



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)

March 31, 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



# LIST OF FIGURES



# LIST OF TABLES



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



## 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 [krizhevskyLearningMultipleLayers] 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

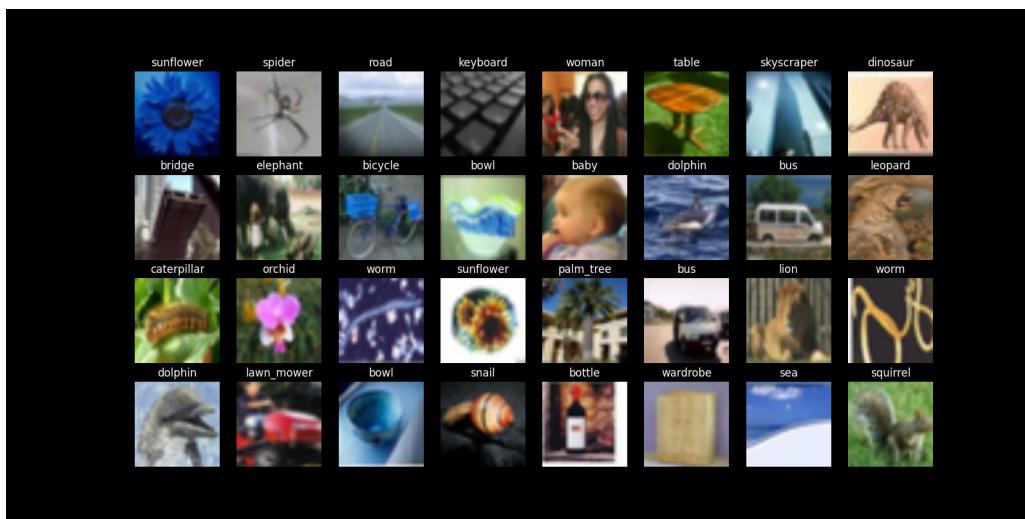


Figure 2.1: CIFAR100 Sample



Figure 2.2: Stanford Dogs Sample

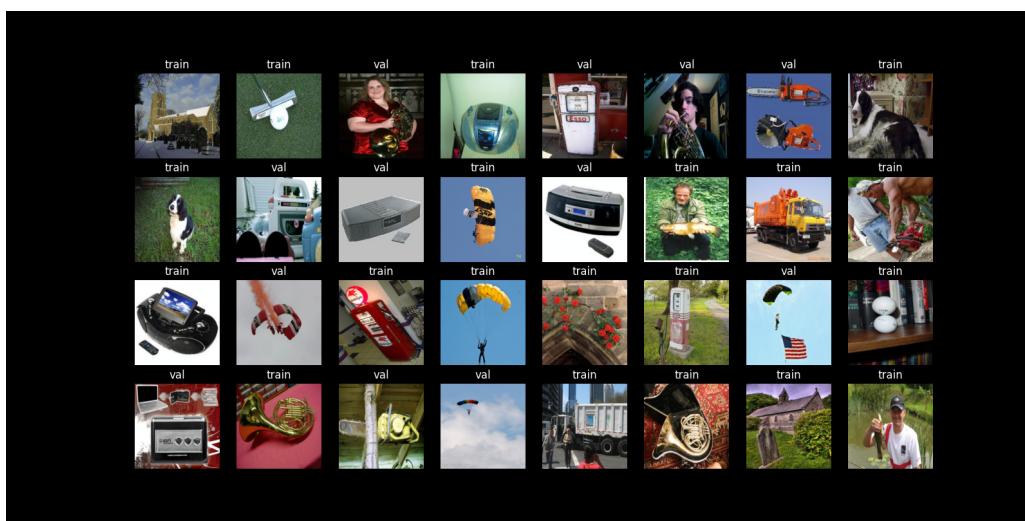


Figure 2.3: Imagenette Sample

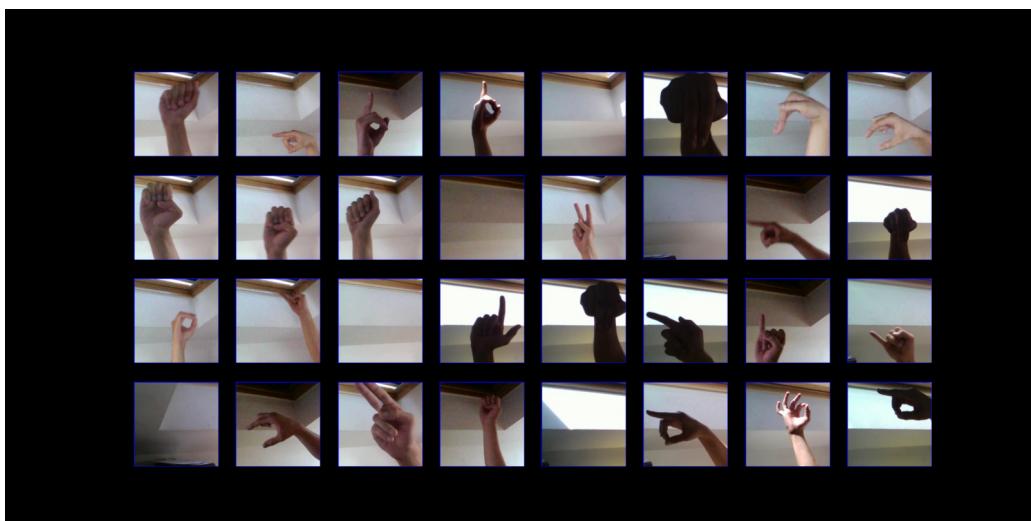


Figure 2.4: ASL Sample



## CHAPTER 3

## STATE OF THE ART

## 3.1 Gradient Based Explanations

Beware Of Inmates

Interpretation Is Fragile

Sanity Checks

The Unreliability Of Saliency Methods

There And Back Again

Influence Of Image Class Acc On Saliency Map Estimation

Deconvnet

Deep Inside Conv Nets

Cam

Gradcam++

Guided Backprop

In another paper, the authors propose a score weighted approach (ScoreCAM) to create saliency maps [wangScoreCAMScore]. 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  $L_{ScoreCAM}^c = ReLU(\sum_k w_k^c A^k)$ , 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.

Guided Gradcam

Salience Map

Noise Tunnel

Integrated Gradients

Sam Resnet

Conductance

Deep Fool

Deep Lift

Generalizing Adversarial Exp With Gradcam

Shap

Smooth Grad

Smooth Grad Square

Lime

Sp Lime



Summit  
Rise  
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 [hendrycksAugMixSimpleData2020] 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 [linDivergenceMeasuresBased] 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 [devriesImprovedRegularizationConvolutional2017] 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 [devriesImprovedRegularizationConvolutional2017], where the chosen patch is replaced with zero pixels, in CutMix [yunCutMixRegularizationStrategy2019] 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  $\lambda$ , a sample from the  $\beta(1, 1)$  distribution is chosen. This combination is quite similar to [zhangMixupEmpiricalRisk2018] but differs in the sense that CutMix focuses on generating locally natural images. Building up on [yunCutMixRegularizationStrategy2019], Walawalkar et al. propose an alternative method of replacing patches in an image they call Attentive CutMix [walawalkarAttentiveCutMixEnhanced2020]. 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 [frenchMilkingCowMaskSemiSupervised2020], 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 [zhongRandomErasingData2020], and another that uses Cut Mix



