# Semantic Labeling in Very High Resolution Images via A Self-Cascaded Convolutional Neural Network

Yongcheng Liu[a,b], Bin Fan[a,*], Lingfeng Wang[a], Jun Bai[c], Shiming Xiang[a], Chunhong Pan[a]

[a] *National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, P.R. China*
[b] *School of Computer and Control Engineering, University of Chinese Academy of Sciences, Beijing, P.R. China*
[c] *Research Center for Brain-inspired Intelligence, Institute of Automation, Chinese Academy of Sciences, Beijing, P.R. China*

## Abstract

Semantic labeling for very high resolution (VHR) images in urban areas, is of significant importance in a wide range of remote sensing applications. However, many confusing manmade objects and intricate fine-structured objects make it very difficult to obtain both coherent and accurate labeling results. For this challenging task, we propose a novel deep model with convolutional neural networks (CNNs), i.e., an end-to-end self-cascaded network (ScasNet). Specifically, for confusing manmade objects, ScasNet improves the labeling coherence with sequential global-to-local contexts aggregation. Technically, multi-scale contexts are captured on the output of a CNN encoder, and then they are successively aggregated in a self-cascaded manner. Meanwhile, for fine-structured objects, ScasNet boosts the labeling accuracy with a coarse-to-fine refinement strategy. It progressively refines the target objects using the low-level features learned by CNN's shallow layers. In addition, to correct the latent fitting residual caused by multi-feature fusion inside Scas-Net, a dedicated residual correction scheme is proposed. It greatly improves the effectiveness of ScasNet. Extensive experimental results on three public datasets, including two challenging benchmarks, show that ScasNet achieves the state-of-the-art performance.

*Keywords:* Semantic labeling, Convolutional neural networks (CNNs), Multi-scale contexts, End-to-end.

## 1. Introduction

Semantic labeling in very high resolution (VHR) images is a long-standing research problem in remote sensing field. It plays a vital role in many important applications, such as infrastructure planning, territorial planning and urban change detection (Lu et al., 2017a; Matikainen and Karila, 2011; Zhang and Seto, 2011). The target of this problem is to assign each pixel to a given object category. Note that it is not just limited to building extraction (Li et al., 2015a), road extraction (Cheng et al., 2017b) and vegetation extraction (Wen et al., 2017) which only consider labeling one single category, semantic labeling usually considers several categories simultaneously (Li et al., 2015b; Xu et al., 2016; Xue et al., 2015). As a result, this task is very challenging, especially for the urban areas, which exhibit high diversity of manmade objects. Specifically, on one hand, many manmade objects (e.g., buildings) show various structures, and they are composed of a large number of different materials. Meanwhile, plenty of different manmade objects (e.g., buildings and roads) present much similar visual characteristics. These confusing manmade objects with high intra-class variance and low inter-class variance bring much difficulty for coherent labeling. On the other hand, fine-structured objects in cities (e.g., cars, trees and low vegetations) are quite small or threadlike, and they also interact with

---

*Corresponding author at: National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, P.R. China. E-mail address: bfan@nlpr.ia.ac.cn (B. Fan).

each other through occlusions and cast shadows. These factors always lead to inaccurate labeling results. Furthermore, it poses additional challenge to simultaneously label all these size-varied objects well.

To accomplish such a challenging task, features at different levels are required. Specifically, abstract high-level features are more suitable for the recognition of confusing manmade objects, while labeling of fine-structured objects could benefit from detailed low-level features. Convolutional neural networks (CNNs) (Lecun et al., 1990) in *deep learning* field are well-known for feature learning (Mas and Flores, 2008). CNNs consist of multiple trainable layers which can extract expressive features of different levels (Lecun et al., 1998). Moreover, recently, CNNs with *deep learning* have demonstrated remarkable learning ability in computer vision field, such as scene recognition (Yuan et al., 2015) and image segmentation (Long et al., 2015). Meanwhile, the development of remote sensing has also been greatly promoted by numerous CNNs-based methods (Cheng et al., 2017a). For example, deconvolution networks (Zeiler et al., 2010) are investigated by Lu et al. (Lu et al., 2017b) for remote sensing scene classification, and Chen et al. (Chen et al., 2016b) perform target classification using CNNs for SAR Images.

Based on CNNs, many patch-classification methods are proposed to perform semantic labeling (Mnih, 2013; Mostajabi et al., 2015; Paisitkriangkrai et al., 2016; Nogueira et al., 2016; Alshehhi et al., 2017; Zhang et al., 2017). These methods determine a pixel's label by using CNNs to classify a small patch around the target pixel. However, they are far from optimal, because they ignore the inherent relationship between patches and their time consumption is huge (Maggiori et al., 2017). Typically, fully convolutional networks (FCNs) have boosted the accuracy of semantic labeling a lot (Long et al., 2015; Sherrah, 2016). FCNs perform pixel-level classification directly and now become the normal framework for semantic labeling. Nevertheless, due to multiple *sub-samplings* in FCNs, the final *feature maps* are much coarser than the input image, resulting in less accurate labeling results.

Accordingly, a tough problem locates on how to perform accurate labeling with the coarse output of FCNs-based methods, especially for fine-structured objects in VHR images. To solve this problem, some researches try to reuse the low-level features learned by CNNs' shallow layers (Zeiler and Fergus, 2014). The aim is to utilize the local details (e.g., corners and edges) captured by the *feature maps* in fine resolution. Technically, they perform operations of multi-level feature fusion (Ronneberger et al., 2015; Long et al., 2015; Hariharan et al., 2015; Pinheiro et al., 2016), *deconvolution* (Noh et al., 2015) or *up-pooling* with recorded *pooling* indices (Badrinarayanan et al., 2015). Most of these methods use the strategy of direct stack-fusion. However, this strategy ignores the inherent semantic gaps in features of different levels. An alternative way is to impose boundary detection (Bertasius et al., 2016; Marmanis et al., 2016). It usually requires extra boundary supervision and leads to extra model complexity despite boosting the accuracy of object localization.

Another tricky problem is the labeling incoherence of confusing objects, especially of the various manmade objects in VHR images. To tackle this problem, some researches concentrate on leveraging the multi-context to improve the recognition ability of those objects. They use multi-scale images (Farabet et al., 2013; Mostajabi et al., 2015; Cheng et al., 2016; Liu et al., 2016b; Chen et al., 2016a; Zhao and Du, 2016) or multi-region images (Gidaris and Komodakis, 2015; Luus et al., 2015) as input to CNNs. However, these methods are usually less efficient due to a lot of repetitive computation. Differently, some other researches are devoted to acquire multi-context from the inside of CNNs. They usually perform operations of multi-scale *dilated convolution* (Chen et al., 2015), multi-scale pooling (He et al., 2015b; Liu et al., 2016a; Bell et al., 2016) or multi-kernel convolution (Audebert et al., 2016), and then fuse the acquired multi-scale contexts in a direct stack manner. Nevertheless, this manner not only ignores the hierarchical dependencies among the objects and scenes in different scales, but also neglects the inherent semantic gaps in contexts of different-level information.

