn low-level features (e.g. edges, lines) while deeper layers learn more structured and complex features. The success of CNNs has stimulated the interest to use them for CV tasks. In addition to CNNs, Vision Transformers (ViT) [30] are also gaining popularity and have been widely used in deep learning for CV tasks.

However, these algorithms are data-intensive and often suffer from the overfitting problem [119] - where the model performs well on training data but poorly on test data (unseen data). The issue is exacerbated when large amounts of data are not available, which can occur due to privacy concerns or the need for time-consuming and expensive human labeling tasks [79], [123]. Despite the existence of large datasets such as ImageNet [25], overfitting remains a challenge because the standard training process only learns the important regions, but fails to learn less important features that are necessary for generalization [162]. Moreover, adversarial attacks [57], [97], [168] pose a threat to the accuracy of CNNs, where small, invisible perturbations added to the input image can fool the network and cause it to fail to identify the correct features in

an image.

To address these challenges, data augmentation is often applied, not just in CV tasks, but also in a range of domains such as audio [3], [12], [72], [80], [102], [106], [126], [140] and text [7], [39], [93], [124]. This survey will specifically focus on the CV domain.

Regularization is an effective method for generalizing Convolutional Neural Network (CNN) models from both architectural and data perspectives. Various forms of regularization have been developed, including Data Augmentation [170], Dropout [132], Batch Normalization [63], Transfer Learning [122], [153], and Pre-training [32]. Among these, Image Data Augmentation [170] has proven to be a useful form of regularization in several studies [76], [125], [170]. This technique expands the dataset by altering the sample's appearance or flavor [170] to provide a more diverse range of views. However, performing Image Data Augmentation directly on image data can increase the risk of overfitting and biases, making it both important and challenging, as discussed further in Section IV.

Generally, Image Data Augmentation addresses two primary problems in CNN models. The first problem is a shortage of data or limited data, which can result in overfitting. Image Data Augmentation provides a solution by feeding the model with various scenarios of an image, making the model more generalized and allowing for the extraction of more information from the original dataset. The second problem relates to labeling, where the original dataset has a label for each sample. Augmenting the sample preserves the label of the original sample and assigns it to the augmented sample.

Numerous surveys have been conducted on the topic of image data augmentation. For example, Wang et al. explored and compared several traditional data augmentation techniques in their work [109], but this study was limited to image classification tasks only. In another study, Wang et al. reviewed the available data augmentation approaches for face recognition [146]. Khosla et al. briefly discussed warping and oversampling-based data augmentation approaches in their work [68]. However, the authors did not provide a comprehensive taxonomy or a thorough evaluation of the techniques they discussed. Shorten et al. presented a comprehensive survey on image data augmentation in their work [123]. The authors proposed a novel taxonomy, discussed future directions, and addressed the challenges associated with data augmentation. However, the survey lacked an evaluation of image data augmentation for various computer vision tasks. Additionally, as the study is three years old, it may not include the latest state-of-the-art augmentation methods such as cutmix and grid mask. . Recently, Yang et al. conducted a survey on data augmentation in computer vision tasks [158]. However, their study only covered a few data augmentation methods and did not provide any code compilation for result reproducibility. Another study by Xu proposed a novel taxonomy for image data augmentations [157], but did not evaluate the techniques discussed. This paper presents an extended taxonomy for data augmentation and reviews state-of-the-art techniques. The

source code used in this study is available for result reproducibility. It should be noted that this survey does not cover data augmentations based on generative adversarial networks (GANs) due to out of the scope of this paper. But we redirect the reader to [133], [161] for more details about GAN-based data augmentations.

The followings are our contributions:

- A comprehensive image data augmentation taxonomy is presented.
- An extensive survey of state-of-the-art data augmentation techniques, complete with visual examples, is provided.
- The performance of state-of-the-art data augmentation techniques is evaluated and compared for several computer vision tasks.
- The challenges of data augmentation are highlighted and future directions are identified.
- The available codes for data augmentations, following the proposed taxonomy, are compiled for result reproducibility and made available at ³.

The above contributions provide the following benefits:

- A better understanding of data augmentation working mechanism to fix the overfitting problem.
- Our comprehensive analysis and comparison between the existing data augmentation techniques will save researchers time searching this field.
- Facilitates result reproducibility by providing the source code for the different data augmentation techniques investigated.
- Future work will spark interest in the research community.

II. TAXONOMY AND BACKGROUND

The proposed taxonomy, presented in Figure 2, classifies data augmentation into two main branches: Basic and Advanced data augmentations. The former encompasses fundamental techniques for data augmentation, while the latter encompasses more complex techniques. The specifics of each data augmentation method are thoroughly discussed in subsequent sections.

A. Basic Image Data Augmentations

This section describes basic image data augmentation methods and their classifications. They are classified as below:

- **Image Manipulation**
 - *Geometric Manipulation*
 - *Non-Geometric Manipulation*
- **Image Erasing**
 - *Erasing*

1) **Image Manipulation:** Image manipulation refers to the changes made in an image with respect to its position or color. Positional manipulation is made by adjusting the position of the pixels while color manipulations are made by altering the pixel values of the image. Image manipulation is further

³<https://github.com/kmr2017/Advanced-Data-augmentation-codes>

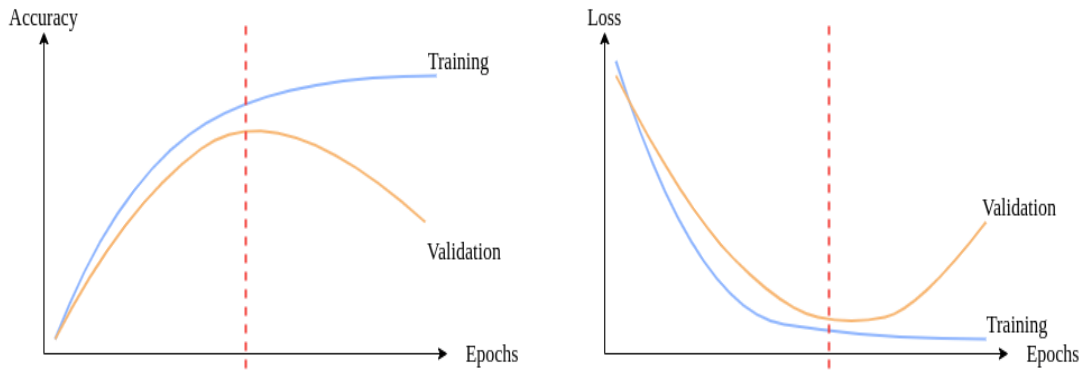


Fig. 1. Overfitting problem: On the left side, overfitting is explained in terms of accuracy, after the inflation point (red dotted line), the training accuracy is increasing but validation accuracy is decreasing. On the right side, alternatively in terms of loss, training loss is decreasing but validation loss is increasing after the red dotted line. The figure is taken from the source ²<https://www.baeldung.com/cs/ml-underfitting-overfitting>

divided into two main categories. Each of them is discussed below:

Geometric Data Augmentation: Geometric data augmentation encompasses modifications to the geometric attributes of an image, including its position, orientation, and aspect ratio. This technique involves transforming the arrangement of pixels within an image through a variety of techniques such as rotation, translation, and shearing. Figure 3 illustrates the most commonly employed geometric augmentations. These methods are widely used in the domain of computer vision to diversify the training data and improve the resilience of models to diverse transformations. The utilization of geometric data augmentation has become a critical component in the development of robust computer vision algorithms. Each of the geometric data augmentations is discussed below:

- (i) **Rotation:** Rotation data augmentation involves rotating an image by a specified angle within the range of 0 to 360 degrees. The precise degree of rotation is a hyperparameter that requires careful consideration based on the nature and characteristics of the dataset. For instance, in the MNIST [27] dataset, rotating all digits by 180 degrees, transforming a right-rotated 6 results into a 9, would not be a meaningful transformation. Therefore, a thorough understanding of the dataset is necessary to determine the optimal degree of rotation and achieve the best results.
- (ii) **Translation:** Translation data augmentation involves shifting an image in any of the upward, downward, right, or left directions, as illustrated in Figure 3, in order to provide a more diverse representation of the data. The magnitude of this type of augmentation must be selected with caution, as an excessive shift can result in a substantial change in the appearance of the image. For example, translating a digit 8 to the left by half the width of the image could result in an augmented image that resembles the digit 3. Hence, it is imperative to consider the nature of the dataset when determining the magnitude of the translation augmentation to ensure its efficacy.

- (iii) **Shearing:** Shearing data augmentation involves shifting one part of an image in one direction, while the other part is shifted in the opposite direction. This technique can provide a new and diverse perspective on the data, thereby improving the robustness of a model. However, excessive shearing can cause significant deformation of the image, making it difficult for the model to accurately recognize the objects within it. It is therefore important to consider the amount of shearing applied to the data carefully in order to avoid over-augmenting the images and introducing unwanted noise. In this way, shearing can be a powerful tool for enhancing the generalization ability of computer vision models, while avoiding the potential drawbacks of over-augmentation. For example, applying excessive shearing on cat image during data augmentation may result in a distorted, stretched appearance, hindering the ability of a model to correctly classify the image as a cat. It is crucial to find a balance between the amount of shearing applied and the desired level of diversity, as excessive shearing can introduce significant noise.

Non-Geometric Data Augmentations The non-geometric data augmentation category focuses on modifications to the visual characteristics of an image, as opposed to its geometric shape. This includes techniques such as noise injection, flipping, cropping, resizing, and color space manipulation, as illustrated in Figure 4. These techniques can help improve the generalization performance of a model by exposing it to a wider variety of image variations during training. However, it is important to consider the trade-off between augmenting the data and preserving the integrity of the underlying information in the image. The following section outlines several classical non-geometric data augmentation approaches.

- (i) **Flipping:** Flipping is a type of image data augmentation technique that involves flipping an image either horizontally or vertically. The efficacy of this method has been demonstrated on various widely-used datasets, including cifar10 and cifar100 [74]. However, care must be taken

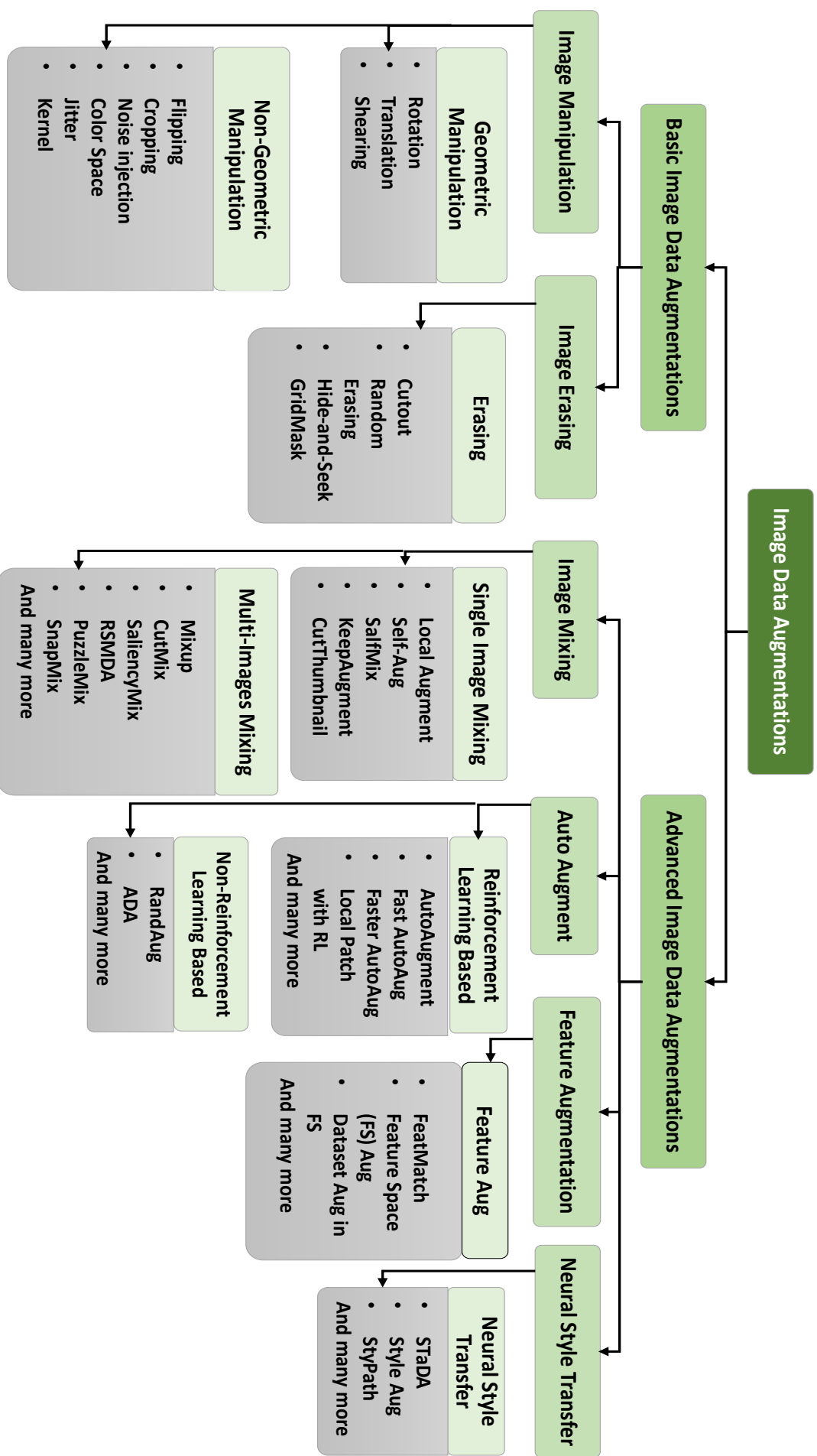


Fig. 2. Image data augmentation taxonomy. Note: All image data augmentation names are not added in this taxonomy due space limit. However, all relevant and remaining image data augmentations are discussed as per taxonomy. The remaining sub-type of categories are discussed in the text.

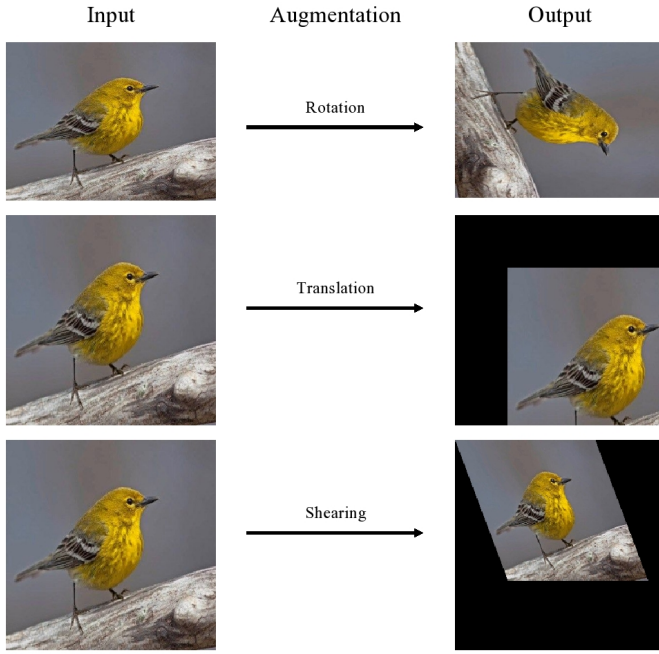


Fig. 3. Overview of the geometric data augmentations.

when applying this technique, as the outcome may depend on the nature of the dataset. For instance, the horizontal flipping of the digit "2" in the Urdu digits dataset [67] may result in the appearance of the digit "6". As such, the choice of flipping must be made carefully to ensure that the desired level of augmentation is achieved without introducing significant noise into the data.

- (ii) **Cropping and resizing:** Cropping is a common pre-processing data augmentation technique that can be applied randomly or to the center of the image. This technique involves trimming the image and then resizing it back to its original size, preserving the original label of the image. However, caution must be exercised when using cropping as a data augmentation method, as it may result in misleading information for the model, such as cropping the upper or lower part of the digit "8" and making it appear as the digit "0".
- (iii) **Noise Injection:** Noise injection is a data augmentation technique that has been demonstrated to enhance the robustness of neural networks in learning features and defending against adversarial attacks. As shown in the survey of nine datasets from the UCI repository [123], the use of noise injection has resulted in impressive performance improvements.
- (iv) **Color Space:** The manipulation of individual channel values within an image, also known as photometric augmentation, is a type of data augmentation that can help control the brightness of the image. Image data typically consists of three channels: Red (R), Green (G), and Blue (B) and has dimensions of Height (H) x Width (W) x

Channels (C). By altering the values of each channel separately, this technique can prevent a model from becoming biased towards specific lighting conditions. The most straightforward approach to perform color space augmentation involves replacing a single channel within the image with a randomly generated channel of the same size, or with a channel filled with either 0 or 255. The utilization of color space manipulation is commonly observed in photo editing applications, where it is used to adjust the brightness or darkness of the image [123].

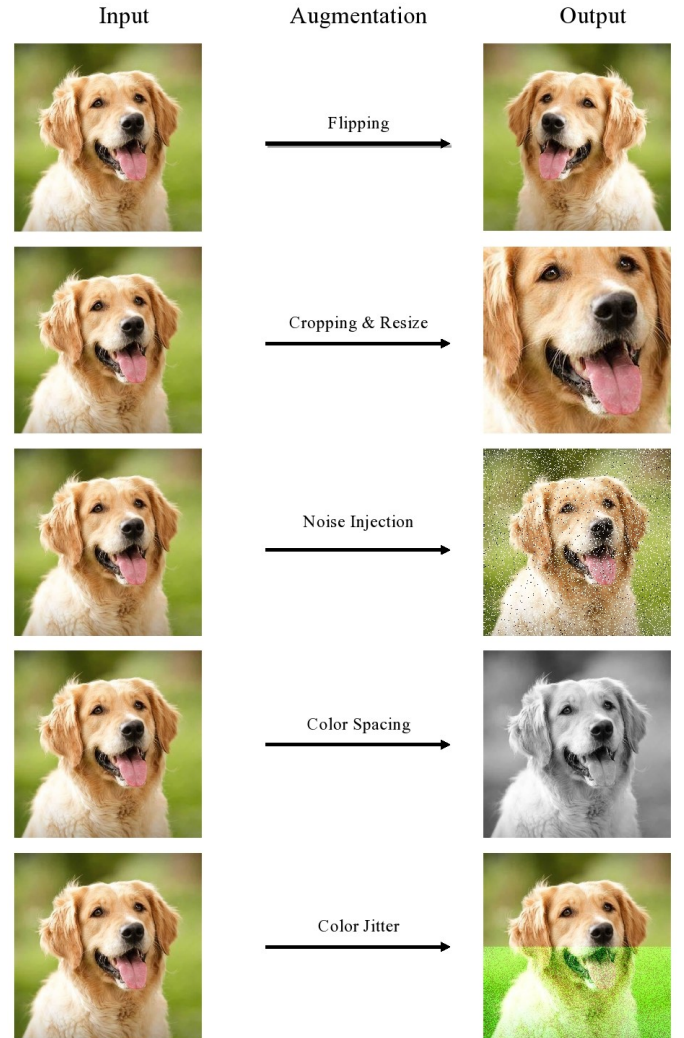


Fig. 4. Overview of the non-geometric data augmentations.

- (v) **Jitter:** Jitter is a technique of data augmentation that involves randomly altering the brightness, contrast, saturation, and hue of an image. The four hyperparameters, i.e., brightness, contrast, saturation, and hue, can be adjusted by specifying their minimum and maximum range. However, it is important to carefully select these ranges as improper adjustments can negatively impact the image's content. For example, increasing the brightness of X-Ray images used for lung disease detection can result

in the whitening and blending of the lungs in the X-Ray, hindering the diagnosis of the disease.

- (vi) **Kernel Filter:** Kernel filtering is a form of data augmentation that enhances or softens the image. This is achieved by applying a window, with a specified size of $n \times n$, containing a Gaussian-blur or an edge filter to the image. The Gaussian-blur filter serves to soften the image, while the edge filter sharpens its edges either horizontally or vertically.

2) **Image Erasing Data Augmentations:** The data augmentation technique of image erasing involves the process of removing specific parts of an image and replacing them with either 0, 255, or the mean of the entire dataset. This type of data augmentation includes various methods such as cutout, random erasing, hide-and-seek, and grid mask, each with their unique implementation and purpose.

- (i) **Cutout :** The Cutout data augmentation method involves the random removal of a sub-region within an image, which is then filled with a constant value such as 0 or 255, during the training phase. This approach has been shown to result in improved performance on widely used datasets [29]. An illustration of the Cutout data augmentation process is provided in Figure 16.
- (ii) **Random erasing :** Random Erasing (RE) [170] randomly erases the sub-region in the image similar to cutout. But the main difference is, it randomly determines whether to mask out region or not and also determines the aspect ratio and size of the masked region. RE demonstration for different tasks is shown in figure 5.



Fig. 5. Random erasing examples for different tasks [170].

