# Use the Detection Transformer as a Data Augmenter

Luping Wang and Bin Liu⋆

Research Center for Applied Mathematics and Machine Intelligence,
Zhejiang Lab, Hangzhou 311121, China
{wangluping,liubin}@zhejianglab.com

**Abstract.** Detection Transformer (DETR) is a Transformer architecture based object detection model. In this paper, we demonstrate that it can also be used as a data augmenter. We term our approach as DETR assisted CutMix, or DeMix for short. DeMix builds on CutMix, a simple yet highly effective data augmentation technique that has gained popularity in recent years. CutMix improves model performance by cutting and pasting a patch from one image onto another, yielding a new image. The corresponding label for this new example is specified as the weighted average of the original labels, where the weight is proportional to the area of the patch. CutMix selects a random patch to be cut. In contrast, DeMix elaborately selects a semantically rich patch, located by a pre-trained DETR. The label of the new image is specified in the same way as in CutMix. Experimental results on benchmark datasets for image classification demonstrate that DeMix significantly outperforms prior art data augmentation methods including CutMix. Oue code is available at https://github.com/ZJLAB-AMMI/DeMix.

**Keywords:** detection transformer · object detection · data augmentation · CutMix · image classification

## 1 Introduction

Data augmentation is a technique used in machine learning where existing data is modified or transformed to create new data. The principle behind data augmentation is that even small changes to existing data can create new useful examples for training. For example, flipping an image horizontally can create a new training example that is still representative of the original object. In general, the goal of data augmentation is to increase the diversity of the training data while preserving its underlying structure. By introducing variations into the training data, models can learn to generalize better and perform well on unseen data. Data augmentation has been around for decades, but recent advances in deep learning have made it more widely used and effective. It has been commonly used in computer vision tasks like image classification [9, 22, 26], object detection [23, 41, 42], and segmentation [5, 11, 33, 39], as well as in natural

---
⋆ Corresponding author

language processing [2, 6, 18]. As a powerful technique that can significantly improve model performance with little additional effort or cost, its widespread adoption and continued development are indicative of its importance in modern machine learning.

There are many typical methods of data augmentation, including flipping, rotating, scaling, cropping, adding noise, and changing color, which can be applied randomly or according to specific rules, depending on the task and the desired results [27]. Recent years also witnessed notable developments in data augmentation techniques, among which CutMix [35], a method that combines parts of multiple images to create new training samples, is our focus in this work.

CutMix is a simple yet highly effective data augmentation technique that has gained popularity in recent years. Given a pair of training examples $A$ and $B$ and their associated labels $y_A$ and $y_B$, CutMix selects a random crop region of $A$ and replaces it with a patch of the same size cut from $B$, yielding a new data example. The corresponding label for this new example is specified as the weighted average of the original labels, where the weight is proportional to the area of the patch. As the patch to be cut is randomly selected, it may totally come from the background or the area of the object or an area that mixes the background and the object, while the contribution of this patch to the label of the resulting new image is deterministic. We argue that this is not reasonable. For example, if the patch selected to be cut is all from the background, then its contribution to the label of the resulting new image should be negligible, while using CutMix, its contribution is proportional to its area.

Motivated by the aforementioned flaw of CutMix, we propose a knowledge-guided CutMix, where the knowledge comes from a pre-trained detection transformer (DETR) model. We term our approach DETR assisted CutMix, or DeMix for short. DETR is an object detection model based on the Transformer architecture [3, 13, 24]. Unlike traditional object detection models, DETR directly models the object detection task as a set matching problem, and uses a Transformer encoder to process input images and a decoder to generate object sets, achieving end-to-end object detection. DeMix takes advantage of the following desirable properties of DETR

- Given an image input to a pre-trained DETR, it can provide an estimate of the class, bounding box position, and corresponding confidence score for each object involved in this image;
- It can detect objects of different numbers and sizes simultaneously.

As the pre-trained DETR model is trained based on datasets that are different from our target dataset, the class labels it provides can be meaningless, while bounding box positions it provides are surprisingly informative for us to use in DeMix. DeMix cuts the image patch associated with one of the bounding boxes, given by DETR, in an image example, then resizes and pastes it onto a random crop region of another example, to create the new example. The label of this new example is specified in the same way as CutMix.

