



university of
groningen

faculty of mathematics
and natural sciences

artificial intelligence

Proxy Attention : Comparing and Combining Augmentation with Attention

Graduation Project
(Computational Intelligence and Robotics)

Subhaditya Mukherjee (s4747925)
April 7, 2023

Internal Supervisor: S.H. Mohades Kasaie, PhD
Second Internal Supervisor: Matias Valdenegro, PhD
(Artificial Intelligence, University of Groningen)

Artificial Intelligence
University of Groningen, The Netherlands



CONTENTS

1	Introduction	7
1.1	Context and Novelty	7
1.2	Motivation	7
1.3	Challenges	7
1.4	Problem Statement	7
1.5	Research Questions	7
1.6	Thesis Outline	7
1.7	Contributions	7
2	Background	8
2.1	Interpretability	8
2.2	Gradient Based Explanations	8
2.3	Augmentation	8
2.4	Datasets	8
2.4.1	CIFAR 100	8
2.4.2	Stanford dogs	8
2.4.3	Imagenette	8
2.4.4	ASL	8
2.4.5	Food-101	8
3	State of the Art	11
3.1	Gradient Based Explanations	11
3.2	Augmentation	13
3.2.1	Summary	15
3.2.2	Limitations	15
3.3	Architectures	15
4	Proposed Approach	16
4.1	Design Decisions	16
4.2	Hyper Parameters	16
4.2.1	Clear Every Step	16
4.2.2	Gradient Threshold Considered	16
4.2.3	Multiply Weight	16
4.2.4	Proxy Steps	16
4.2.5	Subset Of Wrongly Classified	16
4.2.6	Gradient Method	16
4.2.7	Architectures	16



5 Implementation	17
5.1 Overview	17
5.2 Hyper parameters	17
5.2.1 Clear Every Step	17
5.2.2 Gradient Method	17
5.2.3 Gradient Threshold Considered	17
5.2.4 Multiply Weight	17
5.2.5 Proxy Steps	17
5.2.6 Subset Of Wrongly Classified	17
5.3 Data Loading and Pre Processing	17
5.3.1 Directory structure	17
5.3.2 Label function	17
5.3.3 Clearing proxy images	17
5.3.4 Encode, Stratify, Kfold	17
5.3.5 train and test, val separate	17
5.3.6 Augmentations	17
5.4 Training Details	18
5.5 Grid Search	18
5.6 Optimizations	18
5.6.1 Mixed Precision	18
5.6.2 Gradient Scaling	18
5.6.3 No grad	18
5.6.4 Batched Proxy step	18
5.6.5 Trial Resumption	18
5.6.6 Models	18
5.7 Gradient Based Methods	18
5.8 Proxy Attention	18
5.8.1 Callback Mechanism	19
5.9 Tensorboard	19
5.10 Transfer learning	19
5.11 Optimizer	19
5.12 LR scheduler	19
5.13 Loss function	19
5.14 Batch Size Finder	19
5.15 Result Aggregation	19
5.16 Inference	19
6 Evaluation	20
6.1 Metric Based Analysis	20
6.2 Visual Based Analysis	20
6.3 Summary	20
7 Conclusion	21
7.1 Contributions	21
7.2 Lessons Learned	21
7.3 Future Work	21
8 Appendix	22



LIST OF FIGURES

2.1	CIFAR100 Sample	9
2.2	Stanford Dogs Sample	9
2.3	Imagenette Sample	9
2.4	ASL Sample	10



LIST OF TABLES



KEY

1. \odot denotes element-wise multiplication
2. NN denotes Neural Network
3. CNN denotes Convolutional Neural Network



CHAPTER 1

INTRODUCTION

1.1 Context and Novelty

1.2 Motivation

1.3 Challenges

1.4 Problem Statement

1.5 Research Questions

1.6 Thesis Outline

1.7 Contributions



CHAPTER 2

BACKGROUND

2.1 Interpretability

- Need for Interpretability

2.2 Gradient Based Explanations

- Taxonomy

2.3 Augmentation

- Taxonomy

2.4 Datasets

To test Proxy Attention, the following datasets were used. Note: Images are resized to 224x224 pixels for consistency. These batch visualizations are generated by the author using the torchvision and matplotlib libraries.

2.4.1 CIFAR 100

The CIFAR 100 dataset, introduced by [1] is an image dataset with 60000 color images with dimensions 32x32 pixels. As the name suggests, the dataset has 100 unique classes. Each of these classes have 500 training images. Some of the classes are - airplane, bird, truck, ship, deer and dog. This dataset is used as a coarse grained classification dataset in this project.

