# Design of Kernels in Convolutional Neural Networks for Image Classification

Zhun Sun    Mete Ozay    Takayuki Okatani

{sun, mozay, okatani}@vision.is.tohoku.ac.jp

**Abstract.** Despite the effectiveness of convolutional neural networks (CNNs) for image classification, our understanding of the effect of shape of convolution kernels on learned representations is limited. In this work, we explore and employ the relationship between shape of kernels which define receptive fields (RFs) in CNNs for learning of feature representations and image classification. For this purpose, we present a feature visualization method for visualization of pixel-wise classification score maps of learned features. Motivated by our experimental results, and observations reported in the literature for modeling of visual systems, we propose a novel design of shape of kernels for learning of representations in CNNs.

In the experimental results, the proposed models also outperform the state-of-the-art methods employed on the CIFAR-10/100 datasets [1] for image classification. We also achieved an outstanding performance in the classification task, comparing to a base CNN model that introduces more parameters and computational time, using the ILSVRC-2012 dataset [2]. Additionally, we examined the region of interest (ROI) of different models in the classification task and analyzed the robustness of the proposed method to occluded images. Our results indicate the effectiveness of the proposed approach.

**Keywords:** convolutional neural networks, deep learning, convolution kernel, kernel design, image classification.

## 1 Introduction

Following the success of convolutional neural networks (CNNs) for large scale image classification [2,3], remarkable efforts have been made to deliver state-of-the-art performance on this task. Along with more complex and elaborate architectures, lots of techniques concerning parameter initialization, optimization and regularization have also been developed to achieve better performance. Despite the fact that various aspects of CNNs have been investigated, design of the convolution kernels, which can be considered as one of the fundamental problems, has been barely studied. Some studies examined how size of kernels affects performance [4], leading to a recent trend of stacking small kernels (e.g. $3 \times 3$) in deep layers of CNNs. However, analysis of the shapes of kernels is mostly left untouched. Although there seems to be no latitude in designing the shape of convolution kernels intuitively (especially $3 \times 3$ kernels), in this work, we suggest that designing the shapes of kernels is feasible and practical, and we analyze its effect on the performance.
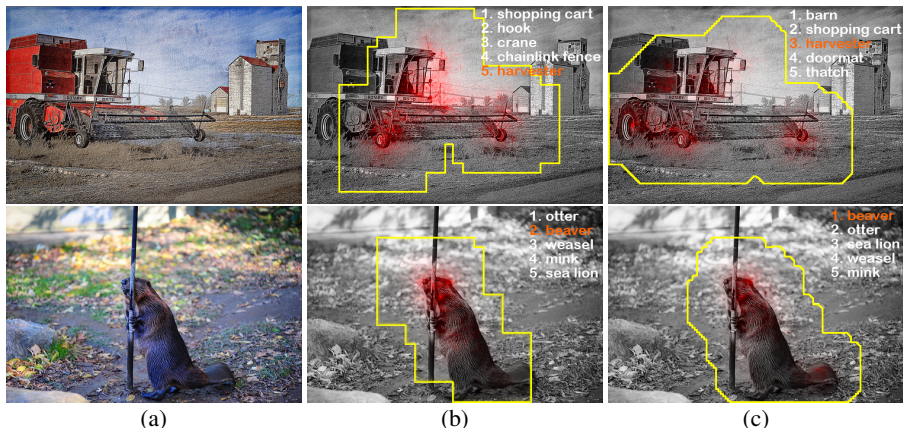
Fig. 1: Examples of visualization of ROI (Sect. C) in two images (a) for CNNs equipped with kernels with (b) square, and (c) our proposed "quasi-hexagonal" shapes (Sect. 2). The pixels marked with red color indicate their maximum contribution for classification scores of the correct classes. For (b), these pixels tend to be concentrated on local, specific parts of the object, whereas for (c), they distribute more across multiple local parts of the object. See texts for more details.

In the early studies of biological vision [5,6,7], it was observed that the receptive fields (RFs) of neurons are arranged in an approximately hexagonal lattice. A recent work reported an interesting result that an irregular lattice with appropriately adjusted asymmetric RFs can be accurate in representation of visual patterns [8]. Intriguingly, hexagonal-shaped filters and lattice structures have been analyzed and employed for solving various problems in computer vision and image processing [9,10]. In this work, motivated by these studies, we propose a method for designing the kernel shapes in CNNs. Specifically, we propose a method to use an asymmetric shape, which simulates hexagonal lattices, for convolution kernels (see Fig. 10 and 4), and then deploy kernels with this shape in different orientations for different layers of CNNs (Sect. 2).

This design of kernel shapes brings multiple advantages. Firstly, as will be shown in the experimental results (Sect. 4.1), CNNs which employ the proposed design method are able to achieve comparable or even better classification performance, compared to CNNs which are constructed using the same architectures (same depth and output channels for each layer) but employing square ($3 \times 3$) kernels. Thus, a notable improvement in computational efficiency (a reduction of 22% parameters and training time) can be achieved as the proposed kernels include fewer weights than $3 \times 3$ kernels. Meanwhile, increasing the number of output channels of our proposed models (to keep the number of parameters same as corresponding models with square shape), leads to a further improvement in performance.

Secondly, CNNs which employ our proposed kernels provide improvement in learning for extraction of discriminative features in a more flexible and robust manner. This results in better robustness to various types of noise in natural images that could make classification erroneous, such as occlusions. Fig. 8 shows examples of visualization of
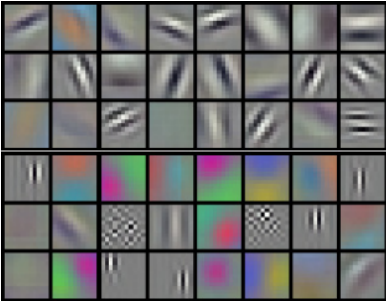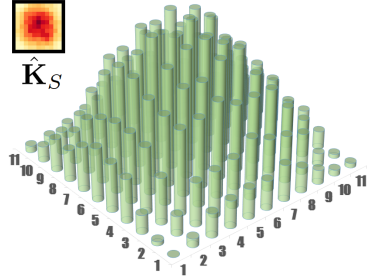
(a) $\{\mathbf{K}_{a,l=1}^{S} \in \mathbb{R}^{K \times K}\}_{a=1}^{48}$.

(b) $\frac{1}{A} \sum_c |\mathbf{K}_{a,1}^{S}(i,j,c)|, \forall i,j,a$

Fig. 2: (a) Visualization of a subset of kernels $\mathbf{K}_{a,l}^{S} \in \mathbb{R}^{K \times K}$, where $K$ is the size of kernel, at the first convolution layer $l = 1$ of AlexNet [3] trained on ImageNet. (b) An average kernel $\hat{\mathbf{K}}_S = \frac{1}{A} \sum_{a=1}^{A} |\mathbf{K}_{a,l}^{S}|$ is depicted at the top-left part. Each bar in the histogram shows a cumulative distribution of values over each channel, $c$.

features extracted using fully-trained CNNs equipped with and without our proposed kernels, which are obtained by the method introduced in Sec. C. These depict the image pixels that have the maximum contribution to the classification score of the correct class (shown in red). It is observed that for CNNs equipped with our proposed kernels, they tend to be less concentrated on local regions and rather distributed across a number of sub-regions, as compared to CNNs with standard square kernels. This property prevents erroneous classification due to occlusions, as will be shown in the experimental results. This also helps to explain the fact that the CNNs equipped with our proposed kernels perform on par with the CNNs equipped with square kernels despite having less number of parameters.

