

MULTI-RESOLUTION DUAL-TREE WAVELET SCATTERING NETWORK FOR SIGNAL CLASSIFICATION

Amarjot Singh and Nick Kingsbury

Signal Processing Group, Department of Engineering, University of Cambridge, U.K.

ABSTRACT

This paper introduces a Deep Scattering network that utilizes Dual-Tree complex wavelets to extract translation invariant representations from an input signal. The computationally efficient Dual-Tree wavelets decompose the input signal into densely spaced representations over scales. Translation invariance is introduced in the representations by applying a non-linearity over a region followed by averaging. The discriminatory information in the densely spaced, locally smooth, signal representations aids the learning of the classifier. The proposed network is shown to outperform Mallat's ScatterNet [1] on four datasets with different modalities on classification accuracy.

Index Terms— DTCWT, Scattering network, Convolutional neural network, USPS dataset, UCI datasets.

1. INTRODUCTION

Signal classification is a difficult problem due to the considerable translation, rotation and scale variations that can hinder the classifier's ability to measure signal similarity [2]. Deep Convolutional Neural Networks (CNNs) [3] have been widely used to eliminate the above-mentioned variabilities and learn invariant as well as discriminative signal representations by using successive kernel operations (linear filters, pooling, and non-linearity). Despite their success, the optimal configuration of these networks is not well understood because of the cascaded nonlinearities.

Scattering convolution network proposed by S. Mallat in [1] provided a mathematical framework to incorporate geometric signal priors to extract discriminative and invariant signal representations. Invariance is introduced in the representations by filtering the input signal with a cascade of multi-scale and multidirectional complex Morlet wavelets followed by pointwise nonlinear modulus and local averaging. The high frequencies lost due to averaging are recovered at the later layers using cascaded wavelet transformations with nonlinearities, justifying the need for a multilayer network. The above class of networks has been widely used in numerous applications as such as object classification [4], audio processing [5], scene understanding [6], biology [7] etc.

This paper proposes an improved Deep Scattering architecture that uses Dual-Tree Complex Wavelet Transform

(DTCWT) [8] to decompose a *multi-resolution input signal* into translation invariant signal representations. The input signal is first decomposed into multi-resolution signal representations that are densely spaced on the scale domain. Translation invariance is then introduced within each representation by applying a non-linearity over a region followed by local averaging. Next, a log non-linearity is used to separate the multiplicative low-frequency illumination components within the representations. Finally, a Support Vector Machine (SVM) is used to create discrimination between different signal classes by learning weights that best summarize the regularities (common coefficients) in the training data and simultaneously ignore the coefficients arising due to the irregularities [9].

The main contributions of the paper and their reasoning are explained below:

- *Multi-scale Input Signal*: The input signal is decomposed into representations that are densely spaced on the scale space using a DTCWT based decimated pyramid of complex values [10]. The multi-scale representations have redundant local regularities that allows the SVM to optimally learn weights that learn discriminatory features (edges between two objects within an image) from fine scale representations while non-discriminator features like the middle of the objects from the coarse features [11].
- *DTCWT Filter Bank*: The proposed network uses DTCWT bank for filtering as opposed to Morlet wavelets [1] due to its perfect reconstruction properties [8]. Perfect reconstruction property allows the DTCWT filter to extract features without any aliasing.
- *Region Non-Linearity*: The extracted representations are cascaded by a *region non-linearity* as opposed to a point non-linearity and then followed by local averaging to produce a regional translation invariant representation. The region non-linearity selects the dominant feature within the region while simultaneously suppressing features with lower magnitudes leading to invariance similar to the max operator in CNNs.

