
Compact Generalized Non-local Network

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 The non-local module [31] is designed for capturing long-range spatio-temporal
2 dependencies in images and videos. Although having shown excellent performance,
3 it lacks the mechanism to model the interactions between positions across channels,
4 which are of vital importance in recognizing fine-grained objects and actions. To
5 address this limitation, we generalize the non-local module and take the correlations
6 between the positions of any two channels into account. This extension unifies
7 second-order feature pooling and achieves state-of-the-arts performance on a variety
8 of fine-grained classification tasks. However, it also leads to an explosion in the
9 computational complexity. To alleviate such an issue, we further propose its
10 compact representation to reduce the high-dimensional feature space and large
11 computation burden involved. Moreover, we try to group the channels and do our
12 generalized non-local method within each group. Experimental results illustrate the
13 significant improvements and practical applicability of the generalized non-local
14 module on both fine-grained object and video classification. Code will be made
15 publicly available to ease the future research.

16

1 Introduction

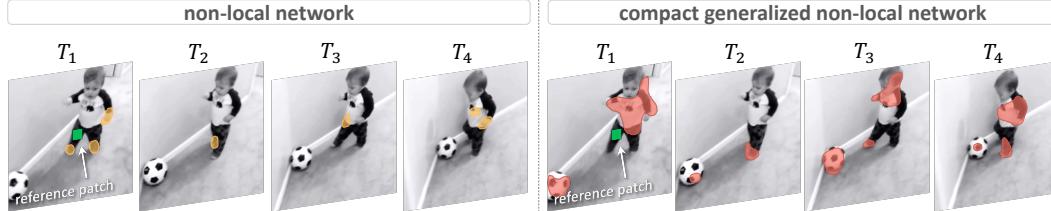


Figure 1: Comparison between non-local (NL) and **compact generalized non-local (CGNL)** on recognizing an action video of kicking the ball. Given the *reference patch* (green rectangle) in the first frame, we visualize for each method the highly related responses in the other frames by thresholding the feature space. CGNL out-performs the original NL in capturing the ball that is not only in long-range distance from the reference patch but also corresponds to different channels in the feature map.

17 Capturing spatio-temporal dependencies between spatial pixels or temporal frames plays key roles in
18 the tasks of fine-grained object and action classification. Modeling such interactions among images
19 and videos is the major topic of various feature extraction techniques, including SIFT, LBP, Dense
20 Trajectory [30], etc. In the past few years, deep neural network automates the feature designing
21 pipeline by stacking multiple end-to-end convolutional or recurrent modules, where each of them
22 processes correlation within spatial or temporal local regions. In general, capturing the long-range
23 dependencies among images or videos still requires multiple stacking of these modules, which greatly
24 hinders the learning and inference efficiency. Recent work [17] also suggests that stacking more
25 layers cannot always increase the effective receptive fields.

26 Inspired by the classical non-local means for image filtering, the recently proposed non-local neural
27 network [31] addresses this challenge by directly modeling the correlation between any two positions
28 in the feature maps in a single module. Without any bells and whistles, the non-local method
29 can greatly improve existing networks on many video classification benchmarks. Despite its great
30 performance, the original non-local network only considers the global spatio-temporal correlation
31 by merging all channels, and it might miss some subtle but important cross-channel clues for
32 discriminating fine-grained objects or actions. For instance, the body, the ball and their interaction
33 are all necessary for describing the action of kicking the ball in Fig. 1, while the original non-local
34 operation learns to focus on the body part relations but neglect the body-ball interactions that usually
35 correspond to different channels of the input features.

36 To improve the effectiveness in fine-grained object and action recognition tasks, this work extends
37 the non-local module by learning explicit correlations among all of the elements across the channels.
38 This extension scale-ups the representation power of the non-local operation to attend the interaction
39 between subtle object parts (*e.g.*, the body and ball in Fig. 1). Second, we propose the compact
40 representation for our generalized non-local module to address the high computation burden issue. We
41 show that as a self-contained module, the compact generalized non-local (CGNL) module provides
42 steady improvements in classification tasks. Third, we also investigate the grouped version of CGNL,
43 which groups the channels and model correlations across channels within each group.

44 We evaluate the proposed CGNL method on the task of fine-grained classification and action recogni-
45 tion. Extensive experimental results show that: 1) The CGNL network are easy to optimize as the
46 original non-local network; 2) Compared with the non-local module, CGNL enjoys capturing richer
47 feature and dense clues for prediction, as shown in Figure 1, which leads to results substantially better
48 than those of the original non-local module.

49 2 Related Works

50 **Channel Correlations:** The mechanism of sharing the same conv kernel among channels of a layer
51 in CNN [12] [11] [23] [26] [7] can be seen as a basic way to capture correlations among channels,
52 which aggregates the channels of feature maps by sum pooling. The SENet [9] may be the first work
53 that explicitly models the interdependencies between the channels of its convolutional features. It
54 aims to select the useful feature maps and suppress the others, and only the global information of
55 each channel is considered. Inspired by [31], we present the generalized non-local (GNL) module,
56 which generalizes the non-local (NL) module to learn the correlations between any two positions
57 across the channels. Compared to the SENet, we model the interdependencies among channels in an
58 explicit and dense manner.