The contributions of the paper are summarized as follows:

1. We propose a method to design convolution kernels in deep layers of CNNs, which is inspired by hexagonal lattice structures employed for solving various problems of computer vision and image processing.
2. We examine classification performance of CNNs equipped with our kernels, and compare the results with state-of-the-art CNNs equipped with square kernels using benchmark datasets, namely ImageNet and CIFAR 10/100. The experimental results show that the proposed method is superior to the state-of-the-art CNN models in terms of computational time and/or classification performance.
3. We introduce a method for visualization of features to qualitatively analyze the effect of kernel design on classification. Additionally, we analyze the robustness of CNNs equipped with and without our kernel design to occlusion by measuring their classification accuracy when some regions on input images are occluded.

## 2   Our approach

We propose a method for designing shape of convolution kernels which will be employed for image classification. The proposed method enables us to reduce the computa-
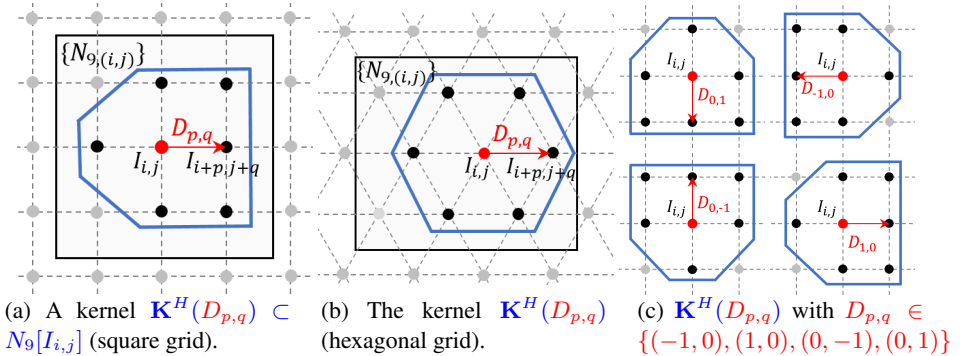
(a) A kernel $\mathbf{K}^H(D_{p,q}) \subset N_9[I_{i,j}]$ (square grid).    (b) The kernel $\mathbf{K}^H(D_{p,q})$ (hexagonal grid).    (c) $\mathbf{K}^H(D_{p,q})$ with $D_{p,q} \in \{(-1,0),(1,0),(0,-1),(0,1)\}$

Fig. 3: (a) Our proposed kernel. (b) It can approximate a hexagonal kernel by shifting through direction $D$. (c) A set of kernel candidates which are denoted as design pattens "U","R", "D", "L" from left to right.
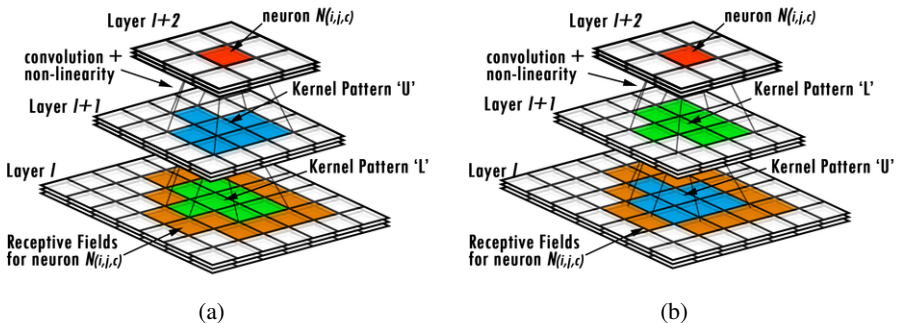


(a)

(b)

Fig. 4: (a) Employment of the proposed method in CNNs by stacking small size "quasi-hexagonal" kernels. (b) The kernels employed at different layers of a two-layer CNN will induce the same pattern of RFs on images observed in (a), if only the kernels designed with the same patterns are used, independent of order of their employment.

tional time of training CNNs providing more compact representations, while preserving the classification performance.

In CNNs [3,11,4], an input image (or feature map) $\mathbf{I} \in \mathbb{R}^{W \times H \times C}$ is convolved with a series of square shaped kernels $\mathbf{K}^S \in \mathbb{R}^{K \times K \times C}$ through its hierarchy. The convolution operation $\mathbf{K}^S * \mathbf{I}$ can be considered as sampling of the image $\mathbf{I}$, and extraction of discriminative information with learned representations. Fig. 9 shows a subset of learned kernels $\mathbf{K}^S$, and the kernel $\hat{\mathbf{K}}_S$ averaged over all the kernels employed at the first layer of AlexNet [3]. Distribution of values of $\hat{\mathbf{K}}_S$ shows that most of the weights at the corner take values close to zero, thus making less contribution for representing features at the higher layers. If a computationally efficient and compressed model is desired, additional methods need to be employed, such as pruning these diluted parameters during fine-tuning [12].

## 2.1  Designing shape of convolution kernels

In this work, we address the aforementioned problems by designing shapes of kernels on a two-dimensional coordinate system. For each channel of a given image $\mathbf{I}$, we associate each pixel $I_{i,j} \in \mathbf{I}$ at each coordinate $(i, j)$ with a lattice point (i.e., a point with integer coordinates) in a square grid (Fig. 3a) [13,14]. If two lattice points in the grid are distinct and each $(i, j)$ differs from the corresponding coordinate of the other by at most 1, then they are called 8-adjacent [13,14]. An 8-neighbor of a lattice point $I_{i,j} \in \mathbf{I}$ is a point that is 8-adjacent to $I_{i,j}$. We define $N_9[I_{i,j}]$ as a set consisting of a pixel $I_{i,j} \in \mathbf{I}$, and its 8 nearest neighbors (Fig. 3a). A shape of a *quasi-hexagonal* kernel $\mathbf{K}^H(D_{p,q}) \subset N_9[I_{i,j}]$ is defined as

$$\mathbf{K}^H(D_{p,q}) = \{I_{i+p,j+q} : I_{i,j} \in N_9[I_{i,j}]\}, \tag{1}$$

where $D_{p,q} \in \mathcal{D}$ is a random variable used as an indicator function employed for designing of shape of $\mathbf{K}^H(D_{p,q})$, and takes values from $\mathcal{D} = \{(-1,0), (1,0), (0,-1), (0,1)\}$ (see Fig. 3c). Then, convolution of the proposed quasi-hexagonal kernel $\mathbf{K}^H(D_{p,q})$ on a neighborhood centered at a pixel located at $(x, y)$ on an image $\mathbf{I}$ is defined as

$$I_{x,y} * \mathbf{K}^H(D_{p,q}) = \sum_{s,t} \mathbf{K}^H_{s,t}(D_{p,q}) I_{x-s,y-t}. \tag{2}$$

## 2.2  Properties of receptive fields and quasi-hexagonal kernels

Aiming at more flexible representation of shapes of natural objects which may diverge from a fixed square shape, we stack "quasi-hexagonal" kernels designed with different shapes, as shown in Fig. 4. For each convolution layers, we randomly select $D_{p,q} \in \mathcal{D}$ according to a uniform distribution to design kernels. Random selection of design patterns of kernels is feasible because the shapes of RFs will not change, independent of the order of employment of kernels if only the kernels designed with the same patterns are used by the corresponding units (see Fig. 4b). Therefore, if a CNN model is deep enough, then RFs with a more stable shape will be induced at the last layer, compared to the RFs of middle layer units.

