

Investigating Convolutional Neural Networks using Spatial Orderness

Rohan Ghosh¹
rghosh92@gmail.com

Anupam K. Gupta¹
anupamgupta1984@gmail.com

Mehul Motani^{1,2}
motani@nus.edu.sg

¹ The N.1 Institute of Health, National University of Singapore

² Department of Electrical and Computer Engineering, National University of Singapore

Abstract

Convolutional Neural Networks (CNN) have been pivotal to the success of many state-of-the-art classification problems, in a wide variety of domains (for e.g. vision, speech, graphs and medical imaging). A commonality within those domains is the presence of hierarchical, spatially agglomerative local-to-global interactions within the data. For two-dimensional images, such interactions may induce an *a priori* relationship between the pixel data and the underlying spatial ordering of the pixels. For instance in natural images, neighboring pixels are more likely contain similar values than non-neighboring pixels which are further apart. To that end, we propose a statistical metric called *spatial orderness*, which quantifies the extent to which the input data (2D) obeys the underlying spatial ordering at various scales. In our experiments, we mainly find that adding convolutional layers to a CNN could be counterproductive for data bereft of spatial order at higher scales. We also observe, quite counter-intuitively, that the spatial orderness of CNN feature maps show a synchronized increase during the initial stages of training, and validation performance only improves after spatial orderness of feature maps start decreasing. Lastly, we present a theoretical analysis (and empirical validation) of the spatial orderness of network weights, where we find that using smaller kernel sizes leads to kernels of greater spatial orderness and vice-versa.

1 Introduction

There has been a large body of theoretical and experimental work exploring various attributes of CNNs which may contribute towards their excellent performance and generalization abilities ([1, 2, 3, 4, 5, 6]). However, the unusual effectiveness of CNNs on a large variety of domains (vision, audio, graphs, medical imaging) is still not entirely comprehended.

Solely from the perspective of a mathematical function, it is intriguing to see a convolutional neural network demonstrate significant performance gains, when compared to a fully connected deep neural network. It is also clear that CNNs and fully-connected neural networks (FC-NNs) showcase very different behaviour. For instance, enough empirical evidence exists ([7, 8, 9]) to suggest that increasing depth in CNNs (by adding convolution layers) almost always leads to considerable performance improvements, on most problems.

However, increasing the depth of a FC-NN quickly leads to worse test performance [12]. Unlike a FC-NN, CNN exhibits translation equivariance across its layers, that enable it to achieve translation equivariant representations deep within the network. This is useful for applications where global translational symmetries exist in the data. However, even in problems where global translational symmetries do not exist (MNIST for instance), CNNs easily outperform FC-NNs [13].

Compared to structure-less FC-NNs, the inductive biases in a CNN are clearly better suited to handle classification problem in various domains. But, instead of a function-based introspection of a CNN, we ask which characteristics of the data itself enable these inductive biases to flourish? Would adding convolution layers still be profitable if one were to synthetically distort these "convolution-conducive" characteristics in the data?

In this work, we systematically explore these questions, based on a hypothesis: *Convolutional structure in a neural network benefits from spatially ordered data*. We define *spatial order* to be the extent to which spatial proximity determines data value proximity. For images, an example of high spatial order in data is when spatially nearby pixels are more likely to have similar values than pixels which are far apart. Similarly, when the pixel intensity differences are independent of the spatial distance between the pixels (e.g. in white-noise images), the data is said to have low spatial order. Spatially ordered data is likely to contain more meaningful spatial structure and hierarchy, and can benefit from locality-preserving feedforward functions; like convolutional operations in CNNs.¹

A simple, novel metric for reliably quantifying spatial order in 2D data is proposed, denoted as *spatial orderness*. This metric can either be computed for a single 2D image, or for an entire dataset of 2D images. First, it is shown that spatial orderness is a reliable quantifier of spatial order in the input: synthetic disruption of spatial structure in the data decreases spatial orderness. Next, we find that adding convolutional depth to a CNN ceases to yield performance improvements, when the data lacks spatial order at higher scales (Figure 2). The remaining experiments and theoretical contributions of this work relates to the computation of spatial orderness of the (a) input, (b) feature maps (during and after training) and (c) kernels (post-training).

2 Spatial Diffusability implies Spatial Orderness

Let us denote a set of 2D data as I_1, I_2, \dots, I_N , all of equal size. We denote underlying data generating probability distribution as D . Consider a spatial location p within the domain of I , such that $I_j(p)$ denotes the value of the image I_j at the spatial location p . Let us denote an image X which is a random variable following the distribution D . Since X is a random variable, it follows that $X(p)$, which represents the value of X at spatial location p , is also a random variable. Proceeding with these definitions, we outline our approach for computing spatial orderness.

For the purpose of quantifying spatial order, we propose a spatial diffusion based generative modelling of X . The key observation is that nearby spatial locations (or regions) must show smaller differences in intensity values (or average intensity), *compared* to spatial locations (or regions) which are further apart. Note that simple correlations between neighboring pixel values is not enough to concretely quantify the above relationship.

¹For graphs, we can extend the definition of spatial order to one of locality. For instance, a graph where every node is connected to every other node would be a counter-example of locality in a graph.

We begin with three random variables extracted from different spatial locations within X : $X(p)$, $X(q)$ and $X(r)$. Importantly, the spatial locations p , q and r are chosen such that $q \in N(p)$ and $r \in N(N(p)) \setminus N(p)$, where $N(p)$ represents the set of neighboring spatial locations of p . Here \setminus is the set subtraction operator, such that r is at a distance of two hops from p . Throughout this paper, we denote this particular spatial arrangement of locations as the *two-hop spatial arrangement*. With this, we can outline the relationship between the random variables $X(p)$, $X(q)$ and $X(r)$ using a normally distributed spatial diffusion process:

$$X(q) = X(p) + \mathcal{N}(0, \sigma) \quad (1)$$

$$X(r) = X(q) + \mathcal{N}(0, \sigma). \quad (2)$$

Combining equations 1 and 2, we have

$$X(r) = X(p) + \mathcal{N}(0, \sqrt{2}\sigma), \quad (3)$$

using which the above equations can be summarized with the relationship:

