

VoronoiPatches: Evaluating A New Data Augmentation Method

Steffen Illium, Gretchen Griffin, Michael Kölle,
Maximilian Zorn, Jonas Nüßlein and Claudia Linnhoff-Popien
Institute of Informatics, LMU Munich, Oettingenstraße 67, Munich, Germany
{*steffen.illum, michael.koelle, jonas.nuesslein, linnhoff*}@ifi.lmu.de

Keywords: voronoi patches, information transport, image classification, data augmentation, deep learning

Abstract: Overfitting is a problem in Convolutional Neural Networks (CNN) that causes poor generalization of models on unseen data. To remediate this problem, many new and diverse data augmentation methods (DA) have been proposed to supplement or generate more training data, and thereby increase its quality. In this work, we propose a new data augmentation algorithm: VoronoiPatches (VP). We primarily utilize non-linear re-combination of information within an image, fragmenting and occluding small information patches. Unlike other DA methods, VP uses small convex polygon-shaped patches in a random layout to transport information around within an image. Sudden transitions created between patches and the original image can, optionally, be smoothed. In our experiments, VP outperformed current DA methods regarding model variance and overfitting tendencies. We demonstrate data augmentation utilizing non-linear re-combination of information within images, and non-orthogonal shapes and structures improves CNN model robustness on unseen data.

1 INTRODUCTION

Fueled by big data and available powerful hardware, Deep Artificial Neural Networks (ANNs) have achieved remarkable performance in computer vision thanks to the recent development steps. But deeper, and wider networks with more and more parameters are data hungry (Aggarwal, 2018) beasts. For a wide variety of problems, a sufficient amount of training data is critical to achieve good performance and avoiding overfitting with modern (often oversized) ANNs. Unfortunately, through mistakes in the process of acquiring, mislabeling, underrepresentation, imbalanced classes, etc., data sources can hard to deal with appropriately. (Illium et al., 2021; Illium et al., 2020). In such cases, straight away learning from such real-world datasets might not be as easy as many common research data sets suggest. To overcome this challenge, Data Augmentation (DA) is commonly used alongside other regularization techniques because of its effectiveness and ease of use (Shorten and Khoshgoftaar, 2019). Over the years, various methods (e.g., occlusion, re-combination, fragmentation) have been developed and reviewed and contrasted in research. In the work at hand, we propose a logical combination of such existing DA-methods: VoronoiPatches (VP).

First, we introduce the concept of DA and Voronoi diagrams in Section 2, then we introduce and discuss

existing and related works in Section 3. In consequence, we propose our approach (VP), the dataset used as well as our experimental setup in Section 4. Finally, we present the results of our experiments in Section 5 just before we conclude in Section 6.

2 PRELIMINARIES

Data Augmentation (DA) is a technique used to reduce overfitting in ANNs, which, in the most extreme cases, is hindered from generalization (the major advantage of ANNs) by perfectly memorizing its training data. The consequence is poor performance on unseen data. A function, learned by an overfit model, exhibits high variance in its output (Shorten and Khoshgoftaar, 2019) by overestimating which ultimately leads to poor overall performance. The amount of a model's variance can be thought of as a function of its size. Assuming finite samples, the variance of a model will increase as its number of parameters increases (Burnham and Anderson, 2002). This is where DA methods can be applied (on training data) to increase the size and diversity of an otherwise limited (e.g., size, balance) data set.

Such approaches have been very successful in the domain of computer vision (and many others) with the advent of deep CNN (Shorten and Khoshgoftaar,

2019) in the past. Image data is especially well suited for augmentation, as one major task of ANNs is to be robust to invariance of objects or image features in general. DA algorithms, on the other hand, are made to generate such invariances. As many other comparable domains exhibit their own challenges and characteristics, in this work, we restrict ourselves to DA methods for images.

In general, DA takes advantage of the assumption that more information can be extracted from the original data to enlarge a data set. It follows that a data set supplemented with augmented data represents a more complete set of all possible data, i.e., closing the real-world gap and promoting the ability to generalize. There are two categories. In a *data warping method*, existing data is transformed to inflate the size of a data set. Whereas, *methods of oversampling* synthesize entirely new data to add to a data set (Shorten and Khoshgoftaar, 2019).

Voronoi diagrams are geometrical structures which define the partition of a space using a finite set of distinct and isolated points (*generator points*). Every other point in the space belongs to the closest generator point. Thus, the points belonging to each generator point form the *regions* of a Voronoi diagram (Okabe et al., 2000). We will further focus the explanations and definitions to 2-dimensional Voronoi diagrams (spanning Euclidean space), as used in this work.

Let S be a set of $n \geq 3$ generator points p, q, r, \dots in Euclidean space \mathbb{R}^2 . The distance d between an arbitrary point $x = (x_1, x_2)$ and the generator point $p = (p_1, p_2)$ is given as:

$$d(p, x) = \sqrt{(p_1 - x_1)^2 + (p_2 - x_2)^2} \quad (1)$$

If we examine generator points p and q , we can define a line which is mutually equidistant as:

$$D(p, q) = \{x | d(p, x) \leq d(q, x)\} \quad (2)$$

A Voronoi region belonging to a generator point $p \in S$, $VR(p, S)$, is the intersection of half-planes $D(p, q)$ where q ranges over all p in S :

$$VR(p, S) = \bigcap_{q \in S, q \neq p} D(p, q) \quad (3)$$

In other words, Voronoi region of p ($VR(p, S)$) is made up of all points $x \in \mathbb{R}^2$ for which p is the nearest neighboring generator point. This results in a convex polygon, which may be bounded or unbounded. The boundaries of regions are called *edges*, which are constrained by their endpoints (*vertices*). An edge belongs to two regions; all points on an edge are clos-

est to exactly two generator points. Vertices are single points that are closest to three or more generator points. Thus, the regions of a Voronoi diagram form a polygonal partition of the plane, $V(P)$ (Aurenhammer, 1991; Aurenhammer et al., 2013).

