

A STRUCTURALLY REGULARIZED CONVOLUTIONAL NEURAL NETWORK FOR IMAGE CLASSIFICATION USING WAVELET-BASED SUBBAND DECOMPOSITION

Pavel Sinha, Ioannis Psaromiligkos, Zeljko Zilic

McGill University, Department of Electrical & Computer Engineering, Montreal, Canada
 pavel.sinha@mail.mcgill.ca, ioannis.psaromiligkos@mcgill.ca, zeljko.zilic@mcgill.ca

ABSTRACT

We propose a convolutional neural network (CNN) architecture for image classification based on subband decomposition of the image using wavelets. The proposed architecture decomposes the input image spectra into multiple critically sampled subbands, extracts features using a single CNN per subband, and finally, performs classification by combining the extracted features using a fully connected layer. Processing each of the subbands by an individual CNN, thereby limiting the learning scope of each CNN to a single subband, imposes a form of structural regularization. This provides better generalization capability as seen by the presented results. The proposed architecture achieves best-in-class performance in terms of total multiply-add-accumulator operations and nearly best-in-class performance in terms of total parameters required, yet it maintains competitive classification performance. We also show the proposed architecture is more robust than the regular full-band CNN to noise caused by weight-and-bias quantization and input quantization.

Index Terms— CNN; wavelet-based subband decomposition; image classification; regularization

1. INTRODUCTION

Deep learning has resulted in state-of-the-art performance in image recognition and vision tasks. Most of these achievements can be attributed to the use of convolutional neural networks (CNNs) [1]. Since then, several other improvements to the CNN architecture have been proposed, including AlexNet [2], VGG [3], GoogleNet [4], ResNet [5], Spatial Pyramid Pooling [6], SqueezeNet [7], and more.

The increasing complexity of CNNs poses challenges to state-of-the-art implementations. There are numerous techniques to reduce the computational cost of CNNs. Pruning

of filters to simplify a CNN was proposed in [8]. Another approach that used sparsity to reduce the number of filters per channel and per stage of a CNN was introduced in [9]. SqueezeNet was introduced in [7] that claimed $50\times$ fewer parameters than AlexNet, by using 1×1 convolutional filters and reducing the overall number of parameters. Another work on model compression was introduced by [10]. A characteristic shared by most of these methods is that they can be reduced architecture-wise to special cases of the base CNN introduced in [2].

Deep networks require several layers of weights to be trained and even with millions of training data samples, overfitting remains inevitable [11]. Some recent techniques to combat overfitting include data augmentation [2], weight regularization [12], dropouts [2], and adaptive regularization of weight vectors [12]. There is also a notion of *structural regularization*, wherein constraints are imposed on the network structure rather than on the weight updates to limit overfitting [13]. Several works focus on this approach. A jointly enforced global wavelet domain sparsity constraint together with a learned analysis sparsity prior was introduced in [14]. A wavelet-regularized semi-supervised learning algorithm using suitably defined spline-like graph wavelets was introduced in [15]. In a recent work, a method of Graph-Spectral-Regularization was introduced in [16].

Multi-resolution analysis using wavelets was introduced by Daubechies [17]. It is well known that decomposing an image into subbands using wavelets is advantageous for image analysis. Not surprisingly, wavelets have been used with CNNs in several works. In [18], a single layer CNN was proposed in which the convolution kernel was wavelet-based. The model could not utilize the subbands from a regularization perspective and did not present an automated learning strategy as in deep learning. Another approach [19] used wavelet decomposition for hierarchical image reconstruction to analyze CT scan images. A similar work on multi-resolution analysis with wavelets and CNNs was presented in [20]. The use of a scattering network as a generic and fixed initialization of the first layer in a CNN achieved similar results compared to learning the weights of the first layer was seen in [21].

In this paper, and in the context of image classification,

This research has benefited from a TechAccel grant, McGill University.

