# Volumetric Transformer Networks

Seungryong Kim[1], Sabine Süsstrunk[2], and Mathieu Salzmann[2]

[1] Department of Computer Science and Engineering, Korea University, Korea
`seungryong_kim@korea.ac.kr`
[2] School of Computer and Communication Sciences, EPFL, Switzerland
`{sabine.susstrunk, mathieu.salzmann}@epfl.ch`

**Abstract.** Existing techniques to encode spatial invariance within deep convolutional neural networks (CNNs) apply the same warping field to all the feature channels. This does not account for the fact that the individual feature channels can represent different semantic parts, which can undergo different spatial transformations w.r.t. a canonical configuration. To overcome this limitation, we introduce a learnable module, the volumetric transformer network (VTN), that predicts channel-wise warping fields so as to reconfigure intermediate CNN features spatially and channel-wisely. We design our VTN as an encoder-decoder network, with modules dedicated to letting the information flow across the feature channels, to account for the dependencies between the semantic parts. We further propose a loss function defined between the warped features of pairs of instances, which improves the localization ability of VTN. Our experiments show that VTN consistently boosts the features' representation power and consequently the networks' accuracy on fine-grained image recognition and instance-level image retrieval.

**Keywords:** Spatial invariance, attention, feature channels, fine-grained image recognition, instance-level image retrieval

## 1 Introduction

Learning discriminative feature representations of semantic object parts is key to the success of computer vision tasks such as fine-grained image recognition [16,72], instance-level image retrieval [42,46], and people re-identification [74,35]. This is mainly because, unlike generic image recognition and retrieval [11,13], solving these tasks requires handling subtle inter-class variations.

A popular approach to extracting object part information consists of exploiting an attention mechanism within a deep convolutional neural network (CNN) [20,67,42,65]. While effective at localizing the discriminative parts, such an approach has limited ability to handle spatial variations due to, e.g., scale, pose and viewpoint changes, or part deformations, which frequently occur across different object instances [14,28,10]. To overcome this, recent methods seek to spatially warp the feature maps of different images to a canonical configuration so as to remove these variations and thus facilitate the subsequent classifiers task. This trend was initiated by the spatial transformer networks (STNs) [28],
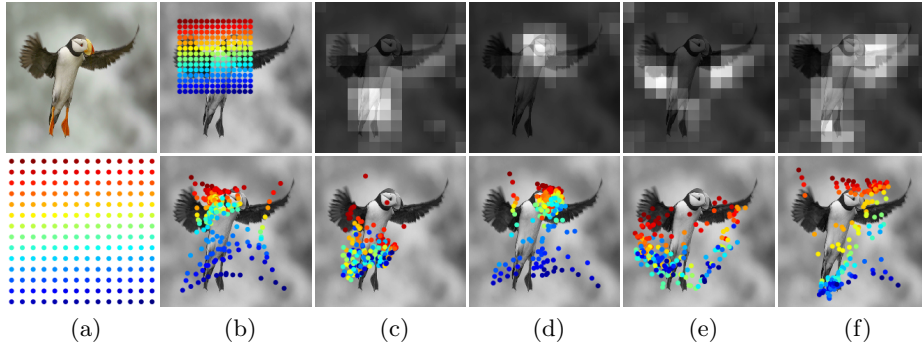
|  (a)  |  (b)  |  (c)  |  (d)  |  (e)  |  (f)  |

**Fig. 1.** Visualization of VTN: (a) input image and target coordinates for warping an intermediate CNN feature map, (b) source coordinates obtained using STNs [28] (top) and SSN [48] (bottom), and (c), (d), (e), and (f) four feature channels and samplers in VTN. Note that the colors in the warping fields represent the corresponding target coordinates. Unlike STNs [28] that applies the same warping field across all the feature channels, VTN maps the individual channels independently to the canonical configuration, by localizing different semantic parts in different channels.

of which many variants were proposed, using a recurrent formalism [37], polar transformations [12], deformable convolutional kernels [10], and attention based samplers [48,73]. All of these methods apply the *same* warping field to *all* the feature channels. This, however, does not account for the findings of [49,5,18], which have shown that the different feature channels of standard image classifiers typically relate to different semantic concepts, such as object parts. Because these semantic parts undergo different transformations w.r.t. the canonical configuration, e.g., the wings of a bird may move while its body remains static, the corresponding feature channels need to be transformed individually.

In this paper, we address this by introducing a learnable module, the volumetric transformer network (VTN), that predicts channel-wise warping fields. As illustrated by Fig. 1, this allows us to correctly account for the different transformations of different semantic parts by reconfiguring the intermediate features of a CNN spatially and channel-wisely. To achieve this while nonetheless accounting for the dependencies between the different semantic parts, we introduce an encoder-decoder network that lets information flow across the original feature channels. Specifically, our encoder relies on a channel-squeeze module that aggregates information across the channels, while our decoder uses a channel-expansion component that distributes it back to the original features.

As shown in previous works [28,69,37,12], training a localization network to achieve spatial invariance is challenging, and most methods [28,69,37,12] rely on indirect supervision via a task-dependent loss function, as supervision for the warping fields is typically unavailable. This, however, does not guarantee that the warped features are consistent across different object instances. To improve the localization ability of the predicted warping fields, we further introduce a loss function defined between the warped features of pairs of instances, so as

to encourage similarity between the representation of same-class instances while pushing that of different-class instances apart.

Our experiments on fine-grained image recognition [64,32,40,22] and instance-level image retrieval [46], performed using several backbone networks and pooling methods, evidence that our VTNs consistently boost the features' representation power and consequently the networks' accuracy.

## 2   Related Work

**Attention mechanisms.** As argued in [76], spatial deformation modeling methods [28,37,10,48], including VTNs, can be viewed as hard attention mechanisms, in that they localize and attend to the discriminative image parts. Attention mechanisms in neural networks have quickly gained popularity in diverse computer vision and natural language processing tasks, such as relational reasoning among objects [4,52], image captioning [67], neural machine translation [3,61], image generation [68,71], and image recognition [23,63]. They draw their inspiration from the human visual system, which understands a scene by capturing a sequence of partial glimpses and selectively focusing on salient regions [27,34].

Unlike methods that consider spatial attention [20,67,42,65], some works [62,70,23,15] have attempted to extract channel-wise attention based on the observation that different feature channels can encode different semantic concepts [49,5,18], so as to capture the correlations among those concepts. In those cases, however, spatial attention was ignored. While some methods [7,65] have tried to learn spatial and channel-wise attention simultaneously, they only predict a fixed spatial attention with different channel attentions. More importantly, attention mechanisms have limited ability to handle spatial variations due to, e.g., scale, pose and viewpoint changes, or part deformations [14,28,10].