$$\mathbb{E} [(X(p) - X(r))^2] = 2 \times \mathbb{E} [(X(p) - X(q))^2], \quad (4)$$

which is independent of σ . Note that equation 4 shows an asymmetric relationship between $X(p)$, $X(q)$ and $X(r)$, due to their spatial positioning. On the other hand, if the $X(p)$, $X(q)$ and $X(r)$ are randomly permuted, the spatial diffusion model doesn't hold. This is because by doing so, we essentially remove any relationship between the pixel values and the pixel locations. In that case, the effect of random permutations results in

$$\mathbb{E} [(X(p) - X(r))^2] = \mathbb{E} [(X(p) - X(q))^2] = \mathbb{E} [(X(q) - X(r))^2], \quad (5)$$

which is a symmetric relationship between $X(p)$, $X(q)$ and $X(r)$. Thus equation 4 and 5 represent opposite extremes of high spatial orderness and low spatial orderness respectively. We now have all the necessary tools to define the spatial orderness metric, and extend it for multiple scales. It is done in the following section.

3 Multi-Scale Spatial Orderness

Given a set of images, I_1, I_2, \dots, I_k , we first extract a fixed number (l) of triples of pixel intensities $(I_{n(1)}(p_1), I_{n(1)}(q_1), I_{n(1)}(r_1)), \dots, (I_{n(l)}(p_l), I_{n(l)}(q_l), I_{n(l)}(r_l))$, such that each triple of spatial locations (p_i, q_i, r_i) follows a 2-hop spatial arrangement. Here $n(i)$ denotes the image from which the i^{th} triple was extracted. Next, we define the spatial orderness measure at the lowest scale as follows,

$$so(I)^1 = \left(\frac{\mathbb{E}_i \left[(I_{n(i)}(p_i) - I_{n(i)}(r_i))^2 \right]}{\mathbb{E}_i \left[(I_{n(i)}(p_i) - I_{n(i)}(q_i))^2 \right]} \right) - 1. \quad (6)$$

Observe that $so(I)^1 = 1$, when the diffusion process $P \rightarrow Q \rightarrow R$ is strictly followed (equation 4), whereas it drops down to zero when all spatial order is removed e.g. by random permutation of pixel values.

With this, we can extend the definition of spatial orderness to multiple spatial scales. For that, a scale-space like decomposition is constructed by averaging $a \times a$ non-overlapping

input regions onto a single pixel value. We let these sets of new mean downsampled images be denoted as $I_1^a, I_2^a, \dots, I_k^a$. For each set of images at each scale, we denote their corresponding spatial orderness values by $so(I)^a$. At the end, we have a set of scalar values $so(I)^1, so(I)^2, \dots, so(I)^p$, which represent the spatial orderness of the data at various scales. We summarize some the ways in which this measure can be interpreted:

- Spatial orderness at the lowest scale is indicative of how much more accurately the value of a pixel can be interpolated from its neighbors *than* its non-neighbors which are at a distance of 2 hops.
- Randomly permuting the spatial locations of pixels (or blocks of pixels) will reduce spatial orderness. Conversely, when a randomly permuted version of an input has an equal likelihood of occurrence to its non-permuted form, the spatial orderness of data is zero at all scales.
- 2D inputs of the form $I[i, j] = a(I) + \mathcal{N}(0, \sigma)$ (a is a constant that can be different for different I), have zero spatial orderness. One can observe that random permutation of pixels does not change the image statistics in this case. Also, note that by controlling variation in $a(I)$ one can arbitrarily increase spatial "correlation" measures, hinting that correlation is not enough to capture spatial order in all cases.

4 Experiments

The experiments reported herewith are conducted on the MNIST [8], Fashion-MNIST [4] and CIFAR-10 [1] datasets. For MNIST and Fashion-MNIST, we perform 2×2 max-pooling after each convolution layer, whereas for CIFAR-10, pooling was only performed after each alternate convolution layer. Additionally, a two layer fully connected network is used for CIFAR-10, whereas a single fc layer is used for MNIST and Fashion-MNIST. For consistency we used 64 units in all hidden layers, with kernels of size 3×3 , except in section 5.2 (variation of kernel size).

4.1 Disrupting Spatial Orderness: Random Block-Swapping

Here we describe a method for disrupting the spatial orderness of the data, by performing block-swapping on the input. First we divide $(N \times N)$ images into blocks of size $k \times k$, such that $N/k \times N/k$ blocks span the entire image. Next, in each iteration of block swapping, a random chosen pair of image blocks are entirely swapped. We then repeat this process for N_s number of iterations. More swaps (larger N_s) will lead to a greater disruption of spatial order, and thus should elicit lower values of spatial orderness, and vice-versa. Furthermore, the block size (k) of the swap is relevant: swapping for a certain k must not greatly impact the spatial orderness at scales less than k , as the spatial arrangement in those scales is not overly affected.²

4.1.1 Random Block-Swapping: Impact on Spatial Orderness

To analyze the effect of block-swapping on spatial orderness measures at various scales, we simply vary the number of block-swap operations on each image of the corresponding datasets. Increasing the number of swaps leads to a steady reduction of spatial order as a

²Block-swapping with larger block-size cannot altogether avoid disrupting the spatial order at lower scales, due to boundary effects of the blocks.

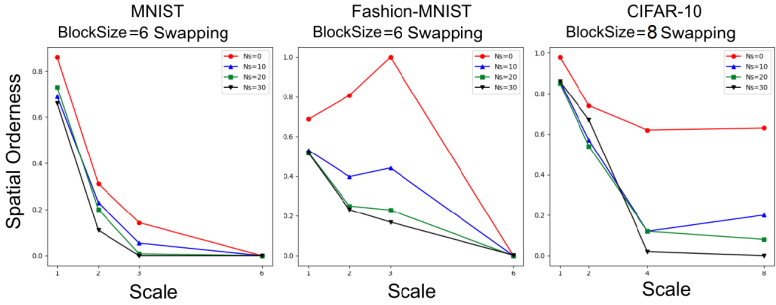


Figure 1: Spatial orderness of the MNIST, Fashion-MNIST and CIFAR-10 datasets at various scales, and their changes with block-swap operations performed on the data. For instance, the plots in red showcase the spatial orderness of the original, unswapped datasets ($N_s = 0$) at specific scales. For each dataset, block swapping was performed for a pre-selected block-size (specified on top).

whole, in the data. Therefore, a metric which measures spatial order must give smaller values when many block-swap operations are performed on the input. For our experiments, we choose four different number of block-swaps ($N_s = (0, 10, 20, 30)$) for the datasets of MNIST (BlockSize=6), Fashion-MNIST (BlockSize=6) and CIFAR-10 (BlockSize=8), generating a total of 12 datasets: MNIST-swap₆(0,10,20,30), CIFAR10-swap₈(0,10,20,30) and Fashion-MNIST-swap₆(0,10,20,30). The results are shown in figure 1.