- (iii) **Hide-and-Seek** The process of hide-and-seek data augmentation [129] involves dividing an image into uniform squares of random size and then randomly removing a specified number of these squares. This technique aims to force neural networks to learn relevant features by hiding important information. A different view of the image is presented at each epoch, as depicted in figure 6. It is important to note that while this technique has been found to be effective in certain applications, it may also result in the removal of important information which could negatively impact the performance of the model.

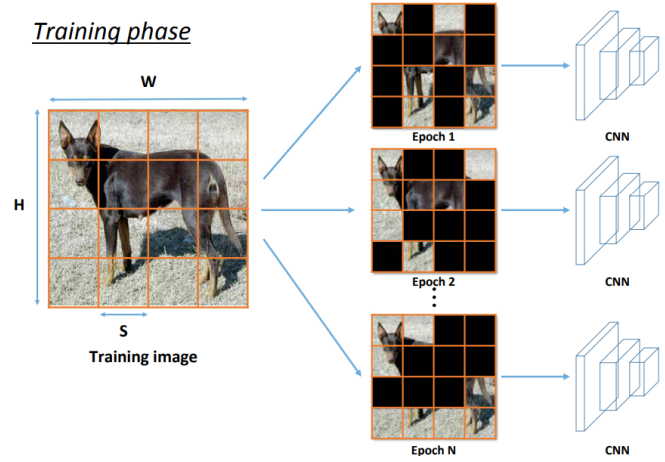


Fig. 6. An example of Hide-and-Seek augmentation [129].

- (iv) **GridMask Data Augmentation** The GridMask data augmentation technique [15] aims to address the challenges associated with randomly removing regions from images. This process, which can completely erase objects or strip away context information, requires a trade-off between the two. To resolve this, GridMask creates a uniform masking pattern and applies it to images as demonstrated in figure 7.

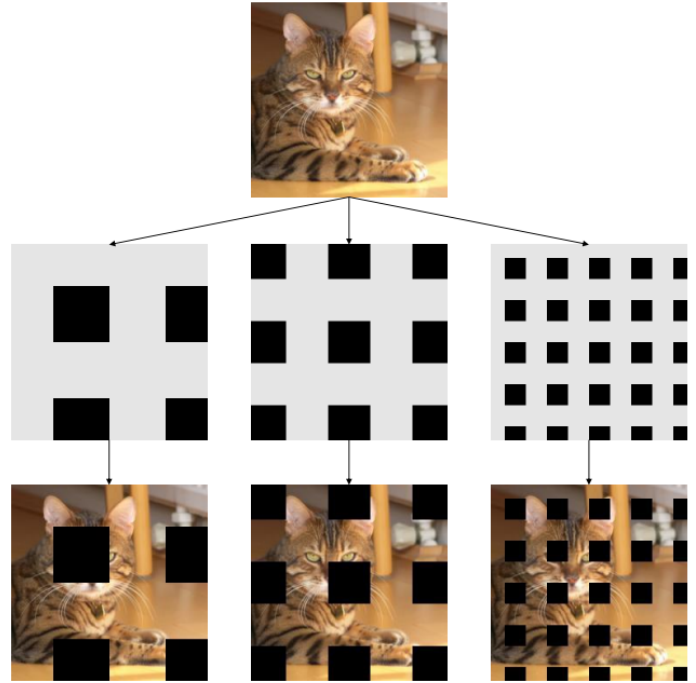


Fig. 7. This figure shows the procedure of GridMask augmentation. They produce a mask and then multiply it with the input image. Image is taken from [15].

B. Advanced Image Data Augmentations

The field of computer vision has seen a surge in interest regarding data augmentation techniques in recent times. This has led to the development of a wide range of innovative methods for augmenting image data, such as mixing images in novel ways, using reinforcement learning, feature-based augmentation, and style-based augmentation. To better understand these advancements, advanced data augmentation techniques have been classified into different major categories. These categories provide a useful framework for surveying the current state of the field and identifying areas for further research and development.

1) **Image Mixing Data Augmentations:** Image mixing data augmentation has gained popularity in computer vision research in recent years. This technique involves blending one or more images, including the same image, resulting in improved deep neural network model accuracy. We categorize image mixing data augmentation into two sub-categories: single image mixing and non-single image mixing. We compare the effectiveness of these sub-categories on benchmark datasets (such as CIFAR10, CIFAR100, ImageNet etc), as shown in table II-B5, II, VII, VIII and II-B5.

Single Image Mixing Data Augmentations: A single-image mixing technique uses only one image and mixes a single image from different mixing points of view. Recently, there has been a lot of work done on single-image augmentation, such as LocalAugment, SelfAugmentation, SelfMix, and many more. The description of each SOTA single image mixing data augmentation has been discussed below.

(i) **Local Augment:** Kim et al. [71] proposed a technique called LocalAugment, which involves dividing an image into smaller patches and applying different types of data augmentation to each patch. The purpose of this technique is to increase diversity in local features, which could help reduce bias and improve generalization performance of neural networks. While this approach does not preserve the global structure of an image, it provides a rich set of local features that can benefit neural network training. Figure 8 and 9 provide visual representations of the LocalAugment technique. Although LocalAugment technique can generate diverse local features of an image, it may not be well-suited for certain types of images that have complex global structures requiring preservation of global spatial relationships. Therefore, it may have some limitations, which should be taken into account while using this technique for image mixing data augmentation.

(ii) **Self-Augmentation:** This work [121] proposes self-augmentation, where a random region of an image is cropped and pasted randomly in the image, improving the generalization capability in few-shot learning. Moreover, the self-augmentation combines regional dropout and knowledge distillation- knowledge from the trained large network is transferred to a small network. The process demonstrated in the figure 10.

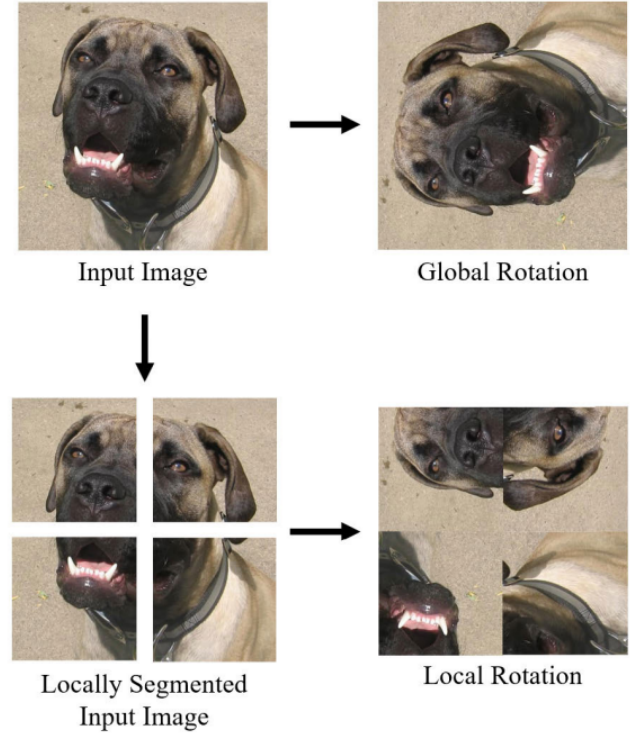


Fig. 8. An example of Global and Local Rotation Image, the example is taken from [71].

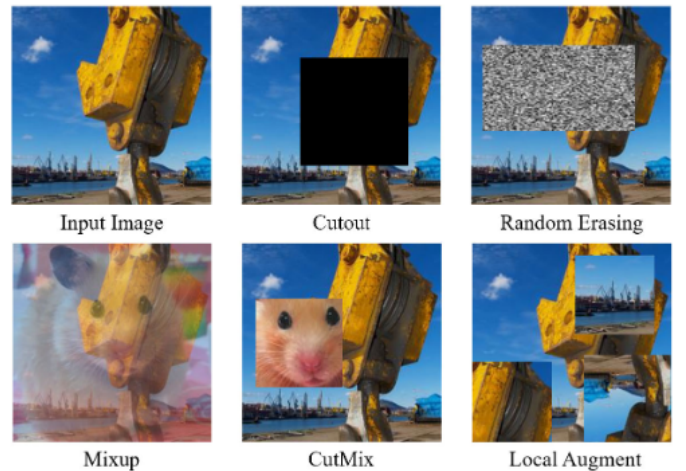


Fig. 9. Comparison of LocalAugment with CutOut, MixUp etc, example is taken from [71].

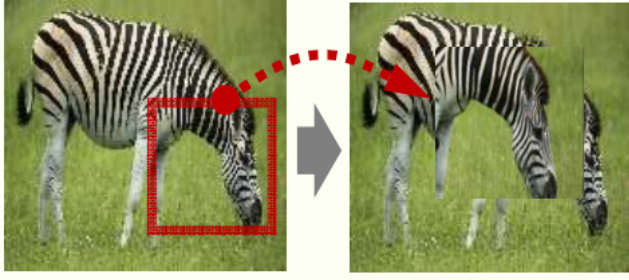


Fig. 10. An example of self augmentation, image is taken from [121]

- (iii) **SalfMix**: This work [20] focuses on whether it is possible to generalize neural networks based on single-image mixed augmentation. For that purpose, it proposes SalfMix, the first salient part of the image is found to decide which part should be removed and which portion should be duplicated. Most salient regions are cropped and placed into non-salient regions. This process is defined and compared with other techniques in figure 11.

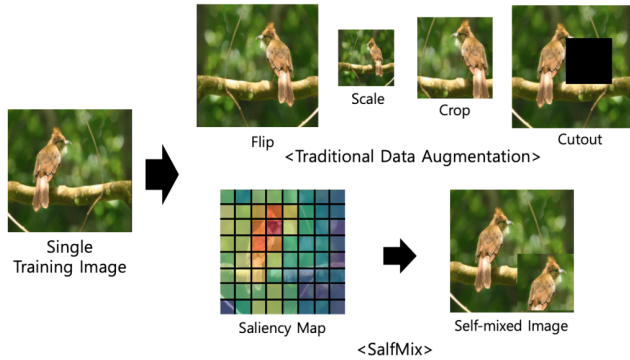


Fig. 11. Conceptual comparison between SalfMix method and other single image-based data augmentation methods, the example is taken from [20].

- (iv) **KeepAugment** KeepAugment [47] is introduced to prevent distribution shift which degrades the performance of neural networks. The idea of KeepAugment is to increase fidelity by preserving the salient features of the image and augmenting the non-salient region. Preserved features help to increase diversity without shifting the distribution. KeepAugment is demonstrated in figure 12.
- (v) **You Only Cut Once** You Only Cut Once (YOCO) [51] is introduced with the aim of recognizing objects from partial information and improving the diversity of augmentation that encourage neural networks to perform better. YOCO makes two pieces of image and augmentation is applied on each piece, then each piece is concatenated for an image and YOCO shows impressive performance and compared with SOTA augmentations, sometimes it outperforms them. It is easy to implement, has no parameters, and is easy to use. The YOCO augmentation

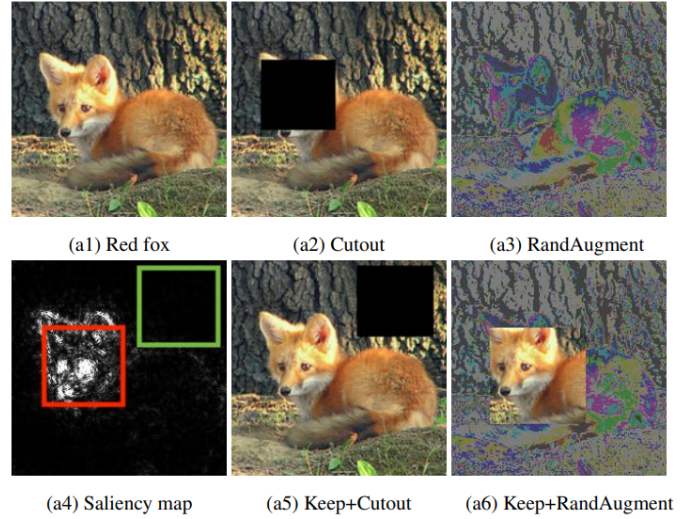


Fig. 12. This image shows the example of KeepAugment with other augmentations, courtesy [47].

process is shown in figure 13.



Fig. 13. An example of YOCO augmentation, image is taken from [51].

- (vi) **Cut-Thumbnail**: Cut-Thumbnail [156] is a novel data augmentation, that resizes the image to a certain small size and then randomly replaces the random region of the image with the resized image, aiming to alleviate the shape bias of the network. The advantage of Cut-thumbnail is, that it not only preserves the original image but also keeps it global in the small resized image. On ImageNet, it shows impressive performance using resnet50. Overall, the cut-thumbnail process and its comparison are shown in figure 15 and figure 14, respectively.

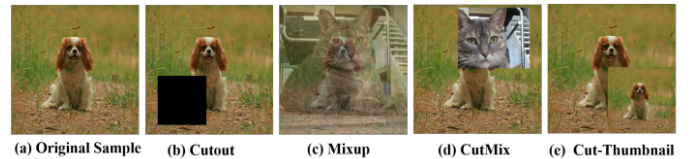


Fig. 14. Comparison between existing data augmentation methods with Cut-Thumbnail, the example is from [156].

Multi-Images Mixing: Multi-Images Mixing data augmentation uses more than one image and applies different mixing strategies. Recently, many researchers have explored a lot of Multi-Images Mixing strategies and still, it is a very attentive topic for many researchers. Recently work has included Mixup, CutMix, SaliencyMix, and many more. Each

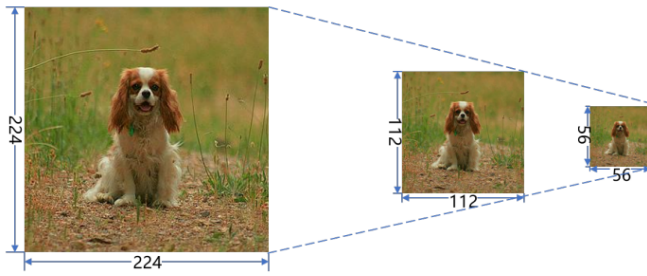


Fig. 15. This image shows an example of reduced images that are called thumbnails. After reducing the image to a certain size of 112×112 or 56×56, The dog is still recognizable even though lots of local details are lost, courtesy [156].

of the relevant non-single image mixing data augmentation techniques is discussed below.

- (i) **Mixup**: Mixup blends two images based on the blending factor (alpha) and the corresponding labels of these images are also mixed in the same way. Mixup data augmentation [163] consistently improved the performance not only in terms of accuracy but also in terms of robustness. Experiments on ImageNet-2012 [119], CIFAR-10, CIFAR-100, Google commands⁴ and UCI datasets⁵ showed impressive results on SOTA methods. Further demonstration and comparison are shown in the figure 16.
- (ii) **CutMix** : CutMix tackles the issues of information loss and region dropout [162]. It is inspired by cutout [29], where any random region is filled with 0 or 255, while in cutmix instead of filling the random region with 0 or 255, the region is filled with a patch from another image. Correspondingly, their labels are also mixed proportionally to the number of pixels mixed. It is compared with other methods and shown in figure 16.



Fig. 16. Overview of the Mixup, Cutout, and CutMix, Example is from [162].

- (iii) **SaliencyMix**: This technique addresses the problem of cutmix and argues that filling a random region of the image with a patch from another will not guarantee that patch has rich information and thereby mixing labels of unguaranteed patches leads the model to learn unnecessary information about the patch [141]. To deal with that issue, saliencyMix first selects the salient part of the image and pastes it to a random region or salient or non-salient of another image. It is shown in figure 17 and figure 18.

⁴<https://research.googleblog.com/2017/08/launching-speech-commands-dataset.html>

⁵<http://archive.ics.uci.edu/ml/index.php>

Target Image	Source Image	Augmented Image	
Mixed label for randomly mixed images		Dog - 80% & Cat 20% ?	Dog - 80% & Cat 20% ?

Fig. 17. An example of SaliencyMix augmentation, image is taken from [141].

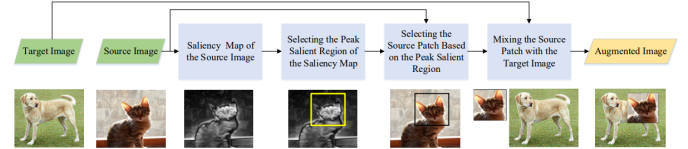


Fig. 18. This image shows the proposed SaliencyMix data augmentation procedure, courtesy [141]

- (iv) **RSMDA: Random Slices Mixing Data Augmentation**: RSMDA [77] addresses issues of feature losing in single image erasing data augmentation. RSMDA gets the slices of one image and mixes them with another image alternatively and the corresponding labels are also mixed accordingly. This work further investigates three different strategies of RSMDA; row-wise slice mixing, column-wise slice mixing and randomness of both. Row-wise slice mixing has shown superior performance. Demonstration of each of the slices mixing strategy is in figure 19.
- (v) **Puzzle Mix**: This article [70] proposes a puzzle mix data augmentation technique that focuses on using explicitly salient information and basic statistics of image wisely with the aim of breaking the misleading supervision of neural networks over existing data augmentations. Furthermore, the demonstration is shown and compared with relevant methods in figure 20.
- (vi) **SnapMix**: The article [60] proposes the Semantically Proportional Mixing (SnapMix) that utilises class activation map (CAM) to reduce the label noise level. SnapMix creates the target label considering the actual salient pixel taking part in the augmented image, which ensures semantic correspondence between the augmented image and mixed labels. The overall process is demonstrated and compared with closely matching augmentations in the figure 21.
- (vii) **FMix**: This article proposes the FMix [52], a kind of mixed sample data augmentation (MSDA), that utilises the random binary masks. These random binary masks are acquired by applying a threshold to low-frequency images that are obtained from the Fourier space. Once the mask is obtained, one color region is applied to input one and another color region is applied to another input. The overall process is shown in figure 22.
- (viii) **MixMo**: This paper [112] focuses on the learning of multi-input multi-output via sub-network. The main mo-

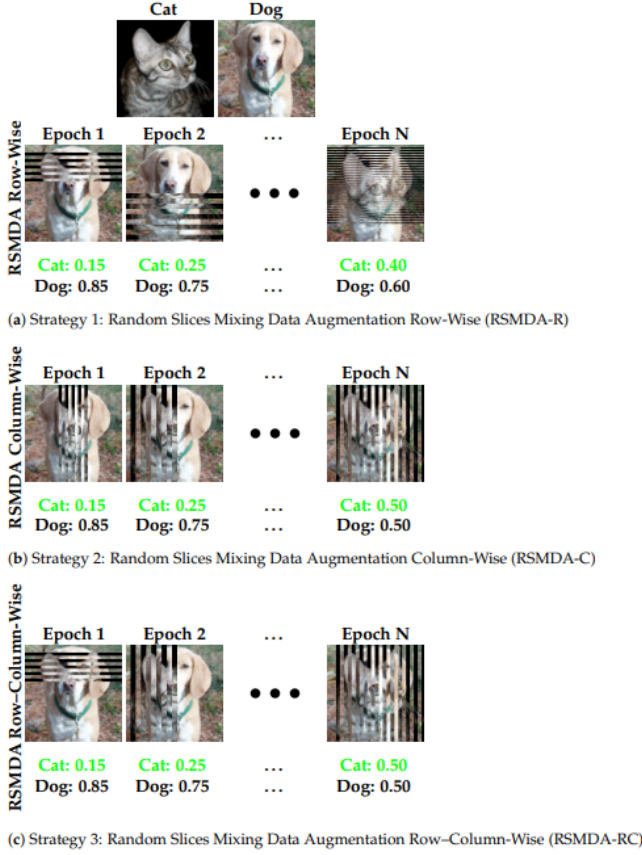


Fig. 19. RSMDA three different strategies. Image is taken from [77].

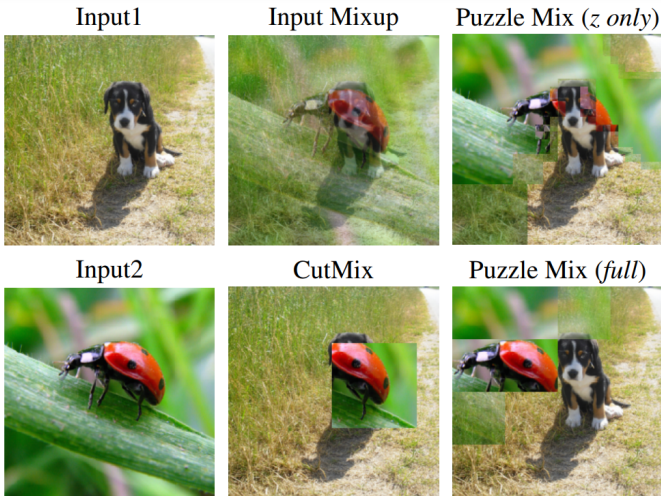


Fig. 20. A visual comparison of the mixup methods. Puzzle Mix ensures to contain sufficient target class information while preserving the local statistics of each in, the example is from [70].

