

# Geometric Neural Phrase Pooling: Modeling the Spatial Co-occurrence of Neurons

Lingxi Xie<sup>1</sup>, Qi Tian<sup>2</sup>, John Flynn<sup>3</sup>, Jingdong Wang<sup>4</sup> and Alan Yuille<sup>1</sup>

<sup>1</sup>The Johns Hopkins University

<sup>2</sup>University of Texas at San Antonio

<sup>3</sup>University of California, Los Angeles

<sup>4</sup>Microsoft Research

## ABSTRACT

Deep Convolutional Neural Networks (CNNs) are playing important roles in state-of-the-art visual recognition. This paper focuses on modeling the spatial co-occurrence of neuron responses, which is less studied in the previous work. For this, we consider the neurons in the hidden layer as neural words, and construct a set of geometric neural phrases on top of them. The idea that grouping neural words into neural phrases is borrowed from the Bag-of-Visual-Words (BoVW) model. Next, the Geometric Neural Phrase Pooling (GNPP) algorithm is proposed to efficiently encode these neural phrases. GNPP acts as a new type of hidden layer, which penalizes the isolated neuron responses after convolution, and can be inserted into a CNN model with little extra computational overhead. Experimental results show that GNPP produces significant and consistent accuracy gain in image classification.

## CONTRIBUTION

In this paper, we present **Geometric Neural Phrase Pooling (GNPP)**, an efficient yet effective algorithm to help CNN training.

GNPP is motivated by one of our previous work, namely Geometric Phrase Pooling (GPP) in the Bag-of-Visual-Words (BoVW) model. GPP works by considering local feature co-occurrence and performing non-linear smoothing. We find that these operations also help network training.

GNPP averages the response of a neuron with the maximal response among its neighboring neurons. This is to penalize the isolated neural responses. We argue that, especially in the high-level layers, isolated responses often correspond to unexpected noise in convolution, because neighboring neurons often share a large part of their receptive fields. Experiments reveal that GNPP often works better on high-level layers, which verifies our motivation.

GNPP can be considered as an intermediate layer in a CNN structure. Since GNPP does not change the geometric shape of data, it can be inserted, at least in theory, to any position in the neural network. In practice, we only insert it between a convolutional layer and a pooling layer because of our motivation.

GNPP produces **consistent accuracy gain** on several state-of-the-art deep networks for object recognition. Moreover, the extra time (+1.29%) and memory (+2.52%) costs of GNPP are almost negligible. Recently, we also observe that GNPP works well in other CNN-based models such as Faster RCNN or DeepLab.

## THE PROPOSED ALGORITHM

### Neural Word and Geometric Neural Phrase

- Defined on a hidden layer  $\mathbf{X}^{(l)}$  of the CNN. For simplicity, denote  $\mathbf{X}^{(l)}$  as  $\mathbf{X}$ .
  - $\mathbf{X}$  is a 3D cube with  $W \times H \times D$  neurons
- We naturally consider the data as a set of  $D$ -dimensional **neural words**:

$$\mathcal{X} = \{\mathbf{x}_{w,h}\}_{w=1,h=1}^{W,H}$$

- A **geometric neural phrase** is a group of neighboring neurons:

$$\mathcal{G}_{w,h} = \{\mathbf{x}_{w,h}^{(k)}\}_{k=0}^K$$

- $\mathbf{x}_{w,h}^{(0)} = \mathbf{x}_{w,h}$ : the **central word**
- $\mathbf{x}_{w,h}^{(k)}$  ( $k > 0$ ): the **side words**, located in a small neighborhood of  $\mathbf{x}_{w,h}$
- For each neural word, there defines a neural phrase.