©2019 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

Published in the *Proceedings of the 2019 IEEE International Conference on Image Processing (ICIP)*. DOI: 10.1109/ICIP.2019.8804202

we leverage subband decomposition to introduce a new structurally regularized CNN architecture wherein multiple CNNs are used to process the input image at different spatial scales as represented by the critically sampled and equally band-limited subbands obtained through a wavelet decomposition. The new architecture represents a departure from the ones used in the above-mentioned methods, where a single CNN was used to process the complete multi-scale wavelet output. As we will see through qualitative arguments and extensive experimental studies the proposed architecture exhibits characteristics that translate to significant computational and classification performance advantages.

2. PROPOSED ARCHITECTURE

The proposed architecture termed Subband Regularized CNN (SRCNN) is presented in Figure 1. In the first stage, the input image is decomposed into subbands through a 2D discrete wavelet transform (2D-DWT). The SRCNN architecture is based on processing each of the subbands separately by individual CNNs. The field of view of each CNN is hence restricted to a dedicated subband, making each CNN indifferent to the rest of the subbands. Importantly, this subband decomposition structure reduces the overall computational cost.

We represent the complete decomposition of the input image X_{in} into K subbands by:

$$(X_0^1, \dots, X_0^K) = \text{DWT}(X_{\text{in}}, K, M) \quad (1)$$

where M is the number of DWT layers, K is the number of subbands, and X_0^k ($k = 1, \dots, K$) are the DWT coefficients for the k^{th} subband. We have chosen the Daubechies (D2) family of basis functions for DWT [22]. This constitutes the simplest Daubechies wavelet basis, with a single vanishing moment. Being symmetric, they offer linear phase characteristics and do not suffer from edge effect characteristics of higher order wavelets [22].

The SRCNN architecture in Figure 1 is a generalized structure. The exact configuration of the architecture implemented in this paper is given in Table 1. The input image is first decomposed into K subbands as described by Equation 1. The subbands are then individually passed through their corresponding CNNs. Finally, the fully connected (FC) layers combine the feature outputs of the subband CNNs and perform image classification. The output of the CNN at the k^{th} subband and i^{th} layer is given by:

$$X_{i+1}^k = \text{Pool}(\text{ReLU}(\text{Conv}(X_i^k, W_i^k), L_i^k), P_i^k) \quad (2)$$

where Conv represents the convolution between the input X_i^k of the i^{th} layer and the weights W_i^k . ReLU(\cdot) indicates the ReLU activation function with L_i^k representing the leakage percentage value [23] which is a real number between 0 and 1. Pool(\cdot) represents the pooling function with pooling parameters P_i^k . The outputs of the subband CNNs are accumulated

to yield X^{FC_0} which is the input to the first FC layer:

$$X^{\text{FC}_0} = (X_I^1, \dots, X_I^K) \quad (3)$$

where I is the number of layers in the subband CNNs. The output at each FC layer is given by:

$$X^{\text{FC}_{n+1}} = \text{ReLU}(W^{\text{FC}_n} \cdot X^{\text{FC}_n}, L^{\text{FC}_n}) \quad (4)$$

where X^{FC_n} denotes the output of the n^{th} FC layer, \cdot indicates matrix multiplication and L^{FC_n} indicates ReLU leakage value. Finally, the output of the last FC layer X^{FC_N} , indexed by N , produces the SRCNN's output Y . Equations 1 to 4 describe the complete input to output relation of the proposed subband based CNN.

3. PROPERTIES OF THE ARCHITECTURE

The proposed architecture emphasizes regularization through its structure, thus it is *structurally regularized*. To enhance regularization effectiveness, the decomposed subbands are critically sampled and band-limited before being processed by individual subband CNNs. Each of the subband CNNs is inhibited from accessing information across the entire spectrum of the input. Overall, each of the CNNs cannot learn sample-specific features present in the entire spectrum of the input. This restriction combined with weight regularization within each CNN, improves regularization, leading to better generalization ability and reduced overfitting, as demonstrated by the accuracy performance comparison in Table 2. Apart from accuracy, the difference in Top-5 and Top-1 accuracy result can be considered as an indicator for generalization effectiveness. A lower difference value indicates better generalization which, in our case, outperforms other state-of-the-art networks.