59 **Compact Representation:** After further investigation, we find that the non-local module contains a
60 second-order feature space (Sect.3.1), which is used widely in previous computer vision tasks, *e.g.*,
61 SIFT [16], Fisher encoding [18], Bilinear model [14] [4] and segmentation task [1]. However, such
62 second-order feature space involves high dimensions heavy computational burdens. In the area of
63 kernel learning [22], there are many prior works such as compact bilinear pooling (CBP) [4] to use
64 Tensor Sketching [19] to address this problem. Fortunately, in mathematics, the whole non-local
65 operation can be viewed as a trilinear formation. It can be fast computed with the associative law of
66 matrix production. To the other types of pairwise function, such as Embedded Gaussian or RBF [20],
67 we propose a tight approximation for them by using Taylor expansion.

68 3 Approach

69 In this section, we introduce a general formulation of the proposed general non-local (GNL) operation.
70 We then show that the original non-local and the bilinear pooling are special cases of this formulation.
71 After that, we illustrate that the general non-local operation can be seen as a modality in the trilinear
72 matrix production and show how to implement our generalized non-local module in a compact
73 representations. In this section, we first review the original non-local and the bilinear pooling and
74 show their connection in formulations. Then a general non-local (GNL) operation is proposed to
75 extend non-local to capture channel-wise feature relations. At last, we show how to implement GNL
76 module in a compact representation.

77 **3.1 Review of Non-local Operation**

78 We begin by briefly reviewing the original non-local operation [31] in matrix form. Suppose that an
 79 image or video is given to the network and let $\mathbf{X} \in \mathbb{R}^{N \times C}$ denote (see notation¹) the input feature
 80 map of the non-local module, where C is the number of channels. For the sake of notation clarity,
 81 we collapse all the spatial (width W and height H) and temporal (video length T) positions in one
 82 dimension, *i.e.*, $N = HW$ or $N = HWT$. To capture long-range dependencies across the whole
 83 feature map, the original non-local operation computes the response $\mathbf{Y} \in \mathbb{R}^{N \times C}$ as the weighted
 84 sum of the features at all positions,

$$\mathbf{Y} = f(\theta(\mathbf{X}), \phi(\mathbf{X}))g(\mathbf{X}), \quad (1)$$

85 where $\theta(\cdot), \phi(\cdot), g(\cdot)$ are learnable transformations on the input. In [31], the authors suggest using
 86 1×1 or $1 \times 1 \times 1$ convolution for simplicity, *i.e.*, the transformations can be written as

$$\theta(\mathbf{X}) = \mathbf{X}\mathbf{W}_\theta \in \mathbb{R}^{N \times C}, \quad \phi(\mathbf{X}) = \mathbf{X}\mathbf{W}_\phi \in \mathbb{R}^{N \times C}, \quad g(\mathbf{X}) = \mathbf{X}\mathbf{W}_g \in \mathbb{R}^{N \times C}, \quad (2)$$

87 parameterized by the weight matrices $\mathbf{W}_\theta, \mathbf{W}_\phi, \mathbf{W}_g \in \mathbb{R}^{C \times C}$ respectively. The pairwise function
 88 $f(\cdot, \cdot) : \mathbb{R}^{N \times C} \times \mathbb{R}^{N \times C} \rightarrow \mathbb{R}^{N \times N}$ computes the affinity between all positions (space or space-time).
 89 There are multiple choices for f , among which dot-product is perhaps the simplest one, *i.e.*,

$$f(\theta(\mathbf{X}), \phi(\mathbf{X})) = \theta(\mathbf{X})\phi(\mathbf{X})^\top. \quad (3)$$

90 Plugging Eq. 2 and Eq. 3 into Eq. 1 yields a trilinear interpretation of the non-local operation,

$$\mathbf{Y} = \mathbf{X}\mathbf{W}_\theta \mathbf{W}_\phi^\top \mathbf{X}^\top \mathbf{X}\mathbf{W}_g, \quad (4)$$

91 where the pairwise matrix $\mathbf{X}\mathbf{W}_\theta \mathbf{W}_\phi^\top \mathbf{X}^\top \in \mathbb{R}^{N \times N}$ encodes the similarity between any locations of
 92 the input feature. The effect of non-local operation can be related to the self-attention module [28]
 93 based on the fact that each position (row) in the result \mathbf{Y} is a linear combination of all the positions
 94 (rows) of $\mathbf{X}\mathbf{W}_g$ weighted by the corresponding row of the pairwise matrix.

95 **3.2 Review of Bilinear Pooling**

96 Analogous to the conventional kernel trick [22], the idea of bilinear pooling [14] has recently been
 97 adopted in CNNs for enhancing the feature representation in various tasks, such as fine-grained
 98 classification, person re-id, action recognition. At a glance, bilinear pooling models pairwise feature
 99 interactions using explicit outer product at the final classification layer:

$$\mathbf{Z} = \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{C \times C}, \quad (5)$$

100 where $\mathbf{X} \in \mathbb{R}^{N \times C}$ is the input feature map generated by the last convolutional layer. Each element
 101 of the final descriptor $z_{c_1 c_2} = \sum_n x_{nc_1} x_{nc_2}$ sum-pools at each location $n = 1, \dots, N$ the bilinear
 102 product $x_{nc_1} x_{nc_2}$ of the corresponding channel pair $c_1, c_2 = 1, \dots, C$.