In summary, although current CNN-based methods have achieved significant breakthroughs in semantic labeling, it is still difficult to label the VHR images in urban areas. The reasons are as follows: 1) Most existing approaches are less efficient to acquire multi-scale contexts for confusing manmade objects recognition; 2) Most existing strategies are less effective to utilize low-level features for accurate labeling, especially for
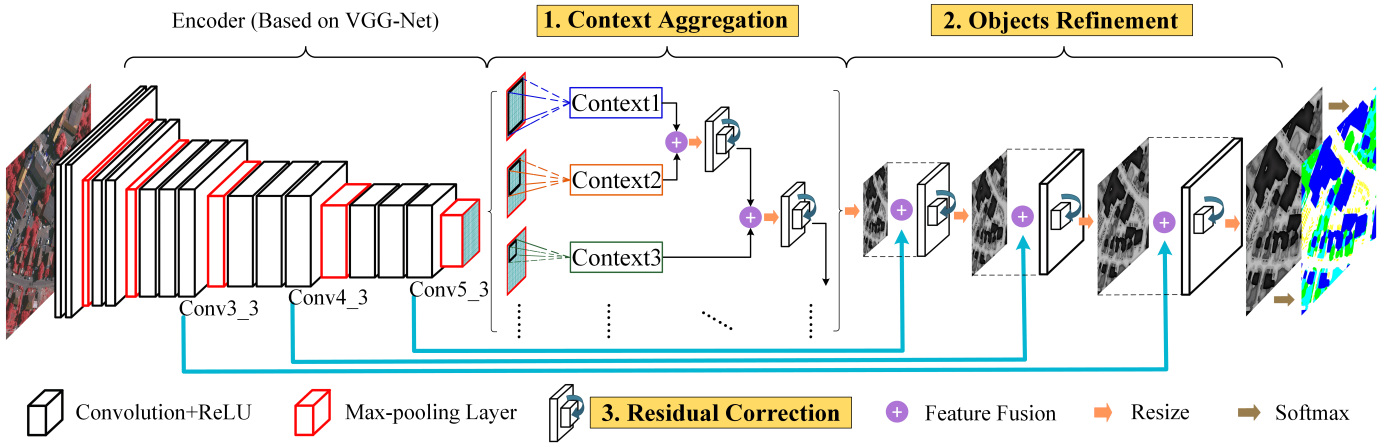
Figure 1: Overview of the proposed ScasNet. (Best viewed in color)

fine-structured objects; 3) Simultaneously fixing the above two issues with a single network is particularly difficult due to a lot of fitting residual in the network, which is caused by semantic gaps in different-level contexts and features.

In this paper, we propose a novel self-cascaded convolutional neural network (ScasNet), as illustrated in Fig. 1. The aim of this work is to further advance the state of the art on semantic labeling in VHR images. To this end, it is focused on three aspects: 1) multi-scale contexts aggregation for distinguishing confusing manmade objects; 2) utilization of low-level features for fine-structured objects refinement; 3) residual correction for more effective multi-feature fusion. Specifically, a conventional CNN is adopted as an encoder to extract features of different levels. On the *feature maps* outputted by the encoder, global-to-local contexts are sequentially aggregated for confusing manmade objects recognition. Technically, multi-scale contexts are first captured by different convolutional operations, and then they are successively aggregated in a self-cascaded manner. With the acquired contextual information, a coarse-to-fine refinement strategy is performed to refine the fine-structured objects. It progressively reutilizes the low-level features learned by CNN's shallow layers with long-span connections. In addition, to correct the latent fitting residual caused by semantic gaps in multi-feature fusion, several residual correction schemes are employed throughout the network. As a result of residual correction, the above two different solutions could work collaboratively and effectively when they are integrated into a single network. Extensive experiments demonstrate the effectiveness of ScasNet. Moreover, the three submodules in ScasNet could not only provide good solutions for semantic labeling, but are also suitable for other tasks such as object detection (Cheng and Han, 2016) and change detection (Zhang et al., 2016; Gong et al., 2017), which will no doubt benefit the development of the remote sensing deep learning techniques.

To sum up, the main contributions of this paper can be highlighted as follows:

- A self-cascaded architecture is proposed to successively aggregate contexts from large scale to small ones. In this way, global-to-local contexts with hierarchical dependencies among the objects and scenes are well retained, resulting in coherent labeling results of confusing manmade objects.

- A coarse-to-fine refinement strategy is proposed, which progressively refines the target objects using the low-level features learned by CNN's shallow layers. Thus, accurate labeling results can be achieved, especially for the fine-structured objects.

- A residual correction scheme is proposed to correct the latent fitting residual caused by semantic gaps in multi-feature fusion. It greatly improves the effectiveness of the above two different solutions.

- All the above contributions constitute a novel end-to-end *deep learning* framework for semantic labelling, as shown in Fig. 1. It achieves the state-of-the-art performance on two challenging bench-

marks by the date of submission: *ISPRS 2D Semantic Labeling Challenge* (ISPRS, 2016) for Vaihingen and Potsdam. Furthermore, these results are obtained using only image data with a single model, without using the elevation data like the Digital Surface Model (DSM), model ensemble strategy or any postprocessing.

A shorter version of this paper appears in (Liu et al., 2017). Apart from extensive qualitative and quantitative evaluations on the original dataset, the main extensions in the current work are:

- More comprehensive and elaborate descriptions about the proposed semantic labeling method.
- Further performance improvement by the modification of network structure in ScasNet.
- Comparative experiments with more state-of-the-art methods on another two challenging datasets for further support the effectiveness of ScasNet.
- More detailed and in-depth analyses, as well as model visualization and complexity analyses of ScasNet.

The remainder of this paper is arranged as follows. The basic modules used in ScasNet are briefly introduced in Section 2. Section 3 presents the details of the proposed semantic labeling method. Experimental evaluations between our method and the state-of-the-art methods, as well as detailed analyses of ScasNet are provided in Section 4. Finally, the conclusion is outlined in Section 5.

## 2. Preliminaries

CNNs (Lecun et al., 1990) are multilayer neural networks that can hierarchically extract powerful low-level and high-level features. The input and output of each layer are sets of arrays called *feature maps*. Commonly, a standard CNN contains three kinds of layers: convolutional layer, nonlinear layer and pooling layer. The convolutional layer offers filter-like function to generate convoluted *feature maps*, while the nonlinear layer simply consists of an elementwise nonlinear activation function applied to each value in the *feature maps*. The pooling layer generalizes the convoluted features into higher level, which makes features more abstract and robust. Meanwhile, in CNNs, the feature extraction module and the classifier module are integrated into one framework, thus the extracted features are more suitable for specific task than hand-crafted features, such as HOG (Dalal and Triggs, 2005), SIFT (Lowe, 2004), and spectral features in remote sensing (Zhang et al., 2012).

