

Appendix

Anonymous ICCV submission

Paper ID 5528

1. Details of the Proving for the relations

We give the details of the relationship between other nonlocal operators in the spectral view discussed in our paper (Sec.3.2). In the following proving, we assume that $X \in \mathbb{R}^{N \times C}$, $Z = g(X) = XW_Z$, $M_{ij} = f(X_i, X_j)$. All the normalized term uses the inverse of the degree $1/d_i$ where $d_i = \sum_j f(X_i, X_j)$. We also merge the output of the operators with the weight kernel $W \in \mathbb{R}^{N \times C}$ and defines it as O for consistency. Thus the target formulations in this section are a bit different with the definition in their own papers.

1.1. Nonlocal Block

The Nonlocal (NL) Block in the spectral view is the same as defining the graph $\mathcal{G} = (\mathbb{V}, D^{-1}M, Z)$ and then using the second term of the Chebyshev Polynomial to approximate the graph filter.

Proof. The NL operator defined in [9] can be formulated as:

$$O_{i,:} = \frac{\sum_j [f(X_{i,:}, X_{j,:})g(X_{j,:})]}{\sum_j f(X_{i,:}, X_{j,:})} W \quad (1)$$

To unify it by our spectral view, we firstly define the graph $\mathcal{G} = (\mathbb{V}, A, Z)$ to represent the graph structure of the NL operator, where the affinity matrix A is calculated by:

$$A = D_M^{-1}M, \quad M = f(X_{i,:}, X_{j,:}) \quad (2)$$

Thus, each element of the affinity matrix A is:

$$A_{ij} = (D_M^{-1}M)_{ij} = \frac{f(X_{i,:}, X_{j,:})}{\sum_j f(X_{i,:}, X_{j,:})} \quad (3)$$

Based on Theorem.1, when using Chebyshev polynomials to approximate the generalized graph filter Ω and only choosing the second term, it becomes :

$$F(A, Z) = AZW \quad (4)$$

Then taking Eq. (3) into this equation, we can get the formulation of the NL operator:

$$F_{i,:}(A, Z) = \frac{\sum_j [f(X_{i,:}, X_{j,:})g(X_{j,:})]}{\sum_j f(X_{i,:}, X_{j,:})} W \quad (5)$$

□

1.2. Nonlocal Stage

The Nonlocal Stage (NS) in the spectral view is the same as defining the graph $\mathcal{G} = (\mathbb{V}, D_M^{-1}M, Z)$ and then using the 1_{st}-order Chebyshev Polynomial to approximate the graph filter with the condition $W_1 = W_2 = -W$.

Proof. The NS operator given defined in [8] can be formulated as:

$$O_{i,:} = \frac{\sum_j [f(X_{i,:}, X_{j,:})(Z_{j,:} - Z_{i,:})]}{\sum_j f(X_{i,:}, X_{j,:})} W \quad (6)$$

Similar with the proof of NL, we can get each element of the affinity matrix A as:

$$A_{ij} = (D_M^{-1}M)_{ij} = \frac{f(X_{i,:}, X_{j,:})}{\sum_j f(X_{i,:}, X_{j,:})} \quad (7)$$

The graph filter Ω on \mathcal{G} is approximated by the Chebyshev polynomial. When using the 1_{st}-order Chebyshev Approximation, it becomes:

$$F(A, Z) = ZW_1 - AZW_2 \quad (8)$$

When sharing the weight for W_1 and W_2 , i.e $W_1 = W_2 = -W$, we get:

$$F(A, Z) = AZW - ZW \quad (9)$$

Then, taking it $Z = g(X) = XW_Z$ and Eq. (7) into this equation, it becomes:

$$F_{i,:}(A, Z) = \frac{\sum_j [f(X_{i,:}, X_{j,:})Z_{j,:}W]}{\sum_j f(X_{i,:}, X_{j,:})} - Z_i W \quad (10)$$