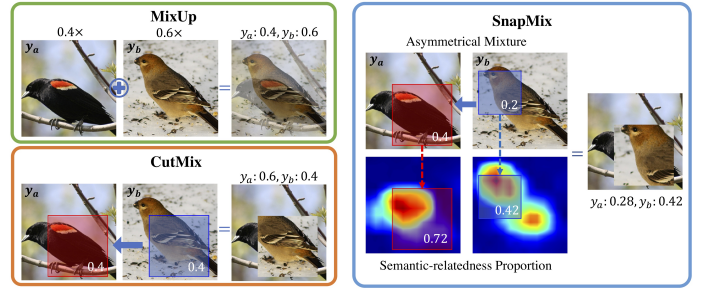


Fig. 21. A visual comparison of Mixup, CutMix, and SnapMix. The figure gives an example where a label generated by SnapMix is visually more consistent with the mixed image semantic structure compared to CutMix and Mixup, courtesy [60].



Fig. 22. Example masks and mixed images from CIFAR-10 for FMix, example is from [52].

tivation of the paper is to replace direct hidden summing operations with more solid mechanisms. For that purpose, it proposes MixMo, which embeds M inputs into the shared space, mixes and passes these to a further layer for classification. Moreover, the overall process is demonstrated in figure 23:

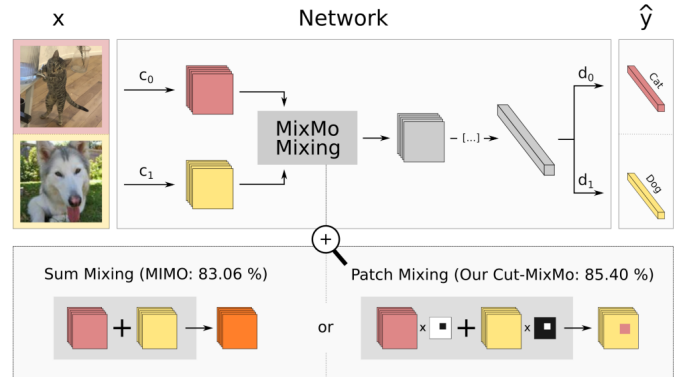


Fig. 23. This image shows the overview of MixMo augmentation, the image is taken from [112].

(ix) **StyleMix:** This paper [59] targets previous approaches that are unable to differentiate between content and style features, approaches such as mixup based data augmentations. To remedy this problem, it proposes two approaches styleMix and StyleCutMix, this is the first

work that separately deals with content and style features of images very carefully and it showed impressive performance on popular benchmark datasets. The overall process is defined and compared with SOTA approaches in figure 24.

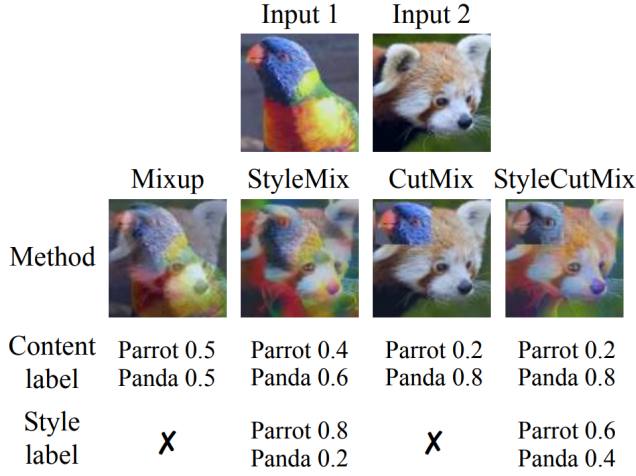


Fig. 24. A Visual comparison of StyleMix [59] and StyleCutMix with Mixup [163] and CutMix [162], example is from [59].

- (x) **RandomMix:** This work [94] improves generalization capability by proposing randomMix, which randomly selects augmentation from a set of image mixing augmentations and applies it to images, enabling the model to look at diverse samples. This method showed impressive results over SOTA image mixing methods. The overall demonstration is shown in figure 25.

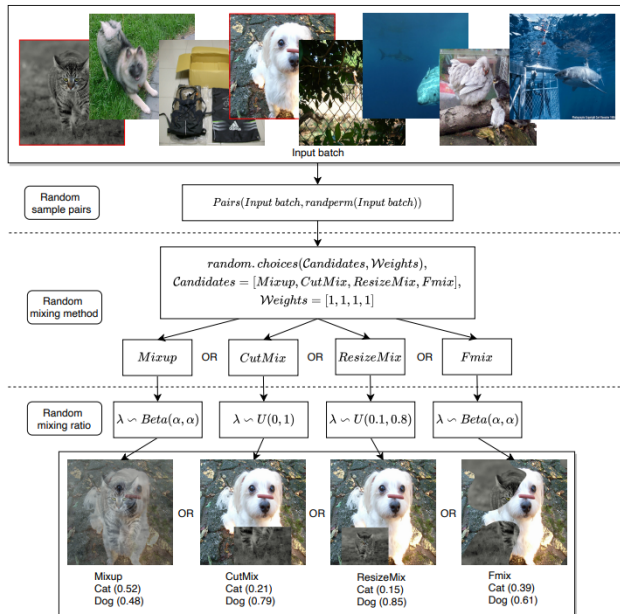


Fig. 25. An illustrative example of RandomMix, image is taken from [94].

- (xi) **MixMatch:** MixMatch data augmentation technique is very useful in semi-supervised learning. MixMatch [10] augments single image K times and passes all K number of images to a classifier, averages their prediction and finally, their predictions are sharpened by adjusting their distribution temperature term. It is demonstrated in figure 26.

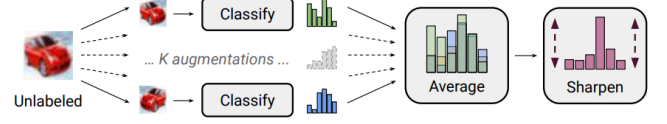


Fig. 26. Diagram of the label guessing process used in MixMatch, courtesy [10].

- (xii) **ReMixMatch:** In this work [9], an extension of MixMatch [10] is proposed to make prior work efficient by introducing distribution alignment and augmentation anchoring. The distribution alignment task aims to minimize the gap between the marginal distribution of predictions on unlabeled data and the marginal distribution of ground truth labels. On the other hand, augmentation anchoring feeds multiple strongly augmented versions of the input into the model and encourages each output to be close to the prediction for a weakly-augmented version of the same input. The process is illustrated in figure 27.

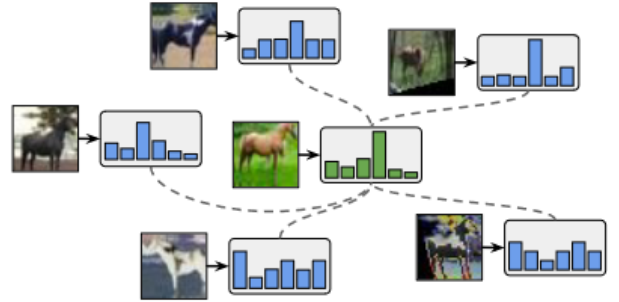


Fig. 27. Anchoring augmentation. It makes predictions on strong augmentations of the same image (blue) using the forecast for a weakly enhanced image (green, centre), courtesy [9].

- (xiii) **FixMatch:** Fixmatch [130] is a method for improving the performance of semi-supervised learning (SSL). It first assigns pseudo-labels to unlabeled images that have a predicted probability above a certain threshold, and then trains the model to match these labels using cross-entropy loss on a strongly augmented version of the image. The process is illustrated in Figure 28.
- (xiv) **AugMix:** The work proposed in [56] presents Augmix, a data augmentation technique that aims to reduce the distribution gap between training and test data. Augmix applies M random augmentations to an input image, each with a random strength, and merges the resulting images

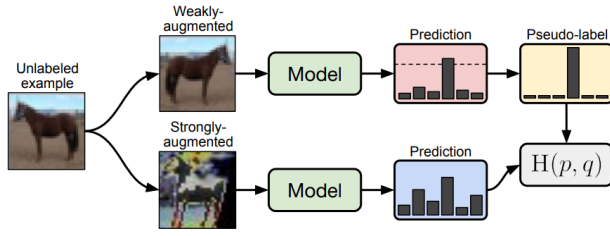


Fig. 28. This image shows the procedure of FixMatch, image is taken from [130].

to produce a new image that spans a wider area of the input space. The process is illustrated in Figure 29, where three branches perform separate augmentations and additional operations are added to increase diversity. The resulting images are then mixed to produce a final augmented image, which is effective in improving model robustness.

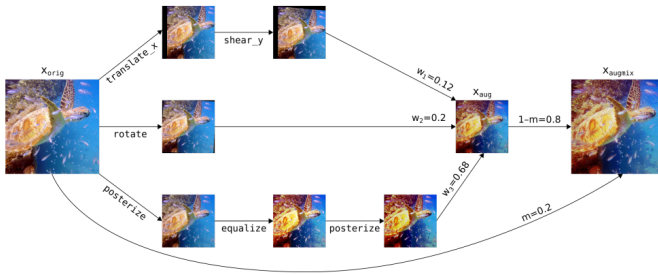


Fig. 29. An overall procedure of AugMix augmentation and example is from [56].

- (xv) **Simple Copy-Paste is a Strong Data Augmentation Method, for Instance Segmentation:** The proposed approach in this work [43] involves copying and pasting instances from one image to another to create an augmented image. This simple technique has shown promising results and is easy to implement. Figure 30 illustrates the process, where instances from two images are pasted onto each other at different scales.



Fig. 30. Image augmentation performed by simple Copy-Paste method, image courtesy [43].

- (xvi) **Improved Mixed-Example Data Augmentation:** Recently, label non-preserving data augmentation techniques

based on linear combinations of two examples have demonstrated promising results. In this paper [135], the authors investigate two research questions: (i) the reasons behind the success of these methods and (ii) the significance of linearity in data augmentations. Figure 31 illustrates the overall process.

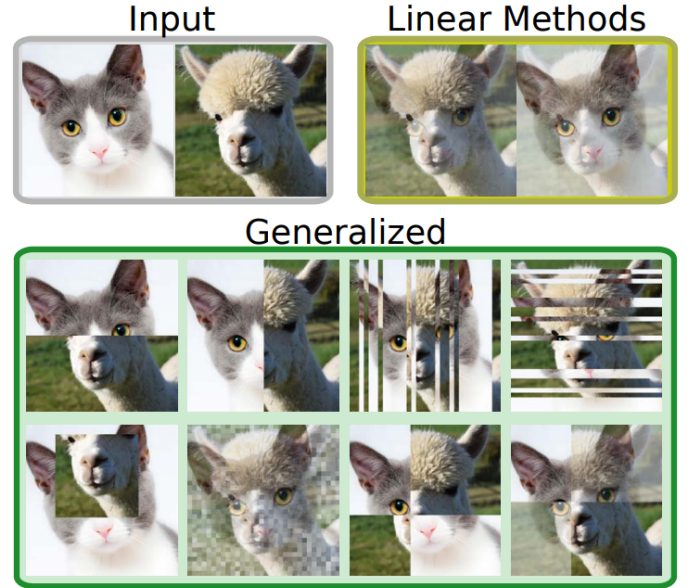


Fig. 31. A visual comparison of linear methods and generalized augmentation performed by Improved Mixed-Example, image is taken from [135].

- (xvii) **Random image cropping and patching (RICAP):** Random Image Cropping and Patching (RICAP) [137] is a new data augmentation technique that cuts and mixes four images rather than two images. The key idea behind RICAP is to crop patch from each of the four images and then mixes these patch to create augmented image. The labels of the images are also mixed in proportion to the area of the patches. This technique showed impressive performance on popular datasets i.e. CIFAR10, CIFAR100, and imageNet. RICAP demonstration is shown in figure 32.
- (xviii) **Cutblur:** This article [159] explores and analyses existing data augmentation techniques for super-resolution and proposes another data augmentation technique for super-resolution, named cutblur that cuts high-resolution image patches and pastes to corresponding low-resolution images and vice-versa. Cutblur shows impressive performance on several super-resolution benchmark datasets. Furthermore, the process is illustrated in figure 33 and 34.
- (xix) **ResizeMix: Mixing Data with Preserved Object Information and True Labels :** The ResizeMix [111] method directly cuts and pastes the source data in four different ways to target the image. These four different ways include salient part, non-part, random part or resized source image to patch, as shown in the figure 35. It addresses two questions:

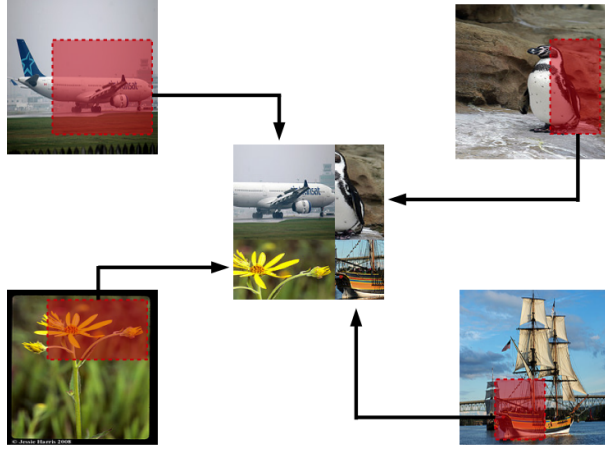


Fig. 32. A conceptual explanation of the RICAP data augmentation, the example is from [137].

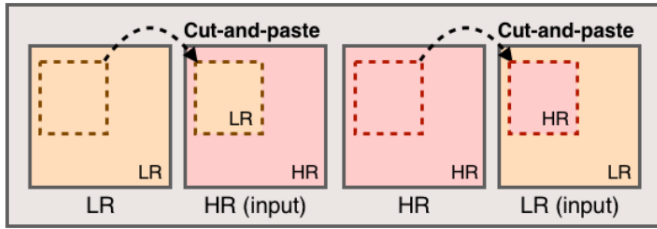


Fig. 33. An Schematic illustration of CutBlur operation, image is taken from [159].

- How to obtain a patch from the source image?
- where to paste the patch from the source image in the target image?

Furthermore, it was found that saliency information is not important to promote mixing data augmentation. ResizeMix is shown in the figure 35.

(xx) **ClassMix: Segmentation-Based Data Augmentation for Semi-Supervised Learning** : This research work [104] proposes novel data augmentation for semi-supervised learning for semantic segmentation task. It showed that traditional data augmentations are not ef-

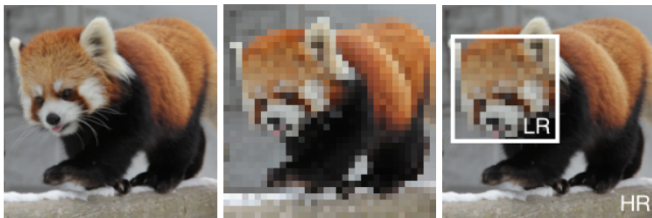


Fig. 34. A visual comparison between High resolution, low resolution and CutBlur, courtesy [159].

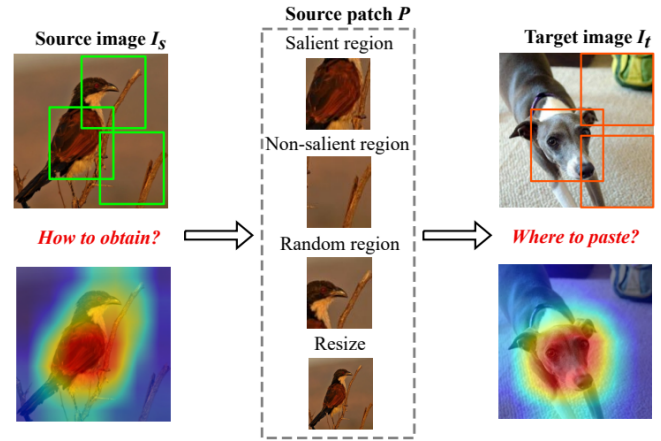


Fig. 35. A visual representation of different cropping manners from the source image and different pasting manners to the target image, image is taken from [111].

fective for image semantic segmentation as they are for image classification. The proposed data augmentation named ClassMix, augments the training sample by mixing unlabeled samples, by exploiting network prediction while taking into account object boundaries. It showed a massive performance gain on two common semantic segmentation datasets for semi-supervised learning. The overall process is shown in the figure 36.

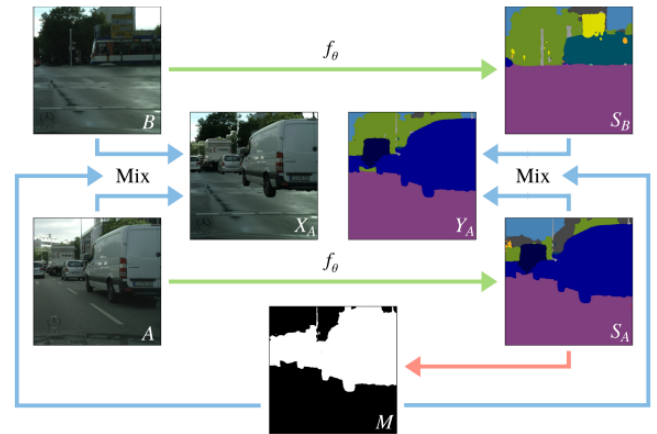


Fig. 36. In a visual representation classMix augmentation, two images are sampled then based on the predictions of each image a binary mask is created. The mask is then used to mix the images and their predictions, the image is taken from [?].

(xxi) **Context Decoupling Augmentation for Weakly Supervised Semantic Segmentation (WSSS)**: This article [134] addresses the problem of traditional data augmentation techniques for WSSS, increasing the same contextual data semantic samples does not add much value in object differentiation, i.e. in image classification, “cat” recognition is due to the cat itself and its surrounding context, these both contexts discourages model to focus

only on the cat. To tackle that issue, this work proposes a novel data augmentation, named Context Decoupling Augmentation (CDA). CDA increases diversity of the specific object and it guides the network to break the dependencies between object and contextual information. In this way, it also provides augmentation and the network focuses on object(s) only rather than object(s) and its contextual information. A comparison of traditional data augmentation and CDA is shown below in the figure 37.

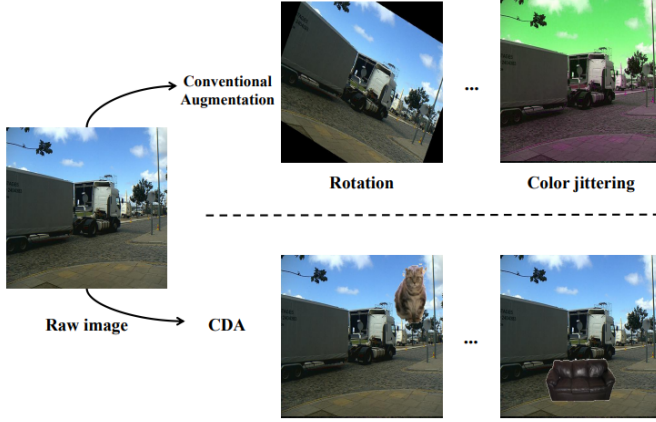


Fig. 37. A visual representation of the difference between the conventional augmentation approach and context decoupling augmentation (CDA), image is taken from [134].

(xxii) **ObjectAug: Object-level Data Augmentation for Semantic Image Segmentation:** This article [164] addresses the problem of mixing image-level data augmentation strategies, which failed to work for segmentation as object and background are coupled and boundaries of objects are not augmented due to their fixed semantic bond with the background. To mitigate this problem, this article [164] proposes a novel approach named ObjectAug, object-level augmentation for semantic segmentation. First, it separates object(s) and backgrounds from an image with the help of semantic labels then each object is augmented using popular data augmentation techniques such as flipping and rotating. Pixel changes due to these data augmentations are restored using image inpainting. Finally, the object(s) and background are coupled to create an augmented image. Experimental results suggest that ObjectAug has shown performance improvement for segmentation tasks. Furthermore, ObjectAug is shown in the figure 38.

2) **AutoAugment:** The goal of this technique is to find the data augmentation policies from training data. It solves the problem of finding the best augmentation policy as a discrete search problem. It consists of a search algorithm and a search space. Furthermore, these techniques are classified into two sub-categories based on reinforcement learning and non-reinforcement learning.

- Reinforcement learning data augmentation

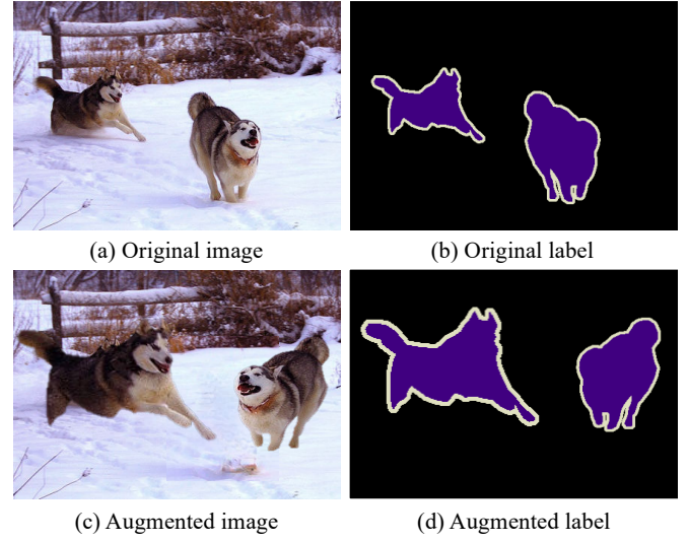


Fig. 38. ObjectAug can perform various augmentation methods for each object to boost the performance of semantic segmentation. The left husky is scaled and shifted, while the right one is flipped and shifted. Thus, the boundaries between objects are extensively augmented to boost their performance, the example is from [164].