First, we look at how spatial orderness changes with scale. As expected, we find that in all three datasets, spatial orderness at the highest scale is significantly lower than in the initial scales. This fact re-affirms the apparent "bag-of-words" like organisation of images at higher scales (objects or patterns are more positionally decorrelated at higher scales) (see [8]).

Next, we note the impact that block swapping has on spatial orderness of the corresponding scales. In all cases, we observe a clear reduction (to zero) of spatial orderness with a greater amount of block swaps, at the corresponding scales. Since block swapping is done with relatively larger block-size, the spatial orderness of the data for scales less than the block-size is not overly affected.

4.1.2 Classification experiments: Is greater convolutional depth always better ?

Here we document CNN classification performance on MNIST-swap₆(0,10,20,30), CIFAR10-swap₈(0,10,20,30) and the Fashion-MNIST-swap₆(0,10,20,30) datasets. The objective of this experiment is to discover if adding convolutional layers to a CNN is still beneficial, when the spatial orderness of the data has been reduced at higher scales. Our primary hypothesis is that convolution layers exploit the spatial orderness of data at multiple scales. Hence, for block-swapped data, we must expect the addition of convolution layers (beyond the scale of the swap) to pay decreasing dividends. Furthermore, because the block-swaps are only done at a higher scale, we should still find that adding initial convolution layers are beneficial, as spatial orderness of initial scales are still preserved (Figure 1).

Results are shown in figure 2. As hypothesized, we find that indeed adding convolution layers lead to decreasing gains, for larger number of block-swaps at the corresponding scales (larger N_s). Also, as anticipated, we observe that initial additions of convolution layers

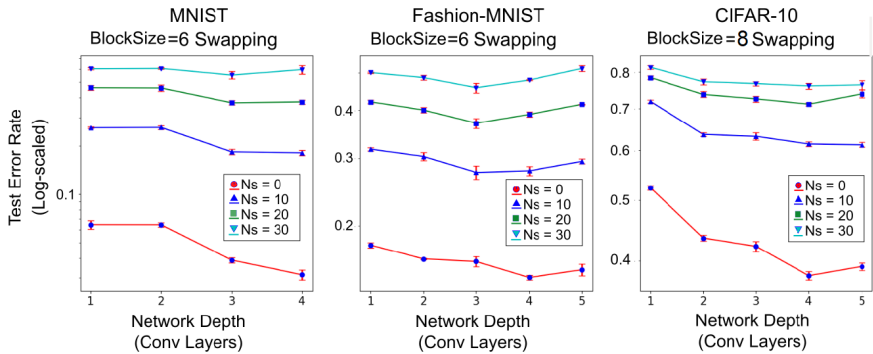


Figure 2: Semilog plots showing the test error rate of networks of different depths, trained on data corrupted by various degrees of spatial block-swapping ($N_s = (0, 10, 20, 30)$) on three different datasets (MNIST, Fashion-MNIST and CIFAR-10). Note that for data lacking in spatial order ($N_s > 0$), depth additions beyond a certain point do not yield improvements. Instead, such additions often significantly increase error rate, for larger N_s .

reduce test errors irrespective of block swapping.

Our findings are consistent with [14], where it was theoretically shown that stacked convolution layers are optimally priored for learning compositional functions. Greater convolutional depth implies greater compositionality which effectively spans a range of scales in the input. By removing spatial order at higher scales by random block-swapping, we are essentially disrupting the compositional structure of the data at higher scales, which leads to diminishing improvements with adding more depth.

4.2 Spatial Orderness of CNN Feature Maps

The previous sections demonstrate that convolutions are more effective when the input data has spatial order at multiple scales. However, note that a convolutional module at a depth of $n + 1$ does not compute on the input data, but rather on the feature map of the n^{th} layer. Since each convolution layer sees the feature map of the previous layers as its input, the computation of the spatial orderness of the feature maps would be a meaningful step. For each feature response map (denoted as a function $f(\cdot)$), computation of spatial orderness proceeds in the same way as for 2D images, treating the feature responses across all training examples as the set of 2D images $f(I_1), f(I_2), \dots, f(I_k)$. As convolution outputs usually have multiple feature maps, the spatial orderness of a layer is estimated as the mean value of the spatial orderness of each individual feature map within that layer.

4.2.1 MNIST and CIFAR-10: Training-time Progression

CNNs are trained on the original datasets of MNIST (3 conv layers), Fashion-MNIST (3 conv layers) and CIFAR-10 (4 conv layers). During training, the validation accuracies and spatial orderness of feature maps at the end of each epoch was monitored. As our initial experiments in Section 4.1.1 demonstrated, spatial orderness of images at highest scale is usually low, which indicates a bag-of-words like, spatially decorrelated organisation of information at higher scales. Hence the feature maps of CNN can be treated in the same way w.r.t their spatial orderness measures, i.e. low spatial orderness must signify higher levels of