**Spatial invariance.** Recent work on spatial deformation modeling seeks to spatially warp the features to a canonical configuration so as to facilitate recognition [28,37,10,12,48]. STNs [28] explicitly allow the spatial manipulation of feature maps within the network while attending to the discriminative parts. Their success inspired many variants that use, e.g., a recurrent formalism [37], polar transformations [12], deformable convolutional kernels [10], and attention based warping [48,73]. These methods typically employ an additional network, called localization network, to predict a warping field, which is then applied to all the feature channels identically. Conceptually, this corresponds to using hard attention [20,67,42,65], but it improves spatial invariance. While effective, this approach concentrates on finding the regions that are most discriminative across all the feature channels. To overcome this, some methods use multi-branches [28,72], coarse-to-fine schemes [16], and recurrent formulations [36], but they remain limited to considering a pre-defined number of discriminative parts, which restricts their effectiveness and flexibility.

**Fine-grained image recognition.** To learn discriminative feature representations of object parts, conventional methods first localize these parts and then

classify the whole image based on the discriminative regions. These two-step methods [6,25] typically require bounding box or keypoint annotations of objects or parts, which are hard to collect. To alleviate this, recent methods aim to automatically localize the discriminative object parts using an attention mechanism [16,72,36,58,51,48,8,73] in an unsupervised manner, without part annotations. However, these methods do not search for semantic part representations in the individual feature channels, which limits their ability to boost the feature representation power. Recently, Chen et al. [8] proposed a destruction and construction learning strategy that injects more discriminative local details into the classification network. However, the problem of explicitly processing the individual feature channels remains untouched.

**Instance-level image retrieval.** While image retrieval was traditionally tackled using local invariant features [39,41] or bag-of-words (BoW) models [56,1], recent methods use deep CNNs [2,59,30,59,42,47] due to their better representation ability. In this context, the main focus has been on improving the feature representation power of pretrained backbone networks [33,55,21], typically by designing pooling mechanisms to construct a global feature, such as max-pooling (MAC) [59], sum-pooling (SPoC) [2], weighted sum-pooling (CroW) [30], regional max-pooling (R-MAC) [59], and generalized mean-pooling (GeM) [47]. These pooling strategies, however, do not explicitly leverage the discriminative parts, and neither do the methods [19,47] that have tried to fine-tune the pretrained backbone networks [33,55,21]. While the approach of [42] does, by learning spatial attention, it ignores the channel-wise variations. Taking such variations into account is the topic of this paper.

## 3    Volumetric Transformer Networks

### 3.1    Preliminaries

Let us denote an intermediate CNN feature map as $U \in \mathbb{R}^{H \times W \times K}$, with height $H$, width $W$, and $K$ channels. To attend to the discriminative object parts and reduce the inter-instance spatial variations in the feature map, recent works [28,69,37,12] predict a warping field to transform the features to a canonical pose. This is achieved via a module that takes $U$ as input and outputs the parameters defining a warping field $G \in \mathbb{R}^{H \times W \times 2}$ to be applied to $U$. The representation in the canonical pose is then obtained via a feature sampling mechanism, which, for every pixel $i$ in the output representation, produces a warped feature such that $V(i) = U(i + G(i))$. As argued above, while this reduces spatial variations and lets the network focus on discriminative image regions, the same warping field is applied across all the channels, without considering the different semantic meanings of these individual channels. Moreover, this does not explicitly constrain the warped features of different instances of the same class to be consistent.

### 3.2    Motivation and Overview

By contrast, our volumetric transformer network (VTN), which we introduce in the remainder of this section, encodes the observation that each channel in an
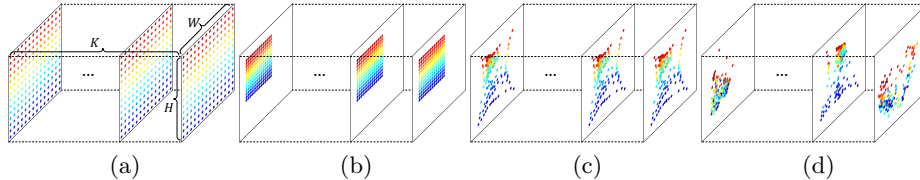
**Fig. 2.** Intuition of VTNs: (a) target coordinates for warping an intermediate CNN feature map and source coordinates obtained using (b) STNs [28], (c) SSN [48], and (d) VTNs, which predict channel-wise warping fields.

intermediate CNN feature map acts as a pattern detector, i.e., high-level channels detect high-level semantic patterns, such as parts and objects [5,7,18,65], and, because these patterns can undergo different transformations, one should separately attend to the discriminative parts represented by the individual channels to more effectively warp them to the canonical pose. To achieve this, unlike existing spatial deformation modeling methods [28,37,10,48], which apply the same warping field to all the feature channels, as in Fig. 2(b), our VTN predicts channel-wise warping fields, shown in Fig. 2(d).

Concretely, a VTN produces a warping field $G_c \in \mathbb{R}^{H \times W \times 2}$ for each channel $c$. Rather than estimating the warping field of each channel independently, to account for dependencies between the different semantic parts, we design two modules, the channel squeeze and expansion modules, that let information flow across the channels. Furthermore, to improve the computational efficiency and localization accuracy, we build a group sampling and normalization module, and a transformation probability inference module at the first and last layer of VTN, respectively. To train the network, instead of relying solely on a task-dependent loss function as in [28,37,10,48], which may yield poorly-localized warping fields, we further introduce a loss function based on the distance between the warped features of pairs of instances, thus explicitly encouraging the warped features to be consistent across different instances of the same class.

### 3.3 Volumetric Transformation Estimator

Perhaps the most straightforward way to estimate channel-wise warping fields is to utilize convolutional layers that take the feature map $U$ as input and output the warping fields $G = \{G_c\} \in \mathbb{R}^{H \times W \times 2 \times K}$. This strategy, however, uses separate convolution kernels for each warping field $G_c$, which might be subject to overfitting because of the large number of parameters involved. As an alternative, one can predict each warping field $G_c$ independently, by taking only the corresponding feature channel $U_c \in \mathbb{R}^{H \times W \times 1}$ as input. This, however, would fail to account for the inter-channel relationships, and may be vulnerable to outlier channels that, on their own, contain uninformative features but can yet be supported by other channels [18,71,29].