2.4.2 Stanford dogs

2.4.3 Imagenette

2.4.4 ASL

2.4.5 Food-101

Plant Village

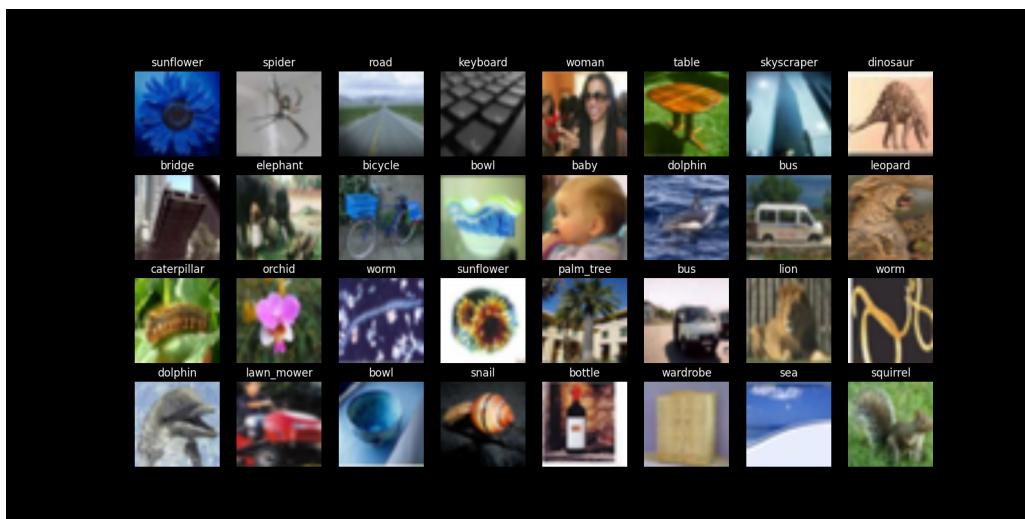


Figure 2.1: CIFAR100 Sample



Figure 2.2: Stanford Dogs Sample

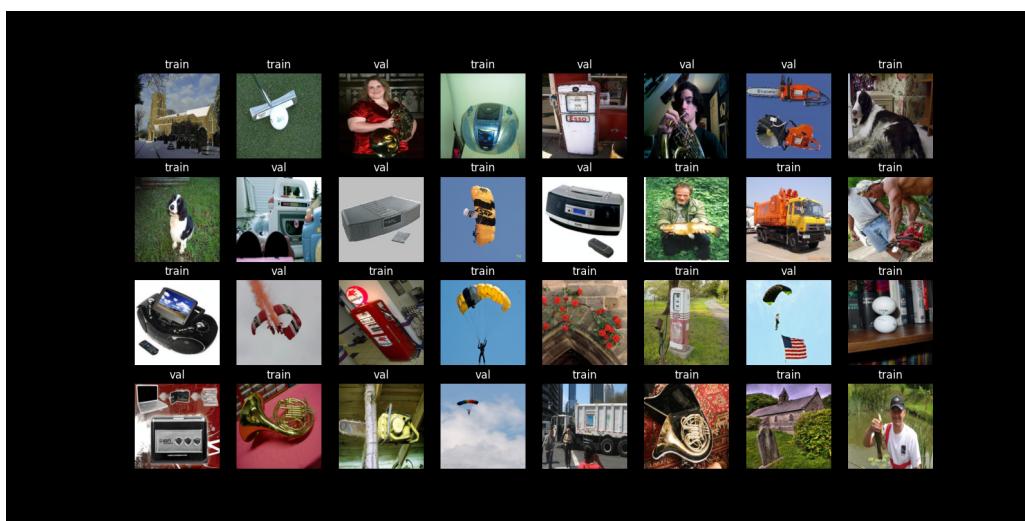


Figure 2.3: Imagenette Sample



Figure 2.4: ASL Sample



CHAPTER 3

STATE OF THE ART

3.1 Gradient Based Explanations

One of the earlier approaches to Saliency maps for CNNs was proposed by Zeiler et al. [2] termed DeconvNet. DeconvNet works by inverting the operations that the network performs in the forward pass. After attaching the DeconvNet layers to the network, propagating through these layers gives a representation of features that the original CNN possessed. For a single class, the relevant reconstruction can be obtained by setting all the activations other than the one corresponding to the class to zero. The resulting image is then used to generate the saliency map. The Conv layer is replaced by a Deconv layer and the ReLU operation has the negative values clamped. While the pooling operation is not strictly invertible, the authors use switch variables that store the position of the maximum value for each pooling operation. While the DeconvNet works to a certain extent, the results are not as accurate as the ones obtained by other methods and are also biased towards the representations of the first layer.

Building on the DeconvNet, Simonyan et al. [3] extrapolate the idea of class visualization to create one of the first approaches to Saliency maps. Their approach, also called Vanilla Gradient ranks the pixels of an image I_0 by how important they are in prediction the Saliency score $S_c(I) \approx w^T I + b$. In this equation, w and b are the weights and biases of the network obtained by back propagating wrt the image itself. The objective to be minimized thus is $\underset{I}{\operatorname{argmax}} S_c(I) - \lambda ||I||_2^2$ where λ is used as a regularization parameter. Using these equations, a saliency map $A \in \mathbb{R}^{m \times n}$ ($m \times n$ stands for *height* \times *width*) can be computed. To find the map, we find the derivative of w , rearrange the elements and then process them according to the number of input channels. If the number of channels is greater than one, the maximum value over the channel is considered $A_{i,j} = \max_{ch} |w_{h(i,j,ch)}|$.

Where, ch is the color channel of the pixel (i, j) , $h(i, j, ch)$ is the index of the w corresponding to that pixel. The Vanilla Gradient method produces an approximate saliency map but has a lot of noise. This leads to issues for more complex images. Many of the issues with Vanilla Gradients and DeconvNets [2] have been addressed by the methods proposed in the following papers.