In the following, each basic layer used in the proposed network will be introduced, and their specific configurations will be presented in Section 3.4.

**Convolutional Layer:** The convolutional (Conv) layer performs a series of convolutional operations on the previous layer with a small kernel (e.g., $3 \times 3$). The output of each convolutional operation is computed by dot product between the weights of the kernel and the corresponding local area (*local receptive field*). A *weight sharing* technique that the parameters (i.e., weights and bias) are shared among each kernel across an entire *feature map*, is adopted to reduce parameters in great deal (Rumelhart et al., 1986).

**Batch Normalization Layer:** Batch normalization (BN) mechanism (Ioffe and Szegedy, 2015) normalizes layer inputs to a Gaussian distribution with zero-mean and unit variance, aiming at addressing the problem of *internal covariate shift*, i.e., the distribution of each layer's inputs changes during training, as the parameters of the previous layers change. Thus, it allows us to use much higher learning rate.

**ReLU Layer:** The rectified linear unit (ReLU) (Glorot et al., 2011; Nair and Hinton, 2010) is usually chosen as the nonlinearity layer. It thresholds the non-positive value as zero and keeps the positive value unchanged, i.e., an elementwise activation as $\max(0, x)$. ReLU can achieve a considerable reduction in training time (Krizhevsky et al., 2012).

**Pooling Layer:** Pooling is a way to perform *sub-sampling* along the spatial dimension. Commonly, there are two kinds of pooling: max-pooling and ave-pooling. Max-pooling samples the maximum in the region to be pooled, while ave-pooling computes the mean value. In our network, we use max-pooling.

**Dropout Layer:** Dropout (Srivastava et al., 2014) is an effective regularization technique to reduce overfitting. It randomly drops units (along with their connections) from the neural network during training, which prevents units from co-adapting too much.

**Interpolation Layer:** Interpolation (Interp) layer performs resizing operation along the spatial dimension. In our network, we use bilinear interpolation.

**Elementwise Layer:** Elementwise (Eltwise) layer performs elementwise operations on two or more previous layers, in which the *feature maps* must be of the same number of channels and the same size. There are three kinds of elementwise operations: *product, sum, max*. In our network, we use *sum* operation.

**Softmax Layer:** The softmax nonlinearity (Bridle, 1989) is applied to the output layer in the case of multiclass classification. It outputs the posterior probabilities over each category.

## 3. Self-cascaded Convolutional Neural Network (ScasNet)

Semantic labeling also called pixel-level classification, is aimed at obtaining all the pixel-level categories in an entire image. For this task, we have to predict the most likely category $\hat{k}$ for a given image $x$ at $j$-th pixel $x^j$, which is given by

$$\hat{k} = \underset{k \in \mathcal{C}}{\mathrm{argmax}}\, p_k(x^j|\boldsymbol{\theta}),\ \forall\, j \in \{1, \cdots, N\}, \tag{1}$$

where $p_k(x^j|\boldsymbol{\theta})$, estimated by a model with parameters $\boldsymbol{\theta}$, denotes the posterior probability of $x^j$ belonging to the $k$-th category in a set of categories $\mathcal{C} = \{1, \cdots, K\}$. $K$ is the number of categories and $N$ is the number of pixels in the given image.

In this work, we perform semantic labeling for VHR images in urban areas by means of a self-cascaded convolutional neural network (ScasNet), which is illustrated in Fig. 1. In the following, we will describe five important aspects of ScasNet, including 1) *Multi-scale contexts Aggregation*, 2) *Fine-structured Objects Refinement*, 3) *Residual Correction*, 4) *ScasNet Configuration*, 5) *Learning and Inference Algorithm*.

### 3.1. Multi-scale contexts Aggregation

Obtaining coherent labeling results for confusing manmade objects in VHR images is not easily accessible, because they are of high intra-class variance and low inter-class variance. To fix this issue, it is insufficient to use only the very local information of the target objects. We need to know the scene information around them, which could provide much wider visual cues to better distinguish the confusing objects. The scene information also means the context, which characterizes the underlying dependencies between an object and its surroundings, is a critical indicator for objects identification. Therefore, we are interested in discussing how to efficiently acquire context with CNNs in this Section.

In CNNs, each unit of deeper layers (*feature maps*) contains more extensive, powerful and abstract information, due to the larger *receptive field* on the input image and higher nonlinearity (Zeiler and Fergus, 2014). Thus, the context acquired from deeper layers can capture wider visual cues and stronger semantics simultaneously. However, only single-scale context may not represent hierarchical dependencies between an object and its surroundings. Naturally, multi-scale contexts are gaining more attention. However, it is very hard to retain the hierarchical dependencies in contexts of different scales using common fusion strategies (e.g., direct stack). To address this issue, we propose a novel self-cascaded architecture, as shown in the middle part of Fig. 1. It is aimed at aggregating global-to-local contexts while well retaining hierarchical dependencies, i.e., the underlying inclusion and location relationship among the objects and scenes in different scales (e.g., the car is more likely on the road, the chimney and skylight is more likely a part of roof and the roof is more likely by the road).
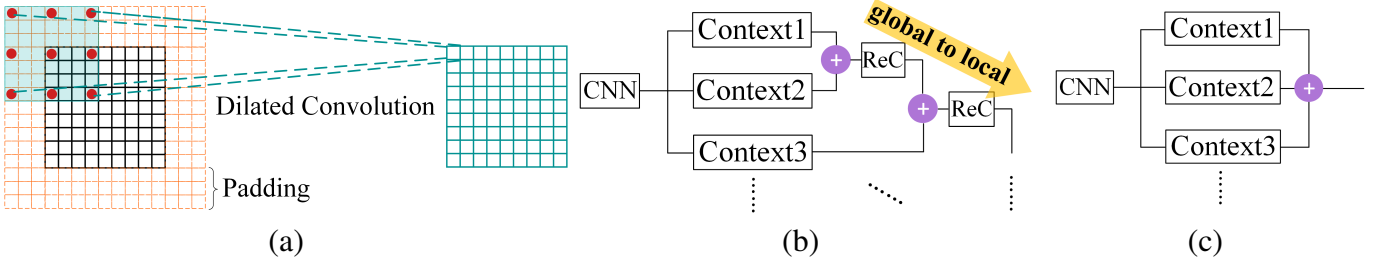
Figure 2: (a) An illustration of *dilated convolution* used in ScasNet to capture context, where the size of *feature map* and *convolution kernel* is $9 \times 9$ and $3 \times 3$, respectively, both the *dilation rate* and the *padding rate* equal 3, and the *padding value* is zero. (b) The proposed multi-context aggregation approach, i.e., performing aggregation sequentially in a self-cascaded manner. (c) Multi-context aggregation in a parallel stack. 'ReC' denotes the proposed residual correction scheme. (Best viewed in color)

Specifically, we perform *dilated convolution* operation on the last layer of the encoder to capture context. The reasons are two-fold. On one hand, dilated convolution expands the *receptive field*, which can capture high-level semantics with wider information. On the other hand, although theoretically, features from high-level layers of a network have very large *receptive fields* on the input image, in practice they are much smaller (Zhou et al., 2015). This problem can be alleviated by *dilated convolution*. Fig. 2(a) illustrates an example of *dilated convolution*. To make the size of *feature map* after *dilated convolution* unchanged, the *padding rate* should be set as the same to the *dilation rate*. More details about *dilated convolution* can be referred in (Yu and Koltun, 2016).