To alleviate these limitations, we introduce the channel squeeze and expansion modules, which yield a trade-off between the two extreme solutions discussed above. We first decompose the input feature map across the channel dimension,
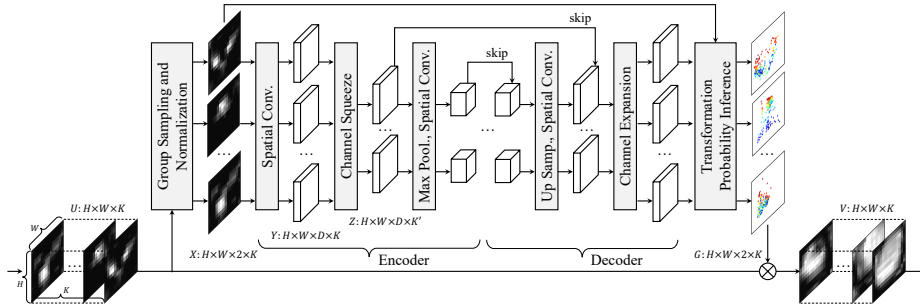
**Fig. 3.** VTN architecture. A VTN consists of a group sampling and normalization module, sequential spatial convolutions, channel squeeze and expansion modules, and a transformation probability inference module.

and apply a shared convolution to each of the $K$ channels. We then combine the original feature channels into $K'$ new channels by a channel-squeeze module, parameterized by a learned matrix $W_{\mathrm{cs}}$, in the encoder and expand these squeezed feature channels into $K$ channels by a channel-expansion module, parameterized by a learned matrix $W_{\mathrm{ce}}$, in the decoder.

Formally, as depicted by Fig. 3, let us define an intermediate CNN feature map after a forward pass through an encoder as $Y = \mathcal{F}(U; W_{\mathrm{s}}) \in \mathbb{R}^{H \times W \times D \times K}$, where each feature channel is processed independently with spatial convolution parameters $W_{\mathrm{s}}$ shared across the channels, introducing an additional dimension of size $D$. We introduce a channel squeeze module, with parameters $W_{\mathrm{cs}} \in \mathbb{R}^{K \times K'}$, $K' < K$, applied to the reshaped $Y \in \mathbb{R}^{HWD \times K}$, whose role is to aggregate the intermediate features so as to output $Z = \mathcal{F}(Y; W_{\mathrm{cs}}) \in \mathbb{R}^{HWD \times K'}$, which can also be reshaped to $\mathbb{R}^{H \times W \times D \times K'}$. In short, this operation allows the network to learn how to combine the initial $K$ channels so as to leverage the inter-channel relationships while keeping the number of trainable parameters reasonable. We then incorporate a channel expansion module, with parameters $W_{\mathrm{ce}} \in \mathbb{R}^{K' \times K}$, which performs the reverse operation, thereby enlarging the feature map $Z \in \mathbb{R}^{H \times W \times D \times K'}$ back into a representation with $K$ channels. This is achieved through a decoder.

We exploit sequential spatial convolution and channel squeeze modules in the encoder, and sequential spatial convolution and channel expansion modules in the decoder. In our experiments, the volumetric transformation estimator consists of an encoder with 4 spatial convolution and channel squeeze modules followed by max-pooling, and a decoder with 4 spatial convolution and channel expansion modules followed by upsampling. Each convolution module follows the architecture Convolution-BatchNorm-ReLU [26].

**Grouping the channels.** In practice, most state-of-the-art networks [55,21,24] extract high-dimensional features, and thus processing all the initial feature channels as described above can be computationally prohibitive. To overcome this, we propose a group sampling and normalization module inspired by group normalization [66] and attention mechanisms [65]. Concretely, a group sampling
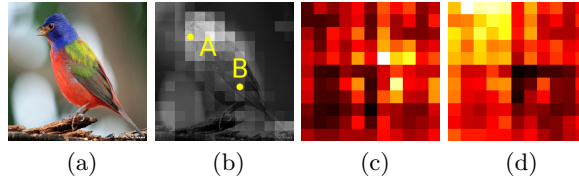
(a)            (b)            (c)            (d)

**Fig. 4.** Visualization of the learned probability of transformation candidates: (a) input image, (b) arbitrary feature channel, and (c) and (d) probabilities of points A and B in (b), respectively. VTNs estimate the probability of transformation candidates, instead of directly estimating the warping fields.

and normalization module takes the feature map $U$ as input and separates it into $C$ groups following the sequential order of the channels. We then aggregate the features $U_c$ in each group $c \in \{1, \ldots, C\}$ by using two pooling operations: $U_c^{\max} \in \mathbb{R}^{H \times W \times 1}$ and $U_c^{\text{avg}} \in \mathbb{R}^{H \times W \times 1}$, and concatenate them as $X_c \in \mathbb{R}^{H \times W \times 2}$, followed by group normalization without per-channel linear transform [66]. We then take the resulting $X = \{X_c\} \in \mathbb{R}^{H \times W \times 2 \times K}$ as input to the volumetric transformation estimator described above, instead of $U$.

**Probabilistic transformation modeling.** Unlike existing spatial deformation modeling methods [28,37] that rely on parametric models, e.g., affine transformation, VTNs estimate non-parametric warping fields, thus having more flexibility. However, regressing the warping fields directly may perform poorly because the mapping from the features to the warping fields adds unnecessary learning complexity. To alleviate this, inspired by [60,57], we design a probabilistic transformation inference module that predicts probabilities for warping candidates, instead of directly estimating the warping field. Specifically, we predict the probability $P_c(i, j)$ of each candidate $j \in N_i$ at each pixel $i$, and compute the warping field $G_c$ by aggregating these probabilities as

$$G_c(i) = \sum_{j \in N_i} P_c(i, j)(j - i). \tag{1}$$

Furthermore, instead of predicting the probability $P_c(i, j)$ directly, we compute a residual probability and then use a softmax layer such that

$$P_c(i, j) = \Psi\left(\left(U_c^{\max}(j) + U_c^{\text{avg}}(j) + E_c(i, j)\right) / \beta\right), \tag{2}$$