In principle, DeMix provides an elaborate improvement to CutMix by borrowing knowledge from a pre-trained DETR. Even if the knowledge borrowed is

totally inaccurate or meaningless, then DeMix reduces to CutMix. Our major contributions can be summarized as follows

- We demonstrates how DETR can be used as a tool for data augmentation, resulting in a new approach DeMix;
- We evaluate the performance of DeMix on several different fine-grained image classification datasets. Experimental results demonstrate that DeMix significantly outperforms all competitor methods, including CutMix;
- Our work sheds light on how to improve a general-purpose technique (data augmentation here) using a special-purpose pre-trained model (DETR here) without fine-tuning.

## 2 Related Works

In this section, we briefly introduce DETR followed by related works on data augmentation techniques, especially CutMix, which are closely related to our DeMix.

### 2.1 DETR

Detection Transformer (DETR) is a type of object detection model based on the Transformer architecture [3, 13, 24]. Unlike traditional object detection models, DETR directly models the object detection task as a set matching problem, and uses a Transformer encoder to process input images and a decoder to generate object sets, achieving end-to-end object detection [3]. In DETR, the image is divided into a set of small patches and each patch is mapped to a feature vector using a convolutional neural network (CNN) encoder. Then, the output of the encoder is passed as input to the Transformer encoder to allow interactions among the features in both spatial and channel dimensions. Finally, the decoder converts the output of the Transformer encoder into an object set that specifies the class, bounding box position, and corresponding confidence score for each object. DETR does not require predefined operations such as non-maximum suppression or anchor boxes. It can detect objects of different numbers and sizes simultaneously. Fig. 1 shows detection results given by a pre-trained DETR model on some image examples in the dataset used in our experiment.

DETR is developed for object detection, while, in this paper, we demonstrate that it could also be used as a data augmenter. Specifically, we leverage DETR to locate a semantically rich patch associated with an object for use in CutMix.

### 2.2 Data Augmentation

Data augmentation is a technique used in machine learning to increase the size of a training dataset by generating new examples from existing ones. The goal of data augmentation is to improve the generalization ability of machine learning models by introducing more variations in the training data.

Fig. 1: Object detection with DETR

There are many existing works on data augmentation, which can be broadly categorized into four groups as follows.

1. Traditional methods: These include commonly used techniques such as random cropping, resizing, flipping, rotating, color jittering, and adding noise [27]. The strength of traditional methods is that they are simple and easy to implement, and can effectively generate new samples from existing ones. However, they may not always work well for complex datasets or tasks, as they do not account for higher-level features and structures.

2. Generative methods: These use generative models, such as variational autoencoders (VAEs) and generative adversarial networks (GANs), to synthesize new data samples [4, 32, 34]. The strength of generative methods is that they can generate high-quality and diverse samples that are similar to the real data distribution, which can help improve model generalization. However, they require large amounts of computational resources and training data to build the generative models, which can be a limitation in some cases.

3. Adversarial methods: These use adversarial attacks to perturb the existing data samples to generate new ones [25, 38, 40]. The strength of adversarial methods is that they can generate realistic and diverse augmented data samples, which can help improve model robustness. However, they require careful design and tuning of hyperparameters to avoid overfitting.

4. Methods that create new training examples by mixing parts of multiple images together such as Mixup [36] and CutMix [35]. Mixup linearly combines data from different examples, while CutMix cuts and pastes patches of images. These methods are much simpler to implement than generative and adversarial methods, while can achieve much better performance than traditional methods.

In summary, each data augmentation method has its own strengths and limitations.

DeMix proposed here is an algorithmic improvement to the aforementioned method CutMix, which has gained popularity in recent years, since it is simple to implement, while can achieve much better performance than traditional

methods. In this paper, we propose DeMix, which is as simple as CutMix, while outperforms it significantly.

In CutMix, the image patch to be cut is randomly selected, while the contribution of this patch to the label of the resulting new image is deterministically proportional to its area. As aforementioned, this is not reasonable. For example, if the patch to be cut is all from the background, then its contribution to the label of the resulting new image should be negligible, while using CutMix, it can be large.

Several work has been proposed to address this flaw of CutMix, where the basic idea is to select a semantically rich instead of a random patch to be cut and paste, see e.g., SaliencyMix [30], the class activation map (CAM) based method [37], Keepaugment [12], and SnapMix [17]. All these advanced methods require image pre-processing, such as computing the saliency map or CAM, prior to data augmentation, while, in contrast, DeMix does not need any pre-processing operation before generating new image examples.

## 3   DeMix: DETR Assisted CutMix

In this section, we describe our DeMix method in detail. Given a pair of image examples, DeMix uses two operations to generate a new image. The first operation employs a pre-trained DETR to identify bounding box positions for each object in the source image denoted as $x_B$. The second operation cuts an image patch associated with one bounding box, resize it, and then pastes it onto a randomly selected crop region of the other image, termed target image and denoted as $x_A$, yielding the new image example. The label of this new image is a weighted average of labels of the original images. Fig. 2 illustrates the operations that make up DeMix.