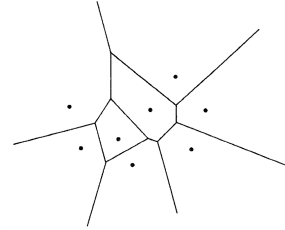


Figure 1: *Voronoi Diagrams*: A simple diagram for eight generator points: $V(P)$ with $|P| = 8$ (Aurenhammer, 1991).

3 RELATED WORK

With the preliminaries introduced, we now survey some existing and related work in the field of Data Augmentation in the context of deep Artificial Neural Networks.

3.1 Occlusion

Methods employing the principle of occlusion mask parts from model input. In consequence, it encounters more varied combinations of an object's features and its context. This forces a stronger recognition of an object by its structure. Two of the earliest methods that use occlusion are *Cutout* (Devries and Taylor, 2017) and *Random Erasing (RE)* (Zhong et al., 2017). Both remove one large, contiguous region from training images (cf. Figure 2, 1A&B). Through *Hide-and-Seek (HaS)* Neural Networks (NNs) learn to focus on the object overall, occluding parts of the input (Singh et al., 2018). This approach removes more varied combinations of smaller regions in a grid pattern from neighboring model input, however, removed regions may also form a larger contiguous region. This, on the same hand, is the major limitation of early DA methods, which can result in the removal of all of an object (or none of it). *GridMask* (Chen et al., 2020) tries to prevent these two extremes by removing uniformly distributed regions (cf. Figure 2, 1C, 2A).

While the use of simple orthogonal shapes and patterns is a common characteristic of occlusion-based methods, Voronoi decomposition-based random region erasing (VDRRE) demonstrates a potential advantage of removing more complex shapes (Abayomi-Alli et al., 2021). For the tasks of facial palsy detection and classification, VDRRE evaluated the use of Voronoi tessellations (cf. Figure 2, 3A&B).

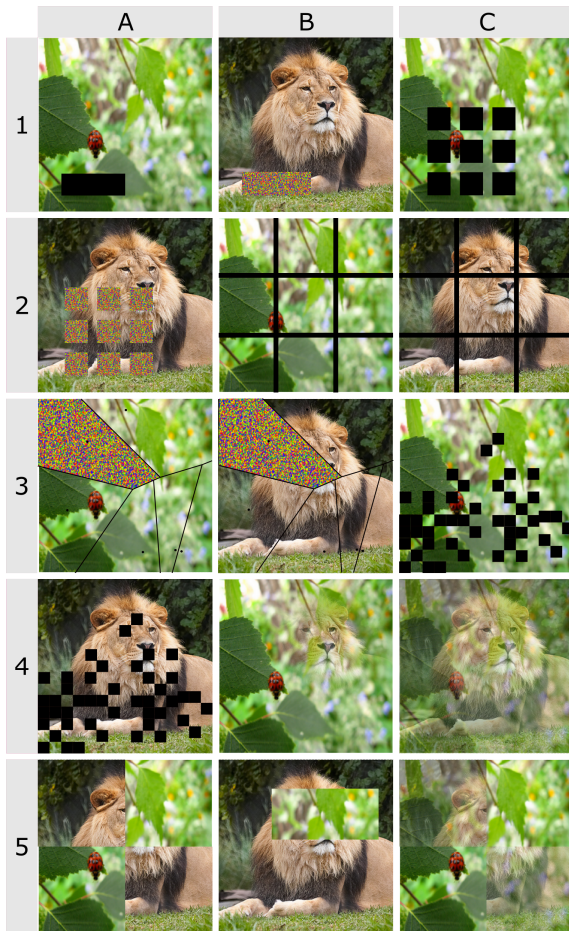


Figure 2: *Combined Augmentations Showcase*: Overview of some available data augmentation methods. **Rows 1-3 + 4A**: Occlusion methods; **4A&B + Row 5**: (Linear) combination methods.

An additional choice for many regional dropout methods is the color of the occlusion mask. Options include random values, the mean pixel value of the data set, or black or white pixels (Yun et al., 2019). Relevant literature includes (Devries and Taylor, 2017; Zhong et al., 2017; Chen et al., 2020; Singh et al., 2018; Abayomi-Alli et al., 2021). Furthermore, regional dropout methods are not always label-preserving depending on the data set (Shorten and Khoshgoftaar, 2019).

3.2 Re-Combination of Data

DA methods which re-combine data, mix training images linearly or non-linearly. This has been found to be an efficient use of training pixels (over regional dropout methods) (Yun et al., 2019), in addition to increasing the variety of data set samples (Takahashi et al., 2018). However, mixed images do not necessarily make sense to a person (Shorten and Khoshgof-

taar, 2019) (e.g., lower section in Figure 2, 4A&B + Row 5) and it is not fully understood why mixing images increases performance (Summers and Dinneen, 2018).

Non-linear mixing methods combine parts of images spatially. *Random image cropping and pasting (RICAP)* (Takahashi et al., 2018) and *CutMix* (Yun et al., 2019) are non-linear methods that combine parts of two or four training images, respectively. Corresponding labels are mixed proportionately to the area of each image used. *RIACP* combines four images in a two by two grid (Takahashi et al., 2018). *CutMix* fills a removed region with a patch cut from the same location in another training image (cf. Figure 2, 5B) (Yun et al., 2019).