- Non-Reinforcement learning data augmentation

Reinforcement Learning data augmentations: Reinforcement learning data augmentation techniques generalize and improve the performance of deep networks in an environment.

- (i) **AutoAugment:** This work [23] automatically finds the best data augmentation rather than manual data augmentation. To address the limitations of manual search-based data augmentation, this article proposes autoaugment, where search space is designed and has policies consisting of many sub-policies. Each sub-policy has two parameters one is the image processing function and the second one is the probability with magnitude. These sub-policies are found using reinforcement learning as a searching algorithm. The overall process is demonstrated in figure 39.

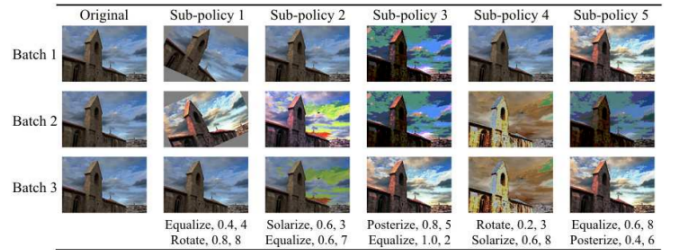


Fig. 39. A visual overview of the sub-policies from ImageNet using AutoAugment, example is from [23].

- (ii) **Fast AutoAugment:** Fast Autoaugment [90] addresses the problem of autoaugment, autoaugment takes a lot of time to find the optimal data augmentation strategy. To reduce the searching time, fast auto augment finds

more optimal data augmentations using an efficient search strategy based on density matching. It reduces the higher order of training time compared to autoaugment. The overall procedure is shown in figure 40.

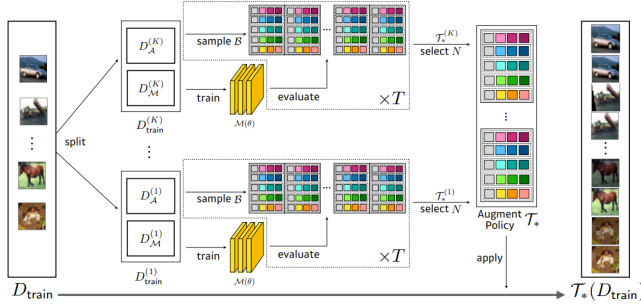


Fig. 40. An overall procedure of augmentation search by Fast AutoAugment algorithm, courtesy [90].

- (iii) **Faster AutoAugment:** This article proposes a faster autoaugment [53] policy intending to find effective data augmentation policies very efficiently. Faster autoaugment is based on a differentiable augmentation searching policy and additionally, it not only estimates gradients for many transformation operations having discrete parameters but also provides a mechanism for choosing operations efficiently. Moreover, it introduces a training objective function with aim of minimising the distance between original and augmented distribution, which is also differentiable. Parameters of augmentations are updated during backpropagation. The Overall process is defined in figure 41:
- (iv) **Reinforcement Learning with Augmented Data:** This paper proposes Reinforcement Learning with Augmented Data (RAD) [84], easily pluggable and enhances the performance of RL algorithms by targeting two issues i) learning data efficiency and ii) generalisation capability for new environments. Furthermore, it shows traditional data augmentation techniques enable RL algorithms to outperform complex SOTA tasks for pixel-based control and state-based control. Overall process is demonstrated in figure 42:
- (v) **Local Patch AutoAugment with Multi-Agent Collaboration:** This is the first work [91] that finds data augmentation policy for patch level using reinforcement learning, named multi-agent reinforcement learning (MARL). MARL starts by dividing images into patches and jointly finds the optimal data augmentation policy for each patch. It shows competitive results on SOTA benchmarks. Overall process is defined in figure 43:
- (vi) **Learning Data Augmentation Strategies for Object Detection:** This work [171] proposes to use autoaugment that learns the best policies for object detection. It finds the best operation and optimal value. Moreover, it addresses two key issues of augmentation for object detection,

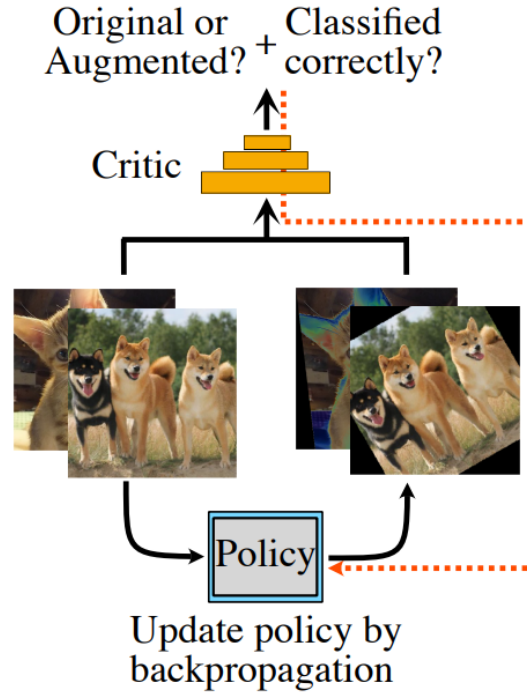


Fig. 41. An Overview of the Faster AutoAugment augmentation, image is taken from [53].

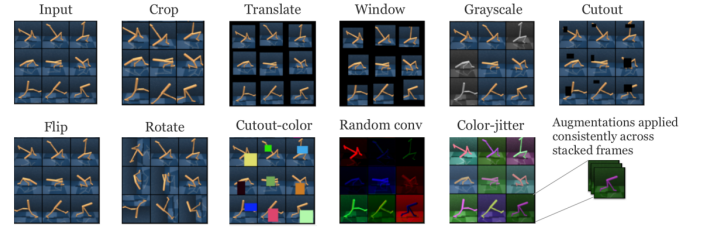


Fig. 42. An overview of different augmentation investigated in RAD, the example is from [84].

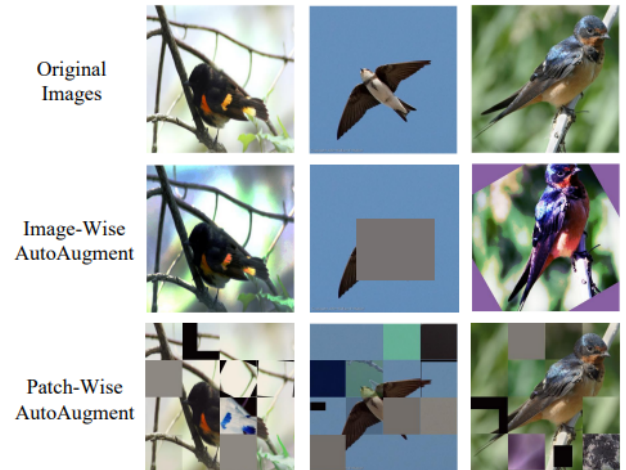


Fig. 43. An illustration of different automated augmentation policies, courtesy [91].

- Classification learned policies can not directly be applied for detection tasks, and it adds more complexity to deal with bounding boxes if geometric augmentations are applied.
- Most researchers think it adds much less value compared to designing new network architecture so gets less attention but augmentation for object detection should be selected carefully.

Some sub-policies for this data augmentation are shown below:

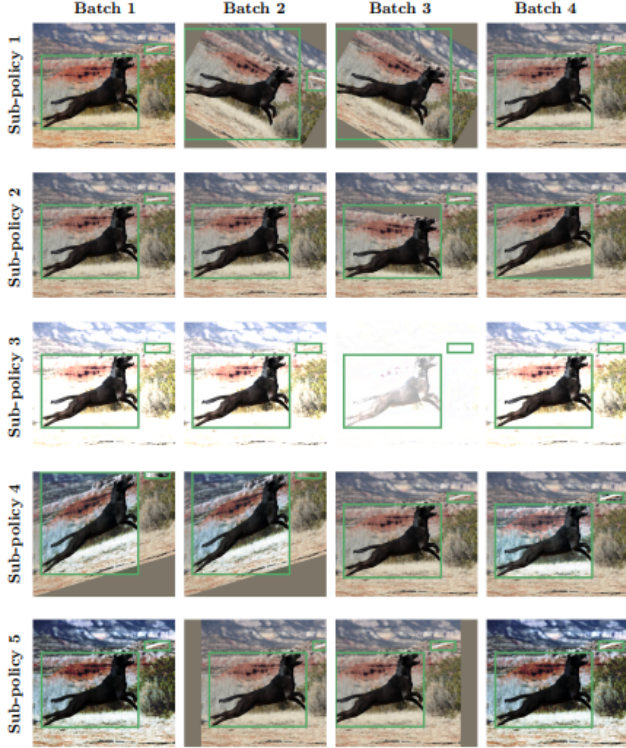


Fig. 44. Different data augmentation sub-policies explored, image is taken from [171].

- Sub-policy 1. (Color, 0.2, 8), (Rotate, 0.8, 10)
- Sub-policy 2. (BBox_Only_ShearY, 0.8, 5)
- Sub-policy 3. (SolarizeAdd, 0.6, 8), (Brightness, 0.8, 10)
- Sub-policy 4. (ShearY, 0.6, 10), (BBox_Only_Equalize, 0.6, 10)
- Sub-policy 5. (Equalize, 0.6, 10), (TranslateX, 0.2, 2)

- Scale-aware Automatic Augmentation for Object Detection:** This work [18] proposes a new data augmentation for object detection named scale aware autoAug, first, it defines a search space where image level and box level data augmentation are prepared for scale invariance, secondly, it also proposes a new search metric named Pareto scale balance for search augmentation effectively and efficiently. Some examples of data augmentation are shown in figure 45.

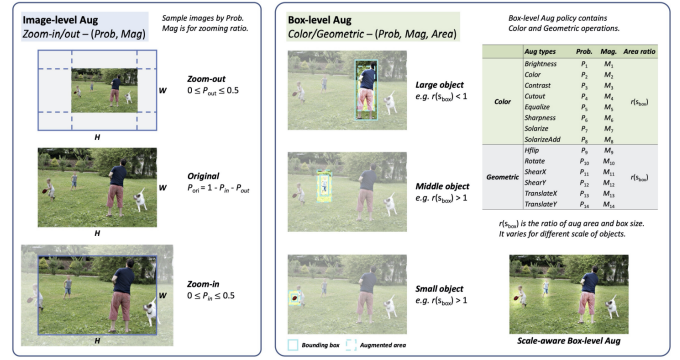


Fig. 45. Example of scale-aware search space which includes image level and box-level augmentation, the example is from, [18].

3) Non-Reinforcement Learning data augmentations:

In auto-augment category, there are some approaches that do not require any reinforcement learning algorithm to find the best data augmentation, we refer to them as non-reinforcement learning data augmentation. We categorise a few of them as discussed below.

- RandAugment:** Previous optimal augmentation finding uses reinforcement or some complex learning strategy that takes a lot of time to find. RandAugment augmentation [24] removes obstacles of a separate searching phase, which makes training more complex and consequently adds computational cost overhead. To break this, randaugment applies randomly N number of data augmentations with M magnitude of all augmentations. Some visualisation is demonstrated in the figure 46:

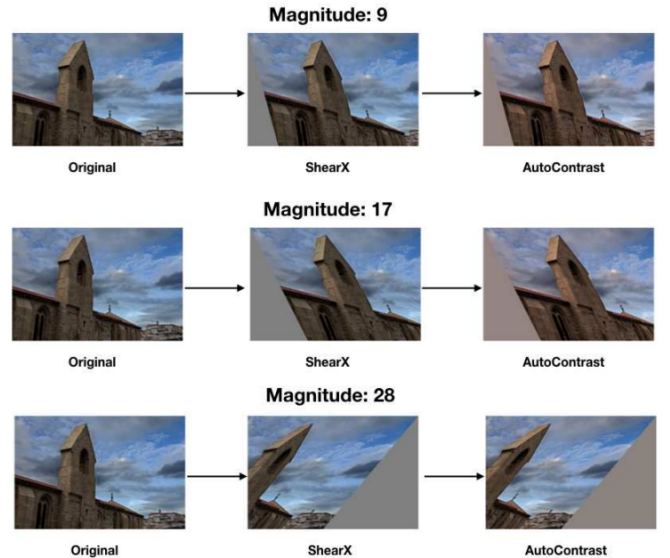


Fig. 46. Example images augmented by RandAugment, image is taken from [24].

- RangeAugment** RangeAugment [98] is a data augmentation technique that aims to improve upon the

shortcomings of existing approaches like AutoAugment and RandAugment. These methods use manually-defined ranges of magnitudes for each type of data augmentation, which can result in sub-optimal policies. In contrast, RangeAugment learns efficient ranges of magnitudes for each augmentation and composite data augmentation by introducing an auxiliary loss based on image similarity. This loss is designed to control the magnitude ranges, resulting in more effective and optimal policies. The process of RangeAugment is illustrated in Figure 47.

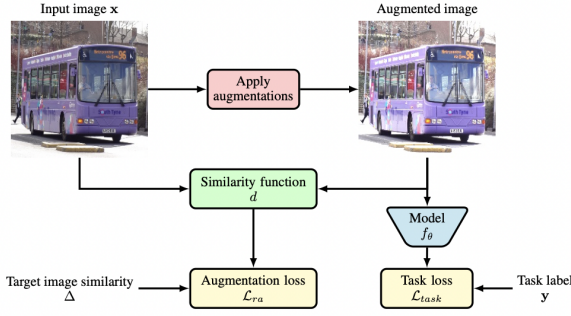


Fig. 47. RangeAugment with neural network training [98].

- (iii) **ADA: Adversarial Data Augmentation for Object Detection:** Data augmentation for object detection has improved performance but it is difficult to understand whether these augmentations are optimal or not. This article [8] provides a systematic way to find optimal adversarial perturbation of data augmentation from an object detection perspective, that is based on game-theoretic interpretation aka Nash equilibrium of data. Nash equilibrium provides the optimal bounding box predictor and optimal design for data augmentation. Optimal adversarial perturbation refers to the worst perturbation of ground truth, that forces the box predictor to learn from the most difficult distribution of samples. An example is shown in figure 48.

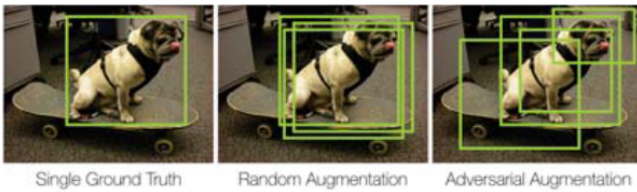


Fig. 48. Annotation distribution types. Adversarial augmentation chooses bounding boxes that are as distinct from the truth as possible while yet containing crucial object characteristics. The example is taken from [8].

- (iv) **Deep CNN Ensemble with Data Augmentation for Object Detection:** This article [49] proposes a new variant of the regions with convolutional neural network (R-CNN) model with two core modifications in training and evaluation. First, it uses several different CNN models as ensembler in R-CNN, secondly, it smartly augments

PASCAL VOC training examples with Microsoft COCO data by selecting a subset from Microsoft COCO datasets that are consistent with PASCAL VOC. Consequently, it increases the dataset size and improves the performance. The schematic diagram is shown in the figure 49.

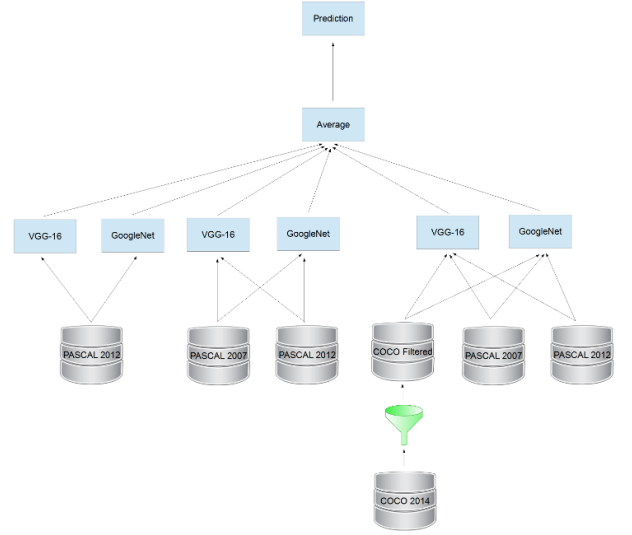


Fig. 49. The proposed schematic diagram. The example is taken from [49].

- (v) **Robust and Accurate Object Detection via Adversarial Learning:** This article [16] first shows classifier performance gain from different data augmentations when it is fine-tuned to object detection tasks and suggests that the performance in terms of accuracy or robustness is not improving. The article provides a unique way of exploring adversarial samples that helps to improve performance. To do so, it augments the example during the fine-tuning stage for object detectors by exploring adversarial samples, which is considered as model-dependent data augmentation. First, it picks the stronger adversarial sample from detector classification and localization layers and ensures the augmentation policy remains consistent. It showed significant performance gain in terms of accuracy and robustness on different object detection tasks. Furthermore, the robustness and accuracy of the proposed method are shown in figure 50.
- (vi) **Perspective Transformation Data Augmentation for Object Detection:** This article [145] proposes a new data augmentation for objection detection named perspective transformation that generates new images captured at different angles. Thus, it mimics images as if they are taken at a certain angle where the camera can not capture those images. This method showed effectiveness on several object detection datasets. An example of the proposed data augmentation is shown in figure 51.
- (vii) **Deep Adversarial Data Augmentation for Extremely Low Data Regimes:** This article [165] addresses the issue of extremely low data regimes-labeled data is very less, no unlabeled data at all. To deal with that problem, it

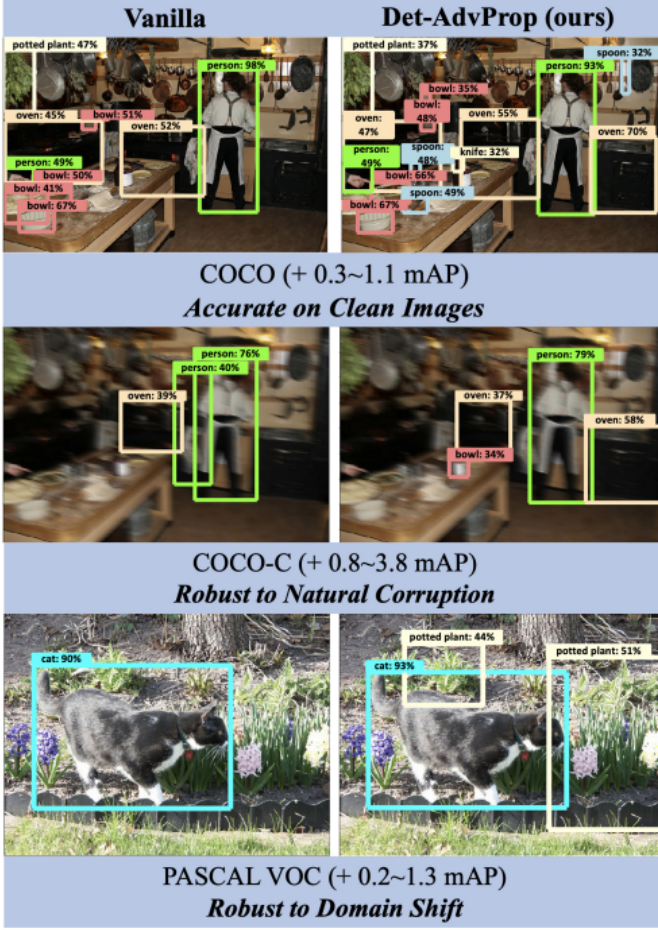


Fig. 50. Overview of Robust and Accurate Object detection via adversarial learning. In the top image, it improves object detector accuracy on clean images. In middle, improves the detector's robustness against natural corruption, and at the bottom, it improves the robustness against cross-dataset domain shift. The image is taken from [18].

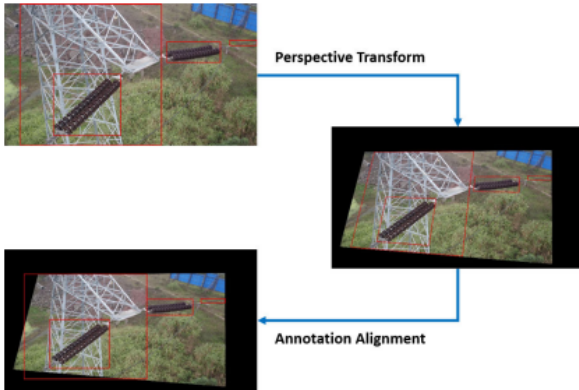


Fig. 51. Perspective transformation data augmentation. An example image is taken from [145]

proposes a deep adversarial data augmentation (DADA), where data augmentation is formulated as a problem of training class conditional and supervised GAN. Furthermore, it also introduces new discriminator loss with aim of fitting data augmentation where real and augmented samples are forced to participate equally and be consistent in finding decision boundaries.