Due to the fact that $\frac{\sum_j f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})}{\sum_j f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})} = 1$, we can get the formulation of the NS operator:

$$F_i(\mathbf{A}, \mathbf{Z}) = \frac{\sum_j \left[f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:}) \mathbf{Z}_{j,:} \mathbf{W} \right]}{\sum_j f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})} - \frac{\sum_j f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})}{\sum_{j,:} f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})} \mathbf{Z}_{i,:} \mathbf{W}$$

$$= \frac{\sum_j \left[f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:}) (\mathbf{Z}_{j,:} - \mathbf{Z}_{i,:}) \right]}{\sum_j f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})} \mathbf{W} \quad (11)$$

□

1.3. Double Attention Block

The Double Attention Block in the spectral view is the same as defining the graph $\mathcal{G} = (\mathbb{V}, \bar{\mathbf{M}}, \mathbf{Z})$ and then using the second term of the Chebyshev Polynomial to approximate the graph filter, i.e $F(\mathbf{A}, \mathbf{Z}) = \bar{\mathbf{M}} \mathbf{Z} \mathbf{W}$:

Proof. The A^2 operator defined in [1] can be formulated as:

$$\mathbf{O} = \sigma(\theta(\mathbf{X})) \sigma(\phi(\mathbf{X})^T) g(\mathbf{X}) = f^a(\mathbf{X}_{i,:}, \mathbf{X}_{j,:}) \mathbf{X} \mathbf{W} \quad (12)$$

The difference between the double A^2 operator and the NL operator is only the kernel function that calculating the affinity matrix [1]. Thus we can use the similar proving strategy to reformulate the A^2 operator into the spectral only by change the affinity matrix as:

$$\mathbf{A} = \bar{\mathbf{M}} = \sigma(\mathbf{X} \mathbf{W}_\phi) \sigma(\mathbf{X} \mathbf{W}_\psi) \quad (13)$$

□

1.4. Compact Generalized Nonlocal Block

When grouping all channels into one group, the Compact Generalized Nonlocal Block in the spectral view is the same as defining the graph $\mathcal{G} = (\mathbb{V}^f, \mathbf{D}_{M^f}^{-1} \mathbf{M}^f, \text{vec}(\mathbf{Z}))$ and then using the second term of the Chebyshev Polynomial to approximate the graph filter, i.e $F(\mathbf{A}, \mathbf{Z}) = \mathbf{D}_{M^f}^{-1} \mathbf{M}^f \text{vec}(\mathbf{Z}) \mathbf{W}$. Note that due to the dimension of the input feature $\text{vec}(\mathbf{Z}) \in \mathbb{R}^{NC \times 1}$ which is different with other nonlocal operators, here we uses $\mathbf{M}^f, \mathbf{A}^f \in \mathbb{R}^{NC \times NC}$ for clarity.

Proof. The CGNL operator defined in [10] can be formulated as:

$$\text{vec}(\mathbf{O}) = f(\text{vec}(\mathbf{X}), \text{vec}(\mathbf{X})) \text{vec}(\mathbf{Z}) \mathbf{W} \quad (14)$$

For simplicity, we use \mathbf{x} to represent $\text{vec}(\mathbf{X})$, thus the target becomes:

$$\mathbf{o} = f(\mathbf{x}, \mathbf{x}) \mathbf{z} \mathbf{W} \quad (15)$$

Then, we define the graph $\mathcal{G} = (\mathbb{V}^f, \mathbf{A}^f, \mathbf{z})$, where the set \mathbb{V}^f contains each index (including position and channel) of the vector \mathbf{x} . The affinity matrix \mathbf{A} is calculated by:

$$\mathbf{A}^f = \mathbf{M}^f, \quad \mathbf{M}^f = f(\mathbf{x}, \mathbf{x}) \quad (16)$$

The graph filter Ω on \mathcal{G} is approximated by the Chebyshev polynomials. When only choosing the second term, we can get the formulation of the CGNL operator:

$$F(\mathbf{A}^f, \mathbf{z}) = \mathbf{A}^f \mathbf{z} \mathbf{W} = f(\mathbf{x}, \mathbf{x}) \mathbf{z} \mathbf{W} \quad (17)$$