The lossless decomposition of the input spectrum into orthogonal subbands allows isolated analysis of the spatial representation of each subband. This is beneficial in the case of corrupted images. Indeed, corruption of the input image by noise, deformities from lens aberration, incorrect exposure, low lighting, etc., does not affect the entire spectrum equally; in reality, some subbands are corrupted more. Isolating the subbands ensures that the corruption of extracted features is limited to the affected subbands, as opposed to a full-band CNN that considers the entire spectrum for feature extraction.

Along similar lines, quantization noise in each weight is confined within the subband and does not affect the entire spectrum. In contrast, in a regular CNN, quantization noise in any weight can potentially corrupt the entire spectrum, since quantization noise can have large bandwidth. Results indicate that compared to full-band CNN, SRCNN proves more robust to weight and input quantization.

The subband decomposition also introduces high degree of sparsity in the subbands, specifically in the non-basebands

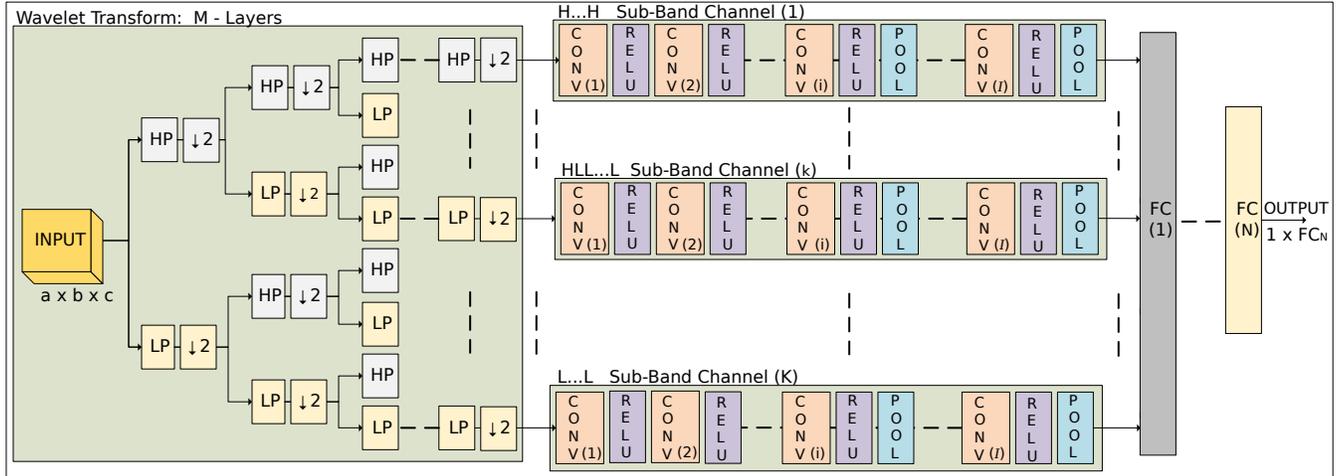


Fig. 1. Architecture of an M -layer SRCNN, parametrized by input dimensions ($a \times b \times c$), number of subbands K , convolutional layers per subband I , FC layers N and output classes FC_N , all open to optimization. The wavelet transform uses High-Pass (HP) and Low-Pass (LP) filters, followed by decimation of 2.

containing mostly edge information. This sparsity is introduced at the very input of the subband CNNs. It is well known that sparse inputs help reduce CNN complexity [9].

Random initialization of weights when training a full-band CNN does not guarantee scanning of the entire spectrum for useful features. In the proposed structure, the CNNs focus only on their corresponding subbands hence the entire spectral decomposed into subbands is covered equally.