We carry out a Monte Carlo simulation to examine this property using different kernel arrangements. Given an image $\mathbf{I} \in \mathbb{R}^{W \times H}$, we first define a stochastic matrix $\mathcal{M} \in \mathbb{R}^{W \times H}$. The elements of the matrix are random variables $\mathcal{M}_{i,j} \in [0, 1]$ whose values represent the probability that a pixel $I_{i,j} \in \mathbf{I}$ is *covered* by an RF. Next, we define $\hat{\mathcal{M}} \triangleq \sum_k \mathcal{M}^k_S$ as an average of RFs for a set of kernel arrangements $\{\mathcal{M}^k_S\}^K_{k=1}$. Then, the difference between $\mathcal{M}^k_S$ and the average $\hat{\mathcal{M}}$ is computed using

$$d(\hat{\mathcal{M}}, \mathcal{M}^k_S) = \|\hat{\mathcal{M}} - \mathcal{M}^k_S\|^2_F / (WH), \tag{3}$$

where $\|\cdot\|^2_F$ is the squared Frobenius norm [15]. Note that, we obtain a better approximation to the average RF as the distance decreases. The results are depicted in Fig. 5. The average $\mathbb{E}[d]$ and variance $\mathbb{V}[d]$ show that a better approximation to the average RF is obtained, if kernels used at different layers are integrated at higher depth.

(a) $Depth = 3$, $\mathbb{E}[d] = 0.075$, $\mathbb{V}[d] = 0.0014$.

(b) $Depth = 5$, $\mathbb{E}[d] = 0.061$, $\mathbb{V}[d] = 0.00092$.

(c) $Depth = 7$, $\mathbb{E}[d] = 0.053$, $\mathbb{V}[d] = 0.00069$.

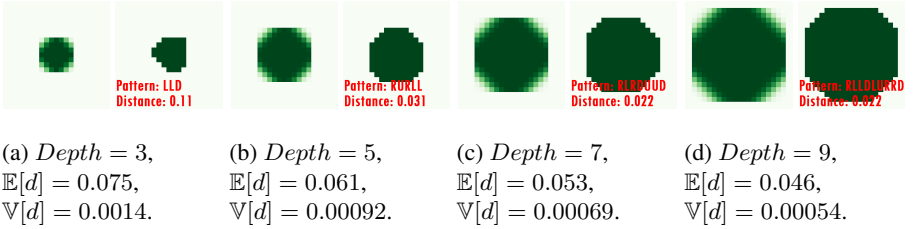(d) $Depth = 9$, $\mathbb{E}[d] = 0.046$, $\mathbb{V}[d] = 0.00054$.

Fig. 5: In (a), (b), (c) and (d), the figures given in left and right show an average shape of kernels emerged from 5000 different shape configurations, and a shape of a kernel designed using a single shape configuration, respectively. It can be seen that the average and variance of $d$ decreases as the kernels are computed at deeper layers. In other words, at deeper layers of CNNs, randomly generated configurations of shapes of kernels can provide better approximations to average shapes of kernels.

## 3 Visualization of regions of interest

We propose a method to visualize the features detected in RFs and the ROI of the image. Following the feature visualization approach suggested in [16], our proposed method provides a saliency map by back-propagating the classification score for a given image and a class. Given a CNN consisting of $L$ layers, the score vector for an input image $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$ is defined as

$$\mathbf{S} = F_1(\mathbf{W}^1, F_2(\mathbf{W}^2, \ldots, F_L(I, \mathbf{W}^L))), \tag{4}$$

where $\mathbf{W}^L$ is the weight vector of the kernel $\mathbf{K}_L$ at the $L^{th}$ layer, and $S^{\mathcal{C}}$ is the $\mathcal{C}^{th}$ element of $\mathbf{S}$ representing the classification score for the $\mathcal{C}^{th}$ class. At the $l^{th}$ layer, we compute a feature map $\mathbf{M}^l$ for each unit $u_{i,j,k}^l \in \mathbf{M}^l$, which takes values from its receptive field $\mathcal{R}(u_{i,j,k}^l)$, and generate a new feature map $\hat{\mathbf{M}}^l$ in which all the units except $u_{i,j,k}^l$ are set to be 0. Then, we feed $\hat{\mathbf{M}}^l$ to the *tail* of the CNN to calculate its score vector as

$$\mathbf{S}(u_{i,j,k}^l) = F_{l+1}(\mathbf{W}^{l+1}, F_{l+2}(\mathbf{W}^{l+2}, \ldots, F_L(\hat{\mathbf{M}}^l, \mathbf{W}^L))). \tag{5}$$

Thereby, we obtain a score map $\mathbb{S}^l$ for all the units of $\mathbf{M}^l$, from which we choose top $N$ most contributed units, i.e. the units with the $N$-highest scores. Then, we back-propagate their score $\mathbf{S}^{\mathcal{C}}(u_{i,j,k}^l)$ for the correct (target) class label towards the forepart of the CNN to rank the contribution of each pixel $p \in \mathbf{I}$ to the score as

$$\mathbb{S}^l(\mathcal{C}, u_{i,j,k}^l) = F_1^{-1}(\mathbf{W}^1, F_2^{-1}(\mathbf{W}^2, \ldots, F_l^{-1}(\mathbf{S}^{\mathcal{C}}(u_{i,j,k}^l), \mathbf{W}^l))), \tag{6}$$

where $\mathbb{S}^l(\mathcal{C}, u_{i,j,k}^l)$ is a score map that has the same dimension with the image $\mathbf{I}$, and that records the contribution of each pixel $p \in \mathbf{I}$ to the $\mathcal{C}^th$ class. Here we choose the top $\Omega$ unit $\{u_\omega^l\}_{\omega=1}^{\Omega}$ with the highest score $\mathbf{S}^{\mathcal{C}}$, where $u_\omega^l$ is the $\omega^{th}$ unit employed at the $l^{th}$ layer. Then, we compute the incorporated saliency map $\mathbf{L}^{\mathcal{C},l} \in \mathbb{R}^{H \times W}$ extracted at the $l^{th}$ layer, for the $\mathcal{C}^{th}$ class as follows

$$\mathbf{L}^{\mathcal{C},l} = \sum_\omega |\mathbb{S}^l(\mathcal{C}, u_\omega^l)|, \tag{7}$$

where $|\cdot|$ is the absolute value function. Finally, the ROI of defined by a set of merged RFs, $\{\mathcal{R}(\mathfrak{u}_\omega^l)\}_{\omega=1}^{\Omega}$ is depicted as a non-zero region in $\mathbf{L}^{\mathcal{C},l}$.

## 4 Experiments

In Sect. 4.1, we examine classification performance of CNNs implemnenting proposed methods using two benchmark datasets, CIFAR-10/100[1] and ILSVRC-2012 (a subset of ImageNet [2]). We first analyze the relationship between shape of kernels, ROI and localization of feature detections on images. Then, we examine the robustness of CNNs for classification of occluded images. Implementation details of the algorithms, and additional results are provided in the supplemental material. We implemented CNN models using the Caffe framework [17], the QH-conv. layer is implemented by utilizing the im2col method to vectorize the inputs and multiply them with corresponding weight matrix.