The performance of the proposed network is tested on four popular datasets selected from different modalities. Such

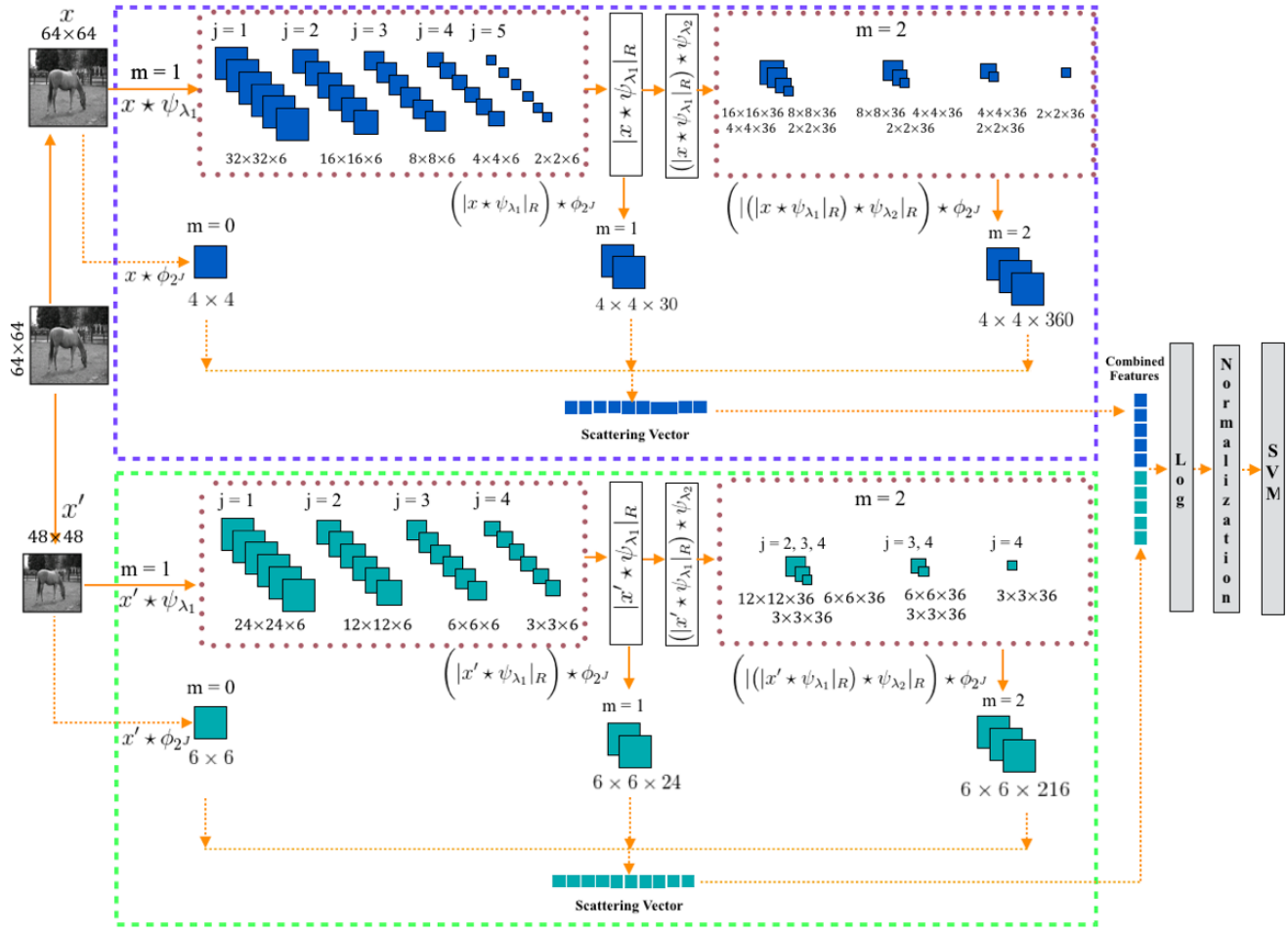


Fig. 1. Illustration shows the input image (x) of size 64×64 resized to images of resolution, x (64×64) and x' (48×48) respectively. Image representations at $m = 1$ are obtained using DTCWT filters at different scales and 6 orientations. Then, L2 region non-linearity is applied on the representations to obtain the regular envelope followed by local averaging to extract the translation invariant coefficients. The information lost due to averaging are recovered by cascaded wavelet filtering at the second layer. Translation invariance is introduced in the recovered frequencies using L2 region non-linearity and local averaging.

diversity of data sets is crucial to verify the generalization of the proposed network to a large variety of problems.