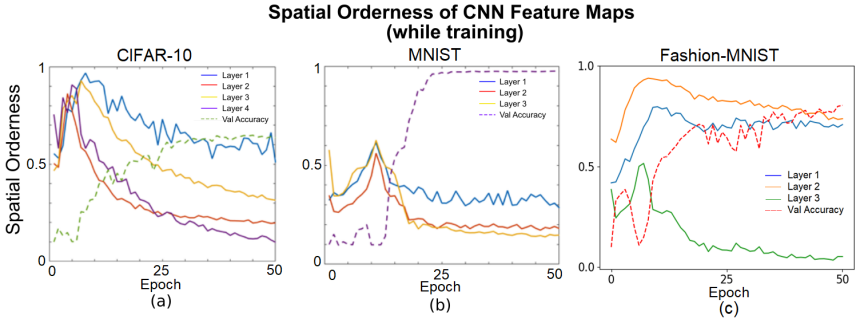


Figure 3: (a), (b) and (c) show the progression of mean spatial orderness of the feature maps (of each layer), while training on CIFAR-10, MNIST and Fashion-MNIST respectively. The validation performance at the end of each epoch is also shown (dotted lines).

abstraction and spatially less redundant information. Concurrently, higher spatial orderness must signify a greater level of spatial redundancy in information, as it essentially implies that spatial gaps in feature values can be effectively "filled in" by means of interpolation using the feature values of the pixel neighbors.

The spatial orderness progression of feature maps of all the CNN layers are shown in figure 7 (additional plots available in supplementary material). For both datasets, we find that the spatial orderness measures show a synchronized "peaking" at the beginning of training. This is an unusual finding, as we should expect the features to get higher in abstraction with more training. Curiously, we do observe that the validation accuracy of the networks *do not* show any steady improvement in this phase. Subsequently, after the peak, we find that the spatial orderness of all layers show a synchronized decrease. The spatial orderness of the last layer shows the greatest reduction after its peak; signifying that spatially de-correlated, abstract concepts are learned mainly in the higher layers. This finding is consistent with our current understanding of CNNs: generalization performance only starts improving when the higher layer representations start capturing increasingly abstract features (see particularly (c) in figure 7). Interestingly, we also notice that the spatial orderness of the other layers decrease with training as well, even including the first layer (blue trajectories). This indicates that while a CNN does eventually captures spatially non-redundant abstract concepts in its deeper layers, all other layers also sequentially strive to capture spatially non-redundant features of low spatial orderness.

5 Spatial Orderness of Kernels

5.1 Theoretical Results

We note that just like the inputs and the feature maps, one can treat the kernels (of size $K \times K$) as 2D images themselves. As such, it is also possible to compute the spatial orderness within the kernels, at the end of training. Convolution is linear in nature, and will elicit larger output responses when the input patches are highly correlated to the kernel form. Thus, kernels with very low spatial order are not likely to extract visually meaningful features, and vice-versa. Hence, from a feature extraction point of view, it is desirable that weights exhibit high spatial

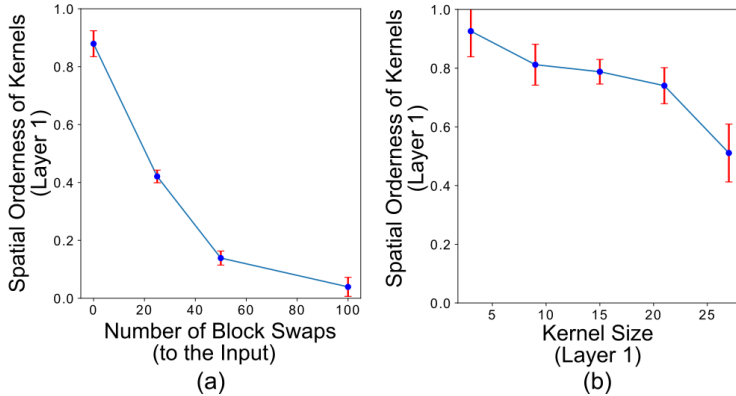


Figure 4: (a) demonstrates that disruption of spatial orderness at the input has an immediate effect on the spatial orderness of the kernels, and (b) shows that the size of kernels affect the spatial orderness of the trained kernels. All experiments were done on MNIST-1000 (1000 training examples used) and each experiment was repeated across six random splits of the data.

orderness.

Here we summarize our theoretical results on the spatial orderness of kernels. Please find our main theoretical results (Theorems 1, Corollaries 1.1 and 1.2) and proofs in the supplementary material. We summarize the theorems as follows.

- **Theorem 1 and Corollary 1.2: How is the spatial orderness of kernels and the spatial orderness of the feature map input related ?** We find that the spatial orderness of the kernels are likely to be higher when the inputs themselves have higher spatial orderness.
- **Corollary 1.1: How is the spatial orderness of kernels related to the choice of kernel size ?** We find that choosing a larger kernel size can lead to kernels with lower spatial orderness³. This shows that the choice of kernel size is quite important w.r.t ensuring spatially ordered kernels.

5.2 Experimental Validation of Results

A 3-layered CNN was trained on six random splits of the MNIST dataset, with random network initializations each time. Only 1000 examples were used for training. Two separate experiments were conducted, guided by the theoretical results: (a) MNIST-swaps₃(0,25,50,100) datasets were generated from each random split on MNIST, on which the CNNs were trained, and (b) kernel size of the first layer was varied in the range (3,27) (for zero-padded convolutions), and the networks were trained on conventional MNIST. The objective of the first experiment was to examine the relationship between the spatial orderness of the input and the kernels; whereas the second experiment’s goal was to discover how the spatial orderness of the kernels was dependent on the choice of kernel size. Note that for the second experiment,

³Note that by "spatial orderness of kernels" we mean the average spatial orderness of post-trained kernel weights (averaged across all kernels within a layer). In the following section, we empirically substantiate the results in the theorems.

the convolution padding of the first layer was zero; i.e. for large kernel sizes the convolution effectively becomes a fully connected layer.

The results are shown in the plots of Figure 4. We find that they conform to the predictions of our theoretical results. We observe a sharp reduction in the spatial orderness of kernels when the input data has less spatial order (using greater block-swaps). Furthermore, we find that the kernel spatial orderness reduces with larger kernel size, almost by a factor of half (from $K=3$ to $K=27$). Note that in the second experiment, no block-swapping was performed on the input.

These results add an interesting perspective on the debate of CNNs versus FC-NNs. Taken together, the results imply that a CNN is more likely to extract visually meaningful (spatially ordered) features, subject to two necessary conditions: (a) the kernel size of the convolutions are small (i.e. more CNN than FC-NN like) and (b) the data on which the network is trained exhibits high spatial orderness.