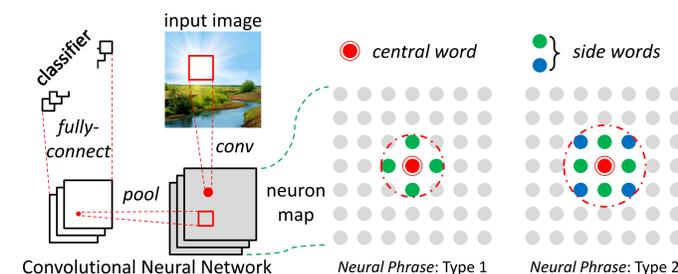
### Geometric Neural Phrase Pooling

- A  $D$ -dimensional vector for each geometric neural phrase individually

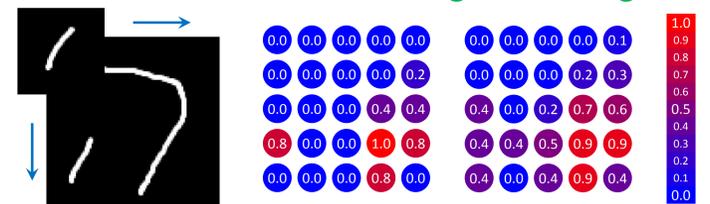
$$\mathbf{z}_{w,h} = \frac{1}{2} \left[ \mathbf{x}_{w,h} + \max_{k>0} \{s_{w,h}^{(k)} \times \mathbf{x}_{w,h}^{(k)}\} \right]$$

- $\max_{k>0} \{ \}$ : dimension-wise maximization
- $s_{w,h}^{(k)} = \sigma$  or  $s_{w,h}^{(k)} = \sigma^2$ , according to the relative position
- $\sigma \in [0,1]$ : the smoothing parameter

### Illustration of Geometric Neural Phrase

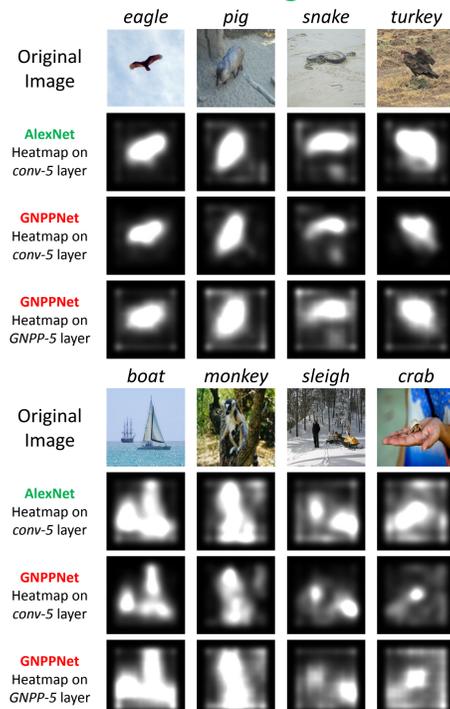


### Geometric Neural Phrase Pooling: Smoothing Effect



GNPP penalizes isolated neural responses, and preserves the clustered responses. In neural networks, especially in high-level layers, the isolated responses often relate to random noise, therefore GNPP works well in these layers.

### Visualization on Neural Responses with or without Adding a GNPP Layer



### LeNet Results on MNIST (UL), SVHN (UR), CIFAR10 (LL) and CIFAR100 (LR)

L1	L2	D	T1 (1.0)	T1 (0.9)	T1 (0.8)	T2 (1.0)	T2 (0.9)	T2 (0.8)
			0.87 ± .02	0.87 ± .02	0.87 ± .02	0.87 ± .02	0.87 ± .02	0.87 ± .02
✓			0.72 ± .04	0.73 ± .03	0.70 ± .05	0.71 ± .06	0.71 ± .06	0.72 ± .04
✓			0.75 ± .03	0.79 ± .02	0.77 ± .05	0.73 ± .04	0.75 ± .04	0.73 ± .05
✓	✓		0.72 ± .03	0.67 ± .04	0.69 ± .04	0.63 ± .03	0.64 ± .03	0.67 ± .03
	✓		0.72 ± .03	0.72 ± .03	0.72 ± .03	0.72 ± .03	0.72 ± .03	0.72 ± .03
✓	✓		0.59 ± .02	0.61 ± .05	0.62 ± .03	0.59 ± .03	0.59 ± .02	0.63 ± .03
✓	✓		0.63 ± .03	0.62 ± .07	0.64 ± .03	0.62 ± .05	0.60 ± .03	0.65 ± .03
✓	✓	✓	0.58 ± .05	0.55 ± .05	0.57 ± .02	0.54 ± .05	0.56 ± .04	0.61 ± .05