Linear mixing methods combine two images by averaging their pixel values (Shorten and Khoshgoftaar, 2019). There are several methods that use this approach including: *Mixup* (Zhang et al., 2017), *Between-class Learning* (Tokozume et al., 2017), and *SamplePairing* (Inoue, 2018) (cf. Figure 2, 4B&C,5C). Interestingly, (Inoue, 2018) found mixing images across the entire training set produced better results than mixing images within classes. Although images may not make semantic sense, they are surprisingly effective at improving model performance (Summers and Dinneen, 2018). A critique on how linear methods introduce “unnatural artifacts”, which may confuse a model, can be found in (Yun et al., 2019).

3.3 Fragmentation

Fragmentation may occur as a side effect of pixel artifacts introduced by occlusions (Lee et al., 2020) or non-linear mixing of images (Takahashi et al., 2018), which create sudden transitions at the edges of removed regions or combined images, respectively. While occluding and mixing information may prevent the network from focusing on ‘easy’ characteristics, deep ANNs may also latch on to created boundaries of these pixel artifacts.

In *SmoothMix* (Lee et al., 2020) the transition between two blended images is softened to deter the network from latching on to them. Alpha-value masks with smooth transitions are applied to two images, which are then combined and blended. Furthermore, *SmoothMix* prevents the network behavior of focusing on ‘strong-edges’ and achieves improved performance while building on earlier image blending methods like *CutMix*.

Cut, Paste and Learn (CPL) (Dwivedi et al., 2017) is a DA method developed for instance detection. To create novel, sufficiently realistic training images, objects are pasted on to random backgrounds while the

object’s edges are either blended or blurred. Although the resulting images look imperfect to the human eye, they perform better in comparison to images created manually (Dwivedi et al., 2017).

While *SmoothMix* and *CPL* aim to minimize artificially introduced boundaries, *MeshCut* (Jiang et al., 2020) uses a grid-shaped mask to fragment information in training images and purposefully introduces boundaries within an image (cf. Figure 2, 2B&C). By doing so, the model learns an object by its many smaller parts; and thereby focuses on broader areas of an object. In contrast to *SmoothMix* and *CPL*, *MeshCut* demonstrates that intentionally created boundaries can improve model performance.

4 METHOD

We noticed the benefits from non-trivial edges (horizontal & vertical) by introducing randomly organized Voronoi diagrams (Abayomi-Alli et al., 2021), as well as the power of occluding multiple small squares (HaS)(Singh et al., 2018), and opted for a conceptual fusion of these approaches. From our perspective, there also is the need for a DA technique which does not occlude wide areas of the image (that may hold the main features), as the occlusion of areas with black (*zero*-values) or Gaussian noise seems not to be fully error-prone. Therefore, we argue in favor of a new category of Data Augmentation methods: *Transport*. The main idea of transport-based approaches is to preserve features and reposition them in their natural context, in contrast to a replacement by non-informative data.

Combined with the concept of nontrivial edges (Voronoi diagram) and occlusion-of-many HaS, we suggest VoronoiPatches (VP). This section now describes our method (VP), the data set used for evaluation, and the experimental setup.

4.1 VoronoiPatches (VP)

We want to supplement training data with novel images generated online through a non-linear transformation. To do so, an image is first partitioned into a set of convex polygons (*patches*) using a Voronoi diagram. A fixed number of bounded patches are randomly chosen and copied from the original image. These are then pasted randomly over the center of a bounded polygonal region (cf. Figure 3). This results in a novel image containing (mostly) the same object features as the original. VP may occlude and duplicate parts of the object in the original image while preserving the original label. The location, number, and visibility of the parts of an object in an image

vary each time VP is applied. By re-combining information within the image instead of replacing them with random values or a single value, information loss is minimized. Due to the random nature of a patch’s shape, size, and final location, patches may overlap or not move at all. This procedure is further described by Algorithm 1.

There are three tunable Hyperparameters (HPs): **Number of generator points:** The approximate size and number of patches generated by a Voronoi partition. Using few generator points results in fewer larger polygons, and vice versa. **Number of patches:** The number of patches to be transported. **Smooth:** The transition style between moved patches and the original image. Which may either be left as they are (sudden) or smoothed. Sudden transitions may create pixel artifacts caused by sudden changes in pixel values, whereas the application of smooth transitions reduces this effect (cf. Figure 4).

Algorithm 1: VoronoiPatches

Input : *sample, generators, patches, smooth*
Output: *aug*

```

begin:
  aug ← copy of sample
  polygons ← Voronoi(generators);
  centroids ← [mean(p) in polygons];
  for range(patches) do
    p ← random(polygons);
    c ← random(centroids);
    Move p such that arith.Mean(p) = c;
    for x, y in p do
      aug[x, y] ← sample[x, y];
      if smooth then
        s ← gauss.Filter(aug);
        borders ← calc_border();
        for xyb in borders do
          | aug[xyb] ← s[xyb];
        end
      end
    end
  end
end
end

```

4.2 Setup & data set

In this section, we give a brief overview of our data pipeline, and introduce the data set used and performance metrics.

Figure 5 summarizes the procedure used in our data pipeline. First, we use an 85:15 split of our training set as an additional validation set for model selec-

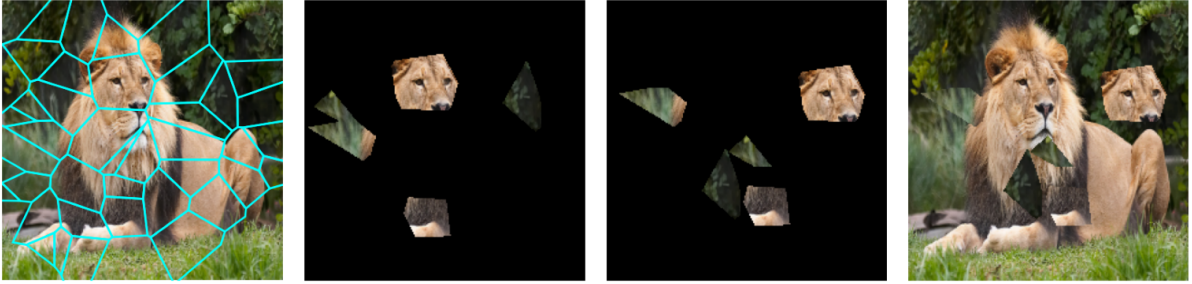


Figure 3: *VoronoiPatches*: An example image and its Voronoi diagram (50 generator points), 5 randomly selected patches *copied*, then *transported* to random locations, and the resulting novel image with sudden transitions (left to right).