where $E_c \in \mathbb{R}^{H \times W \times |N_i|}$ is the output of the volumetric transformation estimator whose the size depends on the number of candidates $|N_i|$. Note that $E_c(i, j)$ is a scalar because $i$ denotes a spatial point over $H \times W$ and $j$ indexes a point among all candidates. $\Psi(\cdot)$ is the softmax operator and $\beta$ is a parameter adjusting the sharpness of the softmax output. At initialization, the network parameters are set to predict zeros, i.e., $E_c(i, j) = 0$, thus the warping fields are determined by candidate feature responses $U_c^{\max} + U_c^{\text{avg}}$, which provide good starting points. As training progresses, the network provides increasingly regularized warping fields. This is used as the last layer of the VTN. Fig. 4 visualizes the learned probability of some transformation candidates.
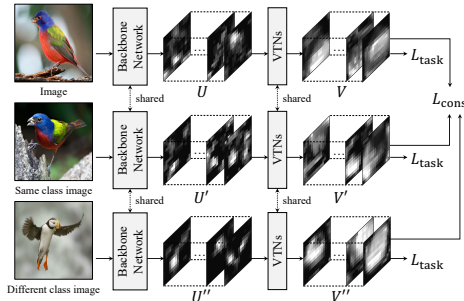
**Fig. 5.** Illustration of training VTNs using our consistency loss function. By simultaneously exploiting a sample from the same class and another sample from a different class, our consistency loss function improves the localization ability and the discriminative power of the intermediate features.

### 3.4   Loss Function

Similarly to existing deformation modeling methods [28,69,37,12], our network can be learned using only the final task-dependent loss function $\mathcal{L}_{\text{task}}$, without using ground-truth warping fields, since all modules are differentiable. This, however, does not explicitly constrain the warped features obtained from the predicted warping fields to be consistent across different object instances of the same class. To overcome this, we draw our inspiration from semantic correspondence works [50,31], and introduce an additional loss function modeling the intuition that the warped features of two instances of the same class should match and be similar. The simplest approach to encoding this consists of using a square loss between such features, which yields

$$\mathcal{L} = \sum_i \|V(i) - V'(i)\|^2, \tag{3}$$

where $V$ and $V'$ are the warped feature maps of two instances of the same class. Minimizing this loss function, however, can induce erroneous solutions, such as constant features at all the pixels. To avoid such trivial solutions, we use a triplet loss function [54,69] simultaneously exploiting a sample $V'$ from the same class as $V$ and another sample $V''$ from a different class. We then express our loss as

$$\mathcal{L}_{\text{cons}} = \sum_i \left[ \|V(i) - V'(i)\|^2 - \|V(i) - V''(i)\|^2 + \alpha \right]_+, \tag{4}$$

where $\alpha > 0$ is a threshold parameter and $[\cdot]_+ = \max(\cdot, 0)$. Our loss function jointly encourages the instances' features from the same class to be similar, and the instances' features from different classes to be dissimilar. Together, this helps to improve the features' discriminative power, which is superior to relying solely on a task-dependent loss function, as in previous methods [28,69,37,12]. Note that our approach constitutes the first attempt at learning warping fields that generate consistent warped features across object instances. To train our VTNs, we then use the total loss $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{task}} + \lambda \mathcal{L}_{\text{cons}}$ with balancing parameter $\lambda$. Fig. 5 depicts the training procedure of VTNs.

### 3.5   Implementation and Training Details

We implemented VTNs using the `Pytorch` library [43]. In our experiments, we use VGGNet [55] and ResNet [21] backbones pretrained on ImageNet [11]. We

| Methods | Backbone | [64] | [32] | [40] |
|---|---|---|---|---|
| Base | VGG-19 | 71.4 | 68.7 | 80.7 |
|  | ResNet-50 | 74.6 | 70.4 | 82.1 |
| Def-Conv [10] | VGG-19 | 74.2 | 70.1 | 82.6 |
|  | ResNet-50 | 76.7 | 72.1 | 83.7 |
| STNs [28] | VGG-19 | 72.1 | 69.2 | 81.1 |
|  | ResNet-50 | 76.5 | 71.0 | 81.2 |
| SSN [48] | VGG-19 | 75.1 | 72.7 | 84.6 |
|  | ResNet-50 | 77.7 | 74.8 | 83.1 |
| ASN [73] | VGG-19 | 76.2 | 74.1 | 82.4 |
|  | ResNet-50 | 78.9 | 75.2 | 85.7 |
| **VTNs** wo/$W_{cs}, W_{ce}$ | VGG-19 | 77.8 | 78.6 | 86.1 |
|  | ResNet-50 | 80.1 | 81.4 | 86.9 |
| **VTNs** wo/Group | VGG-19 | 76.3 | 76.1 | 84.4 |
|  | ResNet-50 | 77.2 | 79.1 | 82.4 |
| **VTNs** wo/T-Probability | VGG-19 | 78.1 | 79.7 | 84.9 |
|  | ResNet-50 | 79.0 | 80.4 | 85.1 |
| **VTNs** wo/$\mathcal{L}_{cons}$ | VGG-19 | 79.2 | 80.2 | 87.1 |
|  | ResNet-50 | 82.4 | 82.1 | 84.9 |
| **VTNs** | VGG-19 | 80.4 | 81.9 | 87.4 |
|  | ResNet-50 | **83.1** | **82.7** | **89.2** |

**Table 1.** Accuracy of VTNs compared to spatial deformation modeling methods on fine-grained image recognition benchmarks (CUB-Birds [64], Stanford-Cars [32], and FGVC-Aircraft [40]).

build VTNs on the last convolution layers of each network. For fine-grained image recognition, we replace the 1000-way softmax layer with a $k$-way softmax layer, where $k$ is number of classes in the dataset [64,32,40,22], and fine-tune the networks on the dataset. The input images were resized to a fixed resolution of $512 \times 512$ and randomly cropped to $448 \times 448$. We apply random rotations and random horizontal flips for data augmentation. For instance-level image retrieval, we utilize the last convolutional features after the VTNs as global representation. To train VTNs, we follow the experimental protocols of [47,19]. We set the hyper-parameters by cross-validation on CUB-Birds [64], and then used the same values for all experiments. We set the size of the transformation candidates $|N_i| = 11 \times 11$, the parameter $\beta = 10$, the number of groups $C = 32$, and the balancing parameter $\lambda = 1$. We also set the threshold parameter $\alpha = 20$ for VGGNet [55] and $\alpha = 30$ for ResNet [21], respectively, because they have different feature distributions. The source code is available online at our project webpage: `http://github.com/seungryong/VTNs/`.