The decomposition reduces the input spatial dimension along rows and columns by 2^M each, where M is the number of decomposition layers. The total reduction of input dimension is effectively on the order of 4^M for two-dimensional input data such as images. The convolution operation accounts for the bulk of computations in a CNN. The total computation cost depends super-linearly on the size of the convolution filters [24] and the sample point counts per dimension, all of which are significantly reduced in our case.

The subband decomposition architecture offers parallel computation along each subband. The parallelism also provides a mechanism to reduce internal memory footprint by sequentially computing each subband and reusing internal scratch memory to compute each subband CNN.

Finally, decomposition of input spectra into subbands is a generalized technique and can be applied to any CNN to improve regularization and thereby improve generalization capacity and improve overall performance.

4. EXPERIMENTAL SETUP AND RESULTS

4.1. Methodology and Training

We use MNIST, CIFAR-10/100, Caltech-101 and ImageNet-2012 datasets (used by [2]) in our experimentation. We compare the proposed architecture against two benchmarks: (i) a full-spectrum base CNN (BCNN) model that closely resem-

bles AlexNet [2] and VGG-16 [3]; (ii) the Transform CNN (TCNN) architecture which shares the same wavelet front-end as SRCNN, except that the subbands are combined and processed by a single CNN with the same number of weights and layers as BCNN. A single layer, subband decomposed TCNN (4 subbands) with 3 input color channels will result in a total input of 12 channels and with $1/2$ the length and width of the original input image. Table 1 shows the parameters of the models used. To study the effect of learning in subband domain, we keep most of the parameters constant across BCNN, SRCNN and TCNN. We compare the number of MAC operations, total number of parameters and accuracy with several well known state-of-the-art CNN architectures.

We train using stochastic gradient descent (SGD) with a mini batch size of 64, batch normalized, randomly picked images per mini batch, momentum of 0.9 and weight decay of 0.0005 [2]. The update equations for W_i^k are:

$$W_i^k(l+1) = W_i^k(l) + V_i^k(l+1) \quad (5)$$

$$V_i^k(l+1) = 0.9 V_i^k(l) - 0.0005 \epsilon W_i^k(l) - \epsilon \left. \frac{\partial L}{\partial w} \right|_{W_i^k(l)} \quad (6)$$

Here, l is the iteration index, $V_i^k(l)$ is the momentum at the l^{th} iteration and k^{th} subband, ϵ the learning rate, and $\left. \frac{\partial L}{\partial w} \right|_{W_i^k(l)}$ is the average over the l^{th} batch of the derivative of the objective function with respect to W_i^k , evaluated at $W_i^k(l)$. We initialize the learning rate to 0.01 and all biases to 1, while we initialize the weights by drawing from a Gaussian distribution with a standard deviation of 0.01.

4.2. Results

Classification Accuracy: Table 3 summarizes the classification accuracy results. As we can see, the 1-Layer SRCNN im-

Table 1. CNN architectural configuration used for BCNN, TCNN and SRCNN. Every convolutional layer is followed by a leaky ReLU [23] with 10% leakage value.