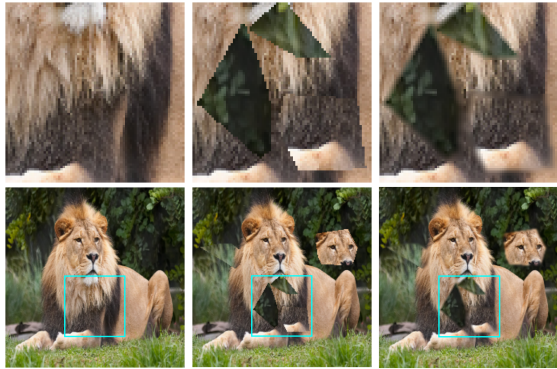


Figure 4: *Voronoi Patches*: The original image (left), the same image augmented with sudden (middle), or smooth transitions (right).

tion. During the training process, performance metrics were calculated and collected each epoch on both the training and validation sets and used to monitor the training process. Measurements were based on the checkpoint with the highest measured accuracy; please note, that this may over- or underestimate true performance (Aggarwal, 2018).

We choose the data set based on two factors: sample size and data set size. By using medium to large resolution images, we ensure more variation in partitions computed and the ability to generate small enough polygons.

The 2012 ImageNet Large-Scale Visual Recognition Challenge (ILSVRC) data set¹, colloquially referred to as ImageNet, consists of 1.2 million training and 60,000 validation full-resolution images categorized and labeled according to a WordNet² based class hierarchy into 1,000 classes (Russakovsky et al., 2015). While the images’ resolution is sufficient, the size of the data set in its entirety is impractical for our experiments.

We used the ‘mixed_10’ data set(Engstrom et al., 2019), which is a subset of the 2012 ImageNet data set

¹<http://www.image-net.org/>

²<https://wordnet.princeton.edu/>

(e.g., in Figure 6). There are 77,237 training images and 3,000 testing images sampled from the 2012 ImageNet training and validation data sets, respectively. These image sets represent ten almost evenly balanced super-classes: dog, bird, insect, monkey, car, cat, truck, fruit, fungus, and boat.

As ‘mixed_10’ is a well-balanced data set, we choose the classification accuracy in % as performance metric. Additionally, we take a look at the variance and entropy introduced by the DA methods.

4.3 Experimental Setup

Baseline: We first established a baseline model based on SqueezeNet 1.0 (through grid search), which provides us with the default performance metrics. Based on this model, we apply different DA methods, in later steps. We chose SqueezeNet 1.0 because it was designed specifically for multi-class image classification using ImageNet and has a small parameter footprint(Iandola et al., 2016). Using 50x fewer parameters, it matched or outperformed the top-1 and top-5 accuracy (Iandola et al., 2016) of the 2012 ILSVRC winner, AlexNet, which has approximately 60 million parameters (Krizhevsky et al., 2017) We only modified the original network architecture by introducing a batch normalization layer as first, to normalize the distribution of our input in a batch-wise fashion. Since DA is performed online, this is our only chance to normalize the distribution of our training data set consistently across all models and DA methods. Initial model HPs were adjusted in a ‘best guess approach’ according to literature. There is no guarantee that we found an optimal configuration; however, our goal was to set up an efficient and well-functioning model for evaluation.

We further chose cross entropy as a common multi-class classification loss (Wang et al., 2022) and ADAM (Kingma and Ba, 2014) as optimizer, which uses adaptive learning rates, which typically require less tuning and converge faster (Ruder, 2016). **Optimizer HPs:** $lr = 0.0001$, $betas = (0.9, 0.999)$,

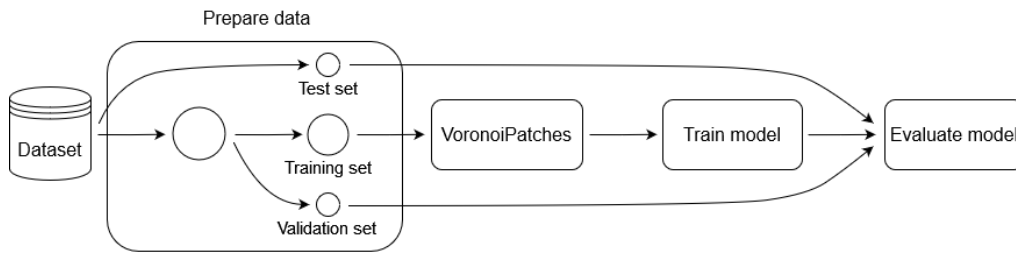


Figure 5: *Data Pipeline*

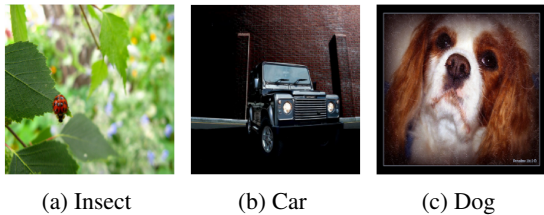


Figure 6: ‘mixed_10’ Data Set: The objects across classes in ‘mixed_10’ images have a wide range of sizes. Examples from classes (a) – (c) illustrate this.

and $eps = 1e - 08$ with $batch_size = 32$. As the SqueezeNet 1.0 architecture (using ReLU) requires images at $224 \times 224 \times 3$ (Iandola et al., 2016), we resized to the fitting resolution and scaled all images channel-wise in the range of $[0, 1]$ using min-max normalization³. By completing this step, the pixel values of our data set have the same scale as the networks’ parameters, which may improve convergence time by stabilizing the training procedure (Aggarwal, 2018). There is no further transformation or augmentation for the baseline model.