## 4    Experiments

### 4.1    Experimental Setup

In this section, we comprehensively analyze and evaluate VTNs on two tasks: fine-grained image recognition and instance-level image retrieval. First, we analyze the influence of the different components of VTNs compared to existing
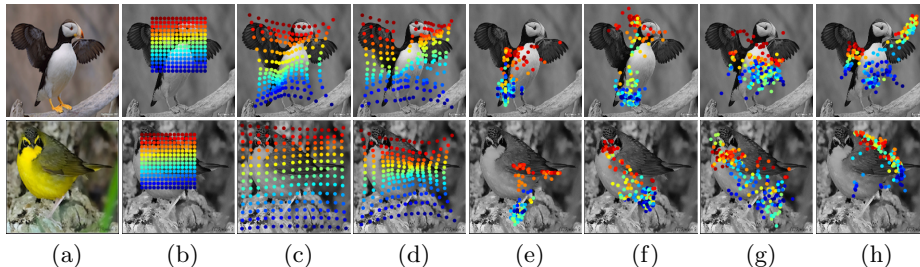
|   (a)   |   (b)   |   (c)   |   (d)   |   (e)   |   (f)   |   (g)   |   (h)   |

**Fig. 6.** Comparison of VTN warping fields with those of existing deformation modeling methods [28,48] on examples from CUB-Birds [64]: (a) input images and source coordinates obtained using (b) STNs [28], (c) SSN [48], (d) ASN [73], and (e), (f), (g), and (h) four feature channel samplers in VTNs. Points with the same color in different images are projected to the same point in the canonical pose. This shows that VTNs not only localize different semantic parts in different channels but identify the same points across different images.

spatial deformation modeling methods [28,10,48,73] and the impact of combining VTNs with different backbone networks [55,21] and second-order pooling strategies [38,17,9,36,35]. Second, we compare VTNs with the state-of-the-art methods on fine-grained image recognition benchmarks [64,32,40,22]. Finally, we evaluate them on instance-level image retrieval benchmarks [46].

### 4.2   Fine-grained Image Recognition

**Analysis of the VTN components.** To validate the different components of our VTNs, we compare them with previous spatial deformation modeling methods, such as STNs [28], deformable convolution (Def-Conv) [10], saliency-based sampler (SSN) [48], and attention-based sampler (ASN) [73] on fine-grained image recognition benchmarks, such as CUB-Birds [64], Stanford-Cars [32], and FGVC-Aircraft [40]. For the comparison to be fair, we apply these methods at the same layer as ours, i.e., the last convolutional layer. In this set of experiments, we utilize VGGNet-19 [55] and ResNet-50 [21] as backbone networks. As an ablation study, we evaluate VTNs without the channel squeeze and expansion modules, denoted by VTNs wo/$W_{cs}, W_{ce}$, without the group sampling and normalization module, denoted by VTNs wo/Group, and without the transformation probability inference module, denoted by VTNs wo/T-Probability. We further report the results of VTNs trained without our consistency loss function, denoted by VTNs wo/$\mathcal{L}_{cons}$.

The results are provided in Table 1 and Fig. 6. Note that all versions of our approach outperform the existing deformation modeling methods [28,10,48,73]. Among these versions, considering jointly spatial and channel-wise deformation fields through our squeeze and expansion modules improves the results. The group sampling and normalization and transformation probability inference modules also boost the results. Using the consistency loss function $\mathcal{L}_{cons}$ further yields higher accuracy by favoring learning a warping to a consistent canonical

| Methods | [64] | [32] | [40] |
|---|---|---|---|
| Base | 74.6 | 70.4 | 82.1 |
| BP [38] | 80.2 | 81.5 | 84.8 |
| CBP [17] | 81.6 | 81.6 | 88.6 |
| KP [9] | 83.2 | 82.9 | 89.9 |
| MPN-COV [36] | 84.2 | 83.1 | 89.7 |
| iSQRT-COV [35] | 88.1 | 90.0 | 92.8 |
| Base+**VTNs** | 83.1 | 82.7 | 89.2 |
| BP [38]+**VTNs** | 84.9 | 84.1 | 90.6 |
| CBP [17]+**VTNs** | 85.2 | 84.2 | 91.2 |
| KP [9]+**VTNs** | 85.1 | 83.2 | 91.7 |
| MPN-COV [36]+**VTNs** | 86.7 | 88.1 | 90.6 |
| iSQRT-COV [35]+**VTNs** | **89.6** | **93.3** | **93.4** |

**Table 2.** Accuracy of VTNs incorporated with second-order pooling methods on fine-grained image recognition benchmarks (CUB-Birds [64], Stanford-Cars [32], and FGVC-Aircraft [40]).

| Methods | Backbone | [64] | [32] | [40] |
|---|---|---|---|---|
| RA-CNN [16] | 3×VGG-19 | 85.3 | 92.5 | 88.2 |
| MA-CNN [72] | 3×VGG-19 | 86.5 | 92.8 | 89.9 |
| DFL-CNN [63] | ResNet-50 | 87.4 | 93.1 | 91.7 |
| DT-RAM [36] | ResNet-50 | 87.4 | 93.1 | 91.7 |
| MAMC [58] | ResNet-50 | 86.5 | 93.0 | 92.9 |
| NTSN [58] | 3×ResNet-50 | 87.5 | 91.4 | 93.1 |
| DCL [8] | VGG-16 | 86.9 | 94.1 | 91.2 |
|  | ResNet-50 | 87.8 | 94.5 | 93.0 |
| TASN [73] | VGG-19 | 86.1 | 93.2 | - |
|  | ResNet-50 | 87.9 | 93.8 | - |
| [35]+TASN [73] | ResNet-50 | 89.1 | - | - |
| DCL [8]+**VTNs** | ResNet-50 | 89.2 | 95.1 | 93.4 |
| [35]+**VTNs** | ResNet-50 | 89.6 | 93.3 | 93.4 |
| [35]+TASN [73]+**VTNs** | ResNet-50 | **91.2** | **95.9** | **94.5** |

**Table 3.** Accuracy of VTNs compared to the state-of-the-art methods on fine-grained image recognition benchmarks (CUB-Birds [64], Stanford-Cars [32], and FGVC-Aircraft [40]).

configuration of instances of the same class and improving the discriminative power of the intermediate features.