Dataset	MNIST Or CIFAR-10/100	Caltech-101 Or ImageNet-2012
Architectures	BCNN / TCNN / SRCNN	BCNN / TCNN / SRCNN
Input Size	28x28x1 Or 32x32x3	224x224x3
SubBand	- / 1-Layer / 1-Layer	- / 1-Layer / 1-Layer
CONV+ReLU	3x3x1x64 / 3x3x4x64 / 3x3x1x16x4	3x3x3x64 / 3x3x12x64 / 3x3x3x16x4
CONV+ReLU	3x3x64x128 / 3x3x64x128 / 3x3x16x32x4	3x3x64x64 / 3x3x64x64 / 3x3x16x16x4
CONV+ReLU	3x3x128x256 / 3x3x128x256 / 3x3x32x64x4	3x3x64x64 / 3x3x64x64 / 3x3x16x16x4
CONV+ReLU	-	3x3x64x64 / 3x3x64x64 / 3x3x16x16x4
CONV+ReLU	-	3x3x64x64 / 3x3x64x64 / 3x3x16x16x4
POOL	2-by-2	2-by-2
CONV+ReLU	3x3x256x512 / 3x3x256x512 / 3x3x64x128x4	3x3x64x128 / 3x3x64x128 / 3x3x16x32x4
CONV+ReLU	3x3x512x128 / 3x3x512x128 / 3x3x128x32x4	3x3x64x128 / 3x3x64x128 / 3x3x16x32x4
CONV+ReLU	-	3x3x64x128 / 3x3x64x128 / 3x3x16x128x4
CONV+ReLU	-	3x3x64x128 / 3x3x64x128 / 3x3x16x32x4
CONV+ReLU	-	3x3x64x128 / 3x3x64x128 / 3x3x16x32x4
POOL	2-by-2	2-by-2
CONV+ReLU	-	3x3x64x128 / 3x3x64x128 / 3x3x16x32x4
CONV+ReLU	-	3x3x128x128
POOL	-	2-by-2
FC-1	4x4x128x4096	4x4x128x4096
DROPOUT [2]	50%	50%
FC-2	4096x1024	4096x1024(C.Tech) Or 4096x4096(Im.Net)
DROPOUT [2]	50%	50%
FC-3	1024x10 / 1024x100	4096x102 / 4096x1000
SOFTMAX	1x10 / 1x100	1x102 / 1x1000

Table 2. Comparison of total MAC operations, parameters used and classification accuracy of 1&2-layer DWT SRCNN architecture with other well established CNN models for the ImageNet-2012 dataset.

Models	MACs	Param. (Million)	Param. (MByte)	Accuracy (Top-1)	Accuracy (Top-5)	Delta Top (5 - 1)
MobileNet V1	569 M	4.24	2	70.9	89.9	19
MobileNet V2	300 M	3.47	1.7	71.8	91	19.2
Google Net	741 M	6.99	3.3	-	92.1	-
AlexNet	724 M	60.95	29.1	62.5	83	20.5
SqueezeNet	451 M	1.24	0.6	57.5	80.3	22.8
ResNet-50	3.9 B	25.6	12.2	75.2	93	17.8
VGG	15.5 B	138	65.8	70.5	91.2	20.7
Inception-V1	1.43 B	7	3.3	69.8	89.3	19.5
SRCNN (1L)	169.5 M	42.05	20.1	65.6	82.17	16.57
SRCNN (2L)	46.34 M	13.64	6.5	-	-	-

proves the state-of-the-art performance for MNIST, CIFAR-10 and CIFAR-100 datasets by a fair margin. Replacing the 1-layer DWT with a 2-layer DWT decomposition, i.e., with a 2-layer subband decomposition or 16 subbands, we achieve an accuracy of 84.37% for TCNN and 88.93% for SRCNN, on the Caltech-101 dataset. With a 1-layer subband decomposition, SRCNN trained on ImageNet-2012, we achieve top-5 and top-1 validation set accuracy [2] of 82.17% and 65.6%, respectively. Table 2 indicates that our proposed architecture achieves accuracy which is competitive with other state-of-the-art CNN networks that are heavily optimized.

Computational Cost: Table 2 compares the total number of multiply-and-accumulate (MAC) operations and parameters used by state-of-the-art CNNs. Both the 1-layer and 2-layer subband decomposed SRCNNs perform best in class in terms of number of total MAC operations needed. On the number of parameters front, the SRCNN architecture performs fairly.

Table 3. Classification accuracy with 1-layer DWT architecture with parameters indicated in Table 1.

Dataset	BCNN	TCNN	SRCNN	State-of-art
MNIST	99.72	99.76	99.83	99.79 [25]
CIFAR-10	95.37	96.59	96.71	96.53 [26]
CIFAR-100	80.72	81.74	82.97	81.70 [27]
CALTECH-101	82.17	83.89	86.93	89.47(ZF-5) [6]

Table 4. Classification accuracy with input quantization for 1-layer DWT architecture.