### 4.1 Classification performance

**Experiments on CIFAR datasets** A list of CNN models used in experiments is given in Table 1a. We used the ConvPool-CNN-C model proposed in [18] as our base model (BASE-A). We employed our method in three different models: i) QH-A retains the structure of the BASE-A by just implementing kernels using the proposed methods, ii) QH-B models a larger number of feature maps compared to QH-A such that QH-B and BASE-A have the same number of parameters, iii) QH-C is a larger model which is used for examination of generalization properties (over/under-fitting) of the proposed QH-models. Following [18] we implement dropout on the input image and at each max pooling layer. We also utilized most of the hyper-parameters suggested in [18] for training the models. We decreased weight decay during the last 100 training epochs to avoid local optima.

Since our proposed kernels have fewer parameters compared to $3\times3$ square shaped kernels, by retaining the same structure as BASE-A, QH-A may benefit from the regularization effects brought by less numbers of total parameters that prevent over-fitting. In order to analyze this regularization property of the proposed method, we implemented a reference model, called BASE-REF with conv-FK (fragmented kernel) layer, which has $3 \times 3$ convolution kernels, and the values of two randomly selected parameters are set to 0 (to keep the number of effective parameters same with quasi-hexagonal kernels). In another reference model (QH-EXT), shape patterns of kernels (Sect. 2) are chosen to be the same ($< R, \ldots, R >$ in this implementation). Moreover, we introduced two additional variants of models using i) different kernel sizes for max pooling (-pool4), and ii) an additional dropout layer before global average pooling (-AD).

Results given in Table 2 show that the proposed QH-A has comparable performance to the base CNN models that employ square shape kernels, despite a smaller number of parameters. Meanwhile, a significant decrement in accuracy appears in the BASE-REF model that employs the same number of parameters as QH-A, which suggests that our proposed model works not only by the employment of a regularization effect but by the utilization of a smaller number of parameters. The inferior performance for QH-EXT

Table 1: CNN configurations. The convolution layer parameters are denoted as <Number of duplication>×conv<kernel>-<number of channels>. A rectified linear unit (ReLU) is followed after each convolution layer. ReLU activation and dropout layer are not shown for brevity. All the conv-3x3/QH/FK layers are set to be stride 1 equipped with pad 1

(a) CNN Configurations - CIFAR

| BASE/BASE-F | QH-A | QH-B/C |
|---|---|---|
| 3×conv-3×3/FK-96 | 3×conv-QH-96 | 3×convH-108/128 |
| maxpool | | |
| 3×conv-3×3/FK-192 | 3×conv-QH-192 | 3×convH-217/256 |
| maxpool | | |
| conv-3×3/FK-192 conv-1×1-192 | convH-192 conv-1×1-192 | conv-QH-217/384 conv-1×1-217/384 |
| conv1-10/100 | | |
| global avepool + soft-max classifier | | |

(b) CNN Configurations - ImageNet

| BASE | QH-BASE | REF-A/B-BASE |
|---|---|---|
| 2×conv-3×3-96 | 2×conv-QH-96 | 2×conv-UB/DIA-96 |
| maxpool | | |
| 2×conv-3×3-192 | 2×conv-QH-192 | 2×conv-UB/DIA-192 |
| maxpool | | |
| 2×conv-3×3-384 | 2×conv-QH-384 | 2×conv-UB/DIA-384 |
| maxpool | | |
| 2×conv-3×3-768 | 2×conv-QH-768 | 2×conv-UB/DIA-768 |
| maxpool | | |
| 2×conv-3×3-1536 | 2×conv-QH-1536 | 2×conv-UB/DIA-1536 |
| maxpool | | |
| conv-3×3-1000 | | |
| conv-1×1-1000 | | |
| global avepool + soft-max classifier | | |

Table 2: Comparison of classification errors using CIFAR-10 dataset (Single models trained without data augmentation)

| Model | Testing Error(%) | Model | Testing Error(%) |
|---|---|---|---|
| BASE-A | 9.02 | BASE-A-pool4 | 8.87 |
| QH-A | 9.10 | QH-A-pool4 | 9.00 |
| BASE-REF | 9.89 | BASE-A-AD | 8.71 |
| QH-EXT | 9.40 | QH-A-AD | 8.79 |

Table 3: Comparison of classification error of models using CIFAR-10/100 datasets (Single models trained without data augmentation)

| Model | Testing Error (%) | | Numbers of Params. |
|---|---|---|---|
| | CIFAR-10 | CIFAR-100 | |
| NIN [19] | 10.41 | 35.68 | $\approx 1M$ |
| DSN [20] | 9.69 | 34.57 | $\approx 1M$ |
| ALL-CNN [18] | 9.08 | 33.71 | $\approx 1.4M$ |
| RCNN [21] | 8.69 | 31.75 | $\approx 1.9M$ |
| Spectral pool [22] | 8.6 | 31.6 | − |
| FMP [23] | − | 31.2 | $\approx 12M$ |
| BASE-A-AD | 8.71 | 31.2 | $\approx 1.4M$ |
| **QH-B-AD** | 8.54 | 30.54 | $\approx 1.4M$ |
| **QH-C-AD** | **8.42** | **29.77** | $\approx 2.4M$ |

model indicates the effectiveness of randomly selecting kernels described in Sect. 2. Moreover, it can also be observed that the implementation of additional dropout and larger size pooling method improves the classification performance of both BASE-A and proposed QH-A in a similar magnitude. Then, the experimental observation implies a general compatibility between the square kernels and the proposed kernels.

Additionally, we compare the proposed methods with state-of-the-art methods for CIFAR-10 and CIFAR-100 datasets. For CIFAR-100, we used the same models implemented for CIFAR-10 with the same hyper-parameters. The results given in Table 4 show that our base model with an additional dropout (BASE-A-AD) provides comparable classification performance for CIFAR-10, and outperforms the state-of-the-art models for CIFAR-100. Moreover, our proposed models (QH-B-AD and QH-C-AD) improve the classification accuracy by adopting more feature maps.

**Experiments on ImageNet** We use an up-scale model of BASE-A model for CIFAR-10/100 as our base model, which stacks 11 convolution layers with kernels that have regular 3×3 square shape, that are followed by a 1×1 convolution layer and a global average pooling layer. Then, we modified the base model with three different types of kernels: i) our proposed quasi-hexagonal kernels (denoted as conv-QH layer), ii) reference kernels where we remove an element located at a corner and one of its adjacent elements located at edge of a standard 3×3 square shape kernel (conv-UB), iii) reference kernels where we remove an element from a corner and an element from a diagonal corner of a standard 3×3 square shape kernel (conv-DIA). Notice that unlike the fragmented kernels we employed in the last experiment, these two reference kernels can also be used to generate aforementioned shapes of RFs. However, unlike the proposed quasi-hexagonal kernels, we cannot assure that these kernels can be used to simulate hexagonal processing. Configurations of the CNN models are given in Table 1a. Dropout [24] is used on an input image at the first layer (with dropout ratio 0.2), and after the last conv-3×3 layer. We employ a simple method for fixing the size of train and test samples to $256 \times 256$ [4], and a patch of $224 \times 224$ is cropped and fed into

Table 4: Comparison of classification performance using validation set of the ILSVRC-2012

| Model | top-1 val.error (%) | top-5 val.error (%) |
|---|---|---|
| BASE | 31.2 | 12.3 |
| QH-BASE | **29.2** | **11.1** |
| REF-A-BASE | 31.4 | 12.4 |
| REF-B-BASE | 31.2 | 12.2 |

network during training. Additional data augmentation methods, such as random color shift [3], are not employed for fast convergence.