□

1.5. Criss-Cross Attention Block

The Criss-Cross Attention Block in the spectral view is the same as defining the graph $\mathcal{G} = (\mathbb{V}, \mathbf{D}_{C \odot M}^{-1} \mathbf{C} \odot \mathbf{M}, \mathbf{X})$ and then using the second term of the Chebyshev Polynomial to approximate the graph filter with node feature \mathbf{X} :

Proof. The criss-cross attention operator defined in [4] can be formulated as:

$$\mathbf{O}_{i,:} = \sum_{j \in \mathbb{V}^i} \mathbf{A}_{ij} \Phi_{j,:} = \frac{\sum_{j \in \mathbb{V}^i} f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:}) \mathbf{X}_{j,:}}{\sum_{j \in \mathbb{V}^i} f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})} \mathbf{W}$$

where the set \mathbb{V}^i is collection of feature vector in \mathbb{V} which are in the same row or column with position u .

Then, we define the graph $\mathcal{G} = (\mathbb{V}, \tilde{\mathbf{A}}, \mathbf{X})$ to represent the criss-cross attention operator in the spectral view. The affinity matrix $\tilde{\mathbf{A}}$ is calculated by:

$$\tilde{\mathbf{A}} = \mathbf{D}_{C \odot M}^{-1} \mathbf{C} \odot \mathbf{M}, \quad \mathbf{M} = f(\mathbf{X}_i, \mathbf{X}_j)$$

$$C_{ij} = \begin{cases} 1 & j \in \mathbb{V}^i \\ 0 & \text{else} \end{cases},$$

We use $\tilde{\mathbf{M}}$ to represent $\mathbf{C} \odot \mathbf{M}$, i.e. $\tilde{\mathbf{M}} = \mathbf{C} \odot \mathbf{M}$. Thus, each element of the affinity matrix $\tilde{\mathbf{M}}$ is:

$$\tilde{M}_{ij} = \begin{cases} M_{ij} & j \in \mathbb{V}^i \\ 0 & \text{else} \end{cases},$$

Thus, we can get the definition of each element in the affinity matrix $\tilde{\mathbf{A}}$:

$$\tilde{A}_{ij} = \begin{cases} \frac{f(\mathbf{X}_i, \mathbf{X}_j)}{\sum_{j \in \mathbb{V}^i} f(\mathbf{X}_i, \mathbf{X}_j)}, & j \in \mathbb{V}^i \\ 0, & \text{else} \end{cases},$$

When using the Chebyshev polynomials to approximate the generalized graph filter Ω on \mathcal{G} and choose the second term, it becomes:

$$F(\tilde{\mathbf{A}}, \mathbf{X}) = \tilde{\mathbf{A}}\mathbf{X}\mathbf{W} \quad (18)$$

When taking Eq.18 into this formulation, we can get the formulation of CC operator:

$$\begin{aligned} F_{i,:}(\tilde{\mathbf{A}}, \mathbf{X}) &= \left(\frac{\sum_{j \in \mathbb{V}^i} f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})\mathbf{X}_{j,:}}{\sum_{j \in \mathbb{V}^i} f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})} + \sum_{j \notin \mathbb{V}^i} 0\mathbf{X}_{j,:} \right) \mathbf{W} \\ &= \frac{\sum_{j \in \mathbb{V}^i} f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})\mathbf{X}_{j,:}}{\sum_{j \in \mathbb{V}^i} f(\mathbf{X}_{i,:}, \mathbf{X}_{j,:})} \mathbf{W} \end{aligned} \quad (19)$$

□

2. Defining Novel Nonlocal with Caylay Filter

As discussed in our paper, the five existing nonlocal-based blocks uses the Chebyshev filter as the generalized graph filter defined in our spectral view. However the Chebyshev filter need to map the eigenvalue to range $[-1, 1]$ in a linear manner, which makes its spectral smooth or lost the interest signal on a small frequency band [6]. This makes the result feature map tend to ignore the large object (having high correlation to a large amount of positions). Cayley filter solves this problem by an additional spectral zoom parameter h , which can help to select interest frequency band.