Datasets	Models	1-bit	2-bits	4-bits	6-bits	8-bits
MNIST	BCNN	97.42	97.53	97.6	97.69	99.72
	TCNN	97.87	98.03	99.45	99.53	99.76
	SRCNN	99.5	99.6	99.63	99.64	99.83
CIFAR-10	BCNN	76.36	76.06	76.39	85.66	95.37
	TCNN	76.66	76.38	79.2	89.56	96.59
	SRCNN	78.13	79.2	80.03	91.32	96.71
CIFAR-100	BCNN	54.89	59.66	68.8	73.85	80.72
	TCNN	58.91	61.03	69.87	74.75	81.74
	SRCNN	60.14	64.79	69.15	74.07	82.97
CALTECH-101	BCNN	69.87	71.96	76.2	79.31	82.17
	TCNN	71.66	76.29	80.91	81.47	83.89
	SRCNN	71.17	74.94	77.61	83.16	86.93

Table 5. Classification accuracy with weight-and-bias quantization for 1-Layer DWT architecture.

Datasets	MNIST			CIFAR-10		
Models	BCNN	TCNN	SRCNN	BCNN	TCNN	SRCNN
8-bits	90.3	90.91	91.65	60.85	61.37	63.54
16-bits	97.53	97.9	99.65	76.13	79.46	79.84
32-bits	99.72	99.76	99.83	95.37	96.59	96.71
Datasets	CIFAR-100			Caltech-101		
Models	BCNN	TCNN	SRCNN	BCNN	TCNN	SRCNN
8-bits	56.17	61.13	63.15	71.13	75.4	82.35
16-bits	53.78	61.13	61.37	79.67	80.01	83.16
32-bits	80.72	81.74	82.97	82.17	83.89	86.93

However, the parameters of all compared CNNs fall between 0.6 to 65.8 MBytes, with the 1-layer and 2-layer SRCNN architectures at 20.1 and 6.5 MBytes, respectively. In practice, the difference between 6.5 and 0.6 MBytes can be ignored, where as a 10 \times reduction in total MAC operations can significantly improve computation time.

Quantization Effects: The effect of input-data quantization and weights-and-biases quantization on classification accuracy is listed in Table 4 and Table 5, respectively, for BCNN, TCNN and SRCNN architectures. MNIST, CIFAR-10/100 and Caltech-101 datasets are native 8 bits per color. We quantize the input image to 1, 2, 4, 6 and 8 bits per color, and the weights and biases to 8, 16 and 32 bits IEEE floating point values. As we can see the SRCNN architecture is more robust than BCNN and TCNN.

5. CONCLUSION

The proposed SRCNN architecture achieves state-of-the-art performance with computation cost less than 10% of an equivalent CNN. Our method owes its performance to structural regularization – the input signal is losslessly decomposed into subbands and the subband CNNs are restrained from learning features in the other subbands, thereby reducing the risk of overfitting. In addition to computational benefits, the distribution of information across different subbands may vary greatly from class to class. As a result, the FC layers of SRCNN have more information compared to FC

layers of a full-band CNN to separate the output space. Further, noise and deformities are isolated to each subband and do not corrupt the rest, making classification robust compared to analyzing the entire spatial representation by a single CNN. Our architecture is also robust to input data quantization and weight-bias quantization error, which is critical in real life CNN applications where quantization is inevitable.