Classification results are given in Table 2. Also, histograms of class-wise accuracy values between BASE and QH-BASE models are given in Fig. 6. The results show that the performance of reference models is slightly better than that of the base model. Notice that since the base model is relatively over-fitted (top5 accuracy for training sets is $\geq 97\%$), these two reference models are more likely to be benefited from the regularization effect brought by less number of parameters. Meanwhile, our proposed QH-BASE outperformed all the reference models, implying the validity of the proposed quasi-hexagonal kernels in approximating hexagonal processing. Detailed analyses concerning compactness of models are provided in the next section.

**Analysis of relationship between compactness of models and classification performance** In this section, we analyze the compactness of learned models for ImageNet and CIFAR-10 datasets. First, we provide a comparison of the number of parameters and computational time of the models in Table 5. The results show that, in the experimental analyses for the CIFAR-10 dataset, QH-A model has a comparable performance to the base model with fewer parameters and computational time. If we keep the same number of parameters (QH-B), then classification accuracy improves for similar computational time. Meanwhile, in the experimental analyses for the ImageNet dataset, our proposed model shows significant improvement in both model size and computational time.

We conducted another set of experiments to analyze the relationship between the classification performance and the number of training samples using CIFAR-10 dataset. The results given in Table 6 show that the QH-A-AD model provides a comparable performance with the base model, and the QH-B-AD model provides a better classification accuracy compared to the base model, as the number of training samples decreases. In an extreme case where only 1000 training samples is selected, QH-A-AD and QH-B-AD outperform the base model by 0.7% and 3.1%, respectively, which indicates the effectiveness of the proposed method.

### 4.2   Visualization of regions of interest

Fig. 6 shows some examples of visualizations depicted using our method proposed in Sect. C. Saliency maps are normalized and image contrast is slightly raised to improve

Table 5: Comparison of number of parameters and computational time of different models

| Model | Num. of params. | Training time (500 samples) | Difference in accuracy |
|---|---|---|---|
| BASE | $\approx 57.3M$ | 51610.5 ms | − |
| QH-BASE | $\approx 44.6M$ | 38815.9 ms | +1.2% |
| BASE-A | $\approx 1.4M$ | 1492 ms | − |
| QH-A | $\approx 1.1M$ | 1227.4 ms | −0.08% |
| QH-B | $\approx 1.4M$ | 1449.9 ms | +0.17% |

Table 6: Comparison of classification error between models BASE-A-AD, QH-A-AD and QH-B-AD with different number of training samples on CIFAR-10 dataset

| Model | Classification Error (%) | | | | |
|---|---|---|---|---|---|
| | Number of Training Samples | | | | |
| | 20K | 10K | 5K | 2K | 1K |
| BASE-AD | 12.6 | 16.8 | 21.8 | 31.0 | 44.9 |
| QH-A-AD | 12.7 | 16.6 | 21.1 | 31.3 | 44.2 |
| QH-B-AD | 12.4 | 16.3 | 20.7 | 30.9 | 41.8 |

visualization of images. We observed that for most of these *correctly* classified testing images, both the BASE model equipped with square kernels and the proposed QH-BASE model equipped with quasi-hexagonal kernels are able to present an ROI that roughly specify the location and some basic shape of the target objects, and vise versa. Since the ROI is directly determined by RFs of neurons with strong reactions toward special features, this observation suggests that the relevance between learned representations and target objects is crucial for recognition and classification using large-scale datasets of natural images such as ImageNet.

However, some obvious difference between the ROI of the base model and the proposed model can be observed: i) ROI of the base model usually involves more background than that of the proposed model. That is, compared to these pixels with strong contributions, the percentage of these pixels that are not essentially contributing to the classification score, is generally higher in the base model. ii) Features learned using the square kernels are more like to be detected within clusters on special parts of the objects. The accumulation of the features located in these clusters results in a superior contribution, compared to the features that are scattered on the images. For instance, in the base model, more neurons have their RFs located in the heads of hare and parrots, thus the heads obtain higher classification scores than other parts of body. iii) As a result of ii), some duplicated important features (e.g, the supporting parts of cart and seats of coach) are overlooked in these top reacted high-level neurons in the base model. Meanwhile, our proposed model with quasi-hexagonal kernels is more likely to obtain *discriminative* features that are spatially distributed on the whole object. In order to
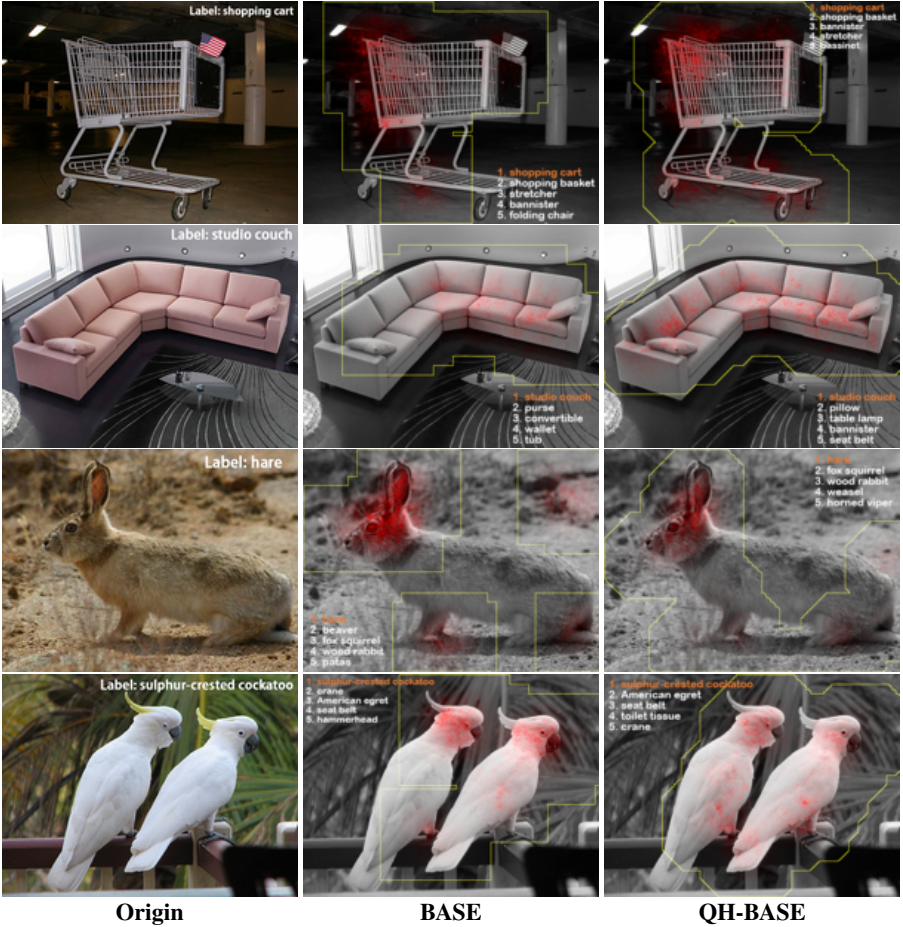
Fig. 6: Examples of visualization of ROI. A ROI demonstrates a union of RFs of the top 40 activated neurons at the last max pooling layer. The pixels marked with red color indicate their contribution to classification score, representing the activated features located at them. Borderlines of ROI are represented using yellow frames. Top 5 class predictions provided by the models are also given, and the correct (target) class is given using orange color.

further analyze the results obtained by employing the square kernel and the proposed kernels for object recognition, we provide a set of experiments using occluded images in the next section.