Our choice of VoronoiPatches HPs was influenced by the characteristics of our data set, as well as reflects the two main goals: To preserve the label of each augmented image and to occlude or repeat the features of an object or its context in a distributed manner. The diversity of samples in our data set creates a challenging situation for choosing an optimal number of generator points. Specifically, the object-to-context ratio can vary greatly (cf. Figure 6a and 6c). With an image, in which the object is tiny, the danger of all features being occluded by VP is still present. Objects like the car (Figure 6b) show an object-to-context ratio, that lies between the extremes of (a) and (c). Given the variety of this ratio in the ‘mixed_10’ data set, it is difficult to estimate what size patch might be too small or large, or how many patches (i.e., total $pixels^2$ moved) are necessary to have a positive impact on model performance. To balance the parameter space of our grid search with the variation in object-to-context ratio in our data set,

³Scikit-learn/Preprocessing: Min-Max-Scale

we chose $generators = \{50, 70, 90\}$, cf. Figure 7. To explore how the size and number of patches might impact performance, we chose $patches = \{5, 10, 15\}$. For all combinations, we also considered both transition styles, $smoothing = \{true, false\}$.

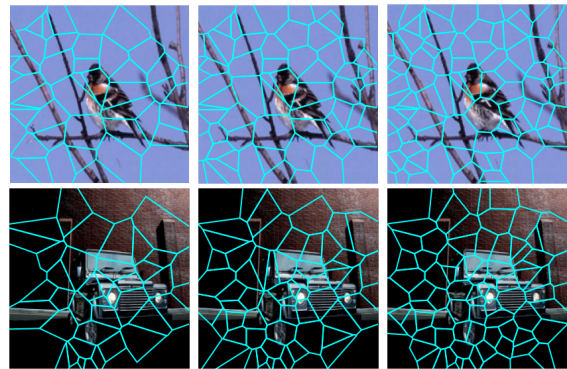


Figure 7: *Voronoi Diagrams*: Computed with (left) 50, (middle) 70, and (right) 90 generator points.

5 EXPERIMENTS

In this section, we describe the course of actions taken to ensure comparable model performance and reproducibility. Finally, we present the observed results of our VP model’s performance against a same-model baseline (no-DA) and other DA methods.

Reproducibility: To ensure the reproducibility of our results, we seeded our interpreter environment, as well as all random number generators involved in the training process. All models (and DA methods) along all grid searches use the same seed to ensure reproducibility and validity. Seeds were chosen at random from $[0, 2^{32} - 1]$. To calculate the avg. expected performance, we re-trained our baseline model with the best performing HP values from a wide seeded grid search. After training to convergence (100 epochs, max. acc. 80.1% at epoch 74), the highest accuracy checkpoints are selected for each seed. To find suitable, VP HPs we performed a grid search over all combinations of $generators$ and $patches$, w/ and

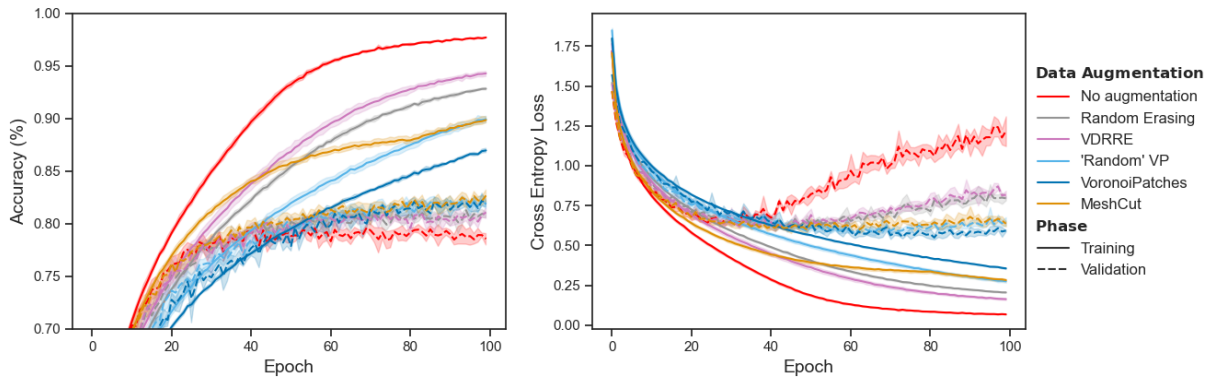


Figure 8: Avg. Model Performance: **left**: Train & val. accuracy over the course of the training (100 epochs). **right**: CE-loss for both train & val. data set. Overfitting clearly visible.

w/o smoothing. Predefined by the baseline model, the training of DA methods is limited to 100 epochs.