In mathematical terms, the process of generating a new image can be explained as follows.:

$$\tilde{x} = (\mathbf{1} - \mathbf{M}_\lambda) \odot x_A + T(\mathbf{M_B} \odot x_B)$$
$$\tilde{y} = (1 - \lambda)y_A + \lambda y_B \tag{1}$$

where $y_A$ and $y_B$ denote labels of $x_A$ and $x_B$, respectively, in the form of one-hot vectors, $\mathbf{M}_\lambda$ and $\mathbf{M_B}$ are binary mask matrices with dimensions $W \times H$, $W$ and $H$ denote the width and height of the images, $0 < \lambda < 1$ is a hyper-parameter defined as the ratio of the area of the randomly selected crop region of $x_A$ to the full area of $x_A$, $\mathbf{M}_\lambda$ denotes the binary mask matrix that defines the aforementioned crop region, $\mathbf{M_B}$ is the binary mask matrix associated to the object bounding box given by DETR, $\odot$ denotes the element-wise multiplication operator, $\mathbf{1}$ represents a matrix of an appropriate size whose elements are all 1, and finally $T(\cdot)$ denotes a linear transformation that aligns the size and position of the cut patch to be consistent with those of the crop region in $x_A$ on which this patch will be pasted.
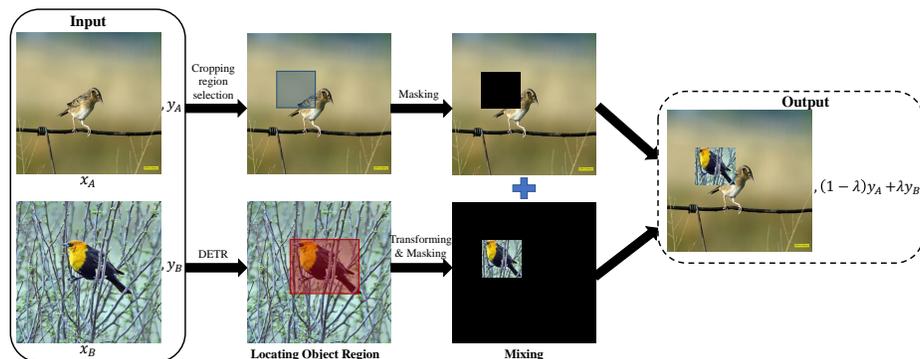
Fig. 2: An example show of the operations of DeMix. Given a pair of a target image and a source image, denoted as $x_A$ and $x_B$, with one-hot labels $y_A$ and $y_B$, respectively, DeMix starts by randomly selecting a cropping region for $x_A$, and locating the object region for $x_B$ based on the object bounding box outputted by a pre-trained DETR. The 'Transforming' operation performed on $x_B$ resizes and relocates the image patch located by DETR to make its size and position consistent with those of the cropping region in $x_A$.

### 3.1   Discussions on the algorithm design of DeMix

In DeMix, we select the target and source images, which are used for generating a new image, in the same way as in CutMix, namely, they are randomly selected from the training set. Since we use a pre-trained DETR for DeMix, which means that we do not need to train the DETR model, the computational overhead of DeMix is comparable to CutMix.

DeMix is based on CutMix, with the same "random cropping" and linear label generation operations. The concept of random cropping has been widely used in deep learning (DL) data augmentation techniques like CutMix [35] and Cutout [7]. The Dropout approach, which is often used for DL regularization [1, 28], is also a form of "random cropping" in essence. However, instead of image patches, it crops neural network weights. Empirically speaking, "random cropping" is a simple yet effective strategy for DL regularization. Its basic mechanism is that DL generally follows "shortcut" learning, making predictions based on "shortcut" features embedded in the training dataset [10]. For instance, if all cows in the training set appear with grass, the DL model could link grass features to cow existence. If a cow appears on a beach without grass in a test image, the model will predict that there is not a cow. By utilizing "random cropping," DL lessens its reliance on shortcut features, reducing the overfitting probability and enhancing the model's generalization ability.