## 4.3   Occlusion and spatially distributed representations

The analyses given in the last section imply that the base CNN models equipped with the square kernel could be vulnerable to recognition of objects in occluded scenes, which is a very common scenario in computer version tasks. In order to analyze the robustness of

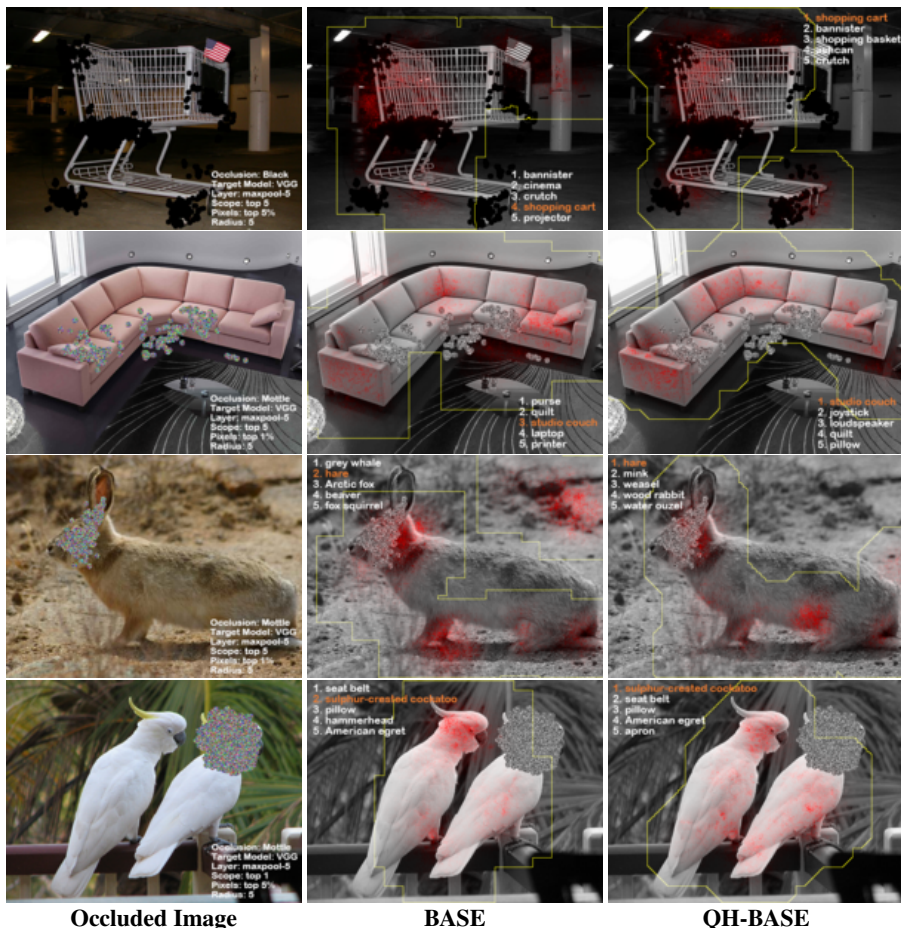| **Occluded Image** | **BASE** | **QH-BASE** |

Fig. 7: Analysis of robustness of different models to occlusion. We use the same proposed method to select neurons and visualize their RFs for each model (see Sect. C). The comparison between the ROI shown in Fig. 6 suggests that the proposed model overcomes the occlusion by detecting features that are spatially distributed on target objects. It can also be seen that, the classification accuracy of the base model is decreased although the ROI of the base model seems to be more adaptive to the shape of objects. This also suggests that the involvement of background may make the CNNs hard to discriminate background from useful features.

the methods to partial occlusion of images, we prepare a set of locally occluded images using the following methods. i) We randomly select 1249 images that are correctly classified by both the base and proposed models using the validation set of ILSVRC-2012 [2]. ii) We select Top1 or Top5 elements with highest classification score at the last maxpool layers of a selected model[1] and calculate the ROI defined by their RFs, as

---

[1] In addition to the BASE and the QH-BASE models, we also employ a "third-party" model, namely VGG [4], to generate the occluded images.

Table 7: Performances on the occlusion datasets. Each column shows the classification accuracy (%) of test models in different occlusion conditions. In the first row, BASE/QH-BASE/VGG indicate the models used for generating occlusion, Top1/Top5 indicate the numbers of selected neurons that control the size of occluded region, Bla./Mot. indicate the patterns of occlusion

| Model | BASE Top1 Bla. | BASE Top1 Mot. | BASE Top5 Bla. | BASE Top5 Mot. | QH-BASE Top1 Bla. | QH-BASE Top1 Mot. | QH-BASE Top5 Bla. | QH-BASE Top5 Mot. | VGG Top1 Bla. | VGG Top1 Mot. | VGG Top5 Bla. | VGG Top5 Mot. | Average accuracy |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BASE | 58.8 | 61.2 | 34.6 | 40.9 | 61.3 | 63.5 | 36.3 | 42.7 | 61.7 | 63.8 | 44.1 | 48.3 | 51.5 |
| QH-BASE | 67.1 | 67.8 | 43.8 | 47.6 | 67.0 | 66.9 | 42.2 | 45.4 | 68.6 | 69.1 | 52.3 | 54.6 | **57.7** |

we described in Sect. C. iii) Within the ROI, we choose 1-10% of pixels that provide the most contribution, and then occlude each of the selected pixels with a small circular occlusion mask (with radius $r = 5$ pixels), which is filled by black (Bla.) or randomly generated colors (Mot.) drawn from a uniform distribution. In total, we generate 120 different occlusion datasets (149880 different occluded images in total), Table 7 shows the classification accuracy on the occluded images. The results show that our proposed quasi-hexagonal kernel model reveal better robustness in this object recognition under targeted occlusion task compared to square kernel model. Some sample images are shown in Fig. 7.

## 5   Conclusion

In this work, we analyze the effects of shapes of convolution kernels on feature representations learned in CNNs and classification performance. We first propose a method to design the shape of kernels in CNNs. We then propose a feature visualization method for visualization of pixel-wise classification score maps of learned features. It is observed that the compact representations obtained using the proposed kernels are beneficial for the classification accuracy. In the experimental analyses, we obtained state-of-the-art performance using ImageNet and CIFAR datasets. Moreover, our proposed methods enable us to implement CNNs with less number of parameters and computational time compared to the base-line CNN models. Additionally, the proposed method improves the robustness of the base-line models to occlusion for classification of partially occluded images. These results confirm the effectiveness of the proposed method for designing of the shape of convolution kernels in CNNs for image classification. In future work, we plan to apply the proposed method to perform other tasks such as object detection and segmentation.

# A   Introduction

In this document, we provide the supplemental material for the paper "Design of Kernels in Convolutional Neural Networks for Image Classification". In the next section, implementation details of the algorithms proposed and employed in the main text are given. In Section C, additional results for visualization of receptive fields are provided.

# B   Implementation detail

## B.1   CNN models implemented for CIFAR

In this subsection, implementation details of the algorithms and models employed in Sect. 4.1.1 of the main text are given.