Results: For all numbers of patches and patch sizes of VP used while training, we observed a performance improvement over baseline performance. While this improvement is clear, our results showed the amount of improvement varies across HP combinations explored. The highest validation accuracy for VP can be reported as 83.6% at HPs: *generators*=70, *patches*=15, and *smooth*=False. We observed an improvement in val. acc. by about 1.3-3.5% over the baseline models for all models trained with VP. On the test set (unseen in training and validation), we evaluated the performance of the resulting best baseline and optimal VP models. The avg. expected performance of our baseline is **80.9%** accuracy, 80.9% macro-recall, and 81% macro-precision. The avg. expected performance of a model trained using VP is **83.3%** accuracy, 83.5% macro-precision, and 83.3% macro-recall. Even though macro-measures are quite similar, the performance of individual classes varies as expected. However, the classes ‘Car’ and ‘Truck’ stand out as the most difficult for all models to separate, as both perform approximately 10-30% worse than the others. Figure 9 shows identified occurrences of noisy ‘Car’ and ‘Truck’ labels among these collected images (VP applied). We cannot rule out that the large performance difference of these classes, on avg., is due to these noisy labels. During our research, we learned, that other researchers had problems with this situation, which is why ImageNet is currently in debate (Beyer et al., 2020).

Smoothing: Unfortunately, we did not observe any clear improvement in the performance of VP using smooth over sudden transitions (*smooth*=false). Although the best performing models use sudden transitions, based on our observations in model training; however, we cannot rule out a possible advantage at this time. Depending on the other HPs, either ap-

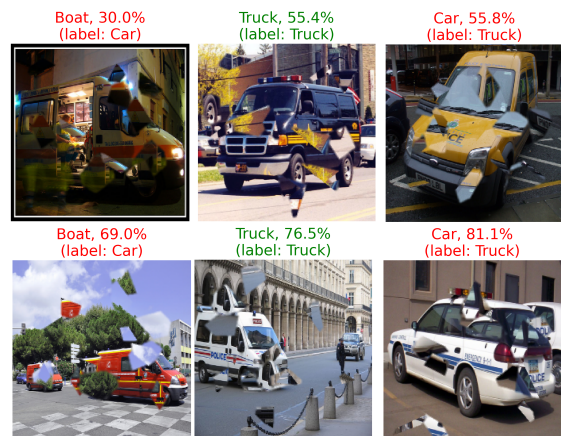


Figure 9: Noisy Labels: Inspection of augmented training images exposed several anomalies. Vans used as ambulances are labeled with ‘Car’ (bottom, left); however, vans used by the police are labeled with ‘Truck’ (middle). A hatchback car used by the police is labeled as ‘Truck’ (bottom, right).

proach can be of advantage. We implemented the width of the transition between patch and image with a fixed factor. In future works, this could be introduced as another HP. Both smoothing and the factor could also be bound to a stochastic process (sample on every VP application).

Comparisons: To verify the validity and utility of VP, we evaluate the performance of the RE, VDRRE and a combination of both, using the same procedure as outlined above. For RE HPs outlined in (Zhong et al., 2017), pixel values ‘black’ (0) and ‘random’ (uniform [0,1)) are chosen. Figure 8 shows an avg. acc. improvement by VP over RE of 0.6%-0.7%. VDRREs parameters are chosen as suggested (Abayomi-Alli et al., 2021). They established six regions, of which one was then occluded with noise. On our data set, VDRRE achieved results, which are close to RE. *MeshCut* (augmentation schedule in training and HPs as stated by the authors (Jiang

et al., 2020)) behaved worse than ‘Random’VP. Inspired by the occlusion with Gaussian noise, we also tested a VP version with randomly filled patches, rather than transporting sections of the image (more or less comparable with VDRRE, but at smaller scale). ‘Random’VP achieved better results than both, VDRRE and RE, while performing close to VP.

Entropy H : To measure the effect of our method on the training data, we employed entropy H (Goodfellow et al., 2016) involving the resulting prob. distribution of each image x (with $P := prob., K := classes$):

$$H(x) = - \sum_{k \in K} P(k) * \log(P(k)) \quad (4)$$

We then aggregated, the entropy scores for all images with the arithmetic mean (cf. Figure 10). The avg. entropy is a measure of how balanced the probability distributions are on avg. (Goodfellow et al., 2016). For a probability distribution of 10 values (for 10 classes) the lower and upper bounds for entropy are 0.00 and 3.32, respectively. The bounds correspond to the case of a perfect classification with 1.0 probability for one class and 0.0 for all others (low entropy) and the case of perfectly balanced probability of 0.1 across all classes (high entropy). As anticipated, the more data transported within an image by VP the more uncertainty is embedded in the model’s output (i.e., the more balanced the model’s probability distributions are). We observe that this effect is greater among all VP HP combinations than RE with either ‘black’ or ‘random’ pixels. We attribute the reduced entropy of RE using ‘random’ pixels over ‘black’ pixels to the model finding patterns in these values.

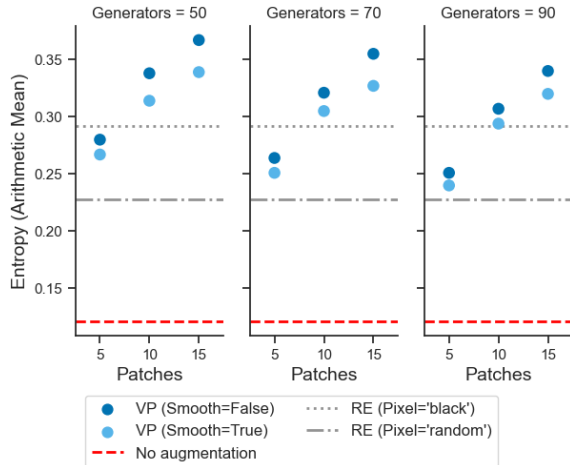


Figure 10: *Entropy*: Avg. effect of Voronoi-Patches (dark blue=sudden, light blue=smooth edges) and RE (gray lines) measured in entropy H (cf. Equation 4).

Main Findings: (1) Figure 11 shows measured vari-

ance for considered DA methods. Over all seeded runs, our proposed method (VP) exhibits the lowest variance, even lower than training on non-augmented training data. By considering the results of the former entropy analysis, especially the last point is quite a surprising find.