6. REFERENCES

- [1] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," in *Proceedings of the IEEE*, 1998, pp. 2278–2324.
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in NIPS*, pp. 1097–1105. 2012.
- [3] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.
- [4] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. E. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," *CoRR*, vol. abs/1409.4842, 2014.
- [5] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CoRR*, vol. abs/1512.03385, 2015.
- [6] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [7] F. N. Iandola, M. W. Moskewicz, K. Ashraf, S. Han, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <1mb model size," *CoRR*, vol. abs/1602.07360, 2016.
- [8] H. Li, A. Kadav, I. Durdanovic, H. Samet, and H. P. Graf, "Pruning filters for efficient convnets," *CoRR*, vol. abs/1608.08710, 2016.
- [9] S. Changpinyo, M. Sandler, and A. Zhmoginov, "The power of sparsity in convolutional neural networks," *CoRR*, vol. abs/1702.06257, 2017.
- [10] V. Sze, Y. H. Chen, T. J. Yang, and J. S. Emer, "Efficient processing of deep neural networks: A tutorial and survey," *CoRR*, vol. abs/1703.09039, 2017.
- [11] T. A. Poggio, H. Mhaskar, L. Rosasco, B. Miranda, and Q. Liao, "Why and when can deep - but not shallow - networks avoid the curse of dimensionality: a review," *CoRR*, vol. abs/1611.00740, 2016.
- [12] K. Crammer, A. Kulesza, and M. Dredze, "Adaptive regularization of weight vectors," in *NIPS*, pp. 414–422. Curran Associates, Inc., 2009.
- [13] X. Sun, "Structure regularization for structured prediction: Theories and experiments," *CoRR*, vol. abs/1411.6243, 2014.
- [14] A. K. Tanc and E. M. Eksioğlu, "Transform learning mri with global wavelet regularization," in *2015 23rd European Signal Processing Conference (EUSIPCO)*, Aug 2015, pp. 1855–1859.
- [15] V. N. Ekambaram, G. Fanti, B. Ayazifar, and K. Ramchandran, "Wavelet-regularized graph semi-supervised learning," in *2013 IEEE Global Conference on Signal and Information Processing*, Dec 2013, pp. 423–426.
- [16] A. Tong, D. V. Dijk, J. S. Stanley, III, M. Amodio, G. Wolf, and S. Krishnaswamy, "Graph Spectral Regularization for Neural Network Interpretability," *ArXiv e-prints*, Sept. 2018.
- [17] I. Daubechies, "Orthonormal bases of compactly supported wavelets," *Communications on Pure and Applied Mathematics*, vol. 41, no. 7, pp. 909–996, Oct 1988.
- [18] S. C. B. Lo, H. Li, J. S. Lin, A. Hasegawa, C. Y. Wu, M. T. Freedman, and S. K. Mun, "Artificial convolution neural network with wavelet kernels for disease pattern recognition," in *Proc. SPIE*, 1995, vol. 2434, pp. 2434 – 2434 – 10.
- [19] E. Kang, J. Min, and J. C. Ye, "Wavenet: a deep convolutional neural network using directional wavelets for low-dose x-ray CT reconstruction," *CoRR*, vol. abs/1610.09736, 2016.
- [20] S. Fujieda, K. Takayama, and T. Hachisuka, "Wavelet convolutional neural networks," *CoRR*, vol. abs/1805.08620, 2018.
- [21] E. Oyallon, E. Belilovsky, and S. Zagoruyko, "Scaling the scattering transform: Deep hybrid networks," *CoRR*, vol. abs/1703.08961, 2017.
- [22] I. Daubechies, "The wavelet transform, time-frequency localization and signal analysis," *IEEE Trans. on Info. Theory*, vol. 36, no. 5, pp. 961–1005, Sept 1990.
- [23] B. Xu, N. Wang, T. Chen, and M. Li, "Empirical evaluation of rectified activations in convolutional network," *CoRR*, vol. abs/1505.00853, 2015.
- [24] K. He and J. Sun, "Convolutional neural networks at constrained time cost," in *CVPR*, 2015, pp. 5353–5360.
- [25] L. Wan, M. Zeiler, S. Zhang, Y. LeCun, and R. Fergus, "Regularization of neural networks using dropconnect," in *Intl. Conf. on Machine Learning*, Atlanta, Georgia, USA, 17-19 Jun 2013, pp. 1058–1066, PMLR.
- [26] B. Graham, "Fractional max-pooling," *CoRR*, vol. abs/1412.6071, 2014.
- [27] S. Zagoruyko and N. Komodakis, "Wide residual networks," *CoRR*, vol. abs/1605.07146, 2016.