We display data augmentation results of MixUp [36], CutMix [35], SaliencyMix [30], and DeMix in Fig. 3. As is shown, DeMix cuts and resizes a patch that covers the whole object region from the source image, while SaliencyMix selects a patch corresponding to the most salient box region that only covers a part of the object. CutMix randomly selects a patch to be cut, which may only con-

tain the background region. MixUp mixes the target and source images through linear combination, which may lead to local ambiguity and unnaturalness in the generated images, as addressed by [35].



target image   source image        MixUp          CutMix        SaliencyMix       DeMix

Fig. 3: Comparison between related data augmentation techniques

As DeMix is a pre-trained DETR assisted data augmentation technique, its performance is strongly connected to the quality of the DETR model. When the DETR model works well to accurately locate the object region in the source image, then DeMix could succeed in selecting a semantically rich patch to be cut and mixed. When the DETR model being used fails to accurately locate the object region in the source image, then, in principle, DeMix reduces to CutMix, since the patch to be cut can be seen as one randomly selected.

## 4    Experiments

In this section, we evaluate the performance of our DeMix method through experiments on image classification tasks that involve different datasets, different neural network architectures of different sizes. Related modern data augmentation techniques, see subsection 4.1, are used for performance comparison.

### 4.1    Experimental Setting

**Datasets** In our experiments, we selected three benchmark fine-grained image datasets for use, namely CUB-200-2011 [31], Stanford-Cars [19], and FGVC-Aircraft [21]. For simplicity, we refer to these three datasets as CUB, Cars, and Aircraft, respectively, in what follows.

**Network Architectures** In order to perform a comprehensive evaluation on the performance of our method, we selected 6 different network architectures in our experiments, including ResNet-(18, 34, 50, 101) [14], InceptionV3 [29], and DenseNet121 [16].

**Comparison Methods** For performance comparison, we selected modern data augmentation techniques in our experiments including CutMix [35], SaliencyMix [30], MixUp [36], and CutOut [8]. We also include a baseline method that refers to a model trained without using any data augmentation technique.

**Futher details on model training** We used the open-source pre-trained DETR model *detr_resnet*50 included in the Pytorch-torchvision package. The initial feature extractor parameter values are set to be equal with those of the ResNet50 feature extractor pre-trained on the ImageNet-1K dataset [15], and the entire DETR model is trained based on the MS-COCO dataset [20]. Fig. 1 demonstrates the detection performance of the aforementioned DETR model on some image examples involved in our experiments.

In our image classification tasks, we followed [17] to specify hyper-parameter values for model training. Specifically, we chose the stochastic gradient descent (SGD) with momentum as the optimizer, and set the momentum factor at 0.9. The initial learning rate for the pre-trained weights was set to 0.001, while that for other weights was set to 0.01. If training from scratch, the initial learning rate for all trainable weights was set to 0.01. When using pre-trained weights, the model was trained for 200 epochs and decayed the learning rate by factor 0.1 at 80, 150, and 180 epoch; otherwise, the model was trained for 300 epochs and decayed the learning rate by factor 0.1 at 150, 225, and 270 epoch.

### 4.2   Experimental Results

We used different data augmentation techniques in image classification and evaluated the performance of each data augmentation technique via its associated classification accuracy. In Tables 1 and 2, we show the image classification performance, in terms of the average top-1 accuracy, with respect to ResNet architectures with different depths. Results of the baseline method, MixUp, CutOut, and CutMix are directly quoted from [17]. The training setting for the other methods, namely SaliencyMix and DeMix, were set as the same as [17] to guarantee that the performance comparison is fair.

We see that DeMix achieved the best performance in almost all the experiments compared to other methods. It also shows that SaliencyMix does not provide a significant performance gain over CutMix on these datasets, while DeMix does. We argue that it is because the discriminative parts of an image are not located in the salient region captured by SaliencyMix, while they are located in the object region detected by the DETR model employed by DeMix.

Furthermore, we found that DeMix performs more stable, compared to other methods, when the network depth varies. For example, on the dataset CUB, CutMix and SaliencyMix perform poorly with shallower network architectures ResNet18, while show significant performance improvement with deeper architectures like ResNet101. This may be because the image samples generated by CutMix and SaliencyMix are more noisy than those generated by DeMix, and deeper networks are better at handling noisy samples. Overall, regardless of the network depth, DeMix outperforms the other methods significantly.

We further conducted experiments to evaluate the performance of DeMix on other neural network architectures, namely InceptionV3 and DenseNet121. The results are shown in table 3. Again, we see a significant performance improvement given by DeMix over the other methods.