(2) By observing the avg. acc. and the CE-loss in Figure 8 we observe a second interesting aspect. All training runs (w/ and w/o DA) came at the cost of still overfitting to the training data of ‘mixed_10’. Our method VP, among all training runs, showed the least amount of overfitting. This can be observed by the growing distance of the CE-loss in Figure 8 (right).

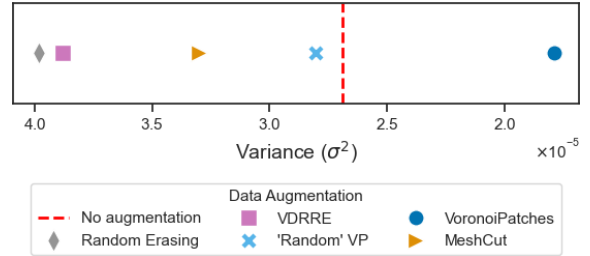


Figure 11: *Model Variance Compare*: Avg. test variance (std) of VP and the other DA methods.

Further Analysis: The number of generator points determines the avg. patch size. Evaluating the impact of this HP, we collected the avg. $pixels^2$ (2D area) for bounded polygons using 50, 70, and 90 generator points. The avg. patch sizes are 932, 673, and 528 $pixels^2$, respectively. We observe that using many smaller patches (15 patches, 528 $pixels^2$), as well as using fewer larger patches results in better model performance (5 patches, 932 $pixels^2$), on avg. As we observe no clear improvement using smooth transitions, we assume, either, that the transition areas’ width is not optimal, or the pixel artifacts have a positive effect similar to the mesh mask used in (Jiang et al., 2020).

To measure how different HP values change our original images, we calculate the Structural SIMilarity Index (SSIM) between orig. and augmented images. SSIM was developed for the measurement ([0-1]) of image degradation by comparing the luminance, contrast, and structures between two images (Wang et al., 2005). $SSIM(original, augmented) = 1.0$ represents the trivial case of a perfect match (‘identical’ images) (Wang et al., 2005). Collecting and aggregating (arithmetic mean) SSIM values for 1,000 iterations with several HP combinations Figure 12 confirms that by increasing the number of patches or decreasing the number of generators, the similarity of original and augmented images decreases. The relationship between the avg. SSIM and

avg. entropy can be reported as approx. linear. As the structural difference in augmented images increases, so does the entropy of the corresponding model output for these images.

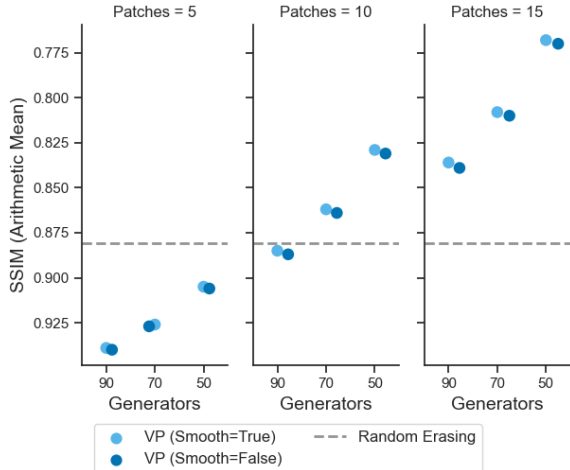


Figure 12: Avg. SSIM for VP and RE.

Further, we observe a drop-off in performance when more than approx. 10,000 total $pixels^2$ are transported. The best performance is achieved by 10,095 total $pixels^2$ ($generators = 70, patches = 15, smooth = False$), a re-combination of 20.1% of an images’ data. There appears to be an optimal avg. patch size and/or number of patches moved similar to the optimal ratio between the grid mask of uniformly distributed squares (Chen et al., 2020), which controls how many squares make up the mask, as well as how much space is left between. As a point of reference, the avg. rectangle size generated by Random Erasing (RE) is 10,176 $pixels^2$. In accordance with the default values (Zhong et al., 2017), RE was applied with a probability of 50%. On avg., RE removes a rectangle equal to 20.3% of an image when applied, whereas VP re-orders patches equal to 20.1%, combined, in each image. We noticed an improvement in performance by re-combining many small VP patches over randomly erasing one large contiguous rectangle with 50% probability (RE).

The variance in patch sizes is also determined by how generator points are distributed across the image. By visual inspection, a trend of smaller patches resulting from more generator points, as well as, larger patches resulting from fewer generator points can be discovered. Furthermore, Figure 13 shows that depending on the size of an object in an image, a patch may occlude or duplicate a larger or smaller feature. Patches generated with 50 generator points in the left-hand column of Figure 13 are large enough to only contain a smaller feature like the face or hand of the

monkey, but large enough to contain a large feature from the insect like its head or thorax.



Figure 13: Comp. of patch sizes using (left) 50, (middle) 70, and (right) 90 gen. points (at 10 patches, smoothing).

6 CONCLUSION

In this work we introduced VoronoiPatches (VP), a novel data augmentation method and category (transport), to solve the problem of overfitting for CNNs. We sought to minimize information loss and pixel artifacts, as well as exploit non-orthogonal shapes and structures in data augmentation. Our method, which employs non-orthogonal shapes and structures to recombine information within an image, outperforms the existing DA methods regarding model variance and overfitting tendencies. Additional experiments analyzed VoronoiPatches’ influence on predictions, as well as the influence of HP values on performance. We show that there are further opportunities to build on our findings and add validity to our initial evaluation of VP. This includes: optimizing smooth transitions, exploring different pixel values or patch shapes (e.g., black pixels, or small squares or rectangles) to better understand the contribution to data augmentation for CNN, mixing images across the training set, and eventually exploring applications of VP in other tasks or fields (e.g., medical image analysis or the field of audio in the form of mel spectrograms). From a practical standpoint, solving the limitation of expensive training would enable more efficient usage of the available data sets. Currently, expensive training is a limitation of VoronoiPatches, which also is a possible future task.