4) **Feature augmentation:** Feature augmentation is another category of data augmentation, where images are transformed into embedding or representation then data augmentation is performed on the embedding of the image. Recently a few works have been done in this area, we selectively highlight the work in a precise way.

- (i) **FeatMatch: Feature-Based Augmentation for Semi-Supervised Learning :** This work [81] presents a novel approach of data augmentation in features space for SSL inspired by an image-based SSL method that uses a combination of augmentations of the images and consistency regularization. Image-based SSL methods are restricted to only conventional data augmentation. To break this end, the feature-based SSL method produced diverse features from complex data augmentations. One key point is, these advanced data augmentations exploit the information from both intra-class and inter-class representations extracted via clustering. The proposed method only showed significant performance gain on mini-Imagenet such as an absolute 17.44% gain on miniImageNet, but also showed robustness on samples that are out-of-distribution. Moreover, the difference between image-level and feature-level augmentation and consistency is shown in figure 52.

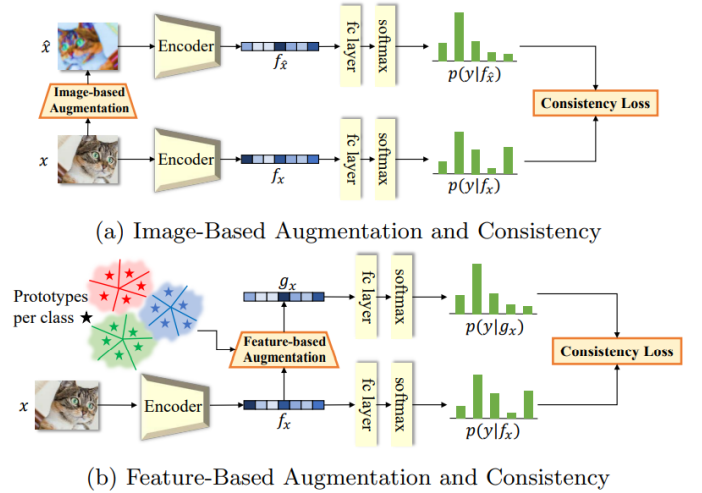


Fig. 52. An overview of featMatch augmentation applied on images and features. Image is taken from [21].

- (ii) **Dataset Augmentation in Feature Space:** This work [28] first used encoder-decoder to learn representation, then on representation apply different transformations such as adding noise, interpolating, or extrapolating. The proposed method has shown performance improve-

ment on both static and sequential data. Moreover, a demonstration of this augmentation is shown in figure 53.

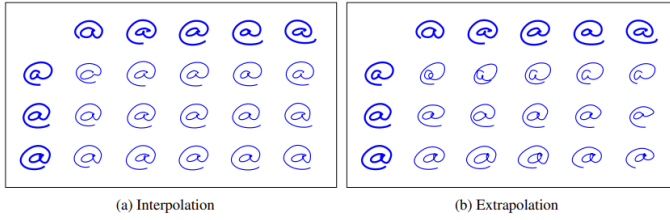


Fig. 53. Overview of interpolation and extrapolation between handwritten characters. Original characters are shown in bold. Image is taken from [28].

- (iii) **Feature Space Augmentation for Long-Tailed Data :** This paper [21] proposed the novel data augmentation in feature space to address the long-tailed issue and uplift the under-represented class samples. The proposed approach first separates class-specific features into generic and specific features with the help of class activation maps. Under-represented class samples are generated by injecting class-specific features of under-represented classes with class-generic features from other confusing classes. It enables diverse data and also deals with the problem of under-represented class samples. It has shown SOTA performance on different datasets. It is demonstrated in figure 54.

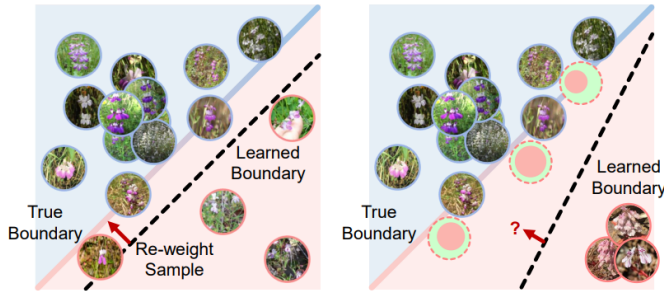


Fig. 54. Left: limited but well-spread data. Right: Without sufficient data. Image is taken from [21].

- (iv) **Adversarial Feature Augmentation for Unsupervised Domain Adaptation:** Generative Adversarial Networks (GANs) showed promising results in unsupervised domain adaptation to learn target domain features indistinguishable from the source domain. This work [143] extends GAN by contributing: i) it forces feature extractor to be domain-invariant ii) To train it via data augmentation in feature space, named feature augmentation. This work explores data augmentation at the feature level with GAN.
- (v) **Understanding data augmentation for classification: when to warp? :** This paper [154] investigates the data augmentation advantages on image space and feature space during training. It proposed two approaches i) data warping which generates extra samples in image space using data augmentations and ii) synthetic over-sampling, which generates samples in feature space. It

also suggests that it is possible to apply general data augmentation techniques in feature space if reasonable data augmentations for data are known.

5) **Neural Style Transfer:** It is another category of data augmentation, which can transfer the artist style of one image to another without changing semantics at a high level. It brings more variety to the training set. The main objective of this neural style transfer is to generate a third image from two images, where one image provides texture content and another provides high-level semantic content. We explore some of the SOTA augmentations for the sub-category.

- (i) **STaDA: Style Transfer as Data Augmentation :** This work [169] thoroughly evaluated different SOTA neural style transfer algorithms as data augmentation for image classification tasks. It shows significant performance gain on Caltech 101 [38] and Caltech 256 [48] datasets. Furthermore, it also combines neural style transfer algorithms with conventional data augmentation methods. A sample of this augmentation is shown in figure 55.



Fig. 55. Overview of the original image and two stylized images by STaDA. Image is taken from [169].

- (ii) **Style Augmentation: Data Augmentation via Style Randomization:** This work [66] proposed a novel data augmentation named style augmentation (SA) based on style neural transfer. SA randomizes the color, contrast, and texture while maintaining the shape and semantic content during the training. This is done by picking an arbitrary style transfer network for randomizing the style and by getting the target style from multivariate normal distribution embedding. It improves performance in three different tasks: classification, regression, and domain adaptation. The style augmentation sample is shown in figure 56.

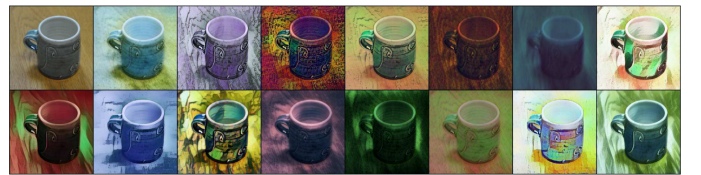


Fig. 56. Overview of Style augmentation applied to an image. The shape is preserved but the style, including color, texture, and contrast is randomized. Image is from [66].

- (iii) **StyPath: Style-Transfer Data Augmentation for Robust Histology Image Classification:** This paper [22] proposes a novel pipeline for Antibody Mediated Rejection (AMR) classification in kidneys based on StyPath data augmentation. StyPath is data augmentation that transfers style intending to reduce bias. The proposed augmentation is much faster than SOTA augmentations for AMR classification. Some samples are shown in figure 57.

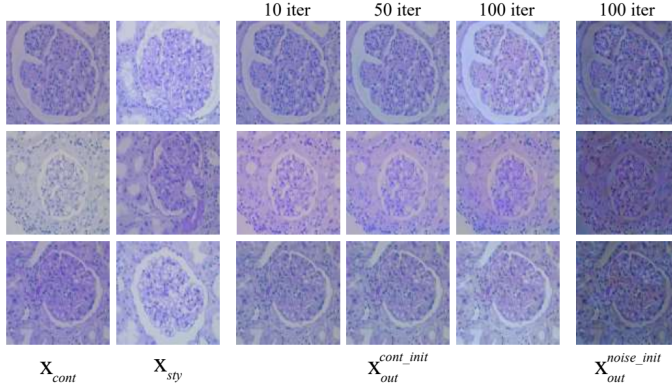


Fig. 57. Comparison of content and random initialization. Authors observe that output images initialized as the noise appeared distorted and discolored and failed to retain the content fidelity. Image is from [22].

- (iv) **A Neural Algorithm of Artistic Style :** This work [42] introduces an artificial system (AS) based on a deep neural network that generates artistic images of high perceptual quality. AS creates neural embedding then it uses the embedding to separate the style and content of the image and then recombines the content and style of target images to generate the artistic image. The sample is shown in figure 59
- (v) **Neural Style Transfer as Data Augmentation for Improving COVID-19 Diagnosis Classification :** This work [58] shows the effectiveness of a cycle GAN, which is mostly used for neural style transfer, augments COVID-19 negative x-ray image to convert into a positive COVID image to balance the dataset and also to increase the diversity of the dataset. It shows that augmenting the images with cycle GAN can improve performance over several different CNN architectures. A sample of this augmentation is shown in figure 58.

III. RESULTS

In this section, we provide the detailed result for various Computer Vision tasks such as image classification, object detection, and semantic segmentation. The main purpose is to show the effect of the data augmentation in CV different tasks and to do so, we compile results from various SOTA data augmentation works.

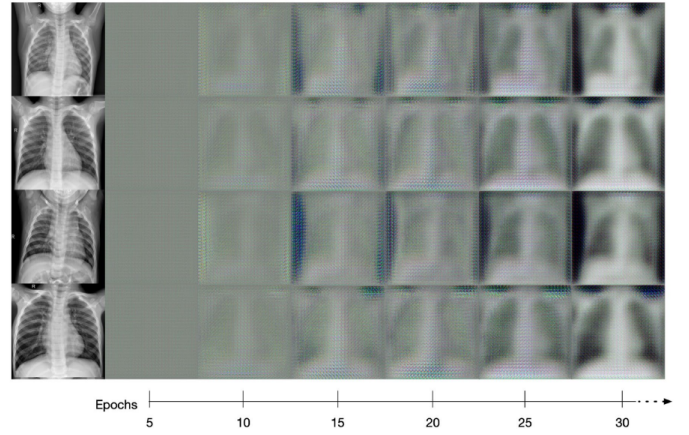


Fig. 58. Overview of generating synthetic COVID images from the healthy category. As the no of epochs grows the quality of the synthetic images improves. An example is from [58].



Fig. 59. Overview of the styled image by the neural algorithm. Image is from [42].

A. Image Classification

In this section, we present the result of several SOTA data augmentation methods for supervised learning and semi-supervised learning. Both are discussed below:

1) *supervised learning results:* In supervised learning, we have data on a large quantity that is fully labeled and we use this data to train the neural network (NN) model. In this section, we compile and compare the results from several SOTA data augmentation methods and put them in two different tables as shown in table II-B5 and table II. In table II-B5 results, ⁺ sign shows traditional data augmentations such as flipping, rotating, and cropping, have been used along with the SOTA augmentation methods. The used datasets are CIFAR10 [74], CIFAR100 [74] and ImageNet [26], and the used networks are wideresnet flavours [55], pyramid network flavours and several popular resnet flavours [55]. Accuracy is the evaluation metric used to compare the different algorithms used. The higher the accuracy, the better. As it can be in table II-B5 and table II, each data augmentation has significantly improved the accuracy.

2) *Semi-supervised learning:* Semi-supervised learning (SSL) is when we have a limited labeled data but unlabeled data is available on the large scale. Labeling the unlabeled data is tedious, time-consuming, and costly [79], [155]. To avoid these issues, SSL is used. There are several techniques

Method	Accuracies			
	CIFAR10	CIFAR10+	CIFAR100	CIFAR100+
ResNet-18 (Baseline)	89.37	95.28	63.32	77.54
ResNet-18 + CutOut	90.69	96.25	65.02	80.58
ResNet-18 + Random Erasing	95.28	95.32	-	-
ResNet-18 + CutMix	90.56	96.22	65.58	80.58
ResNet-18 + SaliencyMix	92.41	96.35	71.27	80.71
ResNet-18 + GridMask	95.28	96.54	-	-
ResNet-50 (Baseline)	87.86	95.02	63.52	78.42
ResNet-50 + CutOut	91.16	96.14	67.03	78.62
ResNet-50 + CutMix	90.84	96.39	68.35	81.28
ResNet-50 + SaliencyMix	93.19	96.54	75.11	81.43
WideResNet-28-10 (Baseline) [141]	93.03	96.13	73.94	81.20
WideResNet-28-10 + CutOut [29]	94.46	96.92	76.06	81.59
WideResNet-28-10 + Random Erasing	96.2	96.92	81.59	82.27
WideResNet-28-10 + GridMask	96.13	97.24	-	-
WideResNet-28-10 + CutMix	94.82	97.13	76.79	83.34
WideResNet-28-10 + PuzzleMix	-	-	-	83.77
WideResNet-28-10 + SaliencyMix	95.96	97.24	80.55	83.44

Note: + sign after dataset name shows that traditional data augmentation methods have been used

TABLE I

BASELINE PERFORMANCE COMPARISON OF VARIOUS AUGMENTATION ON CIFAR10 AND CIFAR100 DATASETS.

of SSL, but recently, data augmentation is employed with the limited labeled data to increase the diversity of the data. Data augmentation with SSL has increased the performance on different datasets and NN architectures. The used dataset are CIFAR10, CIFAR100, SVHN [103] and Mini-ImageNet. Several SSL techniques are used such as pseudoLabel, SSL with memory, label propagation, mean teacher, etc. We compile the results from many SOTA SSL methods with data augmentation and present them in this work. The effect of the data augmentation has also been shown with the different number of samples in SSL as shown in table III, table IV, and table V.

B. Object detection

In this section, we discuss the effectiveness of various image data augmentation techniques on the frequently used COCO2017 [92], PASCAL VOC [35], VOC 2007 [33], and VOC 2012 [34] datasets, which are commonly used for object detection tasks. We compile results from various SOTA data augmentation methods and put them in three different tables as shown in the table II-B5, VII, and VIII. FRCNN along with synthetic data gives the best mAP accuracy on VOC 2007 dataset as shown in table VII. Several classical and automatic data augmentation methods have shown promising performance using different SOTA models on the PASCAL VOC dataset as shown in table II-B5. The DetAdvProp achieves the highest score outperforming AutoAugment [23] on PASCAL VOC 2012 dataset as shown in the table VIII. The scores are in terms of mean average precision (mAP), average precision (AP) at the intersection over union (IOU) of 0.5 (AP50), and AP at IOU of 0.75 (AP75) metrics.

C. Semantic Segmentation

This subsection includes semantic segmentation results on PASCAL VOC and CITYSCAPES datasets, most frequently

used in several research papers. In table IX and table X, we compiled the effectiveness of validation set results on the different datasets with the effect of SOTA data augmentations on the semantic segmentation task. The results are reported in the term of mean intersection over union (mIoU) as the accuracy on the Cityscape dataset and PASCAL VOC dataset as shown in table IX and table X, respectively. We found performance gains on a few metrics such as mIoU and mAP, with several semantic segmentation models: deeplabv3+ [160], DeepLab-v2 [104], Xception-65 [160], ExFuse [166] and Eff-L2 [172]. It has been observed that incorporating data augmentation techniques can enhance the performance of semantic segmentation models. Notably, advanced image data augmentation methods have demonstrated greater improvements in performance compared to traditional techniques. Table IX and table X provide evidence of this improvement. The traditional data augmentations including rotation, scaling, flipping, and shifting [164].

IV. DISCUSSION AND FUTURE DIRECTIONS

A. Current approaches

It is proven that if we provide more data to the model, it improves model performance [50], [136]. A few current tendencies are discussed by Xu et al. [157]. Among these, one way is to collect the data and label it manually, but it is not an efficient way to do this. Another efficient way is to apply data augmentation, the more data augmentations we apply, the better improvement we get in terms of performance but to a certain extent. Currently, image mixing methods and autoaugment methods are successful for image classification tasks, scale aware based auto augment methods are showing promising results in detection tasks and semantic segmentation tasks. But these data augmentation performances can vary with the number of data augmentation applied, as it is known

Augmentation	CIFAR-10		CIFAR-100		ImageNet	
	Accuracy (%)	Model	Accuracy (%)	Model	Accuracy (%)	Model
Cutout [29]	97.04	WRN-28-10	81.59	WRN-28-10	77.1	ResNet-50
Random Erasing [170]	96.92	WRN-28-10	82.27	WRN-28-10	-	-
Hide-and-Seek [129]	95.53	ResNet-110	78.13	ResNet-110	77.20	ResNet-50
GridMask [15]	97.24	WRN-28-10	-	-	77.9	ResNet-50
LocalAugment [71]	-	-	95.92	WRN-22-10	76.87	ResNet-50
SelfMix [20]	96.62	PreActResNet-101	80.11	PreActResNet-101	-	-
KeepAugment [47]	97.8	ResNet-28-10	-	-	80.3	ResNet-101
Cut-Thumbnail [156]	97.8	ResNet-56	95.94	WRN-28-10	79.21	ResNet-50
MixUp [163]	97.3	WRN-28-10	82.5	WRN-28-10	77.9	ResNet-50
CutMix [162]	97.10	WRN-28-10	83.40	WRN-28-10	78.6	ResNet-50
SaliencyMix [141]	97.24	WRN-28-10	83.44	WRN-28-10	78.74	ResNet-50
PuzzleMix [70]	-	-	84.05	WRN-28-10	77.51	ResNet-50
FMix [52]	98.64	Pyramid	83.95	Dense	77.70	ResNet-101
MixMo [112]	96.38	WRN-28-10	82.40	WRN-28-10	-	-
StyleMix [59]	96.44	PyramidNet-200	85.83	PyramidNet-200	77.29	PyramidNet-200
RandomMix [94]	98.02	WRN-28-10	84.84	WRN-28-10	77.88	WRN-28-10
MixMatch [10]	95.05	WRN-28-10	74.12	WRN-28-10	-	-
ReMixMatch [9]	94.71	WRN-28-2	-	-	-	-
FixMatch [130]	95.69	WRN-28-2	77.04	WRN-28-2	-	-
AugMix [56]	-	-	-	-	77.6	ResNet-50
Improved Mixed-Example [135]	96.02	ResNet-18	80.3	ResNet-18	-	-
RICAP [137]	97.18	WRN-28-10	82.56	ResNet-28-10	78.62	WRN-50-2
ResizeMix [111]	97.60	WRN-28-10	84.31	WRN-28-10	79.00	ResNet-50
AutoAugment [23]	97.40	WRN-28-10	82.90	WRN-28-10	83.50	AmoebaNet-C
Fast AutoAugment [90]	98.00	SS(26 2×96d)	85.10	SS(26 2×96d)	80.60	ResNet-200
Faster AutoAugment [53]	98.00	SS(26 2 × 112d)	84.40	SS(26 2×96d)	75.90	ResNet-50
Local Patch AutoAugment [91]	98.10	SS(26 2 × 112d)	85.90	SS(26 2×96d)	81.00	ResNet-200
RandAugment [24]	98.50	PyramidNet	83.30	WRN-28-10	85.00	EfficientNet-B7

TABLE II

PERFORMANCE COMPARISON OF THE VARIOUS IMAGE ERASING AND IMAGE MIXING AUGMENTATIONS FOR IMAGE CLASSIFICATION PROBLEMS. WRN AND SS STAND FOR WIDERESNET AND SHAKE-SHAKE, RESPECTIVELY.

that the combined data augmentation methods show better performance than single one [108], [158].

B. Theoretical aspects

There is no theoretical support available to explain why specific augmentation is improving performance and which sample(s) should be augmented, as the same aspect has been discussed by Yang et al [158] and Shorten et al [123]. Like in random erasing, we randomly erase the region of the image - sometime may erase discriminating features, and the erased image makes no sense to a human. But the reason behind performance improvement is still unknown, which is another open challenge. Most of the time, we find the optimal parameters of the augmentation through an extensive number of experiments or we choose data augmentation based on our experience. But there should be a mechanism for choosing the data augmentation with theoretical support considering model architecture and dataset size. Researching the theoretical aspect is another open challenge for the research community.

C. Optimal number of samples generation

It is a known fact, as we increase data size, it improves the performance [50], [123], [136], [158] but it is not a case - increasing the number of samples will not improve performance after a certain number of samples [78]. What is the optimal number of samples to be generated, depending on the model architecture and dataset size, is a challenging aspect to be explored. Currently, researchers perform many experiments

to find the optimal number of sample generation [78]. But it is not feasible way as it requires time and computational cost. Can we devise a mechanism to find an optimal number of samples, which is an open research challenge?