The paper is divided into the following sections. Section 2 describes the proposed Deep Multi-Resolution DTCWT Scattering Architecture while Section 3 presents the step by step experimental results leading to signal classification. Section 4 draws conclusions about the experimental results.

2. MULTI-RESOLUTION DTCWT SCATTERNET

In this section, the mathematical formulation of the proposed Multi-Resolution DTCWT Scattering Network is presented for 2 layers that decompose a two-dimensional input signal x into translation invariant representations that are densely spaced over scales. However, the formulation can be easily extended to the input signal of other dimensions and network with deeper layers.

First, the multi-resolution representation of the signal x is obtained at interleaved scale (s) [10] using a DTCWT based decimated pyramid of complex values as shown below:

$$[x, x'] = \text{decimated}(x) \quad (1)$$

The proposed scattering network is applied at each resolution (x, x') to produce a translation invariant representation at multiple layers. For the sake of simplicity, the derivation is presented only for the input signal x .

The proposed architecture is realized by arranging DTCWT filters with numerous scales (and orientations (for only 2D signal)) in multiple layers to extract stable and informative signal representations. DTCWT is an advantageous filtering choice due to its perfect reconstruction properties. It allows the extraction of features without any aliasing. DTCWT wavelets are represented by ψ (with a real ψ^a and imaginary ψ^b group). The complex band-pass filter ψ is decomposed

into real and imaginary groups as shown:

$$\psi(t) = \psi^a(t) + i\psi^b(t) \quad t = (t_1, t_2) \quad (2)$$

The signal x is filtered at the first layer (L_1) using a family of DTCWT wavelets $\psi_{\lambda_1}(t)$ at different scales and orientations (λ_1), formulated as:

$$x \star \psi_{\lambda}(t) = x \star \psi_{\lambda}^a(t) + ix \star \psi_{\lambda}^b(t) \quad (3)$$

The wavelet transform response commutes with translations, and is therefore not translation invariant. To build a translation invariant representation, a L_2 smooth non-linearity is first applied over all overlapping regions of size R ($R \times R$ for 2D signal) in feature output, obtained at a particular scale (and six orientation (θ) (for 2D signal)). The non-linearity applied to one of the above-mentioned regions is shown below:

$$|x \star \psi_{\lambda_1}|_R = \sqrt{|G_{real} \star \psi_{\lambda_1}^a(t)|^2 + |G_{imag} \star \psi_{\lambda_1}^b(t)|^2} \quad (4)$$

where R is the size of the region and G is a group of R ($R \times R$ for 2D signal) complex scattering coefficients. L_2 is a non-expansive non-linearity that makes it stable to additive noise and deformations [1]. The region non-linearity selects the dominant feature in the region while simultaneously suppressing features with lower magnitudes. This creates translation invariance in a larger region similar to the max operator in CNNs. The scattering coefficients obtained after applying the region non-linearity to the outputs of every wavelet scales is given by $|x \star \psi_{\lambda_1}|_R$.

Next, the desired translation invariant representation are obtained at the first layer (L_1) by applying a local average on $|x \star \psi_{\lambda_1}|_R$, as shown below:

$$(L_1)_R = \left(|x \star \psi_{\lambda_1}|_R \right) \star \phi_{2^J} \quad (5)$$

The high frequencies coefficients lost by the averaging operator are recovered at the second layer (L_2) by calculating the wavelet coefficients of $|x \star \psi_{\lambda_1}|_R$ by the wavelet at scale and orientation, λ_2 , given as $(|x \star \psi_{\lambda_1}|_R) \star \psi_{\lambda_2}(t)$ [1].