References

- Abayomi-Alli, O. O., Damaševičius, R., Maskeliūnas, R., and Misra, S. (2021). Few-shot learning with a novel voronoi tessellation-based image augmentation method for facial palsy detection. *Electronics*, 10(8).
- Aggarwal, C. C. (2018). *Neural Networks and Deep Learning: A Textbook*. Springer Publishing Company, Incorporated, 1st edition.
- Aurenhammer, F. (1991). Voronoi diagrams—a survey of a fundamental geometric data structure. *ACM Comput. Surv.*, 23(3):345–405.
- Aurenhammer, F., Klein, R., and Lee, D.-T. (2013). *Voronoi Diagrams and Delaunay Triangulations*. World Scientific Publishing Co., Inc., USA, 1st edition.
- Beyer, L., Hénaff, O. J., Kolesnikov, A., Zhai, X., and Oord, A. v. d. (2020). Are we done with imagenet? *arXiv preprint arXiv:2006.07159*.
- Burnham, K. P. and Anderson, D. R. (2002). *Model selection and multimodel inference : a practical information-theoretic approach*. Springer-Verlag, New York, NY, 2nd edition.
- Chen, P., Liu, S., Zhao, H., and Jia, J. (2020). Gridmask data augmentation. *CoRR*, abs/2001.04086.
- Devries, T. and Taylor, G. W. (2017). Improved regularization of convolutional neural networks with cutout. *CoRR*, abs/1708.04552.
- Dwibedi, D., Misra, I., and Hebert, M. (2017). Cut, paste and learn: Surprisingly easy synthesis for instance detection. *CoRR*, abs/1708.01642.
- Engstrom, L., Ilyas, A., Salman, H., Santurkar, S., and Tsipras, D. (2019). Robustness (Python Library).
- Goodfellow, I. J., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press.
- Iandola, F. N., Moskewicz, M. W., Ashraf, K., Han, S., Dally, W. J., and Keutzer, K. (2016). Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <1mb model size. *CoRR*, abs/1602.07360.
- Illum, S., Müller, R., Sedlmeier, A., and Linnhoff-Popien, C. (2020). Surgical mask detection with convolutional neural networks and data augmentations on spectrograms. *Proc. Interspeech 2020*, pages 2052–2056.
- Illum, S., Müller, R., Sedlmeier, A., and Popien, C.-L. (2021). Visual transformers for primates classification and covid detection. In *22nd Annual Conference of the International Speech Communication Association, INTERSPEECH 2021*, pages 4341–4345.
- Inoue, H. (2018). Data augmentation by pairing samples for images classification. *CoRR*, abs/1801.02929.
- Jiang, W., Zhang, K., Wang, N., and Yu, M. (2020). Mesh-cut data augmentation for deep learning in computer vision. *PLoS One*, 15(12):e0243613.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *ArXiv*, abs/1412.6980.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2017). Imagenet classification with deep convolutional neural networks. *Commun. ACM*, 60(6):84–90.
- Lee, J., Zaheer, M., Astrid, M., and Lee, S.-I. (2020). Smoothmix: a simple yet effective data augmentation to train robust classifiers. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 3264–3274.
- Okabe, A., Boots, B., Sugihara, K., and Chiu, S. N. (2000). *Spatial Tessellations: Concepts and Applications of Voronoi Diagrams*. Series in Probability and Statistics. John Wiley and Sons, Inc., 2nd edition.
- Ruder, S. (2016). An overview of gradient descent optimization algorithms. *ArXiv*, abs/1609.04747.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. (2015). ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252.
- Shorten, C. and Khoshgoftaar, T. M. (2019). A survey on image data augmentation for deep learning. *Journal of Big Data*, 6(1):1–48.
- Singh, K. K., Yu, H., Sarmasi, A., Pradeep, G., and Lee, Y. J. (2018). Hide-and-seek: A data augmentation technique for weakly-supervised localization and beyond. *CoRR*, abs/1811.02545.
- Summers, C. and Dinneen, M. J. (2018). Improved mixed-example data augmentation. *CoRR*, abs/1805.11272.
- Takahashi, R., Matsubara, T., and Uehara, K. (2018). Ricap: Random image cropping and patching data augmentation for deep cnns. In Zhu, J. and Takeuchi, I., editors, *Proceedings of The 10th Asian Conference on Machine Learning*, volume 95 of *Proceedings of Machine Learning Research*, pages 786–798. PMLR.
- Tokozume, Y., Ushiku, Y., and Harada, T. (2017). Between-class learning for image classification. *CoRR*, abs/1711.10284.
- Wang, Q., Ma, Y., Zhao, K., and Tian, Y. (2022). A comprehensive survey of loss functions in machine learning. *Annals of Data Science*, 9.
- Wang, Z., Bovik, A., and Sheikh, H. (2005). Structural similarity based image quality assessment. *Digital Video Image Quality and Perceptual Coding, Marcel Dekker Series in Signal Processing and Communications*.
- Yun, S., Han, D., Oh, S. J., Chun, S., Choe, J., and Yoo, Y. (2019). Cutmix: Regularization strategy to train strong classifiers with localizable features. *CoRR*, abs/1905.04899.
- Zhang, H., Cissé, M., Dauphin, Y. N., and Lopez-Paz, D. (2017). mixup: Beyond empirical risk minimization. *CoRR*, abs/1710.09412.
- Zhong, Z., Zheng, L., Kang, G., Li, S., and Yang, Y. (2017). Random erasing data augmentation. *CoRR*, abs/1708.04896.