To defines novel nonlocal block based on Caylay filter with the help of on our proposed framework, we can replace the generalized graph filter $g(\Lambda)$ in Theorem. 1 by the k^{th} -order Caylay filter, which is formulated as:

$$g_{\theta,h}(\Lambda) = \theta_0 + 2\text{Re}\left\{ \sum_{j=1}^k \theta_1 (h\Lambda - i\mathbf{I})^j (h\Lambda + i\mathbf{I})^{-j} \right\} \quad (20)$$

where "Re" means the real part of an imaginary number, i is the imaginary unit, h and θ are parameters that can be learned by SGD. Taking this into Theorem. 1, we get the CaylaySNL operator with the help of Jacobi approximation [6]:

$$F_{cl}(\mathbf{A}, \mathbf{Z}) = \theta_0 \mathbf{Z} + 2 \sum_{j=1}^k \theta_1 (h\mathbf{L} - i\mathbf{I})^j (h\mathbf{L} + i\mathbf{I})^{-j} \mathbf{Z} \quad (21)$$

where \mathbf{L} is the graph Laplacian of \mathbf{A} . Note that, eigenvalues of \mathbf{L} are all real, so we can remove the "Re" and take $\mathbf{L} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ into this polynomial form equation. Finally, we uses the 1st-order for Eq. 21 for high computational efficiency and extend it to multi-channel condition:

$$F_{cl}(\mathbf{A}, \mathbf{Z}) = \mathbf{Z}\mathbf{W}_1 + 2(\mathbf{L}^2\mathbf{Z}\mathbf{W}_{h^2} + \mathbf{Z})\mathbf{W}_2 \quad (22)$$

Similar with other nonlocal-based operator, when adding residual connection, the CaylaySNL operator becomes the CaylaySNL block. Note that the Caylay filter do not require the affinity matrix A is normalized, so we just uses the sigmoid to make the elements of affinity matrix higher than 0. The CaylaySNL block can be formulated as:

$$\begin{aligned} \mathbf{Y} &= \mathbf{X} + \mathbf{Z}\mathbf{W}_1 + 2(\mathbf{L}^2\mathbf{Z}\mathbf{W}_{h^2} + \mathbf{Z})\mathbf{W}_2 \quad (23) \\ \text{s.t. } \mathbf{A} &= \mathbf{D}_M^{-\frac{1}{2}} \hat{\mathbf{M}} \mathbf{D}_M^{-\frac{1}{2}}, \quad \hat{\mathbf{M}} = (\mathbf{M} + \mathbf{M}^\top)/2 \end{aligned}$$

Remark 1. Different with all existing nonlocal-based blocks that derivated from Chebyshev filter, the proposed ChebySNL utilizes the Caylay filter which can better concerns the small frequency band, i.e. object with large scale.

We validate the proposed CaylaySNL on CIFAR-100 shown in Table. 1. It can be see that, benefited form concerning the spectral zoom, the Caylay can generate the highest performance (0.31 higher than the proposed SNL). This shows the efficient that using other type graph filter to define nonlocal-based blocks.