Table 1: Top-1 accuracy (%) of each method for image classification tasks on datasets CUB, Cars, and Aircraft. The classification network is initialized by a pre-trained ResNet18 or ResNet34. The best performance is marked in bold.

| | CUB | | Cars | | Aircraft | |
|---|---|---|---|---|---|---|
| | ResNet18 | ResNet34 | ResNet18 | ResNet34 | ResNet18 | ResNet34 |
| baseline | 82.35 | 84.98 | 91.15 | 92.02 | 87.80 | 89.92 |
| MixUp | **83.17** | 85.22 | 91.57 | 93.28 | 89.82 | 91.02 |
| CutOut | 80.54 | 83.36 | 91.83 | 92.84 | 88.58 | 89.90 |
| CutMix | 80.16 | 85.69 | 92.65 | 93.61 | 89.44 | 91.26 |
| SaliencyMix | 80.69 | 85.17 | 93.17 | 93.94 | **90.61** | 91.72 |
| DeMix | 82.86 | **86.69** | **93.37** | **94.49** | 90.52 | **93.10** |

Table 2: Top-1 accuracy (%) of each method for image classification tasks on datasets CUB, Cars, and Aircraft. The classification network is initialized by a pre-trained ResNet50 or ResNet101. The best performance is marked in bold.

| | CUB | | Cars | | Aircraft | |
|---|---|---|---|---|---|---|
| | ResNet50 | ResNet101 | ResNet50 | ResNet101 | ResNet50 | ResNet101 |
| baseline | 85.49 | 85.62 | 93.04 | 93.09 | 91.07 | 91.59 |
| MixUp | 86.23 | 87.72 | 93.96 | 94.22 | 92.24 | 92.89 |
| CutOut | 83.55 | 84.70 | 93.76 | 94.16 | 91.23 | 91.79 |
| CutMix | 86.15 | 87.92 | 94.18 | 94.27 | 92.23 | 92.29 |
| SaliencyMix | 86.35 | 87.59 | 94.23 | 94.22 | 92.41 | 92.77 |
| DeMix | **86.93** | **88.23** | **94.59** | **94.81** | **93.76** | **94.27** |

Table 3: Top-1 accuracy (%) of each method for image classification tasks on datasets CUB, Cars, and Aircraft. The classification network is initialized by a pre-trained InceptionV3 or DenseNet121. The best performance is marked in bold.

| | CUB | | Cars | | Aircraft | |
|---|---|---|---|---|---|---|
| | InceptionV3 | DenseNet121 | InceptionV3 | DenseNet121 | InceptionV3 | DenseNet121 |
| baseline | 82.22 | 84.23 | 93.22 | 93.16 | 91.81 | 92.08 |
| MixUp | 83.83 | 86.65 | 92.23 | 93.21 | 92.02 | 91.42 |
| CutMix | 84.31 | 86.11 | 93.94 | 94.25 | 92.71 | 93.40 |
| SaliencyMix | 85.07 | 85.26 | **94.18** | 93.65 | 93.58 | 92.95 |
| DeMix | **85.12** | **87.38** | 94.13 | **94.29** | **93.85** | **94.27** |

In previous experiments, we fine-tuned pre-trained classification models using augmented training datasets in the modeling training phase. We also conducted experiments wherein we train the classification models from scratch. In this way, we get a clearer performance evaluation, since it avoids the impact of pre-training on data augmentation performance evaluation. The corresponding results are shown in table 4. We see that DeMix performs best in most cases and comparably in the other ones. In particular, on the CUB dataset, DeMix gives a significant performance improvement compared to CutMix and SaliencyMix. It may be because that, in the CUB dataset, images of different classes have more subtle differences among each other, thus requiring training samples of higher quality; and DeMix can generate samples of higher quality than CutMix and SaliencyMix.

Table 4: Top-1 accuracy (%) of each method for image classification tasks on datasets CUB, Cars, and Aircraft. The classification network architecture is set as ResNet18 or ResNet50, the same as in Table 1, while here the model is trained from scratch, other than pre-trained as shown in Table 1. The best performance is marked in bold.

| | CUB | | Cars | | Aircraft | |
|---|---|---|---|---|---|---|
| | ResNet18 | ResNet50 | ResNet18 | ResNet50 | ResNet18 | ResNet50 |
| baseline | 64.98 | 66.92 | 85.23 | 84.63 | 82.75 | 84.49 |
| MixUp | 67.63 | **72.39** | 89.14 | 89.69 | 86.38 | 86.59 |
| CutMix | 60.03 | 65.28 | 89.11 | 90.13 | 85.60 | 86.95 |
| SaliencyMix | 65.60 | 67.03 | 88.53 | 89.81 | 86.95 | **88.81** |
| DeMix | **70.00** | 71.80 | **89.83** | **91.72** | **88.48** | 88.66 |