103 Despite the distinct design motivation, it is interesting to see that bilinear pooling (Eq. 5) can be
 104 viewed as a special case of the second-order term (Eq. 3) in the non-local operation if we consider,

$$\theta(\mathbf{X}) = \mathbf{X}^\top \in \mathbb{R}^{C \times N}, \quad \phi(\mathbf{X}) = \mathbf{X}^\top \in \mathbb{R}^{C \times N}. \quad (6)$$

105 **3.3 Generalized Non-local Operation**

106 The original non-local operation aims to directly capture long-range dependencies between any two
 107 positions in one layer. However, such dependencies are encoded in a joint location-wise matrix
 108 $f(\theta(\mathbf{X}), \phi(\mathbf{X}))$ by aggregating all channel information together. On the other hand, channel-wise
 109 correlation has been recently explored in both discriminative [14] and generative [27] models through
 110 the covariance analysis across channels. Inspired by these works, we generalize the original non-local
 111 operation to model long-range dependencies between any positions of any channels.

¹Bold capital letters denote a matrix \mathbf{X} , bold lower-case letters a column vector \mathbf{x} . \mathbf{x}_i represents the i^{th} column of the matrix \mathbf{X} . x_{ij} denotes the scalar in the i^{th} row and j^{th} column of the matrix \mathbf{X} . All non-bold letters represent scalars. $\mathbf{1}_m \in \mathbb{R}^m$ is a vector of ones. $\mathbf{I}_n \in \mathbb{R}^{n \times n}$ is an identity matrix. $\text{vec}(\mathbf{X})$ denotes the vectorization of matrix \mathbf{X} . $\mathbf{X} \circ \mathbf{Y}$ and $\mathbf{X} \otimes \mathbf{Y}$ are the Hadamard and Kronecker products of matrices.

112 We first reshape the output of the transformations (Eq. 2) on \mathbf{X} by merging channel into position:

$$\theta(\mathbf{X}) = \text{vec}(\mathbf{XW}_\theta) \in \mathbb{R}^{NC}, \phi(\mathbf{X}) = \text{vec}(\mathbf{XW}_\phi) \in \mathbb{R}^{NC}, g(\mathbf{X}) = \text{vec}(\mathbf{XW}_g) \in \mathbb{R}^{NC}. \quad (7)$$

113 By lifting the row space of the underlying transformations, our generalized non-local (GNL) operation
114 pursues the same goal of Eq. 1 that computes the response $\mathbf{Y} \in \mathbb{R}^{N \times C}$ as:

$$\text{vec}(\mathbf{Y}) = f(\text{vec}(\mathbf{XW}_\theta), \text{vec}(\mathbf{XW}_\phi)) \text{vec}(\mathbf{XW}_g). \quad (8)$$

115 Compared to the original non-local operation (Eq. 4), GNL utilizes a more general pairwise function
116 $f(\cdot, \cdot) : \mathbb{R}^{NC} \times \mathbb{R}^{NC} \rightarrow \mathbb{R}^{NC \times NC}$ that can differentiate between pairs of same location but at
117 different channels. This richer similarity greatly augments the non-local operation in discriminating
118 fine-grained object parts or action snippets that usually correspond to channels of the input feature.
119 Compared to the bilinear pooling (Eq. 5) that can only be used after the last convolutional layer, GNL
120 maintains the input size and can thus be flexibly plugged between any network blocks. In addition,
121 bilinear pooling neglects the spatial correlation which, however, is preserved in GNL.

122 Recently, the idea of dividing channels into groups has been established as a very effective technique
123 in increasing the capacity of CNNs. Well-known examples include Xception [2], MobileNet [8],
124 ShuffleNet [35], ResNeXt [33] and Group Normalization [32]. Given its simplicity and independence,
125 we also realize the channel grouping idea in GNL by grouping all C channels into G groups, each
126 of which contains $C' = C/G$ channels of the input feature. We then perform GNL operation
127 independently for each group to compute \mathbf{Y}' and concatenate the results along the channel dimension
128 to restore the full response \mathbf{Y} .

129 3.4 Compact Representation

130 A straightforward implementation of GNL (Eq. 8) is prohibitive as the quadratic increase with respect
131 to the channel number C in the presence of the $NC \times NC$ pairwise matrix. Although the channel
132 grouping technique can reduce the channel number from C to C/G , the overall computational
133 complexity is still much higher than the original non-local operation. To mitigate this problem, this
134 section proposes a compact representation that leads to an affordable approximation for GNL.