Then, by setting a group of big-to-small *dilation rates* (24, 18, 12 and 6 in the experiment), a series of *feature maps* with global-to-local contexts are generated [1]. That is, multi-scale *dilated convolution* operations correspond to multi-size regions on the last layer of encoder (see Fig. 1). Large region (high-level context) contains more semantics and wider visual cues, while small region (low-level context) otherwise. Meanwhile, the obtained *feature maps* with multi-scale contexts can be aligned automatically due to their equal resolution.

To well retain the hierarchical dependencies in multi-scale contexts, we sequentially aggregate them from global to local in a self-cascaded manner as shown in Fig. 2(b). In this way, high-level context with big *dilation rate* is aggregated first and low-level context with small *dilation rate* next. Formally, it can be described as:

$$\begin{cases} T = \Upsilon\Big[ \cdots \Upsilon\big[ \Upsilon[T_1 \oplus T_2] \oplus T_3 \big] \oplus \cdots \oplus T_n \Big], \\ d_{T_1} > d_{T_2} > d_{T_3} > \cdots > d_{T_n}. \end{cases} \quad (2)$$

Here, $T_1, T_2, \cdots, T_n$ denote $n$-level contexts, $T$ is the final aggregated context and $d_{T_i}$ $(i = 1, \ldots, n)$ is the *dilation rate* set for capturing the context $T_i$. '$\oplus$' denotes the fusion operation. $\Upsilon[\cdot]$ denotes the residual correction process, which will be described in Section 3.3. In fact, the above aggregation rule is consistent with the visual mechanism, i.e., wider visual cues in high-level context could play a guiding role in integrating low-level context. For instance, the visual impression of a whole roof can provide strong guidance for the recognition of chimney and skylight in this roof.

The proposed self-cascaded architecture for multi-scale contexts aggregation has several advantages: 1) The multiple contexts are acquired from deep layers in CNNs, which is more efficient than directly using multiple images as input (Gidaris and Komodakis, 2015); 2) Besides the hierarchical visual cues, the acquired contexts also capture the abstract semantics learned by CNN, which is more powerful for confusing objects recognition; 3) The self-cascaded strategy of sequentially aggregating multi-scale contexts, is more

---

[1]Due to the inherent properties of convolutional operation in each single-scale context (same-scale convolution kernels with large original receptive fields convolve with weight sharing over spatial dimension and summation over channel dimension), the relationship between contexts with same scale can be acquired implicitly.
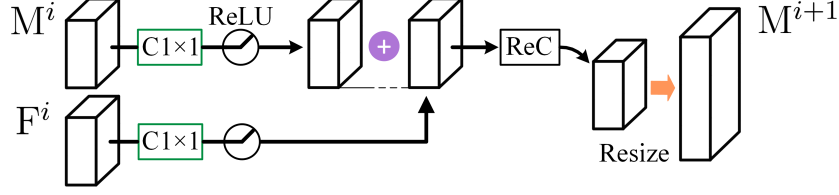
Figure 3: Single process of refinement. 'C1×1' denotes convolutional operation with *kernel size* 1×1, 'ReC' denotes the residual correction scheme. (Best viewed in color)

effective than the parallel stacking strategy (Chen et al., 2015; Liu et al., 2016a), as shown in Fig. 2(c), which potentially loses the hierarchical dependencies in different scales; 4) The more complicated nonlinear operation of Eq. (2) has a stronger capacity to fit the underlying mapping than those stacking operations.

### 3.2. **Fine-structured Objects Refinement**

Besides the complex manmade objects, intricate fine-structured objects also increase the difficulty for accurate labeling in VHR images. Actually, the final *feature maps* outputted by the FCN-based methods is quite coarse due to multiple *sub-samplings*. For example, the size of the last *feature maps* in VGG-Net (Simonyan and Zisserman, 2015) is $1/32$ of input size. Thus, it is very hard to restore the low-level details of objects (e.g., boundary and localization) for accurate labeling, especially for fine-structured objects.

In CNNs, it is found that the low-level features can usually be captured by the shallow layers (Zeiler and Fergus, 2014). Based on this observation, we propose to reutilize the low-level features with a coarse-to-fine refinement strategy, as shown in the rightmost part of Fig. 1. Specifically, the shallow layers with fine resolution are progressively reintroduced into the decoder stream by long-span connections. As a result, the coarse *feature maps* can be refined and the low-level details can be recovered. Each single refinement process is illustrated in Fig. 3, which can be formulated as:

$$M^{i+1} = \Re\left[ \Upsilon\left[ \mathcal{L}(M^i \otimes \mathbf{w}_{M^i}) \oplus \mathcal{L}(F^i \otimes \mathbf{w}_{F^i}) \right] \right], \tag{3}$$

where $M^i$ denotes the refined feature maps of the previous process, and $F^i$ denotes the feature maps to be reutilized in this process coming from a shallower layer. $\mathbf{w}_{M^i}$ and $\mathbf{w}_{F^i}$ are the convolutional weights for $M^i$ and $F^i$ respectively. '$\otimes$' and '$\oplus$' denote the operations of convolution and fusion, respectively. $\mathcal{L}(\cdot)$ is the ReLU activation function. $\Re[\cdot]$ denotes the resize process and $\Upsilon[\cdot]$ denotes the process of residual correction. To fuse finer detail information from the next shallower layer, we resize the current feature maps to the corresponding higher resolution with bilinear interpolation to generate $M^{i+1}$.

It is fairly beneficial to fuse those low-level features using the proposed refinement strategy. On one hand, in fact, the *feature maps* of different resolutions in the encoder (see Fig. 1) represent semantics of different levels (Zeiler and Fergus, 2014). Thus, due to their inherent semantic gaps, stacking all these features directly (Hariharan et al., 2015; Farabet et al., 2013) may not be a good choice. In our method, the influence of semantic gaps is alleviated when a gradual fusion strategy is used. On the other hand, in training stage, the long-span connections allow direct gradient propagation to shallow layers, which helps effective end-to-end training.