The features extracted at the first layer (L_1) are filtered with the DTCWT filter at coarser scales (λ_2) to recover the high frequency components at the second layer (L_2). The recovered frequencies are converted into translation invariant representations by again taking a local average as shown:

$$(L_2)_R = \left(\left(|x \star \psi_{\lambda_1}|_R \right) \star \psi_{\lambda_2}|_R \right) \star \phi_{2^J} \quad (6)$$

The scattering coefficients $S_J x$ for the network at different scales and orientations for two layers at a path p can be obtained using the following:

$$S_J x[p] = \left(\begin{array}{c} x \star \phi_{2^J} \\ \left(|x \star \psi_{\lambda_1}|_R \right) \star \phi_{2^J} \\ \left(\left(|x \star \psi_{\lambda_1}|_R \right) \star \psi_{\lambda_2}|_R \right) \star \phi_{2^J} \end{array} \right)_{\lambda=(2,3,4)} \quad (7)$$

A logarithm non-linearity proposed by Oyallon et al. [4], is applied to the scattering coefficients in order to transform the low-frequency multiplicative components that arise due to illuminations into additive components. These additive coefficients can now be ignored as noise by the classifier. The logarithm applied to the scattering coefficients (Sx) extracted from a dataset with M training signals, where the coefficients computed from a single signal has N dimensions, is given by:

$$\Phi^{M \times N} = \log(Sx^{M \times N} + k) \quad (8)$$

where Sx are the scattering coefficients and k is (small) constant added to reduce the effect of noise magnification at small signal levels. The value of the constant k ($1e^{-6}$) used in [4] is duplicated for all the experiments in this paper.

3. OVERVIEW OF RESULTS

Experiments are conducted on four real-world datasets selected from the image, audio, biology and material modalities to evaluate the performance of the proposed DTCWT multi-resolution scattering network. In order to test the generalization of the proposed network to different problems, a large mixture of data sets from different domains with various sizes and dimensionality are used for experimentation. Please see Table 1 for a detailed description of the datasets used in our experiments.

The proposed multi-resolution DTCWT scattering network is applied on each dataset to extract translation invariant multi-scale representations that are further used for classification. Scattering coefficients in the case of the two-dimensional signals for all the experiments is computed at six orientations ($15^\circ, 45^\circ, 75^\circ, 105^\circ, 135^\circ, 165^\circ$). The discrimination between the signal classes is achieved using a Gaussian SVM. Before the SVM is trained on the training set, each feature is standardized by the mean and standard deviation of the training dataset.

The test set generalization error of the proposed Deep DTCWT multi-resolution scattering network is reported on each Dataset (Table. 1) and compared with the scattering network proposed by Mallat et al [1]. In addition, this error for the proposed network is also compared with the recently proposed machine learning approaches used for classification of the above-mentioned data sets. Only those approaches are considered that don't augment the data sets and apply their algorithm only on the unaltered input signal.

3.1. US Postal Service Dataset

The US postal service dataset consists of two-dimensional structured grayscale image signals with 7291 training observations and 2007 test observations [12]. This dataset was generated by scanning the handwritten digits from envelopes by the U.S. Postal Service. The recorded images are de-slanted and size normalized to 16 x 16 (256) pixels images in the dataset. The objective is to differentiate between 10 different digits between 0 and 9.

Table 1. Classification error (%) on different datasets for each component of the proposed network. DTCWT ScatNet: DSCAT, DTCWT ScatNet + Pooling: DSCATP, Multi-Resolution DTCWT ScatNet: MDSCAT, Multi-Resolution DTCWT ScatNet + Pooling: MDSCATP. The left result in / is without log non-linearity applied while the right is with log applied (*NoLog/Log*).