6 Discussions: Connection to Other Works

Recently it was found that on Imagenet, a bag-of-features based approach with shallow CNNs performs surprisingly close to bigger models which exploit spatial structure at higher scales [9]. Hence, spatial arrangement information beyond a certain scale is not very yielding in terms of improving classification performance. This is consistent with our findings in this paper. As demonstrated in section 4.1.1, the spatial orderness of the image at higher scales is usually lower than the spatial orderness at lower scales, i.e. the information at the higher scales is more "bag-of-words" like than in lower scales. Furthermore, experiments in section 4.2.1 and section 3 of the supplementary material reveal that deeper convolutional feature maps indeed mirror the spatial organization of the input at higher scales.

Another example of testing the generalization abilities of CNNs on the data is in [9]. The authors observe that the CNN fails to generalize well when recognizability-preserving fourier domain filter masks were applied to the input. Three such data distortions were explored: no filtering, low-pass radial mask and random mask. Throughout their experiments, the authors observe that the CNN trained on the low pass filtered radially-masked inputs showed the most consistent performance across datasets, having the smallest generalization gap (among the no data-augmentation schemes). Our analysis on the spatial orderness of kernels in section 5 provides a possible explanation for this observation. Low-pass filtering enhances the spatial orderness of the input, compared to the random mask and unfiltered schemes. By doing so, it also ensures that the kernels within the architecture show greater spatial orderness; which ensures more consistent performance across data distortion variations.

7 Conclusions

A new statistical measure for quantifying spatial order within 2D data at various scales was proposed, called spatial orderness. This measure was shown to be indicative of the spatial organization at various scales, decreasing in value in correlation to the amount of block-swapping performed on the input. The performance gains from adding convolution layers was demonstrated to weaken with greater block-swapping disruption. Interesting bi-phasic trend in the spatial orderness of feature maps was observed during training. Theoretical and empirical results demonstrated the correlation between the spatial orderness of trained

kernels, and the spatial orderness of the input. Additionally, we find that spatial orderness of kernels shows a significant drop with greater kernel-size, as it approaches a FC-NN like configuration.

Acknowledgments

This research was supported by DSO National Laboratories, Singapore (grant no. R-719-000-029-592). We thank Dr. Loo Nin Teow and Dr. How Khee Yin for helpful discussions.

References

- [1] Yoshua Bengio and Yann Lecun. *Scaling learning algorithms towards AI*. MIT Press, 2007.
- [2] Yoshua Bengio et al. Representation learning: A review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, August 2013. ISSN 0162-8828. doi: 10.1109/TPAMI.2013.50.
- [3] Wieland Brendel and Matthias Bethge. Approximating CNNs with bag-of-local-features models works surprisingly well on ImageNet. In *ICLR*, 2019.
- [4] K. He et al. Deep residual learning for image recognition. In *2016 IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, June 2016. doi: 10.1109/CVPR.2016.90.
- [5] Jason Jo and Yoshua Bengio. Measuring the tendency of CNNs to learn surface statistical regularities. *ArXiv*, abs/1711.11561, 2017.
- [6] Kenji Kawaguchi et al. Generalization in deep learning. In *Mathematics of Deep Learning*, Cambridge University Press, to appear., 2018.
- [7] Alex Krizhevsky. Learning multiple layers of features from tiny images. *Technical Report TR-2009*, University of Toronto, Toronto, 2009.
- [8] Y. Lecun. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>.
- [9] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4898–4906. Curran Associates, Inc., 2016.
- [10] H. N. Mhaskar and T. Poggio. Deep vs. shallow networks: An approximation theory perspective. *Analysis and Applications*, 14(06):829–848, 2016.
- [11] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [12] Shizhao Sun et al. On the depth of deep neural networks: A theoretical view. In *AAAI Conf. on Artificial Intelligence*, AAAI’16, pages 2066–2072. AAAI Press, 2016.

- [13] Li Wan et al. Regularization of neural networks using dropconnect. In *ICML*, pages 1058–1066, 2013.
- [14] Han Xiao et al. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *ArXiv*, abs/1708.07747, 2017.
- [15] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *ECCV 2014*, 2014.
- [16] Pan Zhou and Jiashi Feng. Understanding generalization and optimization performance of deep CNNs. *ICML*, 2018.

A Theorems

Please refer to section 5 in the main paper for the implications of the theorems, and empirical validation of the theoretical results.

Theorem 1. *Consider a layer within a Convolutional Neural Network trained by backpropagation, which contains kernel of size $(K \times K)$, s.t. $K \geq 3$. After training, let $w(p), w(q), w(r)$ be extracted tuples of kernel weight values from a particular instance of the kernel W , for locations (p, q, r) within the kernel, following the two-hop spatial arrangement. Let the input feature map for W be denoted as X . Consider all spatial locations within the input feature map a, b and c , which follow the same two hop spatial arrangement as p, q, r . Then we must have,*

$$\begin{aligned} \mathbb{E}_{p,q} [|w(p) - w(q)|^2] &\leq \alpha \mathbb{E}_{a,b,n} [|X_n(a) - X_n(b)|^2], \text{ and} \\ \mathbb{E}_{p,r} [|w(p) - w(r)|^2] &\leq \alpha(1 + so(X)^1) \mathbb{E}_{a,b,n} [|X_n(a) - X_n(b)|^2], \end{aligned}$$

for a certain non-zero valued α . The expectation on the left of the inequality is taken over all possible p, q, r locations within the kernel, which obey the two-hop spatial arrangement. $so(X)^1$ as usual denotes spatial orderness of the input feature map at the first scale, averaged across all examples of the feature map X_1, X_2, \dots, X_N .