Table 1: The Performances of Caylay SNL on CIFAR-100

Models	Top1 (%)	Top5 (%)
PreResNet56	75.33 ^{↑0.00}	93.97 ^{↑0.00}
+ NL	75.29 ^{↓0.04}	94.07 ^{↑0.10}
+ NS	75.39 ^{↑0.06}	93.00 ^{↓0.97}
+ A ²	75.51 ^{↑0.18}	92.90 ^{↓1.07}
+ CGNL	74.71 ^{↓0.62}	93.60 ^{↓0.37}
+ SNL	76.41 ^{↑1.08}	94.38 ^{↑0.41}
+ CaylaySNL	76.72^{↑1.39}	94.54^{↑0.57}

3. External Experiments

3.1. Action Recognition

We also conduct the experiments on UCF-101 dataset with other state-of-the-art action recognition models in our supplementary materials including the P3D [7], the MARS [2], and the VTN [5]. For Pseudo 3D Convolutional Network (P3D) and Motion-augmented RGB Stream (MARS), our SNL block are inserted into the P3D right before the last residual layer of the *res3*. For the Video Transformer Network (VTN), we replace its multi-head self-attention blocks (paralleled-connected NL blocks) into our SNL blocks. We use the model pre-trained on Kinetic dataset and fine-tuning on the UCF-101 dataset. Other setting such as the learning rate and training epochs are the same as the experiment on I3D in our paper. We can see that all the performance are improved when adding our proposed SNL model especially when training end-to-end on the small-scale dataset. In sum, our SNL blocks have shown superior

Table 2: Experiments on Video Person Re-identification

Mars			ILID-SVID			PRID-2011		
Models	Rank1(%)	mAP(%)	Models	Rank1(%)	mAP(%)	Models	Rank1(%)	mAP(%)
ResNet50	82.30 ^{↑0.00}	75.70 ^{↑0.00}	ResNet50	74.70 ^{↑0.00}	81.60 ^{↑0.00}	ResNet50	86.50 ^{↑0.00}	90.50 ^{↑0.00}
+ NL	83.21 ^{↑0.91}	76.54 ^{↑0.84}	+ NL	75.30 ^{↑0.60}	83.00 ^{↑1.40}	+ NL	85.40 ^{↓1.10}	89.70 ^{↓0.80}
+ SNL	83.40 ^{↑1.10}	76.80 ^{↑1.10}	+ SNL	76.30 ^{↑1.60}	84.80 ^{↑3.20}	+ SNL	88.80 ^{↑2.30}	92.40 ^{↑1.90}

results across three SOTAs (the VTN and MARS) in the action recognition tasks (0.30% improvement with VTN, 0.50% improvement with MARS).

Table 3: Experiments with state-of-the-art backbone

Models	Top1(%)
P3D[7]	81.23 ^{↑0.00}
P3D + Ours	82.65 ^{↑1.42}
VTN[5]	90.06 ^{↑0.00}
VTN + Ours	90.34 ^{↑0.30}
MARS[2]	92.29 ^{↑0.00}
MARS + Ours	92.79 ^{↑0.50}

3.2. Person Re-identification

The experiments for video person re-identification are conducted on Mars, ILID-SVID and PRID-2011 datasets. For the backbone, we follow the strategy of [3] that use the temporal pooling to fuse the spatial-temporal features. Note that the models are totally trained on ilidvid and prid2011 rather than fine-tuning the pretrained model on Mars. From Table. 2 (Mars), we can see that our SNL can generate 1.10% improvement both on Rank1 and mAP, which are both higher than the original nonlocal block (0.91% on Rank1, 0.84% on mAP). We also generate experiments on two relatively small datasets: ILID-SVID datasets which contains 300 pedestrians captured by two cameras with 600 tracklets; PRID-2011 dataset which contains 200 pedestrians captured by two cameras with 400 tracklets. In Table. 2 (ILID-SVID), our model can generate 1.60% and 3.20% improvements on the Rank1 and mAP respectively for the ILID-SVID dataset. Moreover, on PRID-2011, we get a significant improvement (2.30% on Rank1, 1.90% on mAP) as shown in Table. 2 (PRID-2011).