Dataset	DSCAT	DSCATP	MDSCAT	MDSCATP
USPS	3.31 / 3.38	3.24 / 3.33	2.89 / 3.11	2.56 / 2.84
Isolet	5.14 / 5.3	5.10 / 5.36	4.75 / 5.02	4.14 / 4.88
Yeast	41.65 / 41.17	45.86 / 45.85	37.04 / 34.62	39.77 / 39.04
Glass	31.78 / 29.16	31.77 / 33.68	27.82 / 24.32	30.05 / 26.06

The proposed scattering network extracts the input signal at 6 resolution (s) (1, 0.85, 0.70, 0.6, 0.5, 0.35) and then extracts scattering coefficients from each resolution at 3 DTCWT scales (J) and 6 orientations (θ). The region non-linearity is applied on overlapping regions (R) of size 2×2 . A cost value (c) of 5 was selected for the linear SVM. All the parameters are selected using 5-fold cross validation. As noted from Table. 1, the proposed network with region non-linearity and without log non-linearity results in the lowest classification error of 2.54%. This increase in error due to the log non-linearity is explained in the previous section. The classification error of the proposed architecture is also compared to ScatterNet [1] and Fuzzy Integral Combination algorithm [3] as presented in Table. 2. The proposed architecture outperforms both the algorithms.

Table 2. Classification error (%) comparison on USPS

Dataset	Proposed	ScatNet [1]	FIC [13]
USPS	2.54	2.6	5.43

3.2. The UCI Isolet Dataset

The Isolet dataset comprises of one-dimensional audio signals collected from 150 speakers uttering all characters in the English alphabet twice. Each speaker contributed 52 training examples with a total of 7797 recordings [14]. The recordings are represented with 617 attributes such as spectral coefficients, contour, sonorant and post-sonorant are provided to classify letter utterance.

The proposed scattering network decomposes the input signal at 4 resolution (s) (1, 0.70, 0.5, 0.35). The translation invariant features are extracted at 6 scattering DTCWT wavelet scales (J) for the input signal at every resolution. Regions (R) of size 1×4 is chosen for the application of the region non-linearity. A cost value (c) of 15 was chosen for the linear SVM. Again, the parameters are selected using 5-fold cross validation. The generalization error is reported on 10-fold cross validation for this dataset. Table. 1 shows that the multi-resolution scattering architecture with region non-linearity and log non-linearity produces the lowest classification error of 4.14%. The proposed method outperformed ScatterNet [1] but was unable to surpass the performance of Extreme entropy machines [15] as shown in Table 3.

Table 3. Classification error (%) comparison on Isolet

Dataset	Proposed	ScatNet [1]	EEM [15]
Isolet	4.14	5.78	2.70

3.3. The UCI Yeast Dataset

This is a highly imbalanced one-dimensional signal dataset that consists of 1484 yeast proteins with 10 cellular binding sites [14]. Each binding site is described with 8 attributes. The aim is to classify the most probable cellular localization site of the proteins.

The proposed scattering network decomposes the input signal at 2 resolution (s) (1, 0.70). The translation invariant features are extracted at 2 scattering DTCWT wavelet scales (J) for the input signal at every resolution. The Region (R) size of 1×2 and a cost value (c) of 15 is chosen using 5-fold cross validation. The generalization error was reported on 10-fold cross validation for this dataset. Table. 1 shows that the multi-resolution scattering architecture with region non-linearity and log non-linearity produces the lowest classification error of 35.02%. The proposed method outperformed ScatterNet [1] but was unable to outrank the instance selection genetic algorithm [16] as shown in Table 4.

Table 4. Classification error (%) comparison on Yeast

Dataset	Proposed	ScatNet [1]	IS [16]
Yeast	35.02	35.89	33.0

3.4. The UCI Glass Dataset

This dataset consists of 214 one-dimensional signals that describe six types of glass based on 9 chemical fractions of the oxide content [14]. This dataset was motivated by a criminological investigation where the correct classification of glass left on the crime scene could be used for evidence. Hence, the aim is to classify between different types of glass.

The proposed scattering network uses the same parameters as mentioned in Section. 3.3 for feature extraction. The generalization error was reported on 10-fold cross validation for this dataset. Table. 1 shows that the multi-resolution scattering architecture with region non-linearity and log non-linearity produces the lowest classification error of 24.32%.