In another paper, the authors propose a score weighted approach (ScoreCAM) to create saliency maps [4]. Like many other methods, the images are first passed through the network and the corresponding activations are obtained from the final convolutional layer. These activation maps are then upsampled and normalized to the range of $[0, 1]$. The portions of the activation maps that were highlighted are then passed through a CNN with a SoftMax layer to obtain the score for each of the current classes. These scores are used to find the relative importance of all the activation maps. Finally the sum of all these maps is computed using a linear combination with the corresponding target score and then passed through a ReLU operation. These operations can be mathematically represented as $L_{ScoreCAM}^c = \text{ReLU}(\sum_k w_k^c A^k)$, where k represents the index considered, c represents the current class and S_k represents the outputs of the aforementioned SoftMax layer. The authors find that the maps obtained using ScoreCAM are less noisy and using this method removes dependency on unstable gradients as compared to other methods.

A variant of GradCAM [5] was proposed by Selvaraju et al. [6] where unlike GradCAM that finds the parts of the image that influence the model's decision, Guided GradCAM takes the positive gradients into account.



These gradients are used to obtain an even more fine-grained representation of the outputs of the saliency map. While GradCAM backpropagates both positive and negative gradients, Guided Backprop only propagates the positive gradients and is defined as a pointwise multiplication of the results of GradCAM and Guided Backpropagation [7].

In combination with attribution methods, Noise Tunnel [8] is an algorithm that improves the accuracy of the masks obtained by these methods. Noise Tunnel was proposed to counter noisy and irrelevant attributions obtained by some of the gradient based methods by adding a Gaussian Noise and then averaging the predictions over sampled attributions. Since all the samples are considered, this method has a significant computational overhead. For Smooth Grad [9], the new attribution is defined as $\hat{M}_c(x) = \frac{1}{n} \sum_1^n M_c(x + \mathcal{N}(0, \sigma^2))$. Where M_c is the attribution calculated by SmoothGrad, $\mathcal{N}(0, 0.01^2)$ is the Gaussian Noise with $\sigma = 0.01$ and n is the number of samples. Similarly for Smooth Grad Square, $\hat{M}_c(x) = \frac{1}{n} \sum_1^n \sqrt{M_c(x + \mathcal{N}(0, \sigma^2))}$. Noise Tunnel can also be used on Var Grad [10] with the equation $\hat{M}_c(x) = \frac{1}{n} \sum_{k=1}^n \{M_c(x + \mathcal{N}(0, \sigma^2))\}^2 - \{\hat{M}_c(x)\}^2$

For a model F , the attribution method Integrated Gradients [11] computes the contribution of each pixel in the image towards the final prediction. The output of the model is used to calculate a pixel wise partial derivative that is then integrated along a path starting from the baseline and ending at the input. Each of the steps are scaled according to the partial derivative obtained in the previous step. For every step k with m total steps over the path, the IG equation is defined as $IntegratedGrads_i^{approx}(x) := (x_i - x'_i) \times \sum_{k=1}^m \frac{\partial F(x' + \frac{k}{m} \times (x - x'))}{\partial x_i} \times \frac{1}{m}$.

Where $(x_i - x'_i)$ is the pixelwise difference between the two images, $\frac{\partial F(x' + \frac{k}{m} \times (x - x'))}{\partial x_i}$ is the partial derivative of the model output F with respect to pixel i at the k -th step of the path and $\frac{1}{m}$ is the scaling factor that ensures that each of the steps taken contribute equally to the final result.

Petsiuk et al. propose RISE [12], a saliency method that randomly alters the input images by applying random noise to each of them. After model predictions are obtained, the saliency map is generated by a combination of the partial maps over each of the modified images. RISE improves accuracy but needs a lot of computation time considering that multiple models need to be trained for each of the random noise samples.

With regard to the relationship between saliency maps and image classification accuracy, Oyama et al. [13] found a strong correlation between the two. The authors found that both the architecture and the initialization strategy influence the final saliency map. By analyzing the generated saliency maps, they find that if the model is randomly initialized and trained for image classification, having limited categories in the original dataset leads to overfitting. On the other hand having a large number of categories suppresses the overfitting for the objects present in the training dataset. On training their proposed network ReadoutNet on a fixation task (a task which requires the network to learn where to focus), they found that the accuracy of estimating the saliency map was linked to the image classification accuracy.

While a large amount of research focuses on interpreting the influence of a single image or neuron, Hohman et al. propose Summit, [14] a novel scalable summarization algorithm. Summit creates an attribution graph that distills the influence of neurons and substructures throughout the network that are used to make the final prediction. The attribution graph is created as a result of combining activation aggregation, a technique to find important neurons and neuron-influence aggregation, a technique to find relationships among the neurons identified in the previous step. To aggregate the activations, after a forward pass through the network, the activation channels maximums are obtained. These are then filtered by class and aggregated by either taking the top k channels or the top k channels by weight. To quantify how much a layer influences the next, the authors aggregate the influences by creating a tensor I^l for all the layers of the network (l). How important channel i of the layer $l - 1$ is determined by the aggregate tensor I_{cij}^l where j represents the output channel and c is the class of the image. Considering the j^{th} kernel of the layer $K^{(j)} \in \mathbb{R}^{H \times W \times C_{l-1}}$, a single channel Y can be represented using the 3D convolution operation by $Y_{::,j} = X * K^{(j)}$. This is equivalent to its representation by the 2D convolution $Y_{::,j} = \sum_{i=1}^{C_{l-1}} X_{::,i} * K_{::,i}^{(j)}$. The value $X_{::,i} * K_{::,i}^{(j)}$ is the contribution of the current channel from the previous layer and the maximum of this value is used to generate the influence map.