The most relevant work with our refinement strategy is proposed in (Pinheiro et al., 2016), however, it is different from ours to a large extent. On one hand, our strategy focuses on performing dedicated refinement considering the specific properties (e.g., small dataset and intricate scenes) of VHR images in urban areas. Specifically, as shown in Fig. 1, only a few specific shallow layers are chosen for the refinement. Those layers that actually contain adverse noise due to intricate scenes are not incorporated. On the other hand, our refinement strategy works with our specially designed residual correction scheme, which will be elaborated in the following Section.
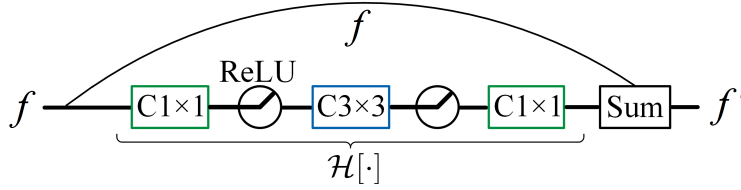
Figure 4: Residual correction scheme. 'C1×1' and 'C3×3' denote convolutional operation with *kernel size* 1×1 and 3×3, respectively. (Best viewed in color)

### 3.3. Residual Correction

It is notable that the proposed two solutions for labeling confusing manmade objects and fine-structured objects are quite different. In order to collaboratively and effectively integrate them into a single network, we have to find a approach to perform effective multi-feature fusion inside the network. This task is very challenging due to two issues. Firstly, as network deepens, it is fairly difficult for CNNs to directly fit a desired underlying mapping (He et al., 2016). Furthermore, this problem is worsened when it comes to fuse features of different levels. Secondly, there exists latent fitting residual when fusing multiple features of different semantics, which could cause the lack of information in the progress of fusion. To address this problem, a residual correction scheme is proposed, as shown in Fig. 4. It is dedicatedly aimed at correcting the latent fitting residual in multi-feature fusion inside ScasNet.

Specifically, building on the idea of deep residual learning (He et al., 2016), we explicitly let the stacked layers fit an inverse residual mapping, instead of directly fitting a desired underlying fusion mapping. Formally, let $f$ denote fused feature and $f'$ denote the desired underlying fusion. We expect the stacked layers to fit another mapping, which we call inverse residual mapping as:

$$\mathcal{H}[\cdot] = f' - f. \tag{4}$$

Actually, the aim of $\mathcal{H}[\cdot]$ is to compensate for the lack of information caused by the latent fitting residual, thus to achieve the desired underlying fusion $f' = f + \mathcal{H}[\cdot]$. Moreover, as demonstrated by (He et al., 2016), the inverse residual learning can be very effective in deep network, because it is easier to fit $\mathcal{H}[\cdot]$ than to directly fit $f'$ when network deepens. As a result, the adverse influence of latent fitting residual in multi-feature fusion can be well counteracted, i.e, the residual is well corrected.

It should be noted that, our residual correction scheme is quite different from the so-called chained residual pooling in RefineNet (Lin et al., 2016) on both function and structure. Functionally, the chained residual pooling in RefineNet aims to capture background context. However, our scheme explicitly focuses on correcting the latent fitting residual, which is caused by semantic gaps in multi-feature fusion. Structurally, the chained residual pooling is fairly complex, while our scheme is simple and efficient. As can be seen in Fig. 4, only one basic residual block is used in our scheme, and it is simply constituted by three convolutional layers and a skip connection.

As shown in Fig. 1, several residual correction modules are elaborately embedded in ScasNet, which can greatly prevent the fitting residual from accumulating. As a result, the proposed two different solutions work collaboratively and effectively, leading to a very valid global-to-local and coarse-to-fine labeling manner. Besides, the skip connection (see Fig. 4) is very beneficial for gradient propagation, resulting in an efficient end-to-end training of ScasNet.

### 3.4. ScasNet Configuration

As depicted in Fig. 1, the encoder network corresponds to a feature extractor that transforms the input image to multi-dimensional shrinking *feature maps*. To achieve this function, any existing CNN structures can be taken as the encoder part. In this paper, we propose two types of ScasNet based on two typical

**Algorithm 1** Learning procedure of the proposed ScasNet

---

**Input:** The image and label data $(\mathbf{x}, \mathbf{y})$.

**Output:** The network parameters $\boldsymbol{\theta}$ of ScasNet.

1: Initialize $\boldsymbol{\theta}$ and the learning rate $\eta$.

2: **Repeat:**

3:   Call the encoder forward pass to obtain feature maps of different levels $\mathrm{F} = \text{FEATUREEXTRACTION}(\mathbf{x}, \boldsymbol{\theta})$.

4:   Aggregate multi-context information $\mathrm{T} = \text{MULTISCALECONTEXTSAGGREGATION}(\mathrm{F}, \boldsymbol{\theta})$ by Eq. (2).

5:   Perform refinement to obtain the refined feature map $f(\mathbf{x}) = \text{REFINEMENT}(\mathrm{T}, \mathrm{F}, \boldsymbol{\theta})$ by Eq. (3).

6:   Calculate $\text{Loss}(\boldsymbol{\theta}) = \text{NORMALIZEDCROSSENTROPYLOSS}(\mathbf{y}, f(\mathbf{x}))$ by Eq. (5) and Eq. (6)

7:   Calculate the back propagation gradient $\frac{\partial \text{Loss}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$ by Eq. (7) and Eq. (8) with chain rule.

8:   Update $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \frac{\partial \text{Loss}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$.

9: **Until:** $\text{Loss}(\boldsymbol{\theta})$ converges

10: **Return:** $\boldsymbol{\theta}$

---

networks, i.e., 16-layer VGG-Net (Simonyan and Zisserman, 2015) and 101-layer ResNet (He et al., 2016). Compared with VGG ScasNet, ResNet ScasNet has better performance while suffering higher complexity.

We supply the trained models of these two CNNs so that the community can directly choose one of them based on different applications which require different trade-off between accuracy and complexity. All codes of the two specific ScasNet are released on the github[*]. For clarity, we briefly introduce their configurations in the following.

**VGG ScasNet:** In VGG ScasNet, the encoder is based on a VGG-Net variant (Chen et al., 2015), which is to obtain finer *feature maps* (about $1/8$ of input size rather than $1/32$). On the last layer of encoder, multi-scale contexts are captured by *dilated convolution* operations with *dilation rates* of 24, 18, 12 and 6. We only choose three shallow layers for refinement as shown in Fig. 1. There are two reasons: 1) shallower layers also carry much adverse noise despite of finer low-level details contained in them; 2) It is very difficult to train a more complex network well with remote sensing datasets, which are usually very small. In the encoder, we always use the last convolutional layer in each stage prior to pooling for refinement, because they contain stronger semantics in that stage. Six residual correction modules are employed for multi-feature fusion. Finally, a softmax classifier is employed to obtain probability maps, which indicate the likelihood of each pixel belonging to a category.

**ResNet ScasNet:** The configuration of ResNet ScasNet is almost the same as VGG ScasNet, except for four aspects: the encoder is based on a ResNet variant (Zhao et al., 2016), four shallow layers are used for refinement, seven residual correction modules are employed for feature fusions and BN layer is used.