L1	L2	L3	T1 (1.0)	T1 (0.9)	T1 (0.8)	T2 (1.0)	T2 (0.9)	T2 (0.8)
			4.63 ± .06	4.63 ± .06	4.63 ± .06	4.63 ± .06	4.63 ± .06	4.63 ± .06
✓			4.46 ± .06	4.47 ± .05	4.42 ± .09	4.42 ± .08	4.42 ± .07	4.43 ± .09
✓			4.15 ± .08	4.18 ± .01	4.17 ± .07	4.08 ± .10	4.19 ± .07	4.20 ± .05
✓	✓		3.76 ± .03	3.72 ± .05	3.77 ± .06	3.53 ± .07	3.64 ± .07	3.65 ± .10
	✓		4.10 ± .05	4.07 ± .03	4.10 ± .05	4.10 ± .07	4.10 ± .03	4.14 ± .07
✓	✓		3.55 ± .10	3.60 ± .03	3.67 ± .06	3.47 ± .05	3.47 ± .02	3.55 ± .09
✓	✓		3.43 ± .06	3.52 ± .07	3.55 ± .04	3.41 ± .03	3.42 ± .04	3.51 ± .05
✓	✓	✓	3.46 ± .07	3.47 ± .06	3.55 ± .06	3.43 ± .05	3.39 ± .01	3.46 ± .03

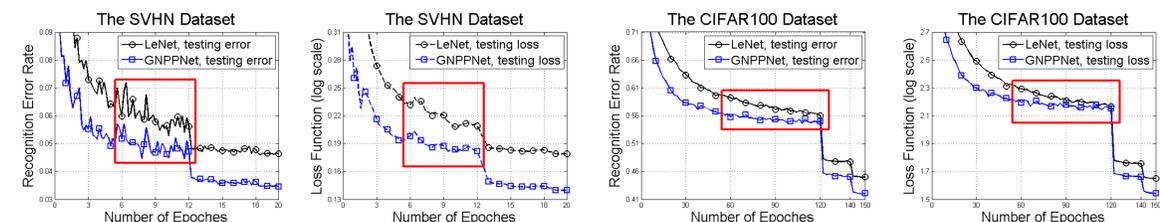
  

L1	L2	L3	T1 (1.0)	T1 (0.9)	T1 (0.8)	T2 (1.0)	T2 (0.9)	T2 (0.8)
			17.07 ± .15	17.07 ± .15	17.07 ± .15	17.07 ± .15	17.07 ± .15	17.07 ± .15
✓			16.67 ± .22	16.80 ± .25	16.84 ± .12	16.65 ± .19	17.03 ± .15	17.04 ± .17
✓			15.79 ± .22	16.09 ± .17	15.95 ± .31	15.69 ± .11	16.07 ± .27	15.90 ± .09
✓	✓		15.49 ± .15	15.31 ± .20	15.51 ± .25	15.27 ± .10	15.29 ± .14	15.28 ± .16
	✓		15.82 ± .23	15.76 ± .18	15.98 ± .14	16.05 ± .29	15.90 ± .25	15.94 ± .09
✓	✓		15.15 ± .20	15.29 ± .12	15.44 ± .19	15.29 ± .32	15.19 ± .35	15.20 ± .35
✓	✓		14.92 ± .18	15.00 ± .18	15.15 ± .15	14.83 ± .25	14.93 ± .20	14.92 ± .16
✓	✓	✓	14.97 ± .17	14.83 ± .23	14.78 ± .17	15.22 ± .16	14.79 ± .26	14.85 ± .26