D. Selection of data augmentation based on model architecture and dataset

Data augmentation selection depends on the nature of the dataset and model architecture. Like on MNIST [27] dataset, geometric transformations are not safe such as rotation on 6 and 9 digits will no longer preserve the label information. For densely parameterized CNN, it is easy to overfit weakly augmented datasets, and for shallow parameterized CNN, it may break generalization capability with data augmentation. It suggests, while selecting the data augmentation, the nature of the dataset and model architecture should be taken into account. Currently, numerous experiments are performed to find model architecture and suitable data augmentation for a specific dataset. Devising a systematic approach to select the data augmentation based on dataset and model architecture is another gap to be filled.

E. Augmentations for spaces

Most of the data augmentation approaches have been explored on the image level - data space. Very few research works have explored data on feature level - feature space. The challenge here arises, in which space should we apply data

TABLE III
COMPARISON ON CIFAR-10 AND SVHN. THE NUMBER REPRESENTS ERROR RATES ACROSS THREE RUNS.

Method	CIFAR-10				SVHN			
	40 labels	250 labels	1,000 labels	4,000 labels	40 labels	250 labels	1,000 labels	4,000 labels
VAT [101]	-	36.03 \pm 2.82	18.64 \pm 0.40	11.05 \pm 0.31	-	8.41 \pm 1.01	5.98 \pm 0.21	4.20 \pm 0.15
Mean Teacher [138]	-	47.32 \pm 4.71	17.32 \pm 4.00	10.36 \pm 0.25	-	6.45 \pm 2.43	3.75 \pm .10	3.39 \pm 0.11
MixMatch [10]	47.54 \pm 11.50	11.08 \pm .87	7.75 \pm .32	6.24 \pm .06	42.55 \pm 14.53	3.78 \pm .26	3.27 \pm .31	2.89 \pm .06
ReMixMatch [9]	19.10 \pm 9.64	6.27 \pm 0.34	5.73 \pm 0.16	5.14 \pm 0.04	3.34 \pm 0.20	3.10 \pm 0.50	2.83 \pm 0.30	2.42 \pm 0.09
UDA	29.05 \pm 5.93	8.76 \pm 0.90	5.87 \pm 0.13	5.29 \pm 0.25	52.63 \pm 20.51	2.76 \pm 0.17	2.55 \pm 0.09	2.47 \pm 0.15
SSL with Memory [17]	-	-	-	11.9 \pm 0.22	-	8.83	4.21	-
Deep Co-Training [110]	-	-	-	8.35 \pm 0.06	-	-	3.29 \pm 0.03	-
Weight Averaging [5]	-	-	15.58 \pm 0.12	9.05 \pm 0.21	-	-	-	-
ICT [142]	-	-	15.48 \pm 0.78	7.29 \pm 0.02	-	4.78 \pm 0.68	3.89 \pm 0.04	-
Label Propagation [64]	-	-	16.93 \pm 0.70	10.61 \pm 0.28	-	-	-	-
SNTG [96]	-	-	18.41 \pm 0.52	9.89 \pm 0.34	-	4.29 \pm 0.23	3.86 \pm 0.27	-
PLCB [4]	-	-	6.85 \pm 0.15	5.97 \pm 0.15	-	-	-	-
II-model [120]	-	53.02 \pm 2.05	31.53 \pm 0.98	17.41 \pm 0.37	-	17.65 \pm 0.27	8.60 \pm 0.18	5.57 \pm 0.14
PseudoLabel [85]	-	49.98 \pm 1.17	30.91 \pm 1.73	16.21 \pm 0.11	-	21.16 \pm 0.88	10.19 \pm 0.41	5.71 \pm 0.07
Mixup [163]	-	47.43 \pm 0.92	25.72 \pm 0.66	13.15 \pm 0.20	-	39.97 \pm 1.89	16.79 \pm 0.63	7.96 \pm 0.14
FeatMatch [81]	-	7.50 \pm 0.64	5.76 \pm 0.07	4.91 \pm 0.18	-	3.34 \pm 0.19	3.10 \pm 0.06	2.62 \pm 0.08
FixMatch [130]	13.81 \pm 3.37	5.07 \pm 0.65	-	4.26 \pm 0.05	3.96 \pm 2.17	2.48 \pm 0.38	2.28 \pm 0.11	-
SelfMatch [69]	93.19 \pm 1.08	95.13 \pm 0.26	-	95.94 \pm 0.08	96.58 \pm 1.02	97.37 \pm 0.43	97.49 \pm 0.07	-

TABLE IV
COMPARISON ON CIFAR-100 AND MINI-IMAGENET. THE NUMBER REPRESENTS ERROR RATES ACROSS TWO RUNS.

Method	CIFAR-100			mini-ImageNet	
	400 labels	4,000 labels	10,000 labels	4,000 labels	10,000 labels
II-model [120]	-	-	39.19 \pm 0.36	-	-
SNTG [96]	-	-	37.97 \pm 0.29	-	-
SSL with Memory [17]	-	-	34.51 \pm 0.61	-	-
Deep Co-Training [110]	-	-	34.63 \pm 0.14	-	-
Weight Averaging [5]	-	-	33.62 \pm 0.54	-	-
Mean Teacher [138]	-	45.36 \pm 0.49	36.08 \pm 0.51	72.51 \pm 0.22	57.55 \pm 1.11
Label Propagation [64]	-	43.73 \pm 0.20	35.92 \pm 0.47	70.29 \pm 0.81	57.58 \pm 1.47
PLCB [4]	-	37.55 \pm 1.09	32.15 \pm 0.50	56.49 \pm 0.51	46.08 \pm 0.11
FeatMatch	-	31.06 \pm 0.41	26.83 \pm 0.04	39.05 \pm 0.06	34.79 \pm 0.22
MixMatch	67.61 \pm 1.32	-	28.31 \pm 0.33	-	-
UDA	59.28 \pm 0.88	-	24.50 \pm 0.25	-	-
ReMixMatch	44.28 \pm 2.06	-	23.03 \pm 0.56	-	-
FixMatch	48.85 \pm 1.75	-	22.60 \pm 0.12	-	-

augmentation, data space, or feature space? It is another interesting aspect that can be explored. For the current approaches, it seems like it depends on the dataset, model architecture, and task. Currently, approaches are conducting experiments in data space and feature space and then selecting the best one [154]. It is not the optimal way to find data augmentation for specific space. It is still an open challenge to be solved.

F. Open research questions

Despite the success of data augmentation techniques in different Computer Vision tasks, it still failed to solve challenges in SOTA data augmentation techniques. After thoroughly reviewing SOTA data augmentation approaches, we found several challenges and difficulties, which are yet to be solved, as it is listed below:

- In image mixing techniques, label smoothing has been used. It makes sense whatever portion of images is mixed, corresponding labels should be mixed accordingly. To the best of our knowledge, none has explored label

smoothing for image manipulation and image erasing subcategories - where the image part is lost. For example, if the image portion is randomly cut out in cutout data augmentation, the corresponding label should be mixed. It is an interesting open research question.

- Currently, data augmentation is performed without considering the importance of an example. All examples may not be difficult for the neural network to learn, but some are. Thus, augmentation should be applied to those difficult examples by measuring the importance of the examples. How neural network behave if data augmentation is applied to those difficult examples?
- In image mixing data augmentations, if we mix more than two images salient parts, that are truly participating in augmentation unlike RICAP [137], what is its effect in terms of accuracy and robustness against adversarial attacks? Note, the corresponding labels of these images will be mixed accordingly.
- In random data augmentation under the auto augmen-

TABLE V
COMPARISON OF TEST ERROR RATES ON CIFAR-10 & SVHN USING WIDERESNET-28 AND CNN-13.

Approach	Method	CIFAR-10 ($N_I=4000$)	SVHN($N_I=1000$)
WideResNet-28			
Pseudo Labeling	Supervised	20.26 \pm 0.38	12.83 \pm 0.47
	PL [85]	17.78 \pm 0.57	7.62 \pm 0.29
	PL-CB [4]	6.28 \pm 0.3	-
Consistency	II Model [83]	16.37 \pm 0.63	7.19 \pm 0.27
	Mean Teacher [138]	15.87 \pm 0.28	5.65 \pm 0.47
	VAT [101]	13.86 \pm 0.27	5.63 \pm 0.20
Regularization	VAT + EntMin [101]	13.13 \pm 0.39	5.35 \pm 0.19
	LGA + VAT [65]	12.06 \pm 0.19	6.58 \pm 0.36
	ICT [142]	7.66 \pm 0.17	3.53 \pm 0.07
Pseudo Labeling	MixMatch [10]	6.24 \pm 0.06	3.27 \pm 0.31
	UDA	5.29 \pm 0.25	2.46 \pm 0.17
	ReMixMatch (Berthelot et al. 2020)	5.14 \pm 0.04	2.42 \pm 0.09
Consistency	FixMatch [130]	4.26 \pm 0.05	2.28 \pm 0.11
	CL	8.92 \pm 0.03	5.65 \pm 0.11
	CL+FA [90]	5.51 \pm 0.14	2.90 \pm 0.19
Regularization	CL+FA [90]+Mixup [163]	5.09 \pm 0.18	2.75 \pm 0.15
	CL+RA+Mixup [163]	5.27 \pm 0.16	2.80 \pm 0.188
CNN-13			
Pseudo Labeling	TSSDL-MT	9.30 \pm 0.55	3.35 \pm 0.27
	LP-MT	10.61 \pm 0.28	-
	Ladder net [117]	12.36 \pm 0.31	-
Consistency	MeanTeacher [138]	12.31 \pm 0.24	3.95 \pm 0.19
	Temporal ensembling [83]	12.16 \pm 0.24	4.42 \pm 0.16
	VAT [101]	11.36 \pm 0.34	5.42
Regularization	NATEntMin [101]	10.55 \pm 0.05	3.86
	SNTG [96]	10.93 \pm 0.14	3.86 \pm 0.27
	ICT [142]	7.29 \pm 0.02	2.89 \pm 0.04
Pseudo Labeling	CL	9.81 \pm 0.22	4.75 \pm 0.28
	CL+RA	5.92 \pm 0.07	3.96 \pm 0.10

tation category, the order of augmentations has not been explored. We believe it has a significant importance. What are the possible ways to explore the order of existing augmentations such as first traditional data augmentations and then image mixing or weight-based?

- Finding an optimal and an ordered number of data augmentation, and the optimal number of samples to be augmented are open challenges. For example, in randAug method, there are N optimal number of augmentations found but it is not known how many, in which order and what samples should be augmented?

V. CONCLUSION

This survey provides a comprehensive overview of state-of-the-art (SOTA) data augmentation techniques for addressing overfitting in computer vision tasks due to limited data. A detailed taxonomy of image data augmentation approaches is presented, along with an overview of each SOTA method and the results of its application to various computer vision tasks such as image classification, object detection, and semantic segmentation. The results for both supervised and semi-supervised learning are also compiled for easy comparison purposes. In addition, the available code for each data augmentation approach is provided to facilitate result reproducibility. The difficulties and challenges of data augmentation are also discussed, along with promising open research questions that

have the potential to further advance the field. This survey is expected to benefit researchers in several ways: (i) a deeper understanding of data augmentation, (ii) the ability to easily compare results, and (iii) the ability to reproduce results with available code.

ACKNOWLEDGMENT

This research was supported by Science Foundation Ireland under grant numbers 18/CRT/6223 (SFI Centre for Research Training in Artificial intelligence), SFI/12/RC/2289/P_2 (Insight SFI Research Centre for Data Analytics), 13/RC/2094/P_2 (Lero SFI Centre for Software) and 13/RC/2106/P_2 (ADAPT SFI Research Centre for AI-Driven Digital Content Technology). For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission.

REFERENCES

- [1] Jiwoon Ahn, Sunghyun Cho, and Suha Kwak. Weakly supervised learning of instance segmentation with inter-pixel relations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2209–2218, 2019.
- [2] Jiwoon Ahn and Suha Kwak. Learning pixel-level semantic affinity with image-level supervision for weakly supervised semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4981–4990, 2018.

Method	Detector	BackBone	AP	AP ₅₀	AP ₇₅	AP _s	AP _m	AP _l
Hand-crafted:								
Dropblock [44]	RetinaNet	ResNet-50	38.4	56.4	41.2	—	—	—
AutoAugment+color Ops [171]	RetinaNet	ResNet-50	37.5	-	-	—	—	—
geometric Ops [171]	RetinaNet	ResNet-50	38.6	-	-	—	—	—
bbox-only Ops [171]	RetinaNet	ResNet-50	39.0	-	-	—	—	—
Mix-up [167]	Faster R-CNN	ResNet-101	41.1	-	-	-	-	-
PSIS* [144]	Faster R-CNN	ResNet-101	40.2	61.1	44.2	22.3	45.7	51.6
Stitcher [19]	Faster R-CNN	ResNet-101	42.1	-	-	26.9	45.5	54.1
GridMask [15]	Faster R-CNN	ResNeXt-101	42.6	65.0	46.5	-	-	-
InstaBoost* [37]	Mask R-CNN	ResNet-101	43.0	64.3	47.2	24.8	45.9	54.6
SNIP (MS test)* [127]	Faster R-CNN	ResNet-101-DCN-C4	44.4	66.2	49.9	27.3	47.4	56.9
SNIPER (MS test)* [128]	Faster R-CNN	ResNet-101-DCN-C4	46.1	67.0	51.6	29.6	48.9	58.1
Traditional Aug [158]	Faster R-CNN	ResNet-101	36.80	58.0	40.0	-	-	-
Traditional Aug* [31]	CenterNet	ResNet-101	41.15	58.01	45.30	-	-	-
Traditional Aug+ [15]	Faster-RCNN	50-FPN (2x)	37.4	58.7	40.5	-	-	-
Traditional Aug+ [15]	Faster-RCNN	50-FPN (2x)+GridMask (p = 0.3)	38.2	60.0	41.4	-	-	-
Traditional Aug+ [15]	Faster-RCNN	50-FPN (2x)+ GridMask (p = 0.5)	38.1	60.1	41.2	-	-	-
Traditional Aug+ [15]	Faster-RCNN	50-FPN (2x)+ GridMask (p = 0.7)	38.3	60.4	41.7	-	-	-
Traditional Aug+ [15]	Faster-RCNN	50-FPN (2x)+ GridMask (p = 0.9)	38.0	60.1	41.2	-	-	-
Traditional Aug+ [15]	Faster-RCNN	50-FPN (4x)	35.7	56.0	38.3	-	-	-
Traditional Aug+ [15]	Faster-RCNN	50-FPN (4x)+ GridMask (p = 0.7)	39.2	60.8	42.2	-	-	-
Traditional Aug+ [15]	Faster-RCNN	X101-FPN (1x))	41.2	63.3	44.8	-	-	-
Traditional Aug+ [15]	Faster-RCNN	X101-FPN (2x))	40.4	62.2	43.8	-	-	-
Traditional Aug+ [15]	Faster-RCNN	X101-FPN (2x)+ GridMask (p = 0.7))	42.6	65.0	46.5	-	-	-
Traditional Aug+ [15]	Faster-RCNN	X101-FPN (2x)+ GridMask (p = 0.7))	42.6	65.0	46.5	-	-	-
KeepAugment: [47]	Faster R-CNN	ResNet50-C4	39.5	—	—	—	—	—
KeepAugment: [47]	Faster R-CNN	ResNet50-FPN	40.7	—	—	—	—	—
KeepAugment: [47]	RetinaNet	ResNet50-FPN	39.1	—	—	—	—	—
KeepAugment: [47]	Faster R-CNN	ResNet101-C4	42.2	—	—	—	—	—
KeepAugment: [47]	Faster R-CNN	ResNet101-FPN	42.9	—	—	—	—	—
KeepAugment: [47]	RetinaNet	ResNet101-FPN	41.2	—	—	—	—	—
DADAAugment: [88]	RetinaNet	ResNet-50	35.9	55.8	38.4	19.9	38.8	45.0
DADAAugment: [88]	RetinaNet	ResNet-50(DADA)	36.6	56.8	39.2	20.2	39.7	46.0
DADAAugment: [88]	Faster R-CNN	ResNet-50	36.6	58.8	39.6	21.6	39.8	45.0
DADAAugment: [88]	Faster R-CNN	ResNet-50 (DADA)	37.2	59.1	40.2	22.2	40.2	45.7
DADAAugment: [88]	Mask R-CNN	ResNet-50	37.4	59.3	40.7	22.2	40.6	46.3
DADAAugment: [88]	Mask R-CNN	ResNet-50(DADA)	37.8	59.6	41.1	22.4	40.9	46.6
AutoAugment: [16]	EfficientDet D0	EfficientNet B0	34.4	52.8	36.7	53.1	40.2	13.9
Det-AdvProp: [16]	EfficientDet D0	EfficientNet B0	34.7	52.9	37.2	54.1	40.6	13.9
AutoAugment: [16]	EfficientDet D1	EfficientNet B1	40.1	59.2	43.2	57.9	45.7	19.9
Det-AdvProp: [16]	EfficientDet D1	EfficientNet B1	40.5	59.2	43.3	58.8	46.2	20.6
AutoAugment: [16]	EfficientDet D2	EfficientNet B2	43.5	62.8	46.6	59.8	48.7	23.9
Det-AdvProp: [16]	EfficientDet D2	EfficientNet B2	43.8	62.6	47.3	61.0	49.6	25.6
AutoAugment: [16]	EfficientDet D3	EfficientNet B3	47.0	66.0	50.8	63.0	51.7	29.8
Det-AdvProp: [16]	EfficientDet D3	EfficientNet B3	47.6	66.3	51.4	64.0	52.2	30.2
AutoAugment: [16]	EfficientDet D4	EfficientNet B4	49.5	68.7	53.7	64.9	54.0	31.9
Det-AdvProp: [16]	EfficientDet D4	EfficientNet B4	49.8	68.6	54.2	65.2	54.2	32.4
AutoAugment: [16]	EfficientDet D5	EfficientNet B5	51.5	70.4	56.0	65.2	56.1	35.4
Det-AdvProp: [16]	EfficientDet D5	EfficientNet B5	51.8	70.7	56.3	66.1	56.2	36.2
Automatic:								
AutoAug-det [171]	RetinaNet	ResNet-50	39.0	-	-	-	-	-
AutoAug-det [171]	RetinaNet	ResNet-101	40.4	-	-	-	-	-
AutoAugment [23]	RetinaNet	ResNet-200	42.1	-	-	-	-	-
AutoAug-det' [171]	RetinaNet	ResNet-50	40.3	60.0	43.0	23.6	43.9	53.8
RandAugmnet* [24]	RetinaNet	ResNet-200	41.9	-	-	-	-	-
AutoAug-det [171]	RetinaNet	ResNet-101	41.8	61.5	44.8	24.4	45.9	55.9
RandAug [24]	RetinaNet	ResNet-101	40.1	-	-	-	-	-
RandAug? [10]	RetinaNet	ResNet-101	41.4	61.4	44.5	25.0	45.4	54.2
Scale-aware AutoAug [18]	RetinaNet	ResNet-50	41.3	61.0	44.1	25.2	44.5	54.6
Scale-aware AutoAug	RetinaNet	ResNet-101	43.1	62.8	46.0	26.2	46.8	56.7
Scale-aware AutoAug	Faster R-CNN	ResNet-101	44.2	65.6	48.6	29.4	47.9	56.7
Scale-aware AutoAug (MS test)	Faster R-CNN	ResNet-101-DCN-C4	47.0	68.6	52.1	32.3	49.3	60.4
Scale-aware AutoAug	FCOS	ResNet-101	44.0	62.7	47.3	28.2	47.8	56.1
Scale-aware AutoAug	FCOS	ResNeXt-32x8d-101-DCN	48.5	67.2	52.8	31.5	51.9	63.0
Scale-aware AutoAug (1200 size)	FCOS	ResNeXt-32x8d-101-DCN	49.6	68.5	54.1	35.7	52.5	62.4
Scale-aware AutoAug (MS Test)	ResNeXt-32x8d-101-DCN	FCOS	51.4	69.6	57.0	37.4	54.2	65.1