In order to understand why DeMix performs better than the other methods involved in our experiments, we investigate the class activation mapping (CAM) associated with each data augmentation technique. CAM is a technique used in deep learning based computer vision to visualize the regions of an image that are most important for a neural network's classification decision. CAM generates a heatmap that highlights the regions of the image that contributed most to the predicted output class, allowing humans to better understand how the model is making its predictions. We checked CAMs given by the classification models trained with aid of different data augmentation techniques. The visualization results on 3 test image examples are shown in Fig. 4. As is shown, using DeMix, the regions of the image that contributed most to the predicted output class match the real objects' regions to a greater extent than using the other data augmentation techniques. For example, for the 2nd test image, the classification model trained with aid of MixUp mainly uses the head of the bird, the model associated with SaliencyMix mainly uses the body of the bird, while the model corresponding to DeMix uses the head and a part of the body together, to generate the predicted label. For the 3rd test image, it is clearer that the model

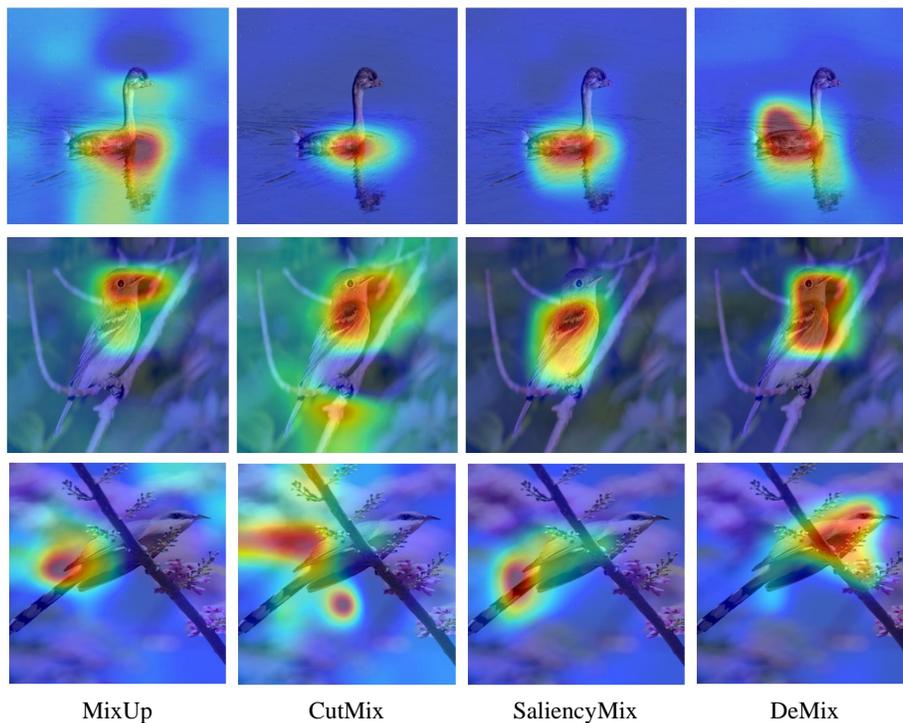associated with DeMix selects a more appropriate region for use in making class predictions.



MixUp          CutMix          SaliencyMix          DeMix

Fig. 4: Visualizations of class activation mapping (CAM) on 3 test image examples

### 4.3   Further Experiments

We conduct an experiment to investigate the influence of the hyperparameter $\lambda$, namely the ratio of the area of the randomly selected crop region, on performance of DeMix. See the result in Table 5, which shows that the performance of DeMix is not very sensitive to the value of $\lambda$. Note that our DeMix is built upon CutMix. It utilizes a random $\lambda$, the same as CutMix, to enhance sample diversity in the augmented dataset, which has been demonstrated to be beneficial for improving the model's performance in terms of generalization.

   We also consider long-tailed recognition tasks on the CUB dataset. Performance comparison results between DeMix with CutMix and SaliencyMix with different imbalance ratios are presented in Table 6. It is shown that, for both architectures ResNet18 and ResNet50, DeMix performs best.