Beware Of Inmates

Interpretation Is Fragile

Sanity Checks

The Unreliability Of Saliency Methods



There And Back Again
Cam
Gradcam++
Guided Backprop
Salience Map
Sam Resnet
Conductance
Deep Fool
Deep Lift
Generalizing Adversarial Exp With Gradcam
Shap
Smooth Grad
Smooth Grad Square
Lime
Sp Lime
Lrp
Var Grad
Visualizing Impact Of Feature Attribution Baselines
Adaptive Whitening Saliency
Bayesian Rule List
Deep Visual Explanations
Dynamic Visual Attention
Embedding Knowledge Into Deep Attention Map
Graph Based Visual Saliency

3.2 Augmentation

Another augmentation strategy proposed by [15] first applies multiple transformations randomly and in parallel chains to each image. These transformations can include combinations of Translation, Rotation, Shearing etc. The outputs of these combinations are then mixed to form a new image, which is then further mixed with the original image to form the new image. This combination is done to improve performance in cases where data shifts are encountered in production. Once the images are mixed, a skip-connection is used to combine the results of the chains. AugMix also uses the Jensen-Shannon Divergence consistency loss [16] to ensure that the images are stable across a range of inputs. Considering KL to be Kullback-Leibler Divergence, the Jensen-Shannon Divergence can be defined as $JS(P_{orig}; P_{augmix1}; P_{augmix2}) = \frac{1}{3}(KL[P_{orig}||M|] + KL[P_{augmix1}||M|] + KL[P_{augmix2}||M|])$, where M is the mean of the three distributions $P_{orig}, P_{augmix1}, P_{augmix2}$.

Devries et al. in their paper [17] propose an augmentation method they call Cutout. In this method, random sized square patches are removed from the images by replacing the corresponding pixels with a constant value (usually 0). Selecting the region involves picking a random pixel value and then creating a uniform sized square around the chosen pixel. The authors also find that Cutout performs better in combination with other methods rather than just being used by itself. Cutout can be expressed as an element-wise multiplication operation $x_{cutout} = x \odot M$, where x is the original image, M is a binary mask of the same size as x with randomly chosen coordinates of a square patch of pixels to be cut out, and \odot denotes element-wise multiplication.

Unlike Cutout [17], where the chosen patch is replaced with zero pixels, in CutMix [18] the chosen patch is replaced with a randomly chosen patch from a different region of the same image. Yun et al. propose this approach as multiple class labels can be learned with a single image. CutMix can be defined by the following operations $\tilde{x} = M \odot x_A + (1 - M) \odot x_B ; \tilde{y} = \lambda y_A + (1 - \lambda) y_B$. where x is an RGB image, y is the respective label, M is a binary mask of the patch of the image that will be dropped and \odot represents element wise multiplication. The new training sample \tilde{x}, \tilde{y} is created by combining two other training samples x_A, y_A and x_B, y_B . To control the combination ratio λ , a sample from the $\beta(1, 1)$ distribution is chosen. This combination is quite similar to [19] but differs in the sense that CutMix focuses on generating locally natural images. Building up on [18],



Walawalkar et al. propose an alternative method of replacing patches in an image they call Attentive CutMix [20]. In this method, instead of randomly pasting patches in the image, a pre-trained network is used to identify attentive regions from the image. Similar to the earlier approach, these patches are then mapped back to the original image. Doing so allows the network to select background regions that are important for the task while also updating the label information.

Many of the algorithms use rectangular or square shaped masks. While they are effective, French et al. propose Cow Mask [21], a new method of masking that uses irregularly shaped masks with a Gaussian filter to reduce noise. The authors also propose two methods of mixing, one that builds up on Random Erasing [22], and another that uses Cut Mix [18]. A pixel wise mixing threshold is also chosen, and either mixing or erasing is applied to the image based on this threshold. This augmentation technique is shown to be effective in semi-supervised learning.

Another approach involving a cut-paste methodology was proposed by [23]. In their paper, the authors propose a new method of augmentation that extracts instances of objects from the images and instead of pasting them on other images, they are pasted on randomly chosen backgrounds. This method leads to pixel artifacts in the images as selecting the objects is a noisy process. To overcome the drop in performance as a result of this, the authors apply a Gaussian blur and poisson blending to the boundaries of the pasted objects. Further augmentation is applied before pasting the objects by rotation, occlusion and truncation. The authors also find that this approach makes the network more robust to artifacts in the images.