Proof. Let X_{rect} be any randomly located rectangular region of size $(K \times K)$ within the i^{th} feature map which is the input to W , and the corresponding scalar output node denoted by o_k , which belongs to the j^{th} feature map in the next layer. The weight update rule for the kernel W , only pertaining to the backpropagation error signal at o_k , at epoch t , for input X_n is

$$\Delta W = -\eta(t) X_{rect} \delta_k^n. \quad (7)$$

Here δ_k^n is the backprop error signal at o_k and $\eta(t)$ is the gradient descent update rate. Let us denote the corresponding nodes within X_{rect} by $(X_n(a), X_n(b), X_n(c))$ which are respectively attached to $(w(p), w(q), w(r))$. Note that the spatial relationship between the feature map locations a, b and c is the same as between p, q and r , i.e. both are 2-hop arrangements in their respective domains. The above update rule only considers a single output node for the update. If one were to consider all updates across all examples and outputs, then the final kernel value for $w(p)$ obtained after T epochs can be simplified as,

$$w(p) = \sum_{t,n,a} -\eta(t) \delta_k^n X_n(a) = \sum_{n,a} C_n X_n(a), \quad (8)$$

where $C_n = \sum_{t,k} \eta(t) \delta_k^n$. Similarly to equation 8, one can obtain forms for $w(q)$ and $w(r)$. Finally we have,

$$\begin{aligned} \mathbb{E}_{p,q} [|w(p) - w(q)|^2] &= \mathbb{E}_{a,b,n} \left[\left(\sum_n C_n (X_n(a) - X_n(b)) \right)^2 \right] \\ &\leq \mathbb{E}_{a,b,n} \left[\sum_n (C_n)^2 \sum_n (X_n(a) - X_n(b))^2 \right] \leq \alpha \mathbb{E}_{a,b,n} [|X_n(a) - X_n(b)|^2], \end{aligned} \quad (9)$$

$$(10)$$

where $\alpha = \sum_n (C_n)^2 N_{total}$. Here N_{total} is the total number of backpropagation driven updates on the kernel values. Similarly, for the locations p and r which are at a 2-hop distance, we have

$$\mathbb{E}_{p,r} [|w(p) - w(r)|^2] \leq \alpha \mathbb{E}_{a,c,n} [|X_n(a) - X_n(c)|^2] \leq \alpha(1 + so(X)^1) \mathbb{E}_{a,b,n} [|X_n(a) - X_n(b)|^2], \quad (11)$$

where $so(X)^1$ is the spatial orderness of the input feature map at the first scale. This completes the proof. \square

Corollary 1.1. *Consider a trained CNN with the same setup as in Theorem 1, with the same symbol definitions. Here, we vary the size of the kernel $K \times K$, with all other network parameters unchanged. For simplicity, we consider the kernels from the first layer. Thus, now each X_i is simply the 2D inputs to the CNN, of size $S \times S$. We consider zero-padded convolutions in the first layer. This ensures that when $K = S$, the convolution layer simply becomes a fully connected layer. We define*

$$D_{ab} = \mathbb{E}_{a,b,n} [|X_n(a) - X_n(b)|^2]. \quad (12)$$

Let us also denote gaussian random variables

$$\epsilon_1 \sim \mathcal{N}\left(0, \frac{\sigma_1^2}{(S-K+1)^2N}\right) \text{ and } \epsilon_2 \sim \mathcal{N}\left(0, \frac{\sigma_2^2}{(S-K+1)^2N}\right), \quad (13)$$

for certain non-zero real constants σ_1 and σ_2 . We wish to compute the uncertainty in the upper bounds of the kernel value differences. It follows that

$$\begin{aligned} |w(p) - w(q)|^2 &\leq \alpha(K)(D_{ab} + \epsilon_1), \text{ and} \\ |w(p) - w(r)|^2 &\leq \alpha(K) \left(1 + so(X)^1 + \frac{\epsilon_2}{D_{ab}} \right) D_{ab}, \end{aligned}$$

for a certain non-zero valued $\alpha(K)$, which is only a function of the kernel size K . Here p, q, r are fixed kernel locations which follow the 2-hop arrangement.

Proof. The main observation required to prove this result is that the number of updates of each kernel value is dependent on the size of the kernel, given that there are a fixed number of training examples. More precisely, for a kernel of size $K \times K$, and an input of size $S \times S$, the total number of updates, N_{total} , to each element of a kernel, is proportional to $(S - K + 1)^2N$ (for zero-padded convolutions), s.t. one can write $N_{total} = \beta(S - K + 1)^2N$. We refer the

reader to equation 10 from Theorem 1. By removing the expectation operator under p, q , we can reformulate the inequality as

$$|w(p) - w(q)|^2 \leq \left(\sum_n (C_n)^2 N_{total} \right) \frac{\sum_{n=1}^{N_{total}} |X_n(a) - X_n(b)|^2}{N_{total}} \quad (14)$$

$$\leq \alpha(K) \left(\mathbb{E}_{a,b,n} [|X_n(a) - X_n(b)|^2] + \mathcal{N} \left(0, \frac{\sigma_1^2}{N_{total}} \right) \right) \quad (15)$$

$$\leq \alpha(K) \left(D_{ab} + \mathcal{N} \left(0, \frac{\sigma_1^2}{(S-K+1)^2 N} \right) \right). \quad (16)$$

Here, σ_1^2 represents the uncertainty involved in the computation of $\mathbb{E}_{a,b,n} [|X_n(a) - X_n(b)|^2]$, for $N_{total} = 1$. Also, $\sigma_1^2 = \sigma_1'^2 / \beta$. Similarly, we can reformulate equation 11, to determine the upper bound on the two-hop kernel value difference as

$$|w(p) - w(r)|^2 \leq \left(\sum_n (C_n)^2 N_{total} \right) \frac{\sum_{n=1}^{N_{total}} |X_n(a) - X_n(c)|^2}{N_{total}} \quad (17)$$

$$\leq \alpha(K) \left(\mathbb{E}_{a,b,n} [|X_n(a) - X_n(c)|^2] + \mathcal{N} \left(0, \frac{\sigma_2^2}{(S-K+1)^2 N} \right) \right) \quad (18)$$

$$\leq \alpha(K) \left((1 + so(X)^1) \mathbb{E}_{a,b,n} [|X_n(a) - X_n(b)|^2] + \mathcal{N} \left(0, \frac{\sigma_2^2}{(S-K+1)^2 N} \right) \right) \quad (19)$$

$$\leq \alpha(K) \left(1 + so(X)^1 + \frac{\varepsilon_2}{D_{ab}} \right) D_{ab}. \quad (20)$$