135 Let us denote $\boldsymbol{\theta} = \text{vec}(\mathbf{XW}_\theta)$, $\boldsymbol{\phi} = \text{vec}(\mathbf{XW}_\phi)$ and $\mathbf{g} = \text{vec}(\mathbf{XW}_g)$, each of which is a NC -D
136 vector column. Without loss of generality, we assume f is a general kernel function (*e.g.*, RBF,
137 bilinear, etc.) that computes a $NC \times NC$ matrix composed by the elements,

$$[f(\boldsymbol{\theta}, \boldsymbol{\phi})]_{ij} \approx \sum_{p=0}^P \alpha_p^2 (\theta_i \phi_j)^p, \quad (9)$$

138 which can be approximated by Taylor series up to certain order P . The coefficient α_p can be computed
139 in closed form once the kernel function is known. Taking RBF kernel for example,

$$[f(\boldsymbol{\theta}, \boldsymbol{\phi})]_{ij} = \exp(-\gamma \|\theta_i - \phi_j\|^2) \approx \sum_{p=0}^P \beta \frac{(2\gamma)^p}{p!} (\theta_i \phi_j)^p, \quad (10)$$

140 where $\alpha_p^2 = \beta \frac{(2\gamma)^p}{p!}$ and $\beta = \exp(-\gamma(\|x\|^2 + \|y\|^2))$ is a constant and $\beta = \exp(-2\gamma)$ if the input
141 vectors $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$ are ℓ_2 -normalized. By introducing two matrices,

$$\boldsymbol{\Theta} = [\alpha_0 \boldsymbol{\theta}^0, \dots, \alpha_P \boldsymbol{\theta}^P] \in \mathbb{R}^{NC \times (P+1)}, \quad \boldsymbol{\Phi} = [\alpha_0 \boldsymbol{\phi}^0, \dots, \alpha_P \boldsymbol{\phi}^P] \in \mathbb{R}^{NC \times (P+1)} \quad (11)$$

142 our compact generalized non-local (CGNL) operation approximates Eq. 8 via a trilinear equation,

$$\text{vec}(\mathbf{Y}) \approx \boldsymbol{\Theta} \boldsymbol{\Phi}^\top \mathbf{g}. \quad (12)$$

143 At first glance, the above approximation still involves the computation of a large pairwise matrix
144 $\boldsymbol{\Theta} \boldsymbol{\Phi}^\top \in \mathbb{R}^{NC \times NC}$. Fortunately, the order of Taylor series is usually relatively small $P \ll NC$.
145 According to the associative law, we could alternatively compute the vector $\mathbf{z} = \boldsymbol{\Phi}^\top \mathbf{g} \in \mathbb{R}^{P+1}$ first
146 and then calculate $\boldsymbol{\Theta} \mathbf{z}$ in a much smaller complexity of $\mathcal{O}(NC(P+1))$.

147 **Complexity analysis:** Table 1 compares the computational complexity of GNL block with CGNL.
148 We cannot afford for directly computing GNL because of its huge complexity of $\mathcal{O}(2(NC)^2)$ in
149 both time and space. Instead, our compact method dramatically eases the heavy calculation to
150 $\mathcal{O}(NC(P+1))$. In Table 2, we investigate the memory cost for each model on one Tesla P40 GPU
151 with 24GB memory. Our models are trained on a server of 8 P40s with 8 clips per batch per GPU (64
152 clips per total mini-batch). We compare on single GPU memory consumption because each GPU is
153 allocated by the same memory during training. This table shows our CGNL operation consumes less
154 GPU memory compared to the NL networks in both 1-block and 5-block settings.

Table 1: Complexity comparison of GNL and CGNL operations, where N and C indicate the number of positions and channels respectively.

| | General NL Method | CGNL Method |
|----------|------------------------|------------------------|
| Strategy | $f(\Theta\Phi^\top)g$ | $\Theta\Phi^\top g$ |
| Time | $\mathcal{O}(2(NC)^2)$ | $\mathcal{O}(NC(P+1))$ |
| Space | $\mathcal{O}(2(NC)^2)$ | $\mathcal{O}(NC(P+1))$ |

Table 2: A summary of memory cost during training. Symbol \times means the models cannot be trained due to the memory explosion. The memory cost is measured by Gigabyte (GB) per GPU.

| model | memory | model | memory |
|---------------|----------|---------------|----------|
| R-101 | 0.86 | R-101 | 0.86 |
| + 1NL block | 0.97 | + 5NL block | 1.09 |
| + 1GNL block | \times | + 5GNL block | \times |
| + 1CGNL block | 0.80 | + 5CGNL block | 0.87 |

155 3.5 Implementation Details

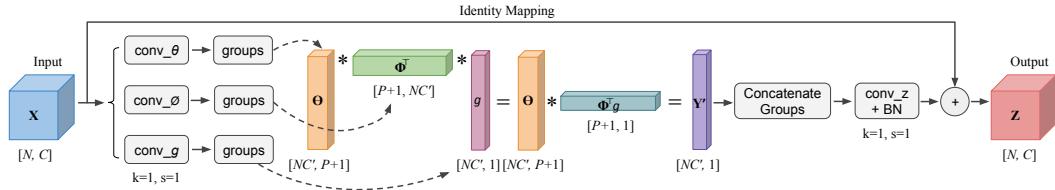


Figure 2: **Compact generalized non-local (CGNL) module with groups.** The feature maps are shown with the shape of their tensors, e.g., $[C, N]$, where $N = THW$ or $N = HW$. The feature maps will be divided along channels into multiple groups after three conv layers whose kernel size and stride both equals 1 ($k = 1, s = 1$). The grouped channels dimension is changed into $C' = C/G$, where G is the group number. The compact representations for generalized non-local module are build within each group. P indicates the order of Taylor expansion for kernel functions.