**Incorporating second-order pooling strategies.** While our VTNs can be used on their own, they can also integrate second-order pooling schemes, such as bilinear pooling (BP) [38], compact BP (CBP) [17], kernel pooling (KP) [9], matrix power normalized covariance pooling (MPN-COV) [36], and iterative matrix square root normalization of covariance pooling (iSQRT-COV) [35], which yield state-of-the-art results on fine-grained image recognition. In this set of experiments, we use ResNet-50 [21] as backbone. As shown in Table 2, our VTNs consistently outperform the corresponding pooling strategy on its own, thus confirming the benefits of using channel-wise warped regions.

| Super Class | ResNet [21] | SSN [48] | TASN [73] | [35]+**VTNs** |
|---|---|---|---|---|
| Plantae | 60.3 | 63.9 | 66.6 | **68.6** |
| Insecta | 69.1 | 74.7 | 77.6 | **79.1** |
| Aves | 59.1 | 68.2 | 72.0 | **72.9** |
| Reptilia | 37.4 | 43.9 | 46.4 | **48.1** |
| Mammalia | 50.2 | 55.3 | 57.7 | **60.6** |
| Fungi | 62.5 | 64.2 | 70.3 | **72.1** |
| Amphibia | 41.8 | 50.2 | 51.6 | **53.9** |
| Mollusca | 56.9 | 61.5 | 64.7 | **66.3** |
| Animalia | 64.8 | 67.8 | 71.0 | **73.2** |
| Arachnida | 64.8 | 73.8 | 75.1 | **78.2** |
| Actinoopterygii | 57.0 | 60.3 | 65.5 | **68.4** |
| Chromista | 57.6 | 57.6 | 62.5 | **64.0** |
| Protozoa | 78.1 | 79.5 | 79.5 | **81.1** |
| Total | 58.4 | 63.1 | 66.2 | **68.2** |

**Table 4.** Accuracy of VTNs compared to the state-of-the-art methods on iNaturalist-2017 [22].

| Methods | top-1 err. | top-5 err. |
|---|---|---|
| GoogLeNet+GAP [75] | 59.00 | - |
| VGGNet+ACoL [70] | 54.08 | 43.49 |
| ResNet+GCAM [53] | 53.42 | 43.12 |
| ResNet+STNs [28]+GCAM [53] | 54.21 | 43.33 |
| ResNet+**VTNs**+GCAM [53] | **52.18** | **41.76** |

**Table 5.** Localization errors on CUB-Birds [64].

**Comparison with the state-of-the-art methods.** We also compare VTNs with the state-of-the-art fine-grained image recognition methods, such as RA-CNN [16], MA-CNN [72], DFL-CNN [63], DT-RAM [36], MAMC [58], NTSN [58], DCL [8], and TASN [73]. Since our VTN is designed as a generic drop-in layer that can be combined with existing backbones and pooling strategies, we report the results of VTNs combined with DCL [8], TASN [73], and iSQRT-COV [35], which are the top-performers on this task. As can be seen in Table 3, our method outperforms the state of the art in most cases. In Table 4, we further evaluate VTNs with iSQRT-COV [35] on iNaturalist-2017 [22], the largest fine-grained recognition dataset, on which we consistently outperform the state of the art.

**Network visualization.** To analyze the feature representation capability of VTNs, we applied Grad-CAM (GCAM) [53], which uses gradients to calculate the importance of the spatial locations, to STN- and VTN-based networks. As shown in Fig. 7 and Table 5, compared to the ResNet-50 [21] backbone, STNs [28] only focus on the most discriminative parts, and thus discard other important parts. Unlike them, VTNs improve the feature representation power by allowing the networks to focus on the most discriminative parts represented by each feature channel.
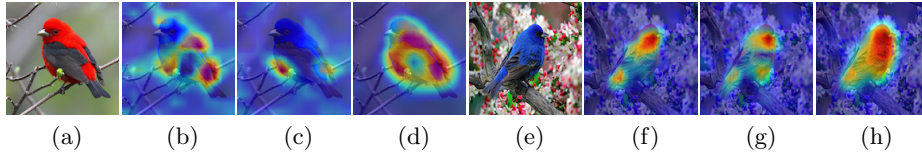
(a)          (b)          (c)          (d)          (e)          (f)          (g)          (h)

**Fig. 7.** Network visualization using Grad-CAM [53]: (a), (e) input images, (b), (f) ResNet-50 [21], (c), (g) ResNet-50 with STNs [28], and (d), (h) ResNet-50 with VTNs.
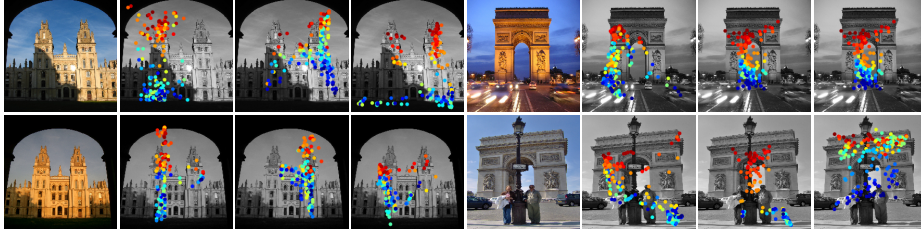


**Fig. 8.** Visualization of warping fields of VTNs at each channel on some instances of the $\mathcal{R}$Oxford and $\mathcal{R}$Paris benchmarks [46]. VTNs not only localize different semantic parts in different channels but identify the same points across different images.

### 4.3  Instance-level Image Retrieval

Finally, we evaluate our VTNs on the task of instance-level image retrieval using the $\mathcal{R}$Oxford and $\mathcal{R}$Paris benchmarks [46], which address some limitations of the standard Oxford-5K [44] and Paris-6K benchmarks [45], such as annotation errors, size of the dataset, and level of difficulty, and comprise 4,993 and 6,322 images, respectively. Following standard practice, we use the mean average precision (mAP) [44] for quantitative evaluation. We follow the evaluation protocol of [46], using two evaluation setups (*Medium* and *Hard*). As baselines, we use a pretrained ResNet-50 [21] backbone, followed by various pooling methods, such as MAC [59], SPoC [2], CroW [30], R-MAC [59], and GeM [47]. We also evaluate deep local attentive features (DELF) [42] with an aggregated selective match kernel [59] and spatial verification [44] that learns spatial attention, and incorporate our VTNs into them. Furthermore, we report the results of end-to-end training techniques [47,19], and incorporate our VTNs on top of them. As evidenced by our significantly better results in Table 6, focusing on the most discriminative parts at each feature channel is one of the key to the success of instance-level image retrieval. Note that the comparison with STNs shows the benefits of our approach, which accounts for different semantic concepts across different feature channels and thus, even for rigid objects, is able to learn more discriminative feature representations than a global warping. Fig. 8 visualizes some warping fields of VTNs.