The definitions of σ_2^2 and $\sigma_2'^2$ are analogous to σ_1^2 and $\sigma_1'^2$ before. This completes the proof. \square

Corollary 1.2. Consider a CNN trained with the same setup and symbol definitions as in Theorem 1. We bound kernel value differences averages across regions of size $d \times d$. Let us denote non-overlapping set of kernel locations p_d, q_d and r_d , each of size $d \times d$, which follow the two hop spatial arrangement. Similarly we denote non-overlapping set of feature map locations in X by a_d, b_d and c_d , each of size $d \times d$, which follow the two hop spatial arrangement. It can then be shown that,

$$\mathbb{E}_{p,q} \left[\left(\mathbb{E}_{p \in p_d} [w(p)] - \mathbb{E}_{q \in q_d} [w(q)] \right)^2 \right] \leq \alpha \mathbb{E}_{a,b,n} \left[\left(\mathbb{E}_{a \in a_d} [X_n(a)] - \mathbb{E}_{b \in b_d} [X_n(b)] \right)^2 \right], \text{ and}$$

$$\mathbb{E}_{p,r} \left[\left(\mathbb{E}_{p \in p_d} [w(p)] - \mathbb{E}_{r \in r_d} [w(r)] \right)^2 \right] \leq \alpha (1 + so(X)^d) \mathbb{E}_{a,b,n} \left[\left(\mathbb{E}_{a \in a_d} [X_n(a)] - \mathbb{E}_{b \in b_d} [X_n(b)] \right)^2 \right].$$

Note that the above is simply a generalization of Theorem 1 ($c = 1$) to arbitrary scales.

Proof. Note the final mathematical expression for $w(p)$ in equation 8. Extending that equation to the average of kernel values across a region of size $d \times d$, $\mathbb{E}_{p \in p_d} [w(q)]$, we have

$$\mathbb{E}_{p \in p_d} [w(q)] = \sum_{t,n,a} -\eta(t) \delta_k^n \mathbb{E}_{a \in a_d} [X_n(a)] = \sum_{n,a} C_n \mathbb{E}_{a \in a_d} [X_n(a)]. \quad (21)$$

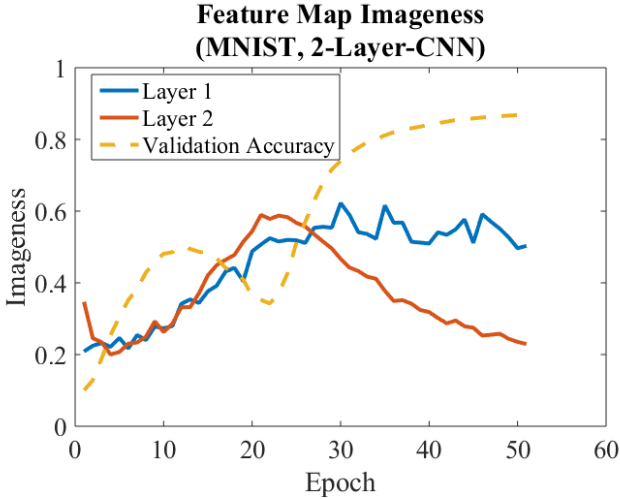


Figure 5: A CNN with two convolutional layers of 3×3 was trained on MNIST, with a small learning rate of 0.0001. Shown are the average spatial orderness of feature maps at the end of each epoch of training, across 50 epochs. Also shown is the validation accuracy of the network at the end of each epoch. As has been observed with previous plots, the validation accuracy starts an unconstrained jump, only after the spatial orderness of the final layer (here Layer 2) starts decreasing.

Subsequently, after a trivial application of Cauchy-Schwarz inequality to the expression $\mathbb{E}_{p,q} \left[\left(\mathbb{E}_{p \in P_d} [w(p)] - \mathbb{E}_{q \in q_d} [w(q)] \right)^2 \right]$, similar to equation 10, the results follow. \square

B Feature Map Spatial Orderness: Progression during Training

Figures 1 and 2 contains plots which depict the progression average spatial orderness of feature maps while training, for different network architecture choices. In Figure 1, the learning rate was decreased 10-fold, such that the co-occurrence of phase (b) and validation accuracy increase can be more precisely observed.

C Correlations between Spatial Orderness of Inputs and Feature Maps

We wish to see whether any correspondence exists between the spatial orderness of the input I_1, I_2, \dots, I_k at various scales, and the spatial orderness of the feature maps at various depths. We train a CNN on the MNIST-swap(0,10,30,60) datasets, in which after every 3×3 convolution, a 2×2 max pooling operation is performed. Note that "effective" receptive field size of each position within a feature map is dependent on the depth of the feature map. Also note that max-pooling operations reduce the overlap between neighboring feature map

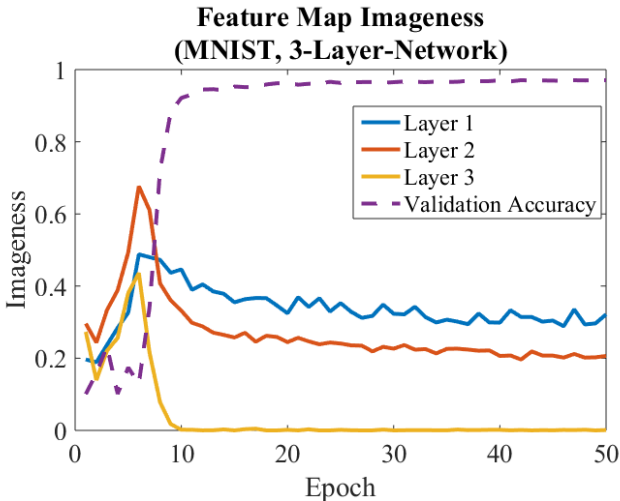


Figure 6: A CNN with three convolutional layers of 3×3 was trained on MNIST with a learning rate of 0.001. The average spatial orderness of feature maps, along with the validation accuracy of various layers are shown across 50 epochs.