In a training phase, we optimise a soft-max loss function at the top layer of a CNN model using stochastic gradient descent with mini-batch 128, and a momentum [**?**] of 0.9 is used. All the models are regularized by weight decay (L2 penalty) with multiplier 0.001 initially. For the QH-models, we decrease the multiplier during the final 100 training epochs to avoid local optima. We also regularize all the models using dropout; dropout with ratio 0.2 is employed for input data, and dropout with ratio 0.5 is employed for each maxpool layer. The learning rate is initially set to $5 \times 10^2$, and then decreased by a factor of 10 after 120, 170, and 220 training epochs. The learning algorithm was stopped after 270 epochs with a final learning rate $5 \times 10^{-5}$.

We apply the global contrast normalization and ZCA whitening which were implemented by Goodfellow *et al.* in the maxout network [**?**], and no further data augmentation is employed for both training and testing images.

## B.2   CNN models implemented for Imagenet

In this subsection, implementation details of the algorithms and models employed in Sect. 4.1.2 of the main text are given.

In a training phase, we optimise a soft-max loss function at the top layer of a CNN model using stochastic gradient descent with mini-batch size 192, and a momentum [**?**] of 0.9 is used. In order to regularize these models, we implement weight decay (L2 penalty) and dropout. For VGG and QH-VGG, weight decay multipliers are set to 0.0005, and dropout with ratio 0.5 is employed for the first two fully connected (fc) layers. For QH-GAP, a weight decay multiplier is set to 0.0001, dropout with ratio 0.5 and 0.2 are employed for the conv3-1000 layer and input data, respectively. The learning rate is initially set to $10^2$, and then decreased by a factor of 10 when the improvement of a validation set accuracy is stopped. The training was stopped after 65 epochs with a final learning rate $10^{-5}$.

Additionally, in a training phase of a CNN, a training image is first resized to a fixed 256×256, then a 224×224 piece is randomly cropped and mirrored as the CNN receives an input image at each iteration. Further augmentation of training data such as random RGB color shift [3] is not employed for a fast convergence. During testing, all the testing images are resized to 256×256, and fed to the CNNs without cropping. Fc

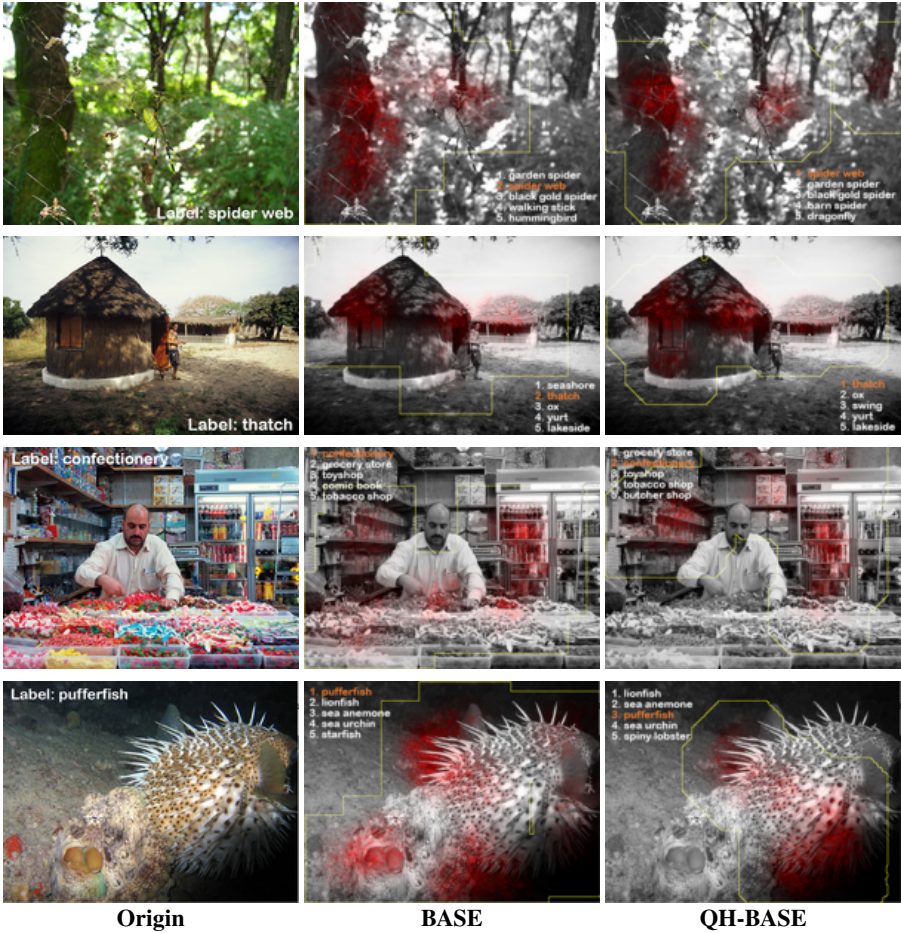**Origin**                    **BASE**                    **QH-BASE**

Fig. 8: Additional results for visualization of RFs. See Sect. 4.2 and Figure 6 in the main text for details.

layers of VGG and QH-VGG models are converted into convolution layers with kernel size $7 \times 7$, hence we could obtain a class score map whose number of channels is equal to the number of classes. Then, the score map is channel-wise average pooled, and fed into the soft-max classifier. We augment the test images by mirroring, and the final score for a test image is averaged from the original and mirrored images.

## C   Visualization of regions of interest

In this subsection, we provide additional results obtained using our proposed visualization method given in Sect. 3 and employed in Sect. 4.2 of the main text.