The proposed method outperformed ScatterNet [1] and Kernelized Vector Quantization [17] as shown in Table 5.

Table 5. Classification error (%) comparison on Glass

Dataset	Proposed	ScatNet [1]	KVQ [17]
Glass	24.32	28.86	31.6

4. CONCLUSION

The paper proposes a ScatterNet that extracts regionally translation invariant features from an input signal that are equally spaced over the scale space. The proposed algorithm was tested on four datasets. It outperformed Mallat’s ScatterNet on all the datasets while was able to outperform the learning based algorithms only on two datasets. Hence, it is necessary to take learning into account. The proposed scattering network can then provide the first two layers of such learning networks. It eliminates translation variability, which can help in learning the next layers. In addition, this network can replace simpler low-level features such as SIFT vectors.

5. REFERENCES

- [1] J. Bruna and S. Mallat, “Invariant scattering convolution networks,” *IEEE Transaction on Pattern Analysis and Machine Intelligence*, vol. 35, pp. 1872–1886, 2013.
- [2] B. Scholkopf and A.J. Smola, “Learning with kernels,” *MIT Press*, 2002.
- [3] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [4] E. Oyallon and S. Mallat, “Deep roto-translation scattering for object classification,” *IEEE CVPR*, pp. 2865–2873, 2015.
- [5] J. Andn, V. Lostanlen, and S. Mallat, “Joint time-frequency scattering for audio classification,” *Proceedings of IEEE MLSP Workshop*, 2015.
- [6] S. Nadella, A. Singh, and SN Omkar, “Aerial scene understanding using deep wavelet scattering network and conditional random field,” in: *Hua G., Jgou H. (eds) Computer Vision — European Conference on Computer Vision (ECCV) Workshops*, vol. 9913, pp. 205–214, 2016.
- [7] V. Miliš and S. Mallat, “Mathematical modeling of lymphocytes selection in the germinal center,” *Journal of Mathematical Biology*, 2016.
- [8] N.G. Kingsbury, “Complex wavelets for shift invariant analysis and filtering of signals,” *Applied and computational harmonic analysis*, vol. 10, pp. 234–253, 2001.
- [9] T. Joachims, “Optimizing search engines using click-through data,” *8th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2002.
- [10] R. Anderson, N.G. Kingsbury, and J. Fauqueur., “Determining multi-scale image feature angles from complex wavelet phases,” *In Proceedings of the Second ICIAR*, pp. 490–498, 2005.
- [11] J.C. Chan, H. Ma, and T.K. Saha, “Partial discharge pattern recognition using multiscale feature extraction and support vector machine,” *2013 IEEE Power and Energy Society General Meeting*, 2013.
- [12] J.J. Hull, “A database for handwritten text recognition research,” *IEEE Transaction on Pattern Analysis and Machine Intelligence*, vol. 16, no. 5, pp. 550–554, 1994.
- [13] B. Hadjadji, Y. Chibani, and H. Nemmour, “Fuzzy integral combination of one-class classifiers designed for multi-class classification,” *Image Analysis and Recognition, LNCS*, vol. 8814, pp. 320–328, 2014.
- [14] D. Newman, S. Hettich, C. Blake, and C. Merz, “UCI repository of machine learning databases,” <http://www.ics.uci.edu/mllearn/MLRepository.html>.
- [15] W.M. Czarnecki and J. Tabor, “Extreme entropy machines: robust information theoretic classification,” *Pattern Analysis and Applications*, pp. 1–18, 2015.
- [16] Z.Y. Chen et al., “Instance selection by genetic-based biological algorithm,” *Soft Computing*, vol. 19, no. 8, pp. 1–18, 2015.
- [17] T. Villmann, S. Haase, and M. Kaden, “Kernelized vector quantization in gradient-descent learning,” *Neurocomputing*, vol. 147, pp. 8395, 2015.