In their paper Singh et al. [24] propose a data augmentation method that takes an image as an input, and divides it into a grid. Each of the sub-grids are then turned off with a given probability. These sub-grids can be connected or independent of each other and the turned off grids are replaced by the average pixel value of all the images in the dataset.

One of the major drawbacks of algorithms that rely on modifying image patches (such as [24, 17, 22]) is that they sometimes delete parts of the image that might be useful to the network. To overcome this problem Chen et al. propose a new method Grid Mask [25] that uses evenly spaced grids to find a balance between the amount of information that is deleted and stored. Using the number of grids and their respective sizes as a hyperparameter, the authors find that Grid Mask is effective in preserving important parts of the image.

Zhang et al. propose another method of data augmentation that uses a CAM [26] to identify the most important regions of an image. These parts are then thresholded, scaled, translated and pasted onto the target image. A similar process is also applied to the target image and the attentive parts of the original image are used to replace the corresponding attentive parts of the target image. Similar to previous methods, the labels are also updated to reflect the changes in the image.

While Cutout augmentation [17] is applied to every image in the dataset, Zhong et al. propose a new method, Random Erasing, that takes a probability of being applied into account [22]. In Random Erasing, contiguous rectangular regions are selected and replaced at random with random upper and lower limits chosen for both region area and aspect ratio. For object detection tasks, a region aware detection algorithm is applied to make the network more robust to occlusion. Note that Cutout removes square patches, while Random Erasing either removes square or rectangular patches.

Many of the augmentation methods that rely on randomly choosing regions to cut and paste from sometimes fail to work well with regions that lack object information. ResizeMix [27] tackles this problem by replacing the patch with a proportional resized version of the selected image. This method is similar to CutMix [18] but differs in the sense that ResizeMix uses a resized version of the entire image instead of a randomly chosen patch.

Another augmentation technique that applies random cropping and pasting is RICAP [28]. In this method, four regions are cropped from different images and then pasted together to form a new image. The created image thus has multiple mixed labels. A uniform distribution is used to determine the area of each cropped region in the final image. The authors propose multiple variants of RICAP that use different points of origin for cropping. They find that the method works best when the cropped regions use the corners as the origin as it allows the network to see more of the image.

While algorithms like Mixup [19] modify the labels of the image proportional to the amount of mixing between the original and the target images, Sample Pairing [29] maintains the same training labels. In their paper, Inoue et al. propose a method that merges images not by cut and paste but by averaging their pixel intensities.



Sample Pairing follows an interval based augmentation policy, where the network is first trained for a 100 epochs normally before being introduced to the mixed images. This process is also repeated cyclically with eight epochs of training with mixed images followed by 2 epochs of training with normal images only.

With the success of mask based approaches for data augmentation, there have been many papers that attempt to fix the flaws of previous research. One such method is SmoothMix [30], which builds up on both CutMix [18] and Cutout [17] but modifies the mask to have softer edges. The intensity of the masked edges gradually decreases and depends on the strength of the mask. The updated pixel values are thus obtained by mixing the mask with the original image according to the formula $\lambda = \frac{\sum_{i=1}^W \sum_{j=1}^H G_{ij}}{WH}$. Where G_{ij} is the pixel value of mask G and H, W are the height and width of the image respectively. The new pixel values are then $(x_{new}, y_{new}) = (G.xa + (1 - G).xb, \lambda.ya + (1 - \lambda).yb)$

One of the older methods of data augmentation is SMOTE [31]. This algorithm is not domain specific but in the context of computer vision, it can be used to balance datasets that suffer from imbalanced labels. SMOTE generates new samples by combining the K-nearest neighbors of the minority class images to form new instances. Although many of the other methods discussed in this paper are more effective, SMOTE is still a useful tool to have.

Huang et al. propose SnapMix [32], where choosing the size of the patch to be cut is determined from the beta distributions of both the original and target images. The extracted patches are then merged with random image regions, each of which are of different sizes. Labels are also updated by taking the composition of the images into account. Cao et al. address the problem of class imbalance by performing data augmentation on images that are part of a minority class. From the labels of the images that were mixed, the final label is chosen as the label of the image with the least representation in the dataset. The authors call this method ReMix [33]. Dvornik et al. propose Visual Context Augmentation [34] that uses a NN to understand the context of objects in the image before pasting them in the target image. The authors generate training data by first generating pairs of context images with the objects masked out. These images are then fed into the NN to learn the difference between objects and backgrounds given the masked pixels. Once the model has learnt this information, instances of the objects are placed into the masked regions of the target image.

While there are many techniques based on Mixup [19], they are mostly focused on generating new samples of images from the existing data. Doing so is useful, but sometimes leads to the generation of examples that confuse the network and are not representative of the actual data. To tackle this issue, Kim et al. [35] propose Pizzle Mix, an algorithm that learns to copy patches of images between each other while taking saliency into account. Puzzle Mix learns to minimize the equation $h(x_0, x_1) = (1 - z) \odot \Pi_0^T x_0 + z \odot \Pi_1^T x_1$ where x_0, x_1 are the two images, z_i is a binary mask, $\lambda = \frac{1}{n} \sum_i z_i$ is the mixing ratio and Π_0, Π_1 represent $n \times n$ grids that denote the amount of mass that is transported during transport of the image patch to another location.