## 5  Conclusion

We have introduced VTNs that predict channel-wise warping fields to boost the representation power of an intermediate CNN feature map by reconfiguring the

| Methods | Medium | | Hard | |
|---|---|---|---|---|
| | $\mathcal{R}$Oxf | $\mathcal{R}$Par | $\mathcal{R}$Oxf | $\mathcal{R}$Par |
| Pretr.+MAC [59] | 41.7 | 66.2 | 18.0 | 44.1 |
| Pretr.+SPoC [2] | 39.8 | 69.2 | 12.4 | 44.7 |
| Pretr.+CroW [30] | 42.4 | 70.4 | 13.3 | 47.2 |
| Pretr.+GeM [47] | 45.0 | 70.7 | 17.7 | 48.7 |
| Pretr.+R-MAC [59] | 49.8 | 74.0 | 18.5 | 52.1 |
| DELF [42,59,44] | 67.8 | 76.9 | 43.1 | 55.4 |
| [47]+GeM | 64.7 | 77.2 | 38.5 | 56.3 |
| [19]+R-MAC | 60.9 | 78.9 | 32.4 | 59.4 |
| [19]+R-MAC+STNs [28] | 61.3 | 79.4 | 36.1 | 59.8 |
| DELF [42,59,44]+**VTNs** | **69.7** | 78.1 | 45.1 | 56.4 |
| [47]+GeM+**VTNs** | 67.4 | 80.5 | **45.5** | 57.1 |
| [19]+R-MAC+**VTNs** | 65.6 | **82.7** | 43.3 | **60.9** |

**Table 6.** Comparison of VTNs with the state-of-the-art methods on the $\mathcal{R}$Oxford and $\mathcal{R}$Paris benchmarks [46].

features spatially and channel-wisely. VTNs account for the fact that the individual feature channels can represent different semantic information and require different spatial transformations. To this end, we have developed an encoder-decoder network that relies on channel squeeze and expansion modules to account for inter-channel relationships. To improve the localization ability of the predicted warping fields, we have further introduced a loss function defined between the warped features of pairs of instances. Our experiments have shown that VTNs consistently boost the features' representation power and consequently the networks' accuracy on fine-grained image recognition and instance-level image retrieval tasks. In the future, we will aim to apply VTNs to other tasks, such as person re-identification and local feature matching.

## Acknowledgments

# References

1. Arandjelovic, R., Zisserman, A.: All about vlad. In: CVPR (2013)
2. Babenko, A., Lempitsky, V.: Aggregating deep convolutional features for image retrieval. In: ICCV (2015)
3. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. In: ICLR (2015)
4. Battaglia, P.W., Pascanu, R., Lai, M., Rezende, D.J., Kavukcuoglu, K.: Interaction networks for learning about objects, relations and physics. In: NeurIPS (2016)
5. Bau, D., Zhou, B., Khosla, A., Oliva, A., Torralba, A.: Network dissection: Quantifying interpretability of deep visual representations. In: CVPR (2017)
6. Berg, T., Liu, J., Lee, S.W., Alexander, M.L., Jacobs, D.W., Belhumeur, P.N.: Birdsnap: Large-scale fine-grained visual categorization of birds. In: CVPR (2014)
7. Chen, L., Zhang, H., Xiao, J., Nie, L., Shao, J., Liu, W., Chua, T.: Sca-cnn: Spatial and channel-wise attention in convolutional networks for image captioning. In: CVPR (2017)
8. Chen, Y., Bai, Y., Zhang, W., Mei, T.: Destruction and construction learning for fine-grained image recognition. In: CVPR (2019)
9. Cui, Y., Zhou, F., Wang, J., Liu, X., Lin, Y., Belongie, S.: Kernel pooling for convolutional neural networks. In: CVPR (2017)
10. Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., Wei, Y.: Deformable convolutional networks. In: ICCV (2017)
11. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A Large-Scale Hierarchical Image Database. In: CVPR (2009)
12. Esteves, C., Allen-Blanchette, C., Zhou, X., Daniilidis, K.: Polar transformer networks. In: ICLR (2018)
13. Everingham, M., Eslami, S.M.A., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The pascal visual object classes challenge: A retrospective. IJCV **111**(1), 98 − 136 (2015)
14. Felzenszwalb, P.F., Girshick, R.B., McAllester, D., Ramanan, D.: Object detection with discriminative trained part based models. IEEE Trans. PAMI **32**(9), 1627 − 1645 (2010)
15. Fu, J., Liu, J., Tian, H., Fang, Z., Lu, H.: Dual attention network for scene segmentation. In: CVPR (2019)
16. Fu, J., Zheng, H., Mei, T.: Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition. In: CVPR (2017)
17. Gao, Y., Beijbom, O., Zhang, N., Darrell, T.: Compact bilinear pooling. In: CVPR (2016)
18. Gonzalez-Garcia, A., Modolo, D., Ferrari, V.: Do semantic parts emerge in convolutional neural networks? IJCV **126**(5), 476–494 (2018)
19. Gordo, A., Almazàn, J., Revaud, J., Larlus, D.: End-to-end learning of deep visual representations for image retrieval. IJCV **124**(2), 237 − 254 (2017)
20. Gregor, K., Danihelka, I., Graves, A., Rezende, D.J., Wierstra, D.: Draw: A recurrent neural network for image generation. In: ICML (2015)
21. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
22. Horn, G.V., Aodha, O.M., Song, Y., Cui, Y., Sun, C., Shepard, A., Adam, H., Perona, P., Belongie, S.: The inaturalist species classification and detection dataset. In: CVPR (2018)
23. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: CVPR (2018)