Change of Network Depth	Error Rate improvement (%)				Spatial Orderness of Final Layer				Correlation
	Ns=0	Ns=10	Ns=30	Ns=60	Ns=0	Ns=10	Ns=30	Ns=60	
1->2	0.42	0.38	0.29	0.18	0.35	0.35	0.24	0.16	0.99
2->3	0.15	0.14	0.06	-0.09	0.08	0.09	0.07	0.04	0.985

Table 1: Table showing the mean reduction of validation error rate (in %), when extra convolution layers are added to a CNN, tested and trained on MNIST. To the mid-right, the mean spatial orderness measures computed at the final convolution layer output of the CNNs are shown, before the layer addition. All block-swapping was done at Scale=3.

units. These two aspects combined hint that the spatial orderness of the feature map at depth k could be related to the spatial orderness of the input at a scale of $2k$. The plot in figure 7 shows the post-training spatial orderness of feature maps (at depths 1,2,3) plotted against the spatial orderness of the input at the corresponding scale (for scales 2,4,6). We can observe that they are correlated. Also observe that reducing the spatial orderness of the input with greater block-swaps, leads to a correlated decrease of the spatial orderness of all layers. As This shows that convolutions operations are intrinsically configured to preserve the spatial order in the input.

D Experiments on MNIST-swap

We wish to narrow down our hypothesis, by conjecturing that a convolution layer directly benefits from the spatial order in the feature maps that it sees. Thus, we compute the spatial orderness of the feature maps of various layers within a CNN. For this experiment, we train CNNs on datasets MNIST-swap(0,10,30,60), consisting of 1,2 and 3 convolution layers each. Next, the spatial orderness of only the final convolution layer’s output of each CNN (after the max-pooling) is recorded, using the approach detailed in the previous paragraph. We

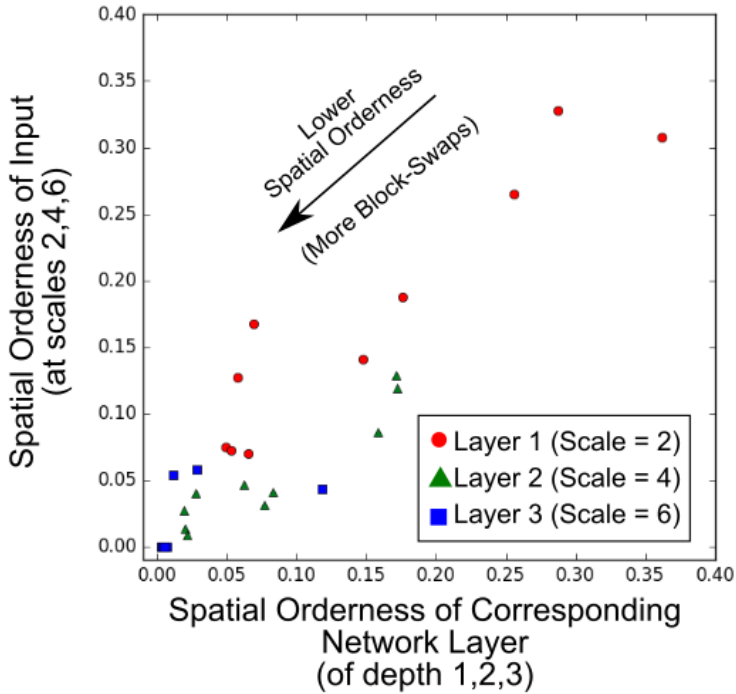


Figure 7: Correlation between spatial orderness of the input data at scales (2,4 and 6), and the mean spatial orderness within the network feature maps at the corresponding depths (1,2 and 3), after 50 epochs of training.

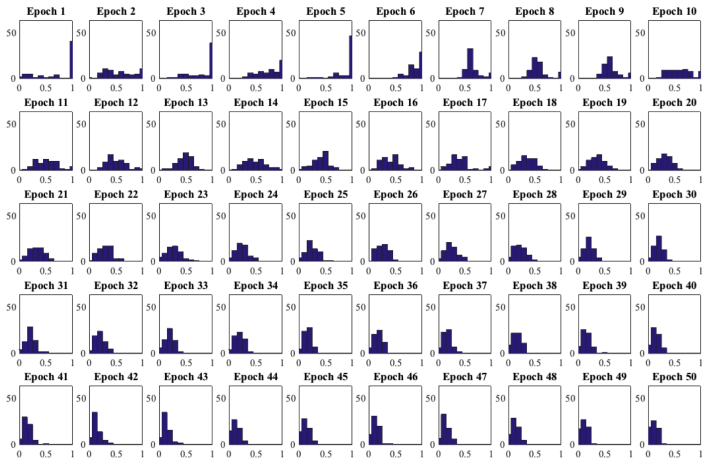


Figure 8: Shown are the histograms of the spatial orderness measures of the feature maps of the final layer of a 3-layer CNN trained on MNIST. The smooth transition from high to low spatial orderness is noticeable.

also note the corresponding improvement in error rates (in %), between CNNs containing different number of conv layers. For accurate testing, training was repeated for six trials for each CNN, and the all recorded measures (spatial orderness and error rate improvements) are averaged across all the trials. To test our narrowed hypothesis, we simply note the correlation between the spatial orderness of the final layer feature maps within a CNN, and the error improvement obtained by adding another convolution layer to the CNNs. Results are shown in Table 1.

The observations are two-fold. First, we observe that reducing spatial orderness of the input leads to reduction of the spatial orderness in the feature maps, which is intuitively sound. Note that more detailed experiments on this correspondence is provided in the supplementary material. Next, we find a significant correlation between the error rate improvement from convolutional layer additions, and the spatial orderness of the final feature maps within a CNN. These observations explain the results in Section 4.1.2.

E Variation of Spatial Orderness within feature maps

In the experiments thus far, all feature maps within each layer were used to compute a single, mean value of spatial orderness. Here, we compute the spatial orderness measure of each feature map within the layers separately, and analyze the changes in their distribution as the network is trained. Figure 8 shows how spatial orderness is distributed among the units of the final layer of a 3-layer CNN being trained on MNIST, and the changes with training. To summarize the observations,