It should be noted that due to the complicated structure, ResNet ScasNet has much difficulty to converge without BN layer. On the contrary, VGG ScasNet can converge well even though the BN layer is not used since it is relatively easy to train. In both of the two types of ScasNet, *sum* fusion operation is performed for efficiency.

### 3.5. Learning and Inference

In the learning stage, original VHR images and their corresponding reference images (i.e., ground truth) are used. Both of them are cropped into a number of patches, which are used as inputs to ScasNet. We use the normalized cross entropy loss as the learning objective, which is defined as

$$\text{Loss}(y, f(x), \boldsymbol{\theta}) = \frac{1}{MN} \sum_{i=1}^{M} \sum_{j=1}^{N} \sum_{k=1}^{K} -\mathrm{I}(y_i^j = k) \log p_k(x_i^j), \tag{5}$$

---

[*]https://github.com/Yochengliu/ScasNet

9

**Algorithm 2** Inference procedure of the proposed ScasNet

---

**Input:** The image data $\mathbf{x}$ and the number of scales $L$.
**Output:** The prediction labeling map $\mathbf{k}$.

 1: Initialize the network parameters $\boldsymbol{\theta}$ outputted by Algorithm 1
    and the average prediction probability map $p_{\mathbf{k}}(\mathbf{x}) = \mathbf{0}$.
 2: **for** $\ell$ in $\{1, \ldots, L\}$ **do**
 3:     Calculate the resized image $\mathbf{x}^\ell = \text{RESIZEWITHBILNEARINTERPOLATION}(\mathbf{x}, \ell)$.
 4:     Obtain the final feature map for the $\ell$-th scale $f^\ell(\mathbf{x}^\ell) = \text{NETWORKFORWARDPASS}(\mathbf{x}^\ell, \boldsymbol{\theta})$.
 5:     Calculate the prediction probability map for the $\ell$-th scale $p_{\mathbf{k}}^\ell(\mathbf{x}^\ell) = \text{SOFTMAXFUNCTION}(f(\mathbf{x}^\ell))$ by Eq. (6).
 6:     Resize $p_{\mathbf{k}}^\ell(\mathbf{x}^\ell)$ to original image size $p_{\mathbf{k}}(\mathbf{x}^\ell) = \text{RESIZEWITHBILNEARINTERPOLATION}(p_{\mathbf{k}}^\ell(\mathbf{x}^\ell), \mathbf{x})$.
 7:     Add $p_{\mathbf{k}}(\mathbf{x}^\ell)$ to the average prediction probability map $p_{\mathbf{k}}(\mathbf{x}) = p_{\mathbf{k}}(\mathbf{x}) + p_{\mathbf{k}}(\mathbf{x}^\ell)$
 8: **end for**
 9: Perform average operation $p_{\mathbf{k}}(\mathbf{x}) = \frac{1}{L} p_{\mathbf{k}}(\mathbf{x})$
10: Calculate the prediction labeling map $\hat{\mathbf{k}} = \underset{\mathbf{k}}{\text{argmax}}\ p_{\mathbf{k}}(\mathbf{x})$ in Eq. (1).

11: **Return:** $\hat{\mathbf{k}}$

---

where $\boldsymbol{\theta}$ represents the parameters of ScasNet; $M$ is the mini-batch size; $N$ is the number of pixels in each patch; $K$ is the number of categories; $\text{I}(y = k)$ is an indicator function, it takes $1$ when $y = k$, and $0$ otherwise; $x_i^j$ is the $j$-th pixel in the $i$-th patch and $y_i^j$ is the ground truth label of $x_i^j$. Let $f(x_i^j)$ denote the output of the layer before softmax (see Fig. 1) at pixel $x_i^j$, the probability of the pixel $x_i^j$ belonging to the $k$-th category $p_k(x_i^j)$ is defined by the softmax function, that is

$$p_k(x_i^j) = \frac{\exp(f_k(x_i^j))}{\sum\limits_{l=1}^{K} \exp(f_l(x_i^j))}. \tag{6}$$

To train ScasNet in the end-to-end manner, $\text{Loss}(\boldsymbol{\theta})$ is minimized w.r.t. the ScasNet parameters $\boldsymbol{\theta}$. We have to first calculate the derivative of the loss in Eq. (5) w.r.t. the parameters of different component layers with chain rule, and then update the parameters layer-by-layer with back propagation. For clarity, we only present the generic derivative of loss to the output of the layer before softmax and other hidden layers. The derivative of $\text{Loss}(\boldsymbol{\theta})$ to the output (i.e., $f_k(x_i^j)$) of the layer before softmax is calculated as:

$$\frac{\partial \text{Loss}(\boldsymbol{\theta})}{\partial f_k(x_i^j)} = \frac{1}{MN} \sum_{i=1}^{M} \sum_{j=1}^{N} \sum_{k=1}^{K} - \text{I}(y_i^j = k) \left(1 - p_k(x_i^j)\right). \tag{7}$$

The specific derivation process can be referred in the Appendix A of supplementary material.

The derivative of $\text{Loss}(\boldsymbol{\theta})$ to each hidden (i.e., $h_k(x_i^j)$) layer can be obtained with the chain rule as:

$$\frac{\partial \text{Loss}(\boldsymbol{\theta})}{\partial h_k(x_i^j)} = \frac{\partial \text{Loss}(\boldsymbol{\theta})}{\partial f_k(x_i^j)} \frac{\partial f_k(x_i^j)}{\partial h_k(x_i^j)}. \tag{8}$$

The first item in Eq. (8) is given in Eq. (7), and the second item also can be obtained by corresponding chain rule.

The pseudo-code of learning procedure of ScasNet is shown in Algorithm 1. In the experiments, we implement ScasNet based on the Caffe framework (Jia et al., 2014). The image patches of size $400 \times 400$ are used as inputs [1]. Due to the limit of GPU memory, we set the mini-batch size as $4$. To train ScasNet,

---

[1]The possibly few number of categories in these patches doesn't influence the high diversity of categories in raw VHR images.