156 Fig. 2 illustrates the workflow of how CGNL processes a feature map \mathbf{X} of the size $N \times C$, where
157 $N = H \times W$ or $N = T \times H \times W$. \mathbf{X} is first fed into three $1 \times 1 \times 1$ convolutional layers that
158 are described by the weights W_θ, W_ϕ, W_g respectively in Eq. 7. To improve the capacity of neural
159 networks, the channel grouping idea [33, 32] is then applied to divide the transformed feature along
160 the channel dimension into G groups. As shown in Fig. 2, we approximate for each group the GNL
161 operation (Eq. 8) using the Taylor series according to Eq. 12. The output \mathbf{Y}' from each group are
162 concatenated together along the channel dimension to restore the full \mathbf{Y} .

163 To achieve generality and compatibility with existing neural network blocks, the CGNL block is
164 implemented by wrapping Eq. 8 in an identity mapping of the input as in residual learning [7]:

$$\mathbf{Z} = BN(\mathbf{Y}\mathbf{W}_z) + \mathbf{X}, \quad (13)$$

165 where $\mathbf{W}_z \in \mathbb{R}^{C \times C}$ denotes a 1×1 or $1 \times 1 \times 1$ convolution layer followed by a Batch Normalization
166 [10].

167 4 Experiments

168 4.1 Datasets

169 Our CGNL module is evaluated on multiple tasks, including fine-grained classification and action
170 recognition. For fine-grained classification, we perform the experiments on the Birds-200-2011
171 (CUB) dataset [29], which contains 11788 images of bird 200 categories. For action recognition,
172 we perform experiments on two challenging datasets, Mini-Kinetics [34] and UCF101 [24]. The
173 Mini-Kinetics dataset contains 200 categories of data. Due to some invalid video links, we use 78265
174 videos for training and 4986 videos for validation. The UCF101 dataset contains 101 categories of
175 action, which are separated into 25 groups with 4-7 videos of each action in a group.

176 4.2 Baselines

177 Given the steady performance and efficiency, the ResNet [7] series (ResNet-50 and ResNet-101) are
178 adopted as our baselines. For video tasks, we keep the same architecture configuration with [31],

Table 3: **Ablations.** Top1 and top5 accuracy (%) on various validation datasets.

- (a) Kernel function results of models with one CGNL module on CUB. The dot production achieves the best result. But the other kernel functions make accuracy be at the edge of baseline performance.
- (c) The results of CGNL networks with channel grouping on CUB. A few groups help to boost the performance. But more groups tend to prevent the CGNL block from capturing the correlations between positions across channels.

| model | top1 | top5 |
|-------------------|-------|-------|
| baseline | 84.06 | 96.19 |
| Dot Production | 85.14 | 96.88 |
| Gaussian RBF | 84.10 | 95.78 |
| Embedded Gaussian | 84.01 | 96.08 |

- (b) Comparison of validation results between NL and CGNL networks on UCF-101. Note that CGNL is not configured with channel grouping.

| model | top1 | top5 |
|----------------|-------|-------|
| R-50 | 81.62 | 94.62 |
| + 1 NL block | 82.88 | 95.74 |
| + 1 CGNL block | 83.38 | 95.42 |

| model | groups | top1 | top5 |
|----------|--------|-------|-------|
| R-101 | - | 85.05 | 96.70 |
| | 1 | 86.17 | 97.82 |
| + 1 CGNL | 4 | 86.24 | 97.05 |
| block | 8 | 86.35 | 97.86 |
| | 16 | 86.13 | 96.75 |
| | 32 | 86.04 | 96.69 |

- (d) The results of CGNL networks with multiple groups on Mini-Kinetics dataset. More groups help the CGNL networks improve top1 accuracy obviously.

| model | gorups | top1 | top5 |
|----------|--------|-------|-------|
| R-50 | - | 75.54 | 92.16 |
| + 1 CGNL | 1 | 77.16 | 93.56 |
| block | 4 | 77.56 | 93.00 |
| | 8 | 77.76 | 93.18 |

where the temporal dimension is trivially addressed by max pooling. Following [31] the convolutional layers in the baselines are implemented as $1 \times k \times k$ kernels, and we insert our CGNL blocks into the network to turn them into compact generalized non-local (CGNL) networks. The configurations of adding 1 and 5 blocks are investigated. [31] suggests that adding 1 block on the 3-th stage is slightly better than the others. So our experiments of adding 1 block all target the 3-th stage of ResNet. The experiments of adding 5 blocks, on the other hand, are configured by inserting 2 blocks on the 3-th stage, and 3 blocks on the 4-th stage, to every other residual block in ResNet-50 and ResNet-101.

Training: For initializing the weights, we use the models pretrained on ImageNet [21]. We extract all the frames of a video in a *dense* manner. Following [31] to generate 32-frames input clips for models, first we randomly crop out 64 consecutive frames from the full-length video and then drop every other frame. The way to choose 32-frames input clips can be viewed as a data augmentation in the temporal domain. We view these 32-frame images as a group, the shorter side of each group is scaled by randomly sampling S from a certain range $[S_{min}, S_{max}]$ (we used $S_{min} = 256$ and $S_{max} = 320$). Then our spatial crop size in each group is $S = 224$. We use a weight decay of 0.0001 and momentum of 0.9. The dropout [25] with ratio 0.5 is inserted between average pooling layer and last fully-connected layer. To keep same with [31], we initialize the BatchNorm (BN) layer in CGNL and NL block with zero value for scale parameter [5]. To train the networks on CUB dataset, we follow the same training strategy above but the random crop method. We random crop an image in size of 448 after rescaling its shorter size into 512.