3.3. Fine-grained Image Classification

The experiments for the fine-grained classification are generated on the Birds-200-2011 (CUB-200) dataset which contains 11,788 images of 200 categories of different birds. We use 5,994 images as the training set and 5,794 images as the testing set [10]. We use the ResNet50 model pre-trained on ImageNet as the backbone and train the models for 110 epochs with the initial learning rate 0.1 which is subsequently divided by 10 at 31, 61, 81 epochs. Table 4

(CUB-200) shows that our model can generate (0.59%) improvement. Compared with the CGNL block concerning channel-wise relations, our SNL is just a bit lower in Top-1 (0.12%). That is because the dependencies among channels play an important role in the fine-grained classification. However, these channel dependencies of CGNL can impede the practical implementations, which needs elaborate preparations for the number of channels per block, the number of blocks and their positions as shown in our main paper. Compared with the other nonlocal block with non-channel concerned, our SNL has improvements with a large margin.

Table 4: Experiments for Nonlocal-based Blocks Added into ResNet50 on CUB-200 Datasets

Models	Top-1 (%)	Top-5 (%)
ResNet50	85.43 ^{↑0.00}	96.70 ^{↑0.00}
+ NL	85.34 ^{↓0.09}	96.77 ^{↑0.07}
+ NS	85.54 ^{↑0.11}	96.56 ^{↓0.14}
+ A ²	85.91 ^{↑0.48}	96.56 ^{↓0.14}
+ CGNL	86.14 ^{↑0.71}	96.34 ^{↓0.36}
+ Ours	86.02 ^{↑0.59}	96.65 ^{↓0.05}

References

- [1] Yunpeng Chen, Yannis Kalantidis, Jianshu Li, Shuicheng Yan, and Jiashi Feng. A²-nets: Double attention networks. In *Neural Information Processing Systems (NeurIPS)*, pages 352–361, 2018. 2
- [2] Nieves Crasto, Philippe Weinzaepfel, Karteek Alahari, and Cordelia Schmid. MARS: Motion-Augmented RGB Stream for Action Recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3, 4
- [3] Jiyang Gao and Ram Nevatia. Revisiting temporal modeling for video-based person reid. *arXiv preprint arXiv:1805.02104*, 2018. 4
- [4] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet: Criss-cross attention for semantic segmentation. In *IEEE International Conference on Computer Vision (ICCV)*, pages 603–612, 2019. 2
- [5] Alexander Kozlov, Vadim Andronov, and Yana Gritsenko. Lightweight network architecture for real-time action recognition. *arXiv preprint arXiv:1905.08711*, 2019. 3, 4
- [6] Ron Levie, Federico Monti, Xavier Bresson, and Michael M Bronstein. Caylennets: Graph convolutional neural networks

432		486
433	with complex rational spectral filters. <i>IEEE Transactions on Signal Processing (TSP)</i> , 67(1):97–109, 2018. 3	487
434	[7] Zhaofan Qiu, Ting Yao, and Tao Mei. Learning spatio-temporal representation with pseudo-3d residual networks. In <i>IEEE International Conference on Computer Vision (ICCV)</i> , pages 5533–5541, 2017. 3, 4	488
435		489
436		490
437		491
438	[8] Yunzhe Tao, Qi Sun, Qiang Du, and Wei Liu. Nonlocal neural networks, nonlocal diffusion and nonlocal modeling. In <i>Neural Information Processing Systems (NeurIPS)</i> , pages 496–506, 2018. 1	492
439		493
440		494
441		495
442	[9] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In <i>IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 7794–7803, 2018. 1	496
443		497
444		498
445		499
446	[10] Kaiyu Yue, Ming Sun, Yuchen Yuan, Feng Zhou, Errui Ding, and Fuxin Xu. Compact generalized non-local network. In <i>Neural Information Processing Systems (NeurIPS)</i> , pages 6510–6519, 2018. 2, 4	500
447		501
448		502
449		503
450		504
451		505
452		506
453		507
454		508
455		509
456		510
457		511
458		512
459		513
460		514
461		515
462		516
463		517
464		518
465		519
466		520
467		521
468		522
469		523
470		524
471		525
472		526
473		527
474		528
475		529
476		530
477		531
478		532
479		533
480		534
481		535
482		536
483		537
484		538
485		539