Table 5: Influence of $\lambda$ on performance of DeMix on the image classification task using the CUB dataset

| $\lambda$ | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 |
|---|---|---|---|---|---|---|---|---|---|
| ResNet18 | 83.83 | 83.28 | 82.93 | 82.78 | 82.67 | 82.36 | 82.02 | 81.43 | 80.83 |
| ResNet50 | 87.33 | 87.02 | 86.90 | 86.90 | 86.87 | 86.61 | 86.49 | 85.90 | 85.40 |

Table 6: Top-1 accuracy (%) comparison on long-tailed CUB dataset with different imbalance ratios

| | ResNet18 | | ResNet50 | |
|---|---|---|---|---|
| Imbalance Ratio | 50% | 10% | 50% | 10% |
| CutMix | 32.15 | 52.50 | 38.54 | 62.63 |
| SaliencyMix | 29.58 | 45.53 | 33.76 | 54.90 |
| DeMix | **32.97** | **52.99** | **38.89** | **62.91** |

## 5   Conclusion

In this paper, we demonstrated that a pre-trained object detection model, namely DETR, can be used as a tool for developing powerful data augmentation techniques. Specifically, we found that a DETR model pre-trained on the MS-COCO dataset [20] can be used to locate semantically rich patches in images of other datasets, such as CUB-200-2011 [31], Stanford-Cars [19], and FGVC-Aircraft [21]. Then we proposed DeMix, a novel data augmentation technique that employs DETR to assist CutMix in locating semantically rich patches to be cut and pasted. Experimental results on several fine-grained image classification tasks that involve different network depths and different network architectures demonstrate that our DeMix performs strikingly better than prior art methods. Our work thus suggests (or confirms) that leveraging the power of a pre-trained (large) model directly, without fine-tuning, is a promising direction for future research to improve task-specific performance.

## References

1. Baldi, P., Sadowski, P.J.: Understanding Dropout. Advances in neural information processing systems **26** (2013)
2. Bayer, M., Kaufhold, M.A., Buchhold, B., Keller, M., Dallmeyer, J., Reuter, C.: Data augmentation in natural language processing: a novel text generation approach for long and short text classifiers. International journal of machine learning and cybernetics **14**(1), 135–150 (2023)
3. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: Computer Vision– ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16. pp. 213–229. Springer (2020)

[4] Chadebec, C., Thibeau-Sutre, E., Burgos, N., Allassonnière, S.: Data augmentation in high dimensional low sample size setting using a geometry-based variational autoencoder. IEEE Trans. on Pattern Analysis and Machine Intelligence (2022)

[5] Chaitanya, K., Karani, N., Baumgartner, C.F., Erdil, E., Becker, A., Donati, O., Konukoglu, E.: Semi-supervised task-driven data augmentation for medical image segmentation. Medical Image Analysis **68**, 101934 (2021)

[6] Chen, J., Shen, D., Chen, W., Yang, D.: Hiddencut: Simple data augmentation for natural language understanding with better generalizability. In: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 4380–4390 (2021)

[7] DeVries, T., Taylor, G.W.: Improved regularization of convolutional neural networks with cutout. arXiv preprint arXiv:1708.04552 (2017)

[8] DeVries, T., Taylor, G.W.: Improved regularization of convolutional neural networks with cutout. arXiv preprint arXiv:1708.04552 (2017)

[9] Fawzi, A., Samulowitz, H., Turaga, D., Frossard, P.: Adaptive data augmentation for image classification. In: 2016 IEEE international conference on image processing (ICIP). pp. 3688–3692. Ieee (2016)

[10] Geirhos, R., Jacobsen, J.H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., Wichmann, F.A.: Shortcut learning in deep neural networks. Nature Machine Intelligence **2**(11), 665–673 (2020)

[11] Ghiasi, G., Cui, Y., Srinivas, A., Qian, R., Lin, T.Y., Cubuk, E.D., Le, Q.V., Zoph, B.: Simple copy-paste is a strong data augmentation method for instance segmentation. In: Proc. of the IEEE/CVF conference on computer vision and pattern recognition (CVPR). pp. 2918–2928 (2021)

[12] Gong, C., Wang, D., Li, M., Chandra, V., Liu, Q.: Keepaugment: A simple information-preserving data augmentation approach. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 1055–1064 (2021)

[13] Han, K., Wang, Y., Chen, H., Chen, X., Guo, J., Liu, Z., Tang, Y., Xiao, A., Xu, C., Xu, Y.: A survey on vision transformer. IEEE transactions on pattern analysis and machine intelligence **45**(1), 87–110 (2022)

[14] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)

[15] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)

[16] Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4700–4708 (2017)

[17] Huang, S., Wang, X., Tao, D.: SnapMix: Semantically proportional mixing for augmenting fine-grained data. In: Proceedings of the AAAI Conference on Artificial Intelligence. pp. 1628–1636 (2021)

[18] Kafle, K., Yousefhussien, M., Kanan, C.: Data augmentation for visual question answering. In: Proc. of the 10th International Conference on Natural Language Generation. pp. 198–202 (2017)