Inference: The models are tested immediately after training is finished. In [31], spatially fully-convolutional [15] inference ² is used for NL networks. For these video clips, the shorter side is resized to 256 pixels and use 3 crops to cover the entire spatial size along the longer side. The final prediction is the averaged softmax scores of all clips. For fine-grained classification, we do 1 center-crop testing with input size 448.

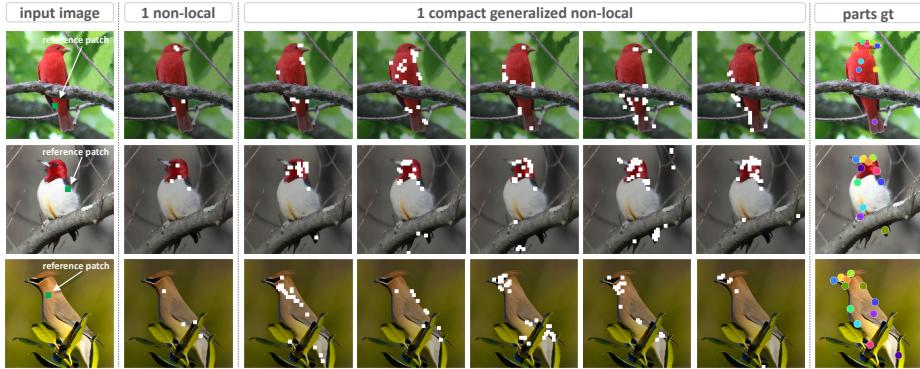
4.3 Ablation Experiments

Kernel Functions: We use three popular kernel functions, namely dot production, embedded Gaussian and Gaussian RBF, in our ablation studies. For dot production, Eq. 12 will be held for direct computation. For embedded Gaussian, the α_p^2 will be $\frac{1}{p!}$ in Eq. 9. And for Gaussian RBF, the corresponding formula is defined as Eq. 10. Following the prior work [3], we expend the Taylor series with fourth order and the hyperparameter γ for RBF is set by $1e-4$. Table 3a suggests that

²<https://github.com/facebookresearch/video-nonlocal-net>

209 dot production is the best kernel functions for CGNL networks. Such experimental observations are
 210 consistent with [31]. The other kernel functions we used, Embedded Gaussian and Gaussian RBF,
 211 has a little improvements for performance. Therefore, we choose the dot production as our main
 212 experimental configuration for other tasks.

213 **Grouping:** The grouping strategy is another important technique. On Mini-Kinetics, Table 3d
 214 suggests that grouping can bring higher accuracy. The improvements brought in by adding groups
 215 are larger than those by reducing the channel reduction ratio. The best top1 accuracy is achieved by
 216 splitting into 8 groups for CGNL networks. On the other hand, however, it is worthwhile to see if
 217 more groups can always improve the results, and Table 3c gives the negative answer, where setting
 218 the group number larger than 8 will decrease the accuracy. This is actually expected, as the affinity in
 219 CGNL block considers the points across channels. When we split the channels into a few groups,
 220 it can facilitate the restricted optimization and ease the training. However, if too many groups are
 adopted, it hinders the affinity to capture the rich correlations between elements across the channels.



221
 222 **Figure 3: Result analysis of the NL block and our CGNL block on CUB.** Column 1: the input images with
 223 a small *reference patch* (green rectangle), which is used to find the highly related patches (white rectangle).
 224 Column 2: the highly related clues for prediction in each feature map found by the NL network. The dimension
 225 of self-attention space in NL block is $N \times N$, where $N = HW$. So its visualization only has one column.
 226 Columns 3 to 7: the most related patches computed by our compact generalized non-local module. We first
 227 pick a reference position in the space of g , then we use the corresponding vectors in Θ and Φ to compute the
 228 attention maps with a threshold. Last column: the ground truth of body parts. The high related areas of CGNL
 229 network can easily cover the standard parts that can provide the prediction clues.

221

222 4.4 Main Results

223 Table 4a shows that although adding 5 NL and CGNL blocks in the baseline networks can both
 224 improve the accuracy, the improvement of CGNL is larger. The same applies to Table 3b and Table 4b.
 225 Similar results are observed in experiments on UCF101 and CUB dataset, where adding 5 CGNL
 226 blocks provides the optimal results both for R-50 and R-101.

227 Table 4a shows the main results on Mini-Kinetics dataset. Compared to the baseline R-50 whose top1
 228 is 75.54%, adding 1 NL block brings improvement by about 1.0%. Similar results can be found in
 229 the experiments based on R-101, where adding 1 CGNL provides about more than 2% improvement,
 230 which is larger than that of adding 1NL block. Table 3b shows the main results on the UCF101
 231 dataset, where adding 1CGNL block achieves higher accuracy than adding 1NL block. And Table 4b
 232 shows the main results on the CUB dataset. To understand the effects brought by CGNL network, we
 233 show the visualization analysis as shown in Fig 3 and Fig 4. Additionally, to investigate the capacity
 234 and the generalization ability of our CGNL network. We test them on the task of object detection and
 235 instance segmentation. We add 1 NL and 1 CGNL block in the R-50 backbone for Mask-RCNN [6].
 236 Table 4c shows the main results on COCO2017 dataset [13] by adopting our 1 CGNL block in the
 237 backbone of Mask-RCNN [6]. It shows that the performance of adding 1 CGNL block is still better
 238 than that of adding 1 NL block.