[**yunCutMixRegularizationStrategy2019**]. 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 [**dwibediCutPasteLearn2017**]. 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. [**singhHideandSeekDataAugmentation2018**] 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 [**singhHideandSeekDataAugmentation2018**, **devriesImprovedRegularizationConvolutional2017**, **zhongRandomErasingData2020**]) 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 [**chenGridMaskDataAugmentation2020**] 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 [**zhouLearningDeepFeatures2016**] 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 [**devriesImprovedRegularizationConvolutional2017**] 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 [**zhongRandomErasingData2020**]. 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 [**qinResizeMixMixingData2020**] tackles this problem by replacing the patch with a proportional resized version of the selected image. This method is similar to CutMix [**yunCutMixRegularizationStrategy2019**] 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 [**takahashiDataAugmentationUsing2018**]. 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 [**zhangMixupEmpiricalRisk2018**] modify the labels of the image proportional to the amount of mixing between the original and the target images, Sample Pairing [**ineoueDataAugmentationPairing2018**] 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 [leeSmoothMixSimpleEffective2020], which builds up on both CutMix [yunCutMixRegularizationStrategy2019] and Cutout [devriesImprovedRegularization2017] 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 [SMOTESyntheticMinority]. 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 [huangSnapMixSemanticallyProportional2021], 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 [caoReMixImagestoImageTranslation2021]. Dvornik et al. propose Visual Context Augmentation [dvornikModelingVisualContext2018] 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.

Attributemix Augmentaiton with curriculum leanring Co mixup Image Mixing and deletion Keep augment Latent space interpo Puzzle mix 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

#### 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 sizer finder

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

### 5.15 Result Aggregation

### 5.16 Inference



---

**Algorithm 1** Batch Size Finder Algorithm

---

**Require:** *dataset\_size, max\_batch\_size, failed*

*batch\_size* = 2

**while** TRUE **do**

**if** *max\_batch\_size* is not None & *batch\_size*  $\geq$  *max\_batch\_size* **then**

*batch\_size*  $\leftarrow$  *max\_batch\_size*

**end if**

**if** *batch\_size*  $\geq$  *dataset\_size* **then**

*batch\_size*  $\leftarrow$  *batch\_size* // 2

**end if**

**if** *failed* is False **then**

**loop**

*inputs*  $\leftarrow$  *random((batch\_size, input\_shape))*

*targets*  $\leftarrow$  *random((batch\_size, output\_shape))*

*outputs*  $\leftarrow$  *model(inputs)*

*loss*  $\leftarrow$  *MSE(outputs, targets)*

*loss.backward()*

*optimizer.step()*

*optimizer.zero\_grad()*

*failed*  $\leftarrow$  True

*batch\_size*  $\leftarrow$  *batch\_size* \* 2

**end loop**

**else if** *failed* is True **then**

*failed*  $\leftarrow$  False

*batch\_size*  $\leftarrow$  *batch\_size* // 2

**end if**

**end while**

---



## 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