24. Huang, G., Liu, Z., Laurens, V.D.M., Weinberger, K.Q.: Densely connected convolutional networks. In: CVPR (2017)
25. Huang, S., Xu, Z., Tao, D., Zhang, Y.: Part-stacked cnn for fine-grained visual categorization. In: CVPR (2016)
26. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: ICML (2015)
27. Itti, L., Koch, C., Niebur, E.: A model of saliency-based visual attention for rapid scene analysis. TPAMI **20**(11), 1254 – 1259 (1998)
28. Jaderberg, M., Simonyan, K., Zisserman, A., Kavukcuoglu, K.: Spatial transformer networks. In: NeurIPS (2015)
29. Jeong, J., Shin, J.: Training cnns with selective allocation of channels. In: ICML (2019)
30. Kalantidis, Y., Mellina, C., Osindero, S.: Cross-dimensional weighting for aggregated deep convolutional features. In: ECCVW (2016)
31. Kim, S., Min, D., Jeong, S., Kim, S., Jeon, S., Sohn, K.: Semantic attribute matching networks. In: CVPR (2019)
32. Krause, J., Stark, M., Deng, J., Fei-Fei, L.: 3d object representations for fine-grained categorization. In: CVPR (2015)
33. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: NeurIPS (2012)
34. Larochelle, H., Hinton, G.E.: Learning to combine foveal glimpses with a third-order boltzmann machine. In: NeurIPS (2010)
35. Li, P., Xie, J., Wang, Q., Zuo, W.: Towards faster training of global covariance pooling networks by iterative matrix square root normalization. In: CVPR (2018)
36. Li, Z., Yang, Y., Liu, X., Zhou, F., Wen, S., Xu, W.: Dynamic computational time for visual attention. In: ICCVW (2017)
37. Lin, C.H., Lucey, S.: Inverse compositional spatial transformer networks. In: CVPR (2017)
38. Lin, T.Y., RoyChowdhury, A., Maji, S.: Bilinear cnn models for fine-grained visual recognition. In: ICCV (2015)
39. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. IJCV **60**(2), 91–110 (2004)
40. Maji, S., Kannala, J., Rahtu, E., Blaschko, M., Vedaldi, A.: Fine-grained visual classification of aircraft. Tech. rep. (2013)
41. Mikolajczyk, K., Schmid, C.: A performance evaluation of local descriptors. TPAMI **27**(10), 1615 – 1630 (2005)
42. Noh, H., Araujo, A., Sim, J., Weyand, T., Han, B.: Large-scale image retrieval with attentive deep local features. In: ICCV (2017)
43. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch (2017)
44. Philbin, J., Chum, O., Isard, M., Sivic, J., Zisserman, A.: Object retrieval with large vocabularies and fast spatial matching. In: CVPR (2007)
45. Philbin, J., Chum, O., Isard, M., Sivic, J., Zisserman, A.: Lost in quantization: Improving particular object retrieval in large scale image databases. In: CVPR (2008)
46. Radenovì, F., Iscen, A., Tolias, G., Avrithis, Y., Chum, O.: Revisiting oxford and paris: Large-scale image retrieval benchmarking. In: CVPR (2018)
47. Radenovì, F., Tolias, G., Chum, O.: Fine-tuning cnn image retrieval with no human annotation. TPAMI **41**(7), 1655–1668 (2019)
48. Recasens, A., Kellnhofer, P., Stent, S., Matusik, W., Torralba, A.: Learning to zoom: a saliency-based sampling layer for neural networks. In: ECCV (2018)

49. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: NeurIPS (2015)
50. Rocco, I., Arandjelović, R., Sivic, J.: End-to-end weakly-supervised semantic alignment. In: CVPR (2018)
51. Rodriguez, P., Gonfaus, J.M., Cucurull, G., Xavierroca, F., Gonzalez, J.: Attend and rectify: A gated attention mechanism for fine-grained recovery. In: ECCV (2018)
52. Santoro, A., Raposo, D., Barrett, D.G., Malinowski, M., Pascanu, R., Battaglia, P., Lillicrap, T.: A simple neural network moduel for relational reasoning. In: NeurIPS (2017)
53. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: ICCV (2017)
54. Simo-Serra, E., Trulls, E., Ferraz, L., Kokkinos, I., Fua, P., Moreno-Noguer, F.: Discriminative learning of deep convolutional feature point descriptors. In: ICCV (2015)
55. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: ICLR (2015)
56. Sivic, J., Zisserman, A.: Video google: A text retrieval approach to object matching in videos. In: ICCV (2003)
57. Sun, K., Xiao, B., Liu, D.: Deep high-resolution representation learning for human pose estimation. In: CVPR (2019)
58. Sun, M., Yuan, Y., Zhou, F., Ding, E.: Multi-attention multi-class constraint for fine-grained image recognition. In: ECCV (2018)
59. Tolias, G., Avrithis, Y., Jgou, H.: Image search with selective match kernels: Aggregation across single and multiple images. IJCV **116**(3), 247–261 (2016)
60. Tompson, J., Goroshin, R., Jain, A., LeCun, Y., Bregler, C.: Efficient object localization using convolutional networks. In: CVPR (2015)
61. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: NeurIPS (2017)
62. Wang, F., Jiang, M., Qian, C., Yang, S., Li, C., Zhang, H., X., W.: Residual attention network for image classification. In: CVPR (2017)
63. Wang, Y., Morariu, V.I., Davis, L.S.: Learning a discriminative filter bank within a cnn for fine-grained recognition. In: CVPR (2018)
64. Welinder, P., Branson, S., Mita, T., Wah, C., Schroff, F., Belongie, S., Perona, P.: Caltech-ucsd birds 200. Tech. Rep. CNS-TR-2010-001, California Institute of Technology (2010)
65. Woo, S., Park, J., Lee, J.Y., Kweon, I.S.: Cbam: Convolutional block attention module. In: ECCV (2018)
66. Wu, Y., He, K.: Group normalization. In: ECCV (2018)
67. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., Bengio, Y.: Show, attend and tell: Neural image caption generation with visual attention. In: ICML (2015)
68. Xu, T., Zhang, P., Huang, Q., Zhang, H., Gan, Z., Huang, X., He, X.: Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In: CVPR (2018)
69. Yi, K.M., Trulls, E., Lepetit, V., Fua, P.: Lift: Learned invariant feature transform. In: ECCV (2016)
70. Zhang, H., Dana, K., Shi, J., Zhang, Z., Wang, X., Tyagi, A., Agrawal, A.: Context encoding for semantic segmentation. In: CVPR (2018)
71. Zhang, H., Goodfellow, I., Metaxas, D., Odena, A.: Self-attention generative adversarial networks. In: ICML (2019)

72. Zheng, H., Fu, J., Mei, T., Luo, J.: Learning multi-attention convolutional neural network for fine-grained image recognition. In: ICCV (2017)
73. Zheng, H., Fu, J., Zha, Z.J., Luo, J.: Looking for the devil in the details: Learning trilinear attention sampling network for fine-grained image recognition. In: CVPR (2019)
74. Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., Tian, Q.: Scalable person re-identification: A benchmark. In: ICCV (2015)
75. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: CVPR (2016)
76. Zhu, X., Cheng, D., Zhang, Z., Lin, S., Dai, J.: An empirical study of spatial attention mechanisms in deep networks. In: ICCV (2019)