Attributemix Augmentaiton with curriculum leanring Co mixup Image Mixing and deletion

Keep augment Latent space interpo Randaugment Random distortion Saliencymix

Spec augment

3.2.1 Summary

3.2.2 Limitations

- Context
- Does not make use of what the network knows
- Does not help the network learn from its mistakes

3.3 Architectures

Resnet 18, 50

VGG

Vision Transformer



CHAPTER 4

PROPOSED APPROACH

4.1 Design Decisions

Efficient Computation Updating Dataloaders Batched Implementation Callbacks Training Resumption Logging

4.2 Hyper Parameters

4.2.1 Clear Every Step

4.2.2 Gradient Threshold Considered

4.2.3 Multiply Weight

4.2.4 Proxy Steps

4.2.5 Subset Of Wrongly Classified

4.2.6 Gradient Method

4.2.7 Architectures

CHAPTER **5**

IMPLEMENTATION

5.1 Overview

5.2 Hyper parameters

5.2.1 Clear Every Step

5.2.2 Gradient Method

5.2.3 Gradient Threshold Considered

5.2.4 Multiply Weight

5.2.5 Proxy Steps

5.2.6 Subset Of Wrongly Classified

5.3 Data Loading and Pre Processing

5.3.1 Directory structure

5.3.2 Label function

5.3.3 Clearing proxy images

5.3.4 Encode, Stratify, Kfold

5.3.5 train and test, val separate

5.3.6 Augmentations

Imagenet Normalize Tensor Num workers



5.4 Training Details

5.5 Grid Search

5.6 Optimizations

5.6.1 Mixed Precision

5.6.2 Gradient Scaling

5.6.3 No grad

5.6.4 Batched Proxy step

5.6.5 Trial Resumption

5.6.6 Models

TIMM

5.7 Gradient Based Methods

5.8 Proxy Attention

Algorithm 1 Single Batch Proxy Attention

Require: *input_wrong*
Require: *CAM*
Require: *proxy_threshold*
Require: *proxy_image_weight*

```
grads ← CAM(input_wrong)
inversed_normalized_inputs ← inverse_normalize(input_wrong)
output ← REPLACE(grads ≥ proxy_threshold, ((1 - proxy_image_weight) * grads) * 
inversed_normalized_inputs, inversed_normalized_inputs)
```

Algorithm 2 Batch Proxy Attention

Require: *input_wrong*
Require: *label_wrong*
Require: *subset_chosen*

```
chosen_inds = CEIL(subset_chosen * LENGTH(input_wrong))
input_wrong_subset = input_wrong[: chosen_inds]
label_wrong_subset = label_wrong[: chosen_inds]
processed_labels ← [], processed_thresholds ← []
for i ← 0 to LENGTH(label_wrong_subset) do
    pass
end for
```



5.8.1 Callback Mechanism

5.9 Tensorboard

5.10 Transfer learning

5.11 Optimizer

5.12 LR scheduler

5.13 Loss function

5.14 Batch Size Finder

To maximize training performance, a batch size finder 3 is used to find the optimal batch size for each of the models.

Algorithm 3 Batch Size Finder Algorithm

```
Require: dataset_size
Require: max_batch_size
batch_size ← 2
while TRUE do
    if max_batch_size is not None & batch_size ≥ max_batch_size then
        batch_size ← max_batch_size
    end if
    if batch_size ≥ dataset_size then
        batch_size ← batch_size//2
    end if
    if failed is False then
        loop
            inputs ← random((batch_size,input_shape))
            targets ← random((batch_size,output_shape))
            outputs ← model(inputs)
            loss ← MSE(outputs,targets)
            loss.backward()
            optimizer.step()
            optimizer.zero_grad()
            failed ← True
            batch_size ← batch_size*2
        end loop
    else if failed is True then
        failed ← False
        batch_size ← batch_size//2
    end if
end while
```