TABLE VI
DATA AUGMENTATION EFFECT ON DIFFERENT OBJECT DETECTION METHODS USING PASCAL VOC DATASET

Method	TSet	mAP	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv
FRCN [45]	7	66.9	74.5	78.3	69.2	53.2	36.6	77.3	78.2	82.0	40.7	72.7	67.9	79.6	79.2	73.0	69.0	30.1	65.4	70.2	75.8	65.8
FRCN* [148]	7	69.1	75.4	80.8	67.3	59.9	37.6	81.9	80.0	84.5	50.0	77.1	68.2	81.0	82.5	74.3	69.9	28.4	71.1	70.2	75.8	66.6
ASDN [148]	7	71.0	74.4	81.3	67.6	57.0	46.6	81.0	79.3	86.0	52.9	75.9	73.7	82.6	83.2	77.7	72.7	37.4	66.3	71.2	78.2	74.3
IRE	7	70.5	75.9	78.9	69.0	57.7	46.4	81.7	79.5	82.9	49.3	76.9	67.9	81.5	83.3	76.7	73.2	40.7	72.8	66.9	75.4	74.2
ORE	7	71.0	75.1	79.8	69.7	60.8	46.0	80.4	79.0	83.8	51.6	76.2	67.8	81.2	83.7	76.8	73.8	43.1	70.8	67.4	78.3	75.6
I+ORE	7	71.5	76.1	81.6	69.5	60.1	45.6	82.2	79.2	84.5	52.5	78.7	71.6	80.4	83.3	76.7	73.9	39.4	68.9	69.8	79.2	77.4
FRCN [45]	7+12	70.0	77.0	78.1	69.3	59.4	38.3	81.6	78.6	86.7	42.8	78.8	68.9	84.7	82.0	76.6	69.9	31.8	70.1	74.8	80.4	70.4
FRCN* [148]	7+12	74.8	78.5	81.0	74.7	67.9	53.4	85.6	84.4	86.2	57.4	80.1	72.2	85.2	84.2	77.6	76.1	45.3	75.7	72.3	81.8	77.3
IRE	7+12	75.6	79.0	84.1	76.3	66.9	52.7	84.5	84.4	88.7	58.0	82.9	71.1	84.8	84.4	78.6	76.7	45.5	77.1	76.3	82.5	76.8
ORE	7+12	75.8	79.4	81.6	75.6	66.5	52.7	85.5	84.7	88.3	58.7	82.9	72.8	85.0	84.3	79.3	76.3	46.3	76.3	74.9	86.0	78.2
I+ORE	7+12	76.2	79.6	82.5	75.7	70.5	55.1	85.2	84.4	88.4	58.6	82.6	73.9	84.2	84.7	78.8	76.3	46.7	77.9	75.9	83.3	79.3
SSD	7+12	77.4	81.7	85.4	75.7	69.6	49.9	84.9	85.8	87.4	61.5	82.3	79.2	86.6	87.1	84.7	78.9	50.0	77.4	79.1	86.2	76.3
SSD+ SD (1x) [145]	7+12	78.1	83.2	84.5	76.1	72.1	50.2	85.2	86.3	87.8	63.7	82.8	80.1	85.2	87.2	84.8	80.0	51.5	77.0	82.0	86.1	76.9
SSD + SD(2x) [145]	7+12	78.3	83.6	85.0	76.2	72.0	51.3	85.1	87.2	87.6	64.2	82.5	81.9	85.5	86.5	85.9	81.2	51.2	72.3	82.8	86.9	78.4
SSD +SD(3x) [145]	7+12	77.8	80.4	85.0	76.3	70.1	50.4	84.8	86.3	88.2	61.0	83.5	79.5	87.2	86.9	85.9	78.8	51.2	76.9	79.4	86.5	77.9
FRCN [45]	7+12	73.2	76.5	79.0	70.9	65.5	52.1	83.1	84.7	86.4	52.0	81.9	65.7	84.8	84.6	77.5	76.7	38.8	73.6	73.9	83.0	72.6
FRCN+SD(1x) [156]	7	79.9	85.1	86.6	78.6	75.7	65.2	83.5	88.4	88.9	65.8	83.6	74.3	86.4	84.7	85.5	88.0	62.0	75.5	75.3	87.7	76.3

TABLE VII

VOC 2007 TEST DETECTION AVERAGE PRECISION (%). FRCN* REFERS TO FRCN WITH TRAINING SCHEDULE IN [148] AND SD REFERS TO SYNTHETIC DATA

Model	mAP	AP50	AP75
EfficientDet-D0	55.6	77.6	61.4
+ AutoAugment	55.7 (+0.1)	77.7 (+0.1)	61.8 (+0.4)
+ Det-AdvProp	55.9 (+0.3)	77.9 (+0.3)	62.0 (+0.6)
EfficientDet-D1	60.8	82.0	66.7
+ AutoAugment	61.0 (+0.2)	82.2 (+0.2)	67.2 (+0.5)
+ Det-AdvProp	61.2 (+0.4)	82.3 (+0.3)	67.4 (+0.7)
EfficientDet-D2	63.3	83.6	69.3
+ AutoAugment	62.7 (-0.6)	83.3 (-0.3)	69.2 (-0.1)
+ Det-AdvProp	63.5 (+0.2)	83.8 (+0.2)	69.7 (+0.4)
EfficientDet-D3	65.7	85.3	71.8
+ AutoAugment	65.2 (-0.5)	85.1 (-0.2)	71.3 (-0.5)
+ Det-AdvProp	66.2 (+0.5)	85.9 (+0.6)	72.5 (+0.7)
EfficientDet-D4	67.0	86.0	73.0
+ AutoAugment	67.0 (+0.0)	86.3 (+0.3)	73.5 (+0.5)
+ Det-AdvProp	67.5 (+0.5)	86.6 (+0.6)	74.0 (+1.0)
EfficientDet-D5	67.4	86.9	73.8
+ AutoAugment	67.6 (+0.2)	87.2 (+0.3)	74.2 (+0.4)
+ Det-AdvProp	68.2 (+0.8)	87.6 (+0.7)	74.7 (+0.9)

TABLE VIII

RESULTS ON PASCAL VOC 2012. THE PROPOSED DETADVPROP GIVES THE HIGHEST SCORE ON EVERY MODEL AND METRIC. IT LARGELY OUTPERFORMS AUTOAUGMENT [23] WHEN FACING DOMAIN SHIFT.

- [3] Sidra Aleem, Teerath Kumar, Suzanne Little, Malika Bendeche, Rob Brennan, and Kevin McGuinness. Random data augmentation based enhancement: A generalized enhancement approach for medical datasets. 2022.
- [4] Eric Arazo, Diego Ortego, Paul Albert, Noel E O'Connor, and Kevin McGuinness. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020.
- [5] Ben Athiwaratkun, Marc Finzi, Pavel Izmailov, and Andrew Gordon Wilson. Improving consistency-based semi-supervised learning with weight averaging. *arXiv preprint arXiv:1806.05594*, 2(9):11, 2018.
- [6] Soroush Baseri Saadi, Nazanin Tataei Sarshar, Soroush Sadeghi, Ramin Ranjbarzadeh, Mersedeh Kooshki Forooshani, and Malika Bendeche. Investigation of effectiveness of shuffled frog-leaping optimizer in training a convolution neural network. *Journal of Healthcare Engineering*, 2022, 2022.
- [7] Markus Bayer, Marc-André Kaufhold, and Christian Reuter. A survey on data augmentation for text classification. *ACM Computing Surveys*, 2021.
- [8] Sima Behpour, Kris M Kitani, and Brian D Ziebart. Ada: Adversarial data augmentation for object detection. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1243–1252. IEEE, 2019.
- [9] David Berthelot, Nicholas Carlini, Ekin D Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. Remixmatch: Semi-supervised learning with distribution alignment and augmentation anchoring. *arXiv preprint arXiv:1911.09785*, 2019.
- [10] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in neural information processing systems*, 32, 2019.
- [11] Aisha Chandio, Gong Gui, Teerath Kumar, Irfan Ullah, Ramin Ranjbarzadeh, Arunabha M Roy, Akhtar Hussain, and Yao Shen. Precise single-stage detector. *arXiv preprint arXiv:2210.04252*, 2022.
- [12] Aisha Chandio, Yao Shen, Malika Bendeche, Irum Inayat, and Teerath Kumar. Audd: audio urdu digits dataset for automatic audio urdu digit recognition. *Applied Sciences*, 11(19):8842, 2021.
- [13] Arslan Chaudhry, Puneet K Dokania, and Philip HS Torr. Discovering class-specific pixels for weakly-supervised semantic segmentation. *arXiv preprint arXiv:1707.05821*, 2017.
- [14] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018.
- [15] Pengguang Chen, Shu Liu, Hengshuang Zhao, and Jiaya Jia. Gridmask data augmentation. *arXiv preprint arXiv:2001.04086*, 2020.
- [16] Xiangning Chen, Cihang Xie, Mingxing Tan, Li Zhang, Cho-Jui Hsieh, and Boqing Gong. Robust and accurate object detection via adversarial learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16622–16631, 2021.
- [17] Yanbei Chen, Xiatian Zhu, and Shaogang Gong. Semi-supervised deep learning with memory. In *Proceedings of the European conference on computer vision (ECCV)*, pages 268–283, 2018.
- [18] Yukang Chen, Yanwei Li, Tao Kong, Lu Qi, Ruihang Chu, Lei Li, and Jiaya Jia. Scale-aware automatic augmentation for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9563–9572, 2021.
- [19] Yukang Chen, Peizhen Zhang, Zeming Li, Yanwei Li, Xiangyu Zhang, Gaofeng Meng, Shiming Xiang, Jian Sun, and Jiaya Jia. Stitcher: Feedback-driven data provider for object detection. *arXiv preprint arXiv:2004.12432*, 2(7):12, 2020.
- [20] Jaehyeop Choi, Chaehyeon Lee, Donggyu Lee, and Heechul Jung. Salfmix: A novel single image-based data augmentation technique using a saliency map. *Sensors*, 21(24):8444, 2021.
- [21] Peng Chu, Xiao Bian, Shaopeng Liu, and Haibin Ling. Feature space augmentation for long-tailed data. In *European Conference on Computer Vision*, pages 694–710. Springer, 2020.
- [22] Pietro Antonio Cicalese, Aryan Mobiny, Pengyu Yuan, Jan Becker, Chandra Mohan, and Hien Van Nguyen. Stypath: Style-transfer data augmentation for robust histology image classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 351–361. Springer, 2020.
- [23] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 113–123, 2019.
- [24] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 702–703, 2020.
- [25] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [26] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- [27] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- [28] Terrance DeVries and Graham W Taylor. Dataset augmentation in feature space. *arXiv preprint arXiv:1702.05538*, 2017.
- [29] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.
- [30] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [31] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6569–6578, 2019.
- [32] Dumitru Erhan, Aaron Courville, Yoshua Bengio, and Pascal Vincent. Why does unsupervised pre-training help deep learning? In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 201–208. JMLR Workshop and Conference Proceedings, 2010.
- [33] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>.
- [34] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>.

Method	Model	1/8	1/4	1/2	7/8	Full
SDA [160]	DeepLabV3Plus	74.1	-	-	-	-
SDA + DSBN [160]	DeepLabV3Plus	69.5	-	-	-	-
SDA [160]	DeepLabV3Plus	-	-	-	-	78.7
SDA + DSBN [160]	DeepLabV3Plus	-	-	-	-	79.2
SDA [160]	DeepLabV3Plus	-	-	-	71.4	-
SDA + DSBN [160]	DeepLabV3Plus	-	-	-	72.5	-
AdvSemi [62]	DeepLabV2	58.8	62.3	65.7	-	66.0
S4GAN + MT [100]	DeepLabV2	59.3	61.9	-	-	65.8
CutMix [41]	DeepLabV2	60.3	63.87	-	-	67.7
DST-CBC [40]	DeepLabV2	60.5	64.4	-	-	66.9
ClassMix [104]	DeepLabV2	61.4	63.6	66.3	-	66.2
ECS [99]	DeepLabV3Plus	67.4	70.7	72.9	-	74.8
DSBN [160]	DeepLabV2	67.6	69.3	70.7	-	70.1
SSBN [160]	DeepLabV3Plus	74.1	77.8	78.7	-	78.7
Adversarial [62]	DeepLab-v2	-	58.8	62.3	65.7	-
s4GAN [100]	DeepLab-v2	-	59.3	61.9	-	65.8
French et al [41]	DeepLab-v2	51.20	60.34	63.87	-	-
DST-CBC [40]	DeepLab-v2	48.7	60.5	64.4	-	-
ClassMix-Seg [104]	DeepLab-v2	54.07	61.35	63.63	66.29	-
DeepLab V3plus [164]	MobileNet	-	-	-	-	73.5
DeepLab V3plus [164]	ResNet-50	-	-	-	-	76.9
DeepLab V3plus [164]	ResNet-101	-	-	-	-	78.5
Baseline+ CutOut (16x16, p = 1) [164]	MobileNet	-	-	-	-	72.8
Baseline+ CutMix (p = 1) [164]	MobileNet	-	-	-	-	72.6
Baseline+ ObjectAug [164]	MobileNet	-	-	-	-	73.5

TABLE IX
RESULTS OF PERFORMANCE (mIoU) ON CITYSCAPES VALIDATION SET

- [35] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88:303–308, 2009.
- [36] Junsong Fan, Zhaoxiang Zhang, Chunfeng Song, and Tieniu Tan. Learning integral objects with intra-class discriminator for weakly-supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4283–4292, 2020.
- [37] Hao-Shu Fang, Jianhua Sun, Runzhong Wang, Minghao Gou, Yong-Lu Li, and Cewu Lu. Instaboost: Boosting instance segmentation via probability map guided copy-pasting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 682–691, 2019.
- [38] Li Fei-Fei, Robert Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE transactions on pattern analysis and machine intelligence*, 28(4):594–611, 2006.
- [39] Steven Y Feng, Varun Gangal, Dongyeop Kang, Teruko Mitamura, and Eduard Hovy. Genaug: Data augmentation for finetuning text generators. *arXiv preprint arXiv:2010.01794*, 2020.
- [40] Zhengyang Feng, Qianyu Zhou, Guangliang Cheng, Xin Tan, Jianping Shi, and Lizhuang Ma. Semi-supervised semantic segmentation via dynamic self-training and classbalanced curriculum. *arXiv preprint arXiv:2004.08514*, 1(2):5, 2020.
- [41] Geoff French, Timo Aila, Samuli Laine, Michal Mackiewicz, and Graham Finlayson. Semi-supervised semantic segmentation needs strong, high-dimensional perturbations. 2019.
- [42] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015.
- [43] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2918–2928, 2021.
- [44] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Dropblock: A regularization method for convolutional networks. *Advances in neural information processing systems*, 31, 2018.
- [45] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.
- [46] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.
- [47] Chengyue Gong, Dilin Wang, Meng Li, Vikas Chandra, and Qiang Liu. Keepaugment: A simple information-preserving data augmentation approach. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1055–1064, 2021.
- [48] Gregory Griffin, Alex Holub, and Pietro Perona. Caltech-256 object category dataset. 2007.
- [49] Jian Guo and Stephen Gould. Deep cnn ensemble with data augmentation for object detection. *arXiv preprint arXiv:1506.07224*, 2015.
- [50] Alon Halevy, Peter Norvig, and Fernando Pereira. The unreasonable effectiveness of data. *IEEE intelligent systems*, 24(2):8–12, 2009.
- [51] Junlin Han, Pengfei Fang, Weihao Li, Jie Hong, Mohammad Ali Armin, Ian Reid, Lars Petersson, and Hongdong Li. You only cut once: Boosting data augmentation with a single cut. *arXiv preprint arXiv:2201.12078*, 2022.
- [52] Ethan Harris, Antonia Marcu, Matthew Painter, Mahesan Niranjan, Adam Prügel-Bennett, and Jonathon Hare. Fmix: Enhancing mixed sample data augmentation. *arXiv preprint arXiv:2002.12047*, 2020.
- [53] Ryuichiro Hataya, Jan Zdenek, Kazuki Yoshizoe, and Hideki Nakayama. Faster autoaugment: Learning augmentation strategies using backpropagation. In *European Conference on Computer Vision*, pages 1–16. Springer, 2020.
- [54] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [55] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [56] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. *arXiv preprint arXiv:1912.02781*, 2019.
- [57] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15262–15271, 2021.
- [58] Netzahualcoyotl Hernandez-Cruz, David Cato, and Jesus Favela. Neural style transfer as data augmentation for improving covid-19 diagnosis classification. *SN Computer Science*, 2(5):1–12, 2021.
- [59] Minui Hong, Jinwoo Choi, and Gunhee Kim. Stylemix: Separating content and style for enhanced data augmentation. In *Proceedings of the*

Method	Model	1/100	1/50	1/20	1/8	1/4	Full
GANSeg [131]	VGG16	-	-	-	-	64.1	
AdvSemSeg [62]	ResNet-101	-	-	-	-	68.4	
CCT [105]	ResNet-50	-	-	-	-	69.4	
PseudoSeg [173]	ResNet-101	-	-	-	-	73.2	
DSBN [160]	ResNet-101	-	-	-	-	75.0	
DSBN [160]	Xception-65	-	-	-	-	79.3	
Fully supervised [160]	ResNet-101	-	-	-	-	78.3	
Fully supervised [160]	Xception-65	-	-	-	-	79.2	
Adversarial [62]	DeepLab-v2	-	57.2	64.7	69.5	72.1	-
s4GAN [100]	DeepLab-v2	-	63.3	67.2	71.4	-	75.6
French et.al [41]	DeepLab-v2	53.79	64.81	66.48	67.60	-	-
DST-CBC [40]	DeepLab-v2	61.6	65.5	69.3	70.7	71.8	-
ClassMix:Seg* [104]	DeepLab-v2	54.18	66.15	67.77	71.00	72.45	-
Mixup [163]	IRNet	-	-	-	-	-	49
CutOut [29]	IRNet	-	-	-	-	-	48.9
CutMix [162]	IRNet	-	-	-	-	-	49.2
Random pasting [134]	IRNet	-	-	-	-	-	49.8
CCNN [107]	VGG16	-	-	-	-	-	35.6
SEC [73]	VGG16	-	-	-	-	-	51.1
STC [151]	VGG16	-	-	-	-	-	51.2
AdvEra [150]	VGG16	-	-	-	-	-	55.7
DCSP [13]	ResNet101	-	-	-	-	-	61.9
MDC [152]	VGG16	-	-	-	-	-	60.8
MCOF [147]	ResNet101	-	-	-	-	-	61.2
DSRG [61]	ResNet101	-	-	-	-	-	63.2
AffinityNet [2]	ResNet-38	-	-	-	-	-	63.7
IRNet [1]	ResNet50	-	-	-	-	-	64.8
FickleNet [86]	ResNet101	-	-	-	-	-	65.3
SEAM [149]	ResNet38	-	-	-	-	-	65.7
ICD [36]	ResNet101	-	-	-	-	-	64.3
IRNet + CDA [134]	ResNet50	-	-	-	-	-	66.4
SEAM + CDA [134]	ResNet38	-	-	-	-	-	66.8
DeepLab V3 [164]	MobileNet	-	-	-	-	-	71.9
DeepLab V3 [164]	ResNet-50	-	-	-	-	-	77.8
DeepLab V3 [164]	ResNet-101	-	-	-	-	-	78.4
DeepLab V3plus [164]	MobileNet	-	-	-	-	-	73.8
DeepLab V3plus [164]	ResNet-50	-	-	-	-	-	78.8
DeepLab V3plus [164]	ResNet-101	-	-	-	-	-	79.6
Baseline+R.Rotation [164]	ObjectAug	-	-	-	-	-	69.5
Baseline +R.Scailing [164]	ObjectAug	-	-	-	-	-	70.3
Baseline + R.Flipping [164]	ObjectAug	-	-	-	-	-	69.6
Baseline + R.Shifting [164]	ObjectAug	-	-	-	-	-	70.7
Baseline + All [164]	ObjectAug	-	-	-	-	-	73.8
Baseline + CutOut (16×16, p = 0.5) [164]	MobileNet	-	-	-	-	-	71.9
Baseline + CutOut (16×16, p = 1) [164]	MobileNet	-	-	-	-	-	72.3
Baseline + CutMix (p = 0.5) [164]	MobileNet	-	-	-	-	-	72.7
Baseline + CutMix (p = 1) [164]	MobileNet	-	-	-	-	-	72.4
Baseline + ObjectAug [164]	MobileNet	-	-	-	-	-	73.8
Baseline + CutOut (16×16, p=0.5) + ObjectAug [164]	MobileNet	-	-	-	-	-	73.9
Baseline + CutMix (p=0.5) + ObjectAug [164]	MobileNet	-	-	-	-	-	74.1
DeepLabv3+ [14]	EfficientNet-B7	-	-	-	-	-	84.6
ExFuse [166]	EfficientNet-B7	-	-	-	-	-	85.8
Eff-B7 [172]	EfficientNet-B7	-	-	-	-	-	85.2
Eff-L2 [172]	EfficientNet-B7	-	-	-	-	-	88.7
Eff-B7 NAS-FPN [43]	EfficientNet-B7	-	-	-	-	-	83.9
Eff-B7 NAS-FPN w/ Copy-Paste pre-training [43]	EfficientNet-B7	-	-	-	-	-	86.6