L1	L2	L3	T1 (1.0)	T1 (0.9)	T1 (0.8)	T2 (1.0)	T2 (0.9)	T2 (0.8)
			44.99 ± .19	44.99 ± .19	44.99 ± .19	44.99 ± .19	44.99 ± .19	44.99 ± .19
✓			44.62 ± .17	44.53 ± .45	44.78 ± .06	44.43 ± .29	44.58 ± .36	44.58 ± .52
✓			43.34 ± .23	43.71 ± .19	43.37 ± .26	43.21 ± .23	43.03 ± .25	43.37 ± .30
✓	✓		43.11 ± .24	42.77 ± .37	42.99 ± .24	42.96 ± .32	42.81 ± .38	43.08 ± .39
	✓		43.99 ± .07	43.63 ± .11	43.50 ± .26	43.38 ± .37	43.34 ± .27	43.46 ± .25
✓	✓		42.85 ± .38	42.81 ± .27	42.82 ± .29	43.08 ± .27	42.79 ± .34	42.93 ± .22
✓	✓		42.35 ± .30	42.34 ± .31	42.04 ± .20	42.92 ± .33	42.72 ± .25	42.54 ± .29
✓	✓	✓	42.97 ± .29	42.77 ± .36	42.36 ± .18	43.31 ± .34	42.85 ± .18	42.60 ± .36

### Testing Error and Loss Curves on SVHN (left) and CIFAR100 (right)



## RESULTS

### Results on some *small* datasets

- We use BigNet [49] and Wide ResNet [50] as baselines. GNPP is inserted before the last pooling layer of each network.

	MNIST	SVHN	CIF10	CIF100
Zeiler [42]	0.45	—	2.80	—
Goodfellow [5]	0.47	—	2.47	—
Lin [21]	0.47	—	2.35	—
Wan [38]	0.52	0.21	—	1.94
Lee [21]	0.39	—	1.92	—
Liang [20]	0.31	—	1.77	—
BigNet, without GNPP	0.86	0.48	3.93	3.48
BigNet, with GNPP	0.68	0.43	3.65	3.25
WRN, without GNPP	0.66	0.45	3.69	3.27
WRN, with GNPP	0.63	0.41	3.61	3.21

### Results on ImageNet

- We use AlexNet [2] as our baseline. GNPP is inserted before the last pooling layer.

ILSVRC2012	top-1	top-5
AlexNet, w/o GNPP	43.19	19.87
AlexNet, w/ GNPP	42.16	19.24

- GNPP builds latent neural connections. With GNPP, the equivalent # of connections between *conv-4* and *conv-5* increases from 149.5M to 348.9M. An alternative way is to increase the # of convolutional kernels, e.g., using 512 kernels at *conv-5* increases the number to 299.0M.

ILSVRC2012	top-1	top-5	Time	Memory
AlexNet, w/ more kernels	42.45	19.47	+9.97%	+5.58%
AlexNet, w/ GNPP	42.16	19.24	+1.29%	+2.52%

GNPP is more effective and more efficient.

## REFERENCES

Key references are numbered as they appear in the paper.

[1] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-scale Hierarchical Image Database. CVPR, 2009.

[2] A. Krizhevsky, I. Sutskever, and G. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. NIPS, 2012.

[26] L. Xie, Q. Tian, M. Wang and B. Zhang. Spatial Pooling of Heterogeneous Features for Image Classification. TIP, 2014.

[49] Nagadomi. The Kaggle CIFAR10 Network. <https://github.com/nagadomi/kaggle-cifar10-torch7/>. 2014.

[50] S. Zagoruyko, N. Komodakis. Wide Residual Networks. arXiv preprint, arXiv: 1605.07146, 2016.

## ACKNOWLEDGEMENTS

This paper is supported by iARPA MICRONS contract D16PC00007, ONR N00014-12-1-0883, ARO grants W911NF-15-1-0290, Faculty Research Gift Awards by NEC Labs of America and Blippar, and NSFC 61429201. We thank Junhua Mao, Cihang Xie and Zhuotun Zhu for discussion.