5.15 Result Aggregation

5.16 Inference



CHAPTER **6**

EVALUATION

6.1 Metric Based Analysis

6.2 Visual Based Analysis

6.3 Summary



CHAPTER 7

CONCLUSION

7.1 Contributions

7.2 Lessons Learned

7.3 Future Work



CHAPTER 8

APPENDIX



BIBLIOGRAPHY

- [1] Alex Krizhevsky. “Learning Multiple Layers of Features from Tiny Images”. In: () .
- [2] Matthew D. Zeiler and Rob Fergus. *Visualizing and Understanding Convolutional Networks*. Nov. 28, 2013. doi: 10.48550/arXiv.1311.2901. arXiv: arXiv:1311.2901. URL: <http://arxiv.org/abs/1311.2901> (visited on 11/28/2022). preprint.
- [3] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. *Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps*. Apr. 19, 2014. arXiv: arXiv:1312.6034. URL: <http://arxiv.org/abs/1312.6034> (visited on 11/18/2022). preprint.
- [4] Haofan Wang et al. *Score-CAM: Score-Weighted Visual Explanations for Convolutional Neural Networks*. Version 2. Apr. 13, 2020. arXiv: arXiv:1910.01279. URL: <http://arxiv.org/abs/1910.01279> (visited on 02/16/2023). preprint.
- [5] Ramprasaath R Selvaraju et al. “Grad-CAM: Visual Explanations From Deep Networks via Gradient-Based Localization”. In: () , p. 9.
- [6] Ramprasaath R. Selvaraju et al. *Grad-CAM: Why Did You Say That?* Jan. 25, 2017. arXiv: arXiv: 1611.07450. URL: <http://arxiv.org/abs/1611.07450> (visited on 02/20/2023). preprint.
- [7] Jost Tobias Springenberg et al. *Striving for Simplicity: The All Convolutional Net*. Apr. 13, 2015. doi: 10.48550/arXiv.1412.6806. arXiv: arXiv:1412.6806. URL: <http://arxiv.org/abs/1412.6806> (visited on 11/18/2022). preprint.
- [8] Narine Kokhlikyan et al. *Captum: A Unified and Generic Model Interpretability Library for PyTorch*. Sept. 16, 2020. doi: 10.48550/arXiv.2009.07896. arXiv: arXiv:2009.07896. URL: <http://arxiv.org/abs/2009.07896> (visited on 04/04/2023). preprint.
- [9] Daniel Smilkov et al. *SmoothGrad: Removing Noise by Adding Noise*. June 12, 2017. doi: 10.48550/arXiv.1706.03825. arXiv: arXiv:1706.03825. URL: <http://arxiv.org/abs/1706.03825> (visited on 11/28/2022). preprint.
- [10] Lorenz Richter et al. “VarGrad: A Low-Variance Gradient Estimator for Variational Inference”. In: *Advances in Neural Information Processing Systems*. Vol. 33. Curran Associates, Inc., 2020, pp. 13481–13492. URL: <https://proceedings.neurips.cc/paper/2020/hash/9c22c0b51b3202246463e986c7e205df-Abstract.html> (visited on 02/20/2023).
- [11] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. *Axiomatic Attribution for Deep Networks*. June 12, 2017. doi: 10.48550/arXiv.1703.01365. arXiv: arXiv:1703.01365. URL: <http://arxiv.org/abs/1703.01365> (visited on 03/24/2023). preprint.
- [12] Vitali Pletschuk, Abir Das, and Kate Saenko. *RISE: Randomized Input Sampling for Explanation of Black-box Models*. Sept. 25, 2018. arXiv: arXiv:1806.07421. URL: <http://arxiv.org/abs/1806.07421> (visited on 02/20/2023). preprint.



- [13] Taiki Oyama and Takao Yamanaka. “Influence of Image Classification Accuracy on Saliency Map Estimation”. In: *CAAI Transactions on Intelligence Technology* 3.3 (2018), pp. 140–152. ISSN: 2468-2322. DOI: 10.1049/trit.2018.1012. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1049/trit.2018.1012> (visited on 10/03/2022).
- [14] Fred Hohman et al. *Summit: Scaling Deep Learning Interpretability by Visualizing Activation and Attribution Summarizations*. Sept. 2, 2019. arXiv: arXiv:1904.02323. URL: <http://arxiv.org/abs/1904.02323> (visited on 02/20/2023). preprint.
- [15] Dan Hendrycks et al. *AugMix: A Simple Data Processing Method to Improve Robustness and Uncertainty*. Feb. 17, 2020. arXiv: arXiv:1912.02781. URL: <http://arxiv.org/abs/1912.02781> (visited on 01/16/2023). preprint.
- [16] Jianhua Lin. “Divergence Measures Based on the Shannon Entropy”. In: () .
- [17] Terrance DeVries and Graham W. Taylor. *Improved Regularization of Convolutional Neural Networks with Cutout*. Nov. 29, 2017. DOI: 10.48550/arXiv.1708.04552. arXiv: arXiv:1708.04552. URL: <http://arxiv.org/abs/1708.04552> (visited on 03/27/2023). preprint.
- [18] Sangdoo Yun et al. “CutMix: Regularization Strategy to Train Strong Classifiers With Localizable Features”. In: *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. 2019 IEEE/CVF International Conference on Computer Vision (ICCV). Seoul, Korea (South): IEEE, Oct. 2019, pp. 6022–6031. ISBN: 978-1-72814-803-8. DOI: 10.1109/ICCV.2019.00612. URL: <https://ieeexplore.ieee.org/document/9008296/> (visited on 02/20/2023).
- [19] Hongyi Zhang et al. *Mixup: Beyond Empirical Risk Minimization*. Apr. 27, 2018. DOI: 10.48550/arXiv.1710.09412. arXiv: arXiv:1710.09412. URL: <http://arxiv.org/abs/1710.09412> (visited on 03/27/2023). preprint.
- [20] Devesh Walawalkar et al. *Attentive CutMix: An Enhanced Data Augmentation Approach for Deep Learning Based Image Classification*. Apr. 5, 2020. DOI: 10.48550/arXiv.2003.13048. arXiv: arXiv:2003.13048. URL: <http://arxiv.org/abs/2003.13048> (visited on 03/29/2023). preprint.
- [21] Geoff French, Avital Oliver, and Tim Salimans. *Milking CowMask for Semi-Supervised Image Classification*. June 5, 2020. DOI: 10.48550/arXiv.2003.12022. arXiv: arXiv:2003.12022. URL: <http://arxiv.org/abs/2003.12022> (visited on 03/31/2023). preprint.
- [22] Zhun Zhong et al. “Random Erasing Data Augmentation”. In: *Proceedings of the AAAI Conference on Artificial Intelligence* 34.07 (Apr. 3, 2020), pp. 13001–13008. ISSN: 2374-3468, 2159-5399. DOI: 10.1609/aaai.v34i07.7000. URL: <https://aaai.org/ojs/index.php/AAAI/article/view/7000> (visited on 10/21/2022).
- [23] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. “Cut, Paste and Learn: Surprisingly Easy Synthesis for Instance Detection”. In: Proceedings of the IEEE International Conference on Computer Vision. 2017, pp. 1301–1310. URL: https://openaccess.thecvf.com/content_iccv_2017/html/Dwibedi_Cut_Paste_and_ICCV_2017_paper.html (visited on 03/31/2023).
- [24] Krishna Kumar Singh et al. *Hide-and-Seek: A Data Augmentation Technique for Weakly-Supervised Localization and Beyond*. Nov. 6, 2018. DOI: 10.48550/arXiv.1811.02545. arXiv: arXiv:1811.02545. URL: <http://arxiv.org/abs/1811.02545> (visited on 03/27/2023). preprint.
- [25] Pengguang Chen et al. *GridMask Data Augmentation*. Jan. 13, 2020. DOI: 10.48550/arXiv.2001.04086. arXiv: arXiv:2001.04086. URL: <http://arxiv.org/abs/2001.04086> (visited on 03/31/2023). preprint.
- [26] Bolei Zhou et al. “Learning Deep Features for Discriminative Localization”. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Las Vegas, NV, USA: IEEE, June 2016, pp. 2921–2929. ISBN: 978-1-4673-8851-1. DOI: 10.1109/CVPR.2016.319. URL: <http://ieeexplore.ieee.org/document/7780688/> (visited on 02/20/2023).