Table 1: The detailed information of experimental setting on the three datasets. 'offline/online' denotes the training set for offline validation and the training set for online test, respectively.

| Dataset | TRAINING SET | | VALIDATION SET | TEST SET |
|---|---|---|---|---|
| | Images | Patches (400×400) | Images | Images |
| Massachusetts Building | 141 | 20727 | 0 | 10 |
| Vaihingen Challenge | offline/online 8/16 | offline/online 12384/24400 | 8 | 17 |
| Potsdam Challenge | offline/online 14/24 | offline/online 16800/28800 | 10 | 14 |



(a) Massachusetts Building  (b) Vaihingen Challenge  (c) Potsdam Challenge
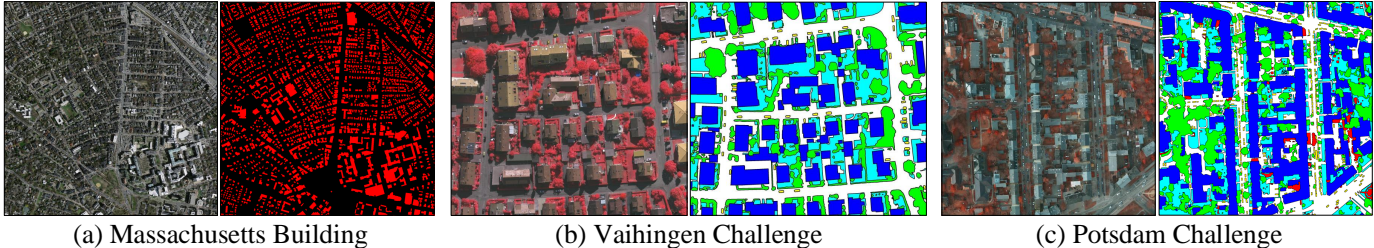
Figure 5: The image samples and corresponding ground truth on the three datasets. The label of Massachusetts building includes two categories: building (red) and background (black). The label of Vaihingen and Potsdam challenge includes six categories: impervious surface (imp surf, white), building (blue), low vegetation (low veg, cyan), tree (green), car (yellow) and clutter/background (red), where the boundary (black) is depicted for visual clarity.

we use stochastic gradient descent (SGD) with initial learning rate of $0.01$, and drop the learning rate by a factor of $0.1$ every $20$ epochs. The momentum and weight decay are set as $0.9$ and $0.0005$, respectively. Experimentally, ScasNet is trained for about $80$ epochs.

The pseudo-code of inference procedure is shown in Algorithm 2. In the inference stage, we perform multi-scale inference of 0.5, 1 and 1.5 times the size of raw images (i.e., $L = 3$ scales), and we average the final outputs at all the three scales. Specifically, we first crop a resized image (i.e., $\mathbf{x}^\ell$) into a series of patches without overlap. Then, the prediction probability maps of these patches are predicted by inputting them into ScasNet with a forward pass. Finally, the entire prediction probability map (i.e., $p_{\mathbf{k}}^\ell(\mathbf{x}^\ell)$) of this image is constituted by the probability maps of all patches. The purpose of multi-scale inference is to mitigate the discontinuity in final labeling map caused by the interrupts between patches.

## 4. Experiments and Evaluations

In this section, dataset description, experimental setting, comparing methods and extensive experiments in both qualitative and quantitative comparisons are first presented. Then, the proposed ScasNet is analyzed in detail by a series of ablation experiments.

### 4.1. Dataset Description

We evaluate the proposed ScasNet on three challenging public datasets for semantic labeling.

**Massachusetts Building Dataset:** This dataset is proposed by Mnih (Mnih, 2013). It consists of 151 aerial images of the Boston area, with each of the images being $1500 \times 1500$ pixels at a GSD (Ground Sampling Distance) of 1m. The ground truth of all these images are available. We randomly split the data into a training set of 141 images, and a test set of 10 images. As Fig. 5(a) shows, it covers mostly urban areas and buildings of all sizes, including houses and garages.

**ISPRS Vaihingen Challenge Dataset:** This is a benchmark dataset for *ISPRS 2D Semantic labeling challenge* in Vaihingen (ISPRS, 2016). It consists of 3-band IRRG (Infrared, Red and Green) image data, and corresponding DSM (Digital Surface Model) and NDSM (Normalized Digital Surface Model) data. Overall, there are 33 images of $\approx 2500 \times 2000$ pixels at a GSD of $\approx$ 9cm in image data. Among them,

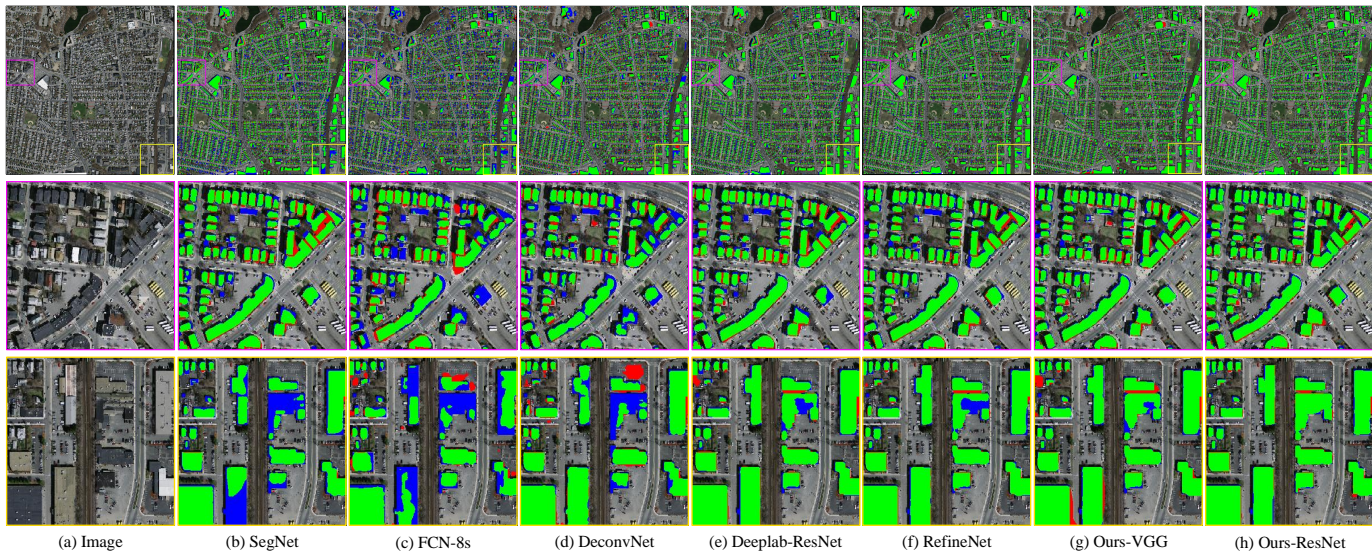| (a) Image | (b) SegNet | (c) FCN-8s | (d) DeconvNet | (e) Deeplab-ResNet | (f) RefineNet | (g) Ours-VGG | (h) Ours-ResNet |

Figure 6: Qualitative comparison with the state-of-the-art deep models on Massachusetts building TEST SET. The 1st row illustrates the overall results of one image sample, and the last two rows show the close-ups of the corresponding regions in the 1st row. In the colored figures, true positive (tp) is marked in green, false positive (fp) in red and false negative (fn) in blue. (Best viewed in color)

the ground truth of only 16 images are available, and those of the remaining 17 images are withheld by the challenge organizer for online test. For offline validation, we randomly split the 16 images with ground truth available into a training set of 8 images, and a validation set of 8 images. For online test, we use all the 16 images as training set. Note that DSM and NDSM data in all the experiments on this dataset are not used.