[19] Krause, J., Stark, M., Deng, J., Fei-Fei, L.: 3d object representations for fine-grained categorization. In: Proceedings of the IEEE international conference on computer vision workshops. pp. 554–561 (2013)

[20] Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: 13th European Conference on Computer Vision (ECCV). pp. 740–755. Springer (2014)

[21] Maji, S., Rahtu, E., Kannala, J., Blaschko, M., Vedaldi, A.: Fine-grained visual classification of aircraft. arXiv preprint arXiv:1306.5151 (2013)

[22] Mikołajczyk, A., Grochowski, M.: Data augmentation for improving deep learning in image classification problem. In: 2018 international interdisciplinary PhD workshop (IIPhDW). pp. 117–122. IEEE (2018)

[23] Montserrat, D.M., Lin, Q., Allebach, J., Delp, E.J.: Training object detection and recognition cnn models using data augmentation. Electronic Imaging **2017**(10), 27–36 (2017)

[24] Parmar, N., Vaswani, A., Uszkoreit, J., Kaiser, L., Shazeer, N., Ku, A., Tran, D.: Image transformer. In: International conference on machine learning. pp. 4055–4064. PMLR (2018)

[25] Peng, X., Tang, Z., Yang, F., Feris, R.S., Metaxas, D.: Jointly optimize data augmentation and network training: Adversarial data augmentation in human pose estimation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2226–2234 (2018)

[26] Perez, L., Wang, J.: The effectiveness of data augmentation in image classification using deep learning. arXiv preprint arXiv:1712.04621 (2017)

[27] Shorten, C., Khoshgoftaar, T.M.: A survey on image data augmentation for deep learning. Journal of big data **6**(1), 1–48 (2019)

[28] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: a simple way to prevent neural networks from overfitting. The journal of machine learning research **15**(1), 1929–1958 (2014)

[29] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2818–2826 (2016)

[30] Uddin, A., Monira, M., Shin, W., Chung, T., Bae, S.H., et al.: Saliencymix: A saliency guided data augmentation strategy for better regularization. arXiv preprint arXiv:2006.01791 (2020)

[31] Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The caltech-ucsd birds-200-2011 dataset (2011)

[32] Wu, Z., Wang, S., Qian, Y., Yu, K.: Data augmentation using variational autoencoder for embedding based speaker verification. In: INTERSPEECH. pp. 1163–1167 (2019)

[33] Xu, J., Li, M., Zhu, Z.: Automatic data augmentation for 3d medical image segmentation. In: 23rd International Conference on Medical Image Computing and Computer Assisted Intervention (MICCAI). pp. 378–387. Springer (2020)

[34] Yang, H., Zhou, Y.: IDA-GAN: A novel imbalanced data augmentation gan. In: 2020 25th International Conference on Pattern Recognition (ICPR). pp. 8299–8305. IEEE (2021)

[35] Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y.: Cutmix: Regularization strategy to train strong classifiers with localizable features. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 6023–6032 (2019)

[36] Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. arXiv preprint arXiv:1710.09412 (2017)

[37] Zhang, W., Cao, Y.: A new data augmentation method of remote sensing dataset based on class activation map. Journal of Physics: Conference Series **1961**, 012023 (2021)

[38] Zhang, X., Wang, Z., Liu, D., Ling, Q.: Dada: Deep adversarial data augmentation for extremely low data regime classification. In: IEEE international conference on acoustics, speech and signal processing (icassp). pp. 2807–2811. IEEE (2019)

[39] Zhao, A., Balakrishnan, G., Durand, F., Guttag, J.V., Dalca, A.V.: Data augmentation using learned transformations for one-shot medical image segmentation. In: Proc. of the IEEE/CVF conference on computer vision and pattern recognition (CVPR). pp. 8543–8553 (2019)

[40] Zhao, L., Liu, T., Peng, X., Metaxas, D.: Maximum-entropy adversarial data augmentation for improved generalization and robustness. Advances in Neural Information Processing Systems **33**, 14435–14447 (2020)

[41] Zhong, Z., Zheng, L., Kang, G., Li, S., Yang, Y.: Random erasing data augmentation. In: Proc. of the AAAI conference on artificial intelligence. pp. 13001–13008 (2020)

[42] Zoph, B., Cubuk, E.D., Ghiasi, G., Lin, T.Y., Shlens, J., Le, Q.V.: Learning data augmentation strategies for object detection. In: Proceedings of European Conference on Computer Vision (ECCV). pp. 566–583. Springer (2020)