TABLE X

RESULTS OF PERFORMANCE MEAN INTERSECTION OVER UNION (MIOU) ON THE PASCAL VOC 2012 VALIDATION SET

- IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14862–14870, 2021.
- [60] Shaoli Huang, Xinchao Wang, and Dacheng Tao. Snapmix: Semantically proportional mixing for augmenting fine-grained data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1628–1636, 2021.
 - [61] Zilong Huang, Xinggang Wang, Jiasi Wang, Wenyu Liu, and Jingdong Wang. Weakly-supervised semantic segmentation network with deep seeded region growing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7014–7023, 2018.
 - [62] Wei-Chih Hung, Yi-Hsuan Tsai, Yan-Ting Liou, Yen-Yu Lin, and Ming-Hsuan Yang. Adversarial learning for semi-supervised semantic segmentation. *arXiv preprint arXiv:1802.07934*, 2018.
 - [63] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
 - [64] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Label propagation for deep semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5070–5079, 2019.
 - [65] Jacob Jackson and John Schulman. Semi-supervised learning by label gradient alignment. *arXiv preprint arXiv:1902.02336*, 2019.
 - [66] Philip TG Jackson, Amir Atapour Abarghouei, Stephen Bonner, Toby P Breckon, and Boguslaw Obara. Style augmentation: data augmentation via style randomization. In *CVPR workshops*, volume 6, pages 10–11, 2019.
 - [67] Wisal Khan, Kislay Raj, Teerath Kumar, Arunabha M Roy, and Bin Luo. Introducing urdu digits dataset with demonstration of an efficient and robust noisy decoder-based pseudo example generator. *Symmetry*, 14(10):1976, 2022.
 - [68] Cherry Khosla and Baljit Singh Saini. Enhancing performance of deep learning models with different data augmentation techniques: A survey. In *2020 International Conference on Intelligent Engineering and Management (ICIEM)*, pages 79–85. IEEE, 2020.
 - [69] Byoungjip Kim, Jinho Choo, Yeong-Dae Kwon, Seongho Joe, Seungjai Min, and Youngjune Gwon. Selfmatch: Combining contrastive self-supervision and consistency for semi-supervised learning. *arXiv preprint arXiv:2101.06480*, 2021.
 - [70] Jang-Hyun Kim, Wonho Choo, and Hyun Oh Song. Puzzle mix: Exploiting saliency and local statistics for optimal mixup. In *International Conference on Machine Learning*, pages 5275–5285. PMLR, 2020.
 - [71] Youmin Kim, AFM Shahab Uddin, and Sung-Ho Bae. Local augment: Utilizing local bias property of convolutional neural networks for data augmentation. *IEEE Access*, 9:15191–15199, 2021.
 - [72] Tom Ko, Vijayaditya Peddinti, Daniel Povey, and Sanjeev Khudanpur. Audio augmentation for speech recognition. In *Sixteenth annual conference of the international speech communication association*, 2015.
 - [73] Alexander Kolesnikov and Christoph H Lampert. Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In *European conference on computer vision*, pages 695–711. Springer, 2016.
 - [74] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
 - [75] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
 - [76] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
 - [77] Teerath Kumar, Alessandra Mileo, Rob Brennan, and Malika Ben-dechache. Rsmda: Random slices mixing data augmentation. *Applied Sciences*, 13(3):1711, 2023.
 - [78] Teerath Kumar, Jinbae Park, Muhammad Salman Ali, AFM Uddin, and Sung-Ho Bae. Class specific autoencoders enhance sample diversity. *Journal of Broadcast Engineering*, 26(7):844–854, 2021.
 - [79] Teerath Kumar, Jinbae Park, Muhammad Salman Ali, AFM Shahab Uddin, Jong Hwan Ko, and Sung-Ho Bae. Binary-classifiers-enabled filters for semi-supervised learning. *IEEE Access*, 9:167663–167673, 2021.
 - [80] Teerath Kumar, Jinbae Park, and Sung-Ho Bae. Intra-class random erasing (icre) augmentation for audio classification. In *Proceedings Of The Korean Society Of Broadcast Engineers Conference*, pages 244–247. The Korean Institute of Broadcast and Media Engineers, 2020.
 - [81] Chia-Wen Kuo, Chih-Yao Ma, Jia-Bin Huang, and Zsolt Kira. Feat-match: Feature-based augmentation for semi-supervised learning. In *European Conference on Computer Vision*, pages 479–495. Springer, 2020.
 - [82] Jiss Kuruvilla, Dhanya Sukumaran, Anjali Sankar, and Siji P Joy. A review on image processing and image segmentation. In *2016 international conference on data mining and advanced computing (SAPIENCE)*, pages 198–203. IEEE, 2016.
 - [83] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *arXiv preprint arXiv:1610.02242*, 2016.
 - [84] Misha Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Reinforcement learning with augmented data. *Advances in Neural Information Processing Systems*, 33:19884–19895, 2020.
 - [85] Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3, page 896, 2013.
 - [86] Jungbeom Lee, Eunji Kim, Sungmin Lee, Jangho Lee, and Sungroh Yoon. Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5267–5276, 2019.
 - [87] Victor Lempitsky, Pushmeet Kohli, Carsten Rother, and Toby Sharp. Image segmentation with a bounding box prior. In *2009 IEEE 12th international conference on computer vision*, pages 277–284. IEEE, 2009.
 - [88] Yonggang Li, Guosheng Hu, Yongtao Wang, Timothy Hospedales, Neil M Robertson, and Yongxin Yang. Dada: Differentiable automatic data augmentation. *arXiv preprint arXiv:2003.03780*, 2020.
 - [89] Junhao Liew, Yunchao Wei, Wei Xiong, Sim-Heng Ong, and Jiashi Feng. Regional interactive image segmentation networks. In *2017 IEEE international conference on computer vision (ICCV)*, pages 2746–2754. IEEE Computer Society, 2017.
 - [90] Sungbin Lim, Ildoo Kim, Taesup Kim, Chiheon Kim, and Sungwoong Kim. Fast autoaugment. *Advances in Neural Information Processing Systems*, 32, 2019.
 - [91] Shiqi Lin, Tao Yu, Ruoyu Feng, Xin Li, Xin Jin, and Zhibo Chen. Local patch autoaugment with multi-agent collaboration. *arXiv preprint arXiv:2103.11099*, 2021.
 - [92] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.
 - [93] Pei Liu, Xuemin Wang, Chao Xiang, and Weiye Meng. A survey of text data augmentation. In *2020 International Conference on Computer Communication and Network Security (CCNS)*, pages 191–195. IEEE, 2020.
 - [94] Xiaoliang Liu, Furao Shen, Jian Zhao, and Changhai Nie. Randommix: A mixed sample data augmentation method with multiple mixed modes. *arXiv preprint arXiv:2205.08728*, 2022.
 - [95] Xiaolong Liu, Zhidong Deng, and Yuhang Yang. Recent progress in semantic image segmentation. *Artificial Intelligence Review*, 52(2):1089–1106, 2019.
 - [96] Yucen Luo, Jun Zhu, Mengxi Li, Yong Ren, and Bo Zhang. Smooth neighbors on teacher graphs for semi-supervised learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8896–8905, 2018.
 - [97] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
 - [98] Sachin Mehta, Saeid Naderiparizi, Fartash Faghri, Maxwell Horton, Lailin Chen, Ali Farhadi, Oncel Tuzel, and Mohammad Rastegari. Rangeaugment: Efficient online augmentation with range learning. *arXiv preprint arXiv:2212.10553*, 2022.
 - [99] Robert Mendel, Luis Antonio de Souza, David Rauber, Joao Paulo Papa, and Christoph Palm. Semi-supervised segmentation based on error-correcting supervision. In *European Conference on Computer Vision*, pages 141–157. Springer, 2020.
 - [100] Sudhanshu Mittal, Maxim Tatarchenko, and Thomas Brox. Semi-supervised semantic segmentation with high-and low-level consistency. *IEEE transactions on pattern analysis and machine intelligence*, 43(4):1369–1379, 2019.

- [101] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1979–1993, 2018.
- [102] Loris Nanni, Gianluca Maguolo, and Michelangelo Paci. Data augmentation approaches for improving animal audio classification. *Ecological Informatics*, 57:101084, 2020.
- [103] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [104] Viktor Olsson, Wilhelm Tranheden, Juliano Pinto, and Lennart Svensson. Classmix: Segmentation-based data augmentation for semi-supervised learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1369–1378, 2021.
- [105] Yassine Ouali, Céline Hudelot, and Myriam Tami. Semi-supervised semantic segmentation with cross-consistency training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12674–12684, 2020.
- [106] Jinbae Park, Teerath Kumar, and Sung-Ho Bae. Search of an optimal sound augmentation policy for environmental sound classification with deep neural networks. In *Proceedings Of The Korean Society Of Broadcast Engineers Conference*, pages 18–21. The Korean Institute of Broadcast and Media Engineers, 2020.
- [107] Deepak Pathak, Philipp Krahenbuhl, and Trevor Darrell. Constrained convolutional neural networks for weakly supervised segmentation. In *Proceedings of the IEEE international conference on computer vision*, pages 1796–1804, 2015.
- [108] Porntiwa Pawara, Emmanuel Okafor, Lambert Schomaker, and Marco Wiering. Data augmentation for plant classification. In *International conference on advanced concepts for intelligent vision systems*, pages 615–626. Springer, 2017.
- [109] Luis Perez and Jason Wang. The effectiveness of data augmentation in image classification using deep learning. *arXiv preprint arXiv:1712.04621*, 2017.
- [110] Siyuan Qiao, Wei Shen, Zhishuai Zhang, Bo Wang, and Alan Yuille. Deep co-training for semi-supervised image recognition. In *Proceedings of the european conference on computer vision (eccv)*, pages 135–152, 2018.
- [111] Jie Qin, Jiemin Fang, Qian Zhang, Wenyu Liu, Xingang Wang, and Xinggang Wang. Resizemix: Mixing data with preserved object information and true labels. *arXiv preprint arXiv:2012.11101*, 2020.
- [112] Alexandre Ramé, Rémy Sun, and Matthieu Cord. Mixmo: Mixing multiple inputs for multiple outputs via deep subnetworks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 823–833, 2021.
- [113] Ramin Ranjbarzadeh, Shadi Dorosti, Saeid Jafarzadeh Ghouschi, Annalina Caputo, Erfan Babaee Tirkolaee, Sadia Samar Ali, Zahra Arshadi, and Malika Bendeche. Breast tumor localization and segmentation using machine learning techniques: Overview of datasets, findings, and methods. *Computers in Biology and Medicine*, page 106443, 2022.
- [114] Ramin Ranjbarzadeh, Saeid Jafarzadeh Ghouschi, Nazanin Tataei Sarshar, Erfan Babaee Tirkolaee, Sadia Samar Ali, Teerath Kumar, and Malika Bendeche. Me-ccnn: Multi-encoded images and a cascade convolutional neural network for breast tumor segmentation and recognition. *Artificial Intelligence Review*, pages 1–38, 2023.
- [115] Ramin Ranjbarzadeh, Nazanin Tataei Sarshar, Saeid Jafarzadeh Ghouschi, Mohammad Saleh Esfahani, Mahboub Parhizkar, Yaghoob Pourasad, Shokofeh Anari, and Malika Bendeche. Mrfe-cnn: multi-route feature extraction model for breast tumor segmentation in mammograms using a convolutional neural network. *Annals of Operations Research*, pages 1–22, 2022.
- [116] Ramin Ranjbarzadeh, Payam Zarbakhsh, Annalina Caputo, Erfan Babaee Tirkolaee, and Malika Bendeche. Brain tumor segmentation based on an optimized convolutional neural network and an improved chimp optimization algorithm. *Available at SSRN 4295236*, 2022.
- [117] Antti Rasmus, Mathias Berglund, Mikko Honkala, Harri Valpola, and Tapani Raiko. Semi-supervised learning with ladder networks. *Advances in neural information processing systems*, 28, 2015.
- [118] Arunabha M Roy, Jayabrata Bhaduri, Teerath Kumar, and Kislay Raj. Wildect-yolo: An efficient and robust computer vision-based accurate object localization model for automated endangered wildlife detection. *Ecological Informatics*, page 101919, 2022.
- [119] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- [120] Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. *Advances in neural information processing systems*, 29, 2016.
- [121] Jin-Woo Seo, Hong-Gyu Jung, and Seong-Whan Lee. Self-augmentation: Generalizing deep networks to unseen classes for few-shot learning. *Neural Networks*, 138:140–149, 2021.
- [122] Ling Shao, Fan Zhu, and Xuelong Li. Transfer learning for visual categorization: A survey. *IEEE transactions on neural networks and learning systems*, 26(5):1019–1034, 2014.
- [123] Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of big data*, 6(1):1–48, 2019.
- [124] Connor Shorten, Taghi M Khoshgoftaar, and Borko Furht. Text data augmentation for deep learning. *Journal of big Data*, 8(1):1–34, 2021.
- [125] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [126] Aditya Singh, Ramin Ranjbarzadeh, Kislay Raj, Teerath Kumar, and Arunabha M Roy. Understanding eeg signals for subject-wise definition of armoni activities. *arXiv preprint arXiv:2301.00948*, 2023.
- [127] Bharat Singh and Larry S Davis. An analysis of scale invariance in object detection snip. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3578–3587, 2018.
- [128] Bharat Singh, Mahyar Najibi, and Larry S Davis. Sniper: Efficient multi-scale training. *Advances in neural information processing systems*, 31, 2018.
- [129] Krishna Kumar Singh, Hao Yu, Aron Sarmasi, Gautam Pradeep, and Yong Jae Lee. Hide-and-seek: A data augmentation technique for weakly-supervised localization and beyond. *arXiv preprint arXiv:1811.02545*, 2018.
- [130] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in Neural Information Processing Systems*, 33:596–608, 2020.
- [131] Nasim Souly, Concetto Spampinato, and Mubarak Shah. Semi supervised semantic segmentation using generative adversarial network. In *Proceedings of the IEEE international conference on computer vision*, pages 5688–5696, 2017.
- [132] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- [133] Xingzhe Su. A survey on data augmentation methods based on gan in computer vision. In *The International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery*, pages 852–865. Springer, 2020.
- [134] Yukun Su, Ruizhou Sun, Guosheng Lin, and Qingyao Wu. Context decoupling augmentation for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7004–7014, 2021.
- [135] Cecilia Summers and Michael J Dinneen. Improved mixed-example data augmentation. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1262–1270. IEEE, 2019.
- [136] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*, pages 843–852, 2017.
- [137] Ryo Takahashi, Takashi Matsubara, and Kuniaki Uehara. Ricap: Random image cropping and patching data augmentation for deep cnns. In *Asian conference on machine learning*, pages 786–798. PMLR, 2018.
- [138] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017.
- [139] Nazanin Tataei Sarshar, Ramin Ranjbarzadeh, Saeid Jafarzadeh Ghouschi, Gabriel Gomes de Oliveira, Shokofeh Anari, Mahboub Parhizkar, and Malika Bendeche. Glioma brain

- tumor segmentation in four mri modalities using a convolutional neural network and based on a transfer learning method. In *Proceedings of the 7th Brazilian Technology Symposium (BTSym'21) Emerging Trends in Human Smart and Sustainable Future of Cities (Volume 1)*, pages 386–402. Springer, 2022.
- [140] Muhammad Turab, Teerath Kumar, Malika Bendeche, and Takfari-nas Saber. Investigating multi-feature selection and ensembling for audio classification. *arXiv preprint arXiv:2206.07511*, 2022.
- [141] AFM Uddin, Mst Monira, Wheemyung Shin, TaeChoong Chung, Sung-Ho Bae, et al. Saliencymix: A saliency guided data augmentation strategy for better regularization. *arXiv preprint arXiv:2006.01791*, 2020.
- [142] Vikas Verma, Kenji Kawaguchi, Alex Lamb, Juho Kannala, Yoshua Bengio, and David Lopez-Paz. Interpolation consistency training for semi-supervised learning. *arXiv preprint arXiv:1903.03825*, 2019.
- [143] Riccardo Volpi, Pietro Morerio, Silvio Savarese, and Vittorio Murino. Adversarial feature augmentation for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5495–5504, 2018.
- [144] Hao Wang, Qilong Wang, Fan Yang, Weiqi Zhang, and Wangmeng Zuo. Data augmentation for object detection via progressive and selective instance-switching. *arXiv preprint arXiv:1906.00358*, 2019.
- [145] Ke Wang, Bin Fang, Jiye Qian, Su Yang, Xin Zhou, and Jie Zhou. Perspective transformation data augmentation for object detection. *IEEE Access*, 8:4935–4943, 2019.
- [146] Xiang Wang, Kai Wang, and Shiguo Lian. A survey on face data augmentation for the training of deep neural networks. *Neural computing and applications*, 32(19):15503–15531, 2020.
- [147] Xiang Wang, Shaodi You, Xi Li, and Huimin Ma. Weakly-supervised semantic segmentation by iteratively mining common object features. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1354–1362, 2018.
- [148] Xiaolong Wang, Abhinav Shrivastava, and Abhinav Gupta. A-fast-rcnn: Hard positive generation via adversary for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2606–2615, 2017.
- [149] Yude Wang, Jie Zhang, Meina Kan, Shiguang Shan, and Xilin Chen. Self-supervised equivariant attention mechanism for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12275–12284, 2020.
- [150] Yunchao Wei, Jiashi Feng, Xiaodan Liang, Ming-Ming Cheng, Yao Zhao, and Shuicheng Yan. Object region mining with adversarial erasing: A simple classification to semantic segmentation approach. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1568–1576, 2017.
- [151] Yunchao Wei, Xiaodan Liang, Yunpeng Chen, Xiaohui Shen, Ming-Ming Cheng, Jiashi Feng, Yao Zhao, and Shuicheng Yan. Stc: A simple to complex framework for weakly-supervised semantic segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(11):2314–2320, 2016.
- [152] Yunchao Wei, Huaxin Xiao, Honghui Shi, Zequn Jie, Jiashi Feng, and Thomas S Huang. Revisiting dilated convolution: A simple approach for weakly- and semi-supervised semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7268–7277, 2018.
- [153] Karl Weiss, Taghi M Khoshgoufar, and DingDing Wang. A survey of transfer learning. *Journal of Big data*, 3(1):1–40, 2016.
- [154] Sebastien C Wong, Adam Gatt, Victor Stamatescu, and Mark D McDonnell. Understanding data augmentation for classification: when to warp? In *2016 international conference on digital image computing: techniques and applications (DICTA)*, pages 1–6. IEEE, 2016.
- [155] Shasha Xie, Hui Lin, and Yang Liu. Semi-supervised extractive speech summarization via co-training algorithm. In *Eleventh Annual Conference of the International Speech Communication Association*, 2010.
- [156] Tianshu Xie, Xuan Cheng, Xiaomin Wang, Minghui Liu, Jiali Deng, Tao Zhou, and Ming Liu. Cut-thumbnail: A novel data augmentation for convolutional neural network. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 1627–1635, 2021.
- [157] Mingle Xu, Sook Yoon, Alvaro Fuentes, and Dong Sun Park. A comprehensive survey of image augmentation techniques for deep learning. *arXiv preprint arXiv:2205.01491*, 2022.
- [158] Suorong Yang, Weikang Xiao, Mengcheng Zhang, Suhan Guo, Jian Zhao, and Furao Shen. Image data augmentation for deep learning: A survey. *arXiv preprint arXiv:2204.08610*, 2022.
- [159] Jaeyun Yoo, Namhyuk Ahn, and Kyung-Ah Sohn. Rethinking data augmentation for image super-resolution: A comprehensive analysis and a new strategy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8375–8384, 2020.
- [160] Jianlong Yuan, Yifan Liu, Chunhua Shen, Zhibin Wang, and Hao Li. A simple baseline for semi-supervised semantic segmentation with strong data augmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8229–8238, 2021.
- [161] Fei Yue, Chao Zhang, MingYang Yuan, Chen Xu, and YaLin Song. Survey of image augmentation based on generative adversarial network. In *Journal of Physics: Conference Series*, volume 2203, page 012052. IOP Publishing, 2022.
- [162] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019.
- [163] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- [164] Jiawei Zhang, Yanchun Zhang, and Xiaowei Xu. Objectaug: object-level data augmentation for semantic image segmentation. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2021.
- [165] Xiaofeng Zhang, Zhangyang Wang, Dong Liu, Qifeng Lin, and Qing Ling. Deep adversarial data augmentation for extremely low data regimes. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(1):15–28, 2020.
- [166] Zhenli Zhang, Xiangyu Zhang, Chao Peng, Xiangyang Xue, and Jian Sun. Exfuse: Enhancing feature fusion for semantic segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 269–284, 2018.
- [167] Zhi Zhang, Tong He, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of freebies for training object detection neural networks. *arXiv preprint arXiv:1902.04103*, 2019.
- [168] Zhengli Zhao, Dheeru Dua, and Sameer Singh. Generating natural adversarial examples. *arXiv preprint arXiv:1710.11342*, 2017.
- [169] Xu Zheng, Tejo Chalasani, Koustav Ghosal, Sebastian Lutz, and Aljosa Smolic. Stada: Style transfer as data augmentation. *arXiv preprint arXiv:1909.01056*, 2019.
- [170] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 13001–13008, 2020.
- [171] Barret Zoph, Ekin D Cubuk, Golnaz Ghiasi, Tsung-Yi Lin, Jonathon Shlens, and Quoc V Le. Learning data augmentation strategies for object detection. In *European conference on computer vision*, pages 566–583. Springer, 2020.
- [172] Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin Dogus Cubuk, and Quoc Le. Rethinking pre-training and self-training. *Advances in neural information processing systems*, 33:3833–3845, 2020.
- [173] Yuliang Zou, Zizhao Zhang, Han Zhang, Chun-Liang Li, Xiao Bian, Jia-Bin Huang, and Tomas Pfister. Pseudoseg: Designing pseudo labels for semantic segmentation. *arXiv preprint arXiv:2010.09713*, 2020.