-
- [27] Jie Qin et al. *ResizeMix: Mixing Data with Preserved Object Information and True Labels*. Dec. 20, 2020. arXiv: arXiv:2012.11101. URL: <http://arxiv.org/abs/2012.11101> (visited on 03/29/2023). preprint.
 - [28] Ryo Takahashi, Takashi Matsubara, and Kuniaki Uehara. “Data Augmentation Using Random Image Cropping and Patching for Deep CNNs”. In: *IEEE Transactions on Circuits and Systems for Video Technology* 30.9 (Sept. 2020), pp. 2917–2931. ISSN: 1051-8215, 1558-2205. DOI: 10.1109/TCSVT.2019.2935128. arXiv: 1811.09030 [cs]. URL: <http://arxiv.org/abs/1811.09030> (visited on 03/30/2023).
 - [29] Hiroshi Inoue. *Data Augmentation by Pairing Samples for Images Classification*. Apr. 11, 2018. arXiv: arXiv:1801.02929. URL: <http://arxiv.org/abs/1801.02929> (visited on 03/30/2023). preprint.
 - [30] Jin-Ha Lee et al. “SmoothMix: A Simple Yet Effective Data Augmentation to Train Robust Classifiers”. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. 2020, pp. 756–757. URL: https://openaccess.thecvf.com/content_CVPRW_2020/html/w45/Lee_SmoothMix_A_Simple_Yet_Effective_Data_Augmentation_to_Train_Robust_CVPRW_2020_paper.html (visited on 03/29/2023).
 - [31] *SMOTE: Synthetic Minority Over-sampling Technique* | *Journal of Artificial Intelligence Research*. URL: <https://www.jair.org/index.php/jair/article/view/10302> (visited on 03/31/2023).
 - [32] 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* 35.2 (2 May 18, 2021), pp. 1628–1636. ISSN: 2374-3468. DOI: 10.1609/aaai.v35i2.16255. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/16255> (visited on 03/31/2023).
 - [33] Jie Cao et al. “ReMix: Towards Image-to-Image Translation With Limited Data”. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021, pp. 15018–15027. URL: https://openaccess.thecvf.com/content/CVPR2021/html/Cao_ReMix_Towards_Image-to-Image_Translation_With_Limited_Data_CVPR_2021_paper.html (visited on 03/31/2023).
 - [34] Nikita Dvornik, Julien Mairal, and Cordelia Schmid. “Modeling Visual Context Is Key to Augmenting Object Detection Datasets”. In: Proceedings of the European Conference on Computer Vision (ECCV). 2018, pp. 364–380. URL: https://openaccess.thecvf.com/content_ECCV_2018/html/NIKITA_DVORNIK_Modeling_Visual_Context_ECCV_2018_paper.html (visited on 10/21/2022).
 - [35] Jang-Hyun Kim, Wonho Choo, and Hyun Oh Song. “Puzzle Mix: Exploiting Saliency and Local Statistics for Optimal Mixup”. In: *Proceedings of the 37th International Conference on Machine Learning*. International Conference on Machine Learning. PMLR, Nov. 21, 2020, pp. 5275–5285. URL: <https://proceedings.mlr.press/v119/kim20b.html> (visited on 04/04/2023).