239 We observe that adding CGNL block can always obtain better results than adding the NL block with
 240 the same blocks number. These experiments suggest that considering the correlations between any

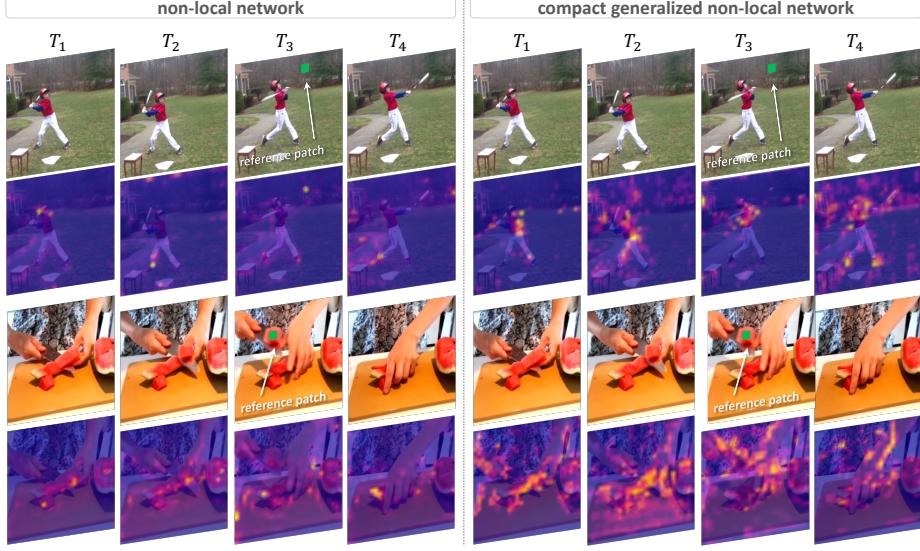


Figure 4: **Visualization with feature heatmaps.** We select a *reference patch* (green rectangle) in one frame, then visualize the high related areas by heatmaps. The CGNL network enjoys capturing dense relationships in feature space than NL networks.

Table 4: **Main results.** Top1 and top5 accuracy (%) on val datasets. And the box and mask AP on COCO.

(a) Main validation results of NL and CGNL networks on Mini-Kinetics. The CGNL networks is established with 8 groups.

| model | top1 | top5 |
|----------------|-------|-------|
| R-50 | 75.54 | 92.16 |
| + 1 NL block | 76.53 | 92.90 |
| + 1 CGNL block | 77.76 | 93.18 |
| + 5 NL block | 77.53 | 94.00 |
| + 5 CGNL block | 78.79 | 94.37 |
| R-101 | 77.44 | 93.18 |
| + 1 NL block | 78.02 | 93.86 |
| + 1 CGNL block | 79.54 | 93.84 |
| + 5 NL block | 79.21 | 93.21 |
| + 5 CGNL block | 79.88 | 93.37 |

(b) Classification results of NL and CGNL networks on CUB. The CGNL models are trained with 8 channel grouping.

| model | top1 | top5 | model | top1 | top5 |
|----------------|-------|-------|----------------|-------|-------|
| R-50 | 84.06 | 96.19 | R-101 | 85.05 | 96.70 |
| + 1 NL block | 84.79 | 96.76 | + 1 NL block | 85.49 | 97.04 |
| + 1 CGNL block | 85.14 | 96.88 | + 1 CGNL block | 86.35 | 97.86 |
| + 5 NL block | 85.10 | 96.18 | + 5 NL block | 86.10 | 96.35 |
| + 5 CGNL block | 85.68 | 96.69 | + 5 CGNL block | 86.24 | 97.23 |

(c) Main results on COCO2017. We add 1 NL block and 1 CGNL block to Mask R-CNN. We observe that the CGNL model achieve better results than NL one under the same training strategy.

| model | AP ^{box} | AP ^{box} ₅₀ | AP ^{box} ₇₅ | AP ^{mask} | AP ^{mask} ₅₀ | AP ^{mask} ₇₅ |
|---------------|-------------------|---------------------------------|---------------------------------|--------------------|----------------------------------|----------------------------------|
| baseline | 34.47 | 54.87 | 36.58 | 30.44 | 51.55 | 31.95 |
| + 1NL block | 35.02 | 55.79 | 37.54 | 30.23 | 52.40 | 32.77 |
| + 1CGNL block | 35.70 | 56.07 | 38.69 | 31.22 | 52.44 | 32.67 |

241 two positions across the channels can significantly improve the performance than that of original
242 non-local methods.