Additional visualization results of receptive fields are shown in Figure 8, Figure 9, Figure 10.
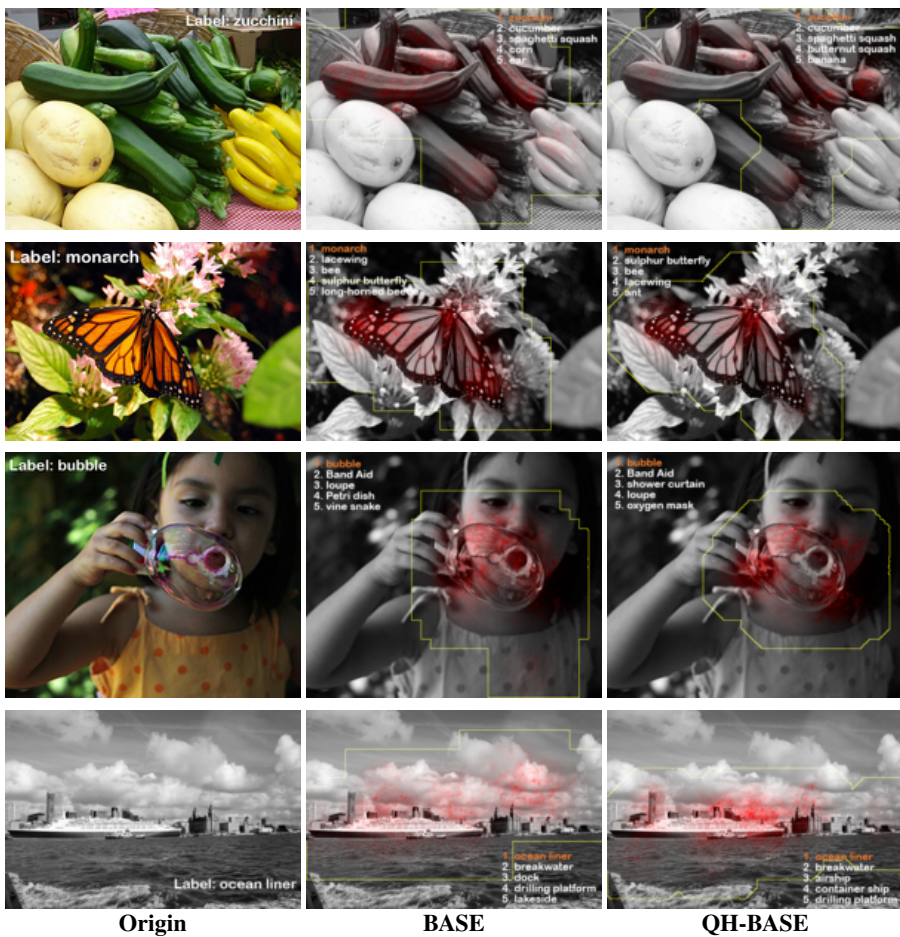
Fig. 9: Additional results for visualization of RFs. See Sect. 4.2 and Figure 6 in the main text for details.
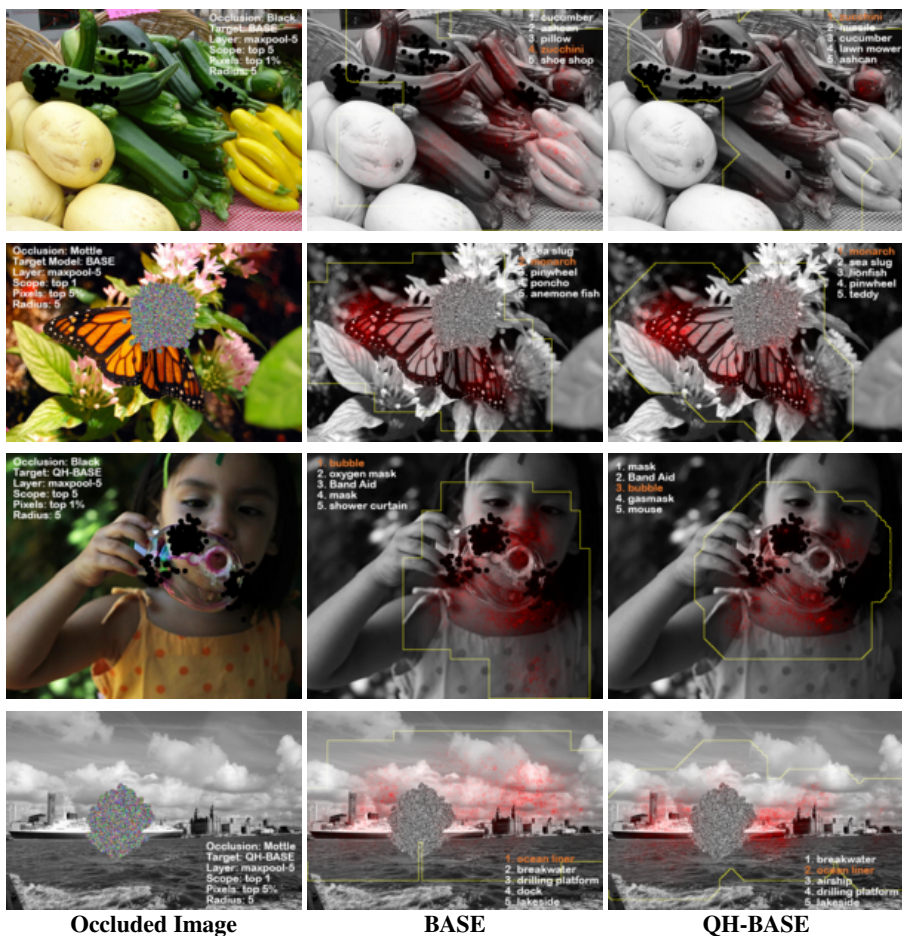
Fig. 10: Additional results for visualization of occluded images. See Sect. 4.3 and Figure 7 in the main text for details.

# References

1. Krizhevsky, A.: Learning multiple layers of features from tiny images (2009) 1, 7
2. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. International Journal of Computer Vision (IJCV) (April 2015) 1–42 1, 7, 13
3. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In Pereira, F., Burges, C., Bottou, L., Weinberger, K., eds.: Advances in Neural Information Processing Systems 25. Curran Associates, Inc. (2012) 1097–1105 1, 3, 4, 10, 15
4. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: Proceedings of ICLR. (2015) 1, 4, 9, 13
5. Hubel, D.H., Wiesel, T.N.: Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. The Journal of Physiology **160**(1) (1962) 106–154 2
6. Mutch, J., Lowe, D.: Object class recognition and localization using sparse features with limited receptive fields. International Journal of Computer Vision **80**(1) (2008) 45–57 2
7. Simoncelli, E.P., Olshausen, B.A.: Natural image statistics and neural representation. Annu Rev Neurosci **24** (2001) 1193–1216 2
8. Liu, Y.S., Stevens, C.F., Sharpee, T.O.: Predictable irregularities in retinal receptive fields. PNAS **106**(38) (2009) 16499–16504 2
9. Kerr, D., Coleman, S., McGinnity, T., Wu, Q., Clogenson, M.: A novel approach to robot vision using a hexagonal grid and spiking neural networks. In: The International Joint Conference on Neural Networks (IJCNN). (June 2012) 1–7 2
10. Mersereau, R.: The processing of hexagonally sampled two-dimensional signals. Proceedings of the IEEE **67**(6) (June 1979) 930–949 2
11. Lecun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. (1998) 2278–2324 4
12. Han, S., Pool, J., Tran, J., Dally, W.J.: Learning both weights and connections for efficient neural networks. CoRR **abs/1506.02626** (2015) 4
13. Gonzalez, R.C., Woods, R.E.: Digital Image Processing (3rd Edition). Prentice-Hall, Inc., Upper Saddle River, NJ, USA (2006) 5
14. Reinhard Klette, A.R.: Digital Geometry: Geometric Methods for Digital Picture Analysis. Morgan Kaufmann, San Francisco (2004) 5
15. Hartley, R., Zisserman, A.: Multiple View Geometry in Computer Vision. Cambridge University Press (2004) 5
16. Simonyan, K., Vedaldi, A., Zisserman, A.: Deep inside convolutional networks: Visualising image classification models and saliency maps. In: Proceedings of the International Conference on Learning Representations (ICLR). (2014) 6
17. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding. arXiv preprint arXiv:1408.5093 (2014) 7
18. Springenberg, J., Dosovitskiy, A., Brox, T., Riedmiller, M.: Striving for simplicity: The all convolutional net. In: ICLR (workshop track). (2015) 7, 9
19. Lin, M., Chen, Q., Yan, S.: Network in network. In: Proceedings of ICLR. (2014) 9
20. Lee, C.Y., Xie, S., Gallagher, P., Zhang, Z., Tu, Z.: Deeply-Supervised Nets. ArXiv e-prints (September 2014) 9
21. Liang, M., Hu, X.: Recurrent convolutional neural network for object recognition. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (June 2015) 9
22. Rippel, O., Snoek, J., Adams, R.P.: Spectral Representations for Convolutional Neural Networks. ArXiv e-prints (June 2015) 9

23. Graham, B.: Fractional max-pooling. CoRR **abs/1412.6071** (2014) 9
24. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. Journal of Machine Learning Research **15** (2014) 1929–1958 9