243 5 Conclusion

244 We have introduced a simple approximated formulation of compact generalized non-local operation,
245 and illustrate it on the task of fine-grained classification and action recognition from RGB images.
246 Our formulation allows for explicit modeling of rich interdependencies between any positions across
247 channels in the feature space. To ease the heavy computation of generalized non-local operation, we
248 propose a compact representation with the simple matrix production by using Taylor expansion for
249 multiple kernel functions. It is easy to implement and requires little additional parameters, making it
250 an attractive alternative to the original non-local block, which only considers the correlations between
251 two positions along the specific channel. Our model produces competitive or state-of-the-art results
252 on various benchmarked datasets.

253 **References**

- 254 [1] J. Carreira, R. Caseiro, J. Batista, and C. Sminchisescu. Semantic segmentation with second-order pooling.
255 In *European Conference on Computer Vision*, pages 430–443. Springer, 2012.
- 256 [2] F. Chollet. Xception: Deep learning with depthwise separable convolutions. *arXiv preprint*, 2016.
- 257 [3] Y. Cui, F. Zhou, J. Wang, X. Liu, Y. Lin, and S. Belongie. Kernel pooling for convolutional neural networks.
258 In *Computer Vision and Pattern Recognition (CVPR)*, 2017.
- 259 [4] Y. Gao, O. Beijbom, N. Zhang, and T. Darrell. Compact bilinear pooling. In *Proceedings of the IEEE*
260 *Conference on Computer Vision and Pattern Recognition*, pages 317–326, 2016.
- 261 [5] P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He.
262 Accurate, large minibatch sgd: training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- 263 [6] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Computer Vision (ICCV), 2017 IEEE*
264 *International Conference on*, pages 2980–2988. IEEE, 2017.
- 265 [7] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the*
266 *IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 267 [8] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam.
268 Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint*
269 *arXiv:1704.04861*, 2017.
- 270 [9] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. *arXiv preprint arXiv:1709.01507*, 2017.
- 271 [10] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal
272 covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- 273 [11] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural
274 networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- 275 [12] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition.
276 *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 277 [13] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft
278 coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer,
279 2014.
- 280 [14] T.-Y. Lin, A. RoyChowdhury, and S. Maji. Bilinear cnn models for fine-grained visual recognition. In
281 *Proceedings of the IEEE International Conference on Computer Vision*, pages 1449–1457, 2015.
- 282 [15] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In
283 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.
- 284 [16] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer*
285 *vision*, 60(2):91–110, 2004.
- 286 [17] W. Luo, Y. Li, R. Urtasun, and R. Zemel. Understanding the effective receptive field in deep convolutional
287 neural networks. In *Advances in Neural Information Processing Systems*, pages 4898–4906, 2016.
- 288 [18] F. Perronnin, J. Sánchez, and T. Mensink. Improving the fisher kernel for large-scale image classification.
289 In *European conference on computer vision*, pages 143–156. Springer, 2010.
- 290 [19] N. Pham and R. Pagh. Fast and scalable polynomial kernels via explicit feature maps. In *Proceedings of*
291 *the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 239–247.
292 ACM, 2013.
- 293 [20] T. Poggio and F. Girosi. Networks for approximation and learning. *Proceedings of the IEEE*, 78(9):1481–
294 1497, 1990.
- 295 [21] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla,
296 M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer*
297 *Vision*, 115(3):211–252, 2015.
- 298 [22] B. Scholkopf and A. J. Smola. *Learning with kernels: support vector machines, regularization, optimization,*
299 *and beyond*. MIT press, 2001.
- 300 [23] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition.
301 *arXiv preprint arXiv:1409.1556*, 2014.
- 302 [24] K. Soomro, A. R. Zamir, and M. Shah. Ucf101: A dataset of 101 human actions classes from videos in the
303 wild. *arXiv preprint arXiv:1212.0402*, 2012.
- 304 [25] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to
305 prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958,
306 2014.
- 307 [26] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, A. Rabinovich,
308 et al. Going deeper with convolutions. Cvpr, 2015.
- 309 [27] I. Ustyuzhaninov, W. Brendel, L. A. Gatys, and M. Bethge. What does it take to generate natural textures?
310 In *ICLR*, 2017.
- 311 [28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin.
312 Attention is all you need. In *NIPS*, 2017.
- 313 [29] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD birds-200-2011 dataset.
314 Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.

- 315 [30] H. Wang, A. Kläser, C. Schmid, and C.-L. Liu. Action recognition by dense trajectories. In *Computer*
316 *Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 3169–3176. IEEE, 2011.
- 317 [31] X. Wang, R. Girshick, A. Gupta, and K. He. Non-local neural networks. *arXiv preprint arXiv:1711.07971*,
318 2017.
- 319 [32] Y. Wu and K. He. Group normalization. *arXiv preprint arXiv:1803.08494*, 2018.
- 320 [33] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural
321 networks. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages
322 5987–5995. IEEE, 2017.
- 323 [34] S. Xie, C. Sun, J. Huang, Z. Tu, and K. Murphy. Rethinking spatiotemporal feature learning for video
324 understanding. *arXiv preprint arXiv:1712.04851*, 2017.
- 325 [35] X. Zhang, X. Zhou, M. Lin, and J. Sun. Shufflenet: An extremely efficient convolutional neural network
326 for mobile devices. *arXiv preprint arXiv:1707.01083*, 2017.