**ISPRS Potsdam Challenge Dataset:** This is a benchmark dataset for *ISPRS 2D Semantic labeling challenge* in Potsdam (ISPRS, 2016). It consists of 4-band IRRGB (Infrared, Red, Green, Blue) image data, and corresponding DSM and NDSM data. Overall, there are 38 images of $6000 \times 6000$ pixels at a GSD of $\approx 5$cm. Among them, the ground truth of only 24 images are available, and those of the remaining 14 images are withheld by the challenge organizer for online test. For offline validation, we randomly split the 24 images with ground truth available into a training set of 14 images, a validation set of 10 images. For online test, we use all the 24 images as training set. Note that only the 3-band IRRG images extracted from raw 4-band data are used, and DSM and NDSM data in all the experiments on this dataset are not used.

Table 1 summarizes the detailed information of all the above datasets. Fig. 5 shows some image samples and the ground truth on the three datasets. As it shows, there are many confusing manmade objects and intricate fine-structured objects in these VHR images, which poses much challenge for achieving both coherent and accurate semantic labeling.

### 4.2. Experimental Setting

The remote sensing datasets are relatively small to train the proposed deep ScasNet. To reduce overfitting and train an effective model, data augmentation, *transfer learning* (Yosinski et al., 2014; Penatti et al., 2015; Hu et al., 2015; Xie et al., 2015) and regularization techniques are applied.

In the experiments, $400 \times 400$ patches cropped from raw images are employed to train ScasNet. For the training sets, we use a two-stage method to perform data augmentation. In the first stage, given an image, we crop it to generate a series of $400 \times 400$ patches with the overlap of 100 pixels. In the second stage, for each patch, we flip it in horizontal and vertical reflections and rotate it counterclockwise at the step of $90°$. The detailed number of patches in the augmented data is presented in Tabel 1.

In the experiments, the parameters of the encoder part (see Fig. 1) in our models are initialized with the models pre-trained on PASCAL VOC 2012 (Everingham et al., 2015). All the other parameters in our

Figure 7: Quantitative comparison (%) with the state-of-the-art deep models on Massachusetts building TEST SET, where the values in bold are the best and the values underlined are the second best.

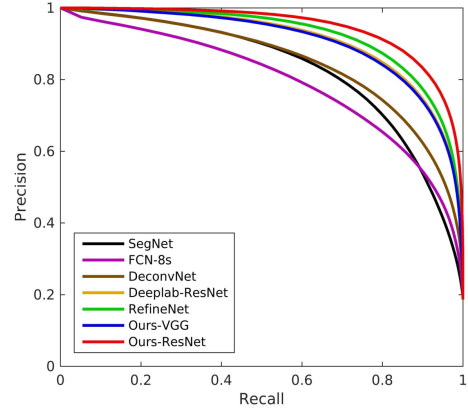| Method | IoU | F1 |
|---|---|---|
| SegNet | 56.38 | 72.11 |
| FCN-8s | 50.94 | 67.96 |
| DeconvNet | 60.74 | 75.57 |
| Deeplab-ResNet | 69.50 | 82.01 |
| RefineNet | 71.92 | 83.67 |
| Ours-VGG | 69.22 | 81.81 |
| Ours-ResNet | **74.34** | **85.58** |



Figure 8: Precision-recall (PR) curves of all the comparing deep models on Massachusetts building TEST SET. (Best viewed in color)



(a) Image    (b) Ground Truth    (c) SegNet    (d) FCN-8s    (e) DeconvNet    (f) Deeplab-ResNet    (g) RefineNet    (h) Ours-VGG    (i) Ours-ResNet
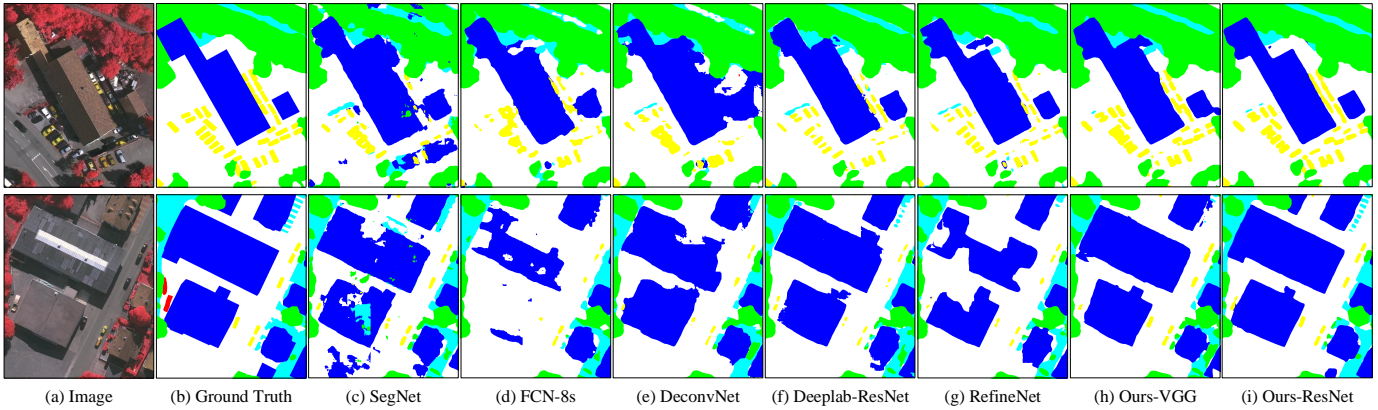
Figure 9: Qualitative comparison with the state-of-the-art deep models on *ISPRS Vaihingen challenge* OFFLINE VALIDATION SET. The label includes six categories: impervious surface (imp surf, white), building (blue), low vegetation (low veg, cyan), tree (green), car (yellow) and clutter/background (red).

models are initialized using the techniques introduced by He et al. (He et al., 2015a).

To avoid overfitting, dropout technique (Srivastava et al., 2014) with ratio of $50\%$ is used in ScasNet, which provides a computationally inexpensive yet powerful regularization to the network.

### *4.3.* **Comparing methods**

To verify the performance, the proposed ScasNet is compared with extensive state-of-the-art methods on two aspects: deep models comparison and benchmark test comparison.

**Comparing Deep Models:** ScasNet is compared with five state-of-the-art deep models on the three datasets. The main information of these models (including our models) is summarized as follows:

1) Ours-VGG: The self-cascaded network with the encoder based on a variant of 16-layer VGG-Net (Chen et al., 2015).
2) Ours-ResNet: The self-cascaded network with the encoder based on a variant of 101-layer ResNet (Zhao et al., 2016).
3) FCN-8s: Long et al. (Long et al., 2015) propose FCN for semantic segmentation, which achieves the state-of-the-art performance on three benchmarks (Everingham et al., 2015; Silberman et al., 2012; Liu et al., 2008). There are three versions of FCN models: FCN-32s, FCN-16s and FCN-8s. We use the best performance model FCN-8s as comparison.

Table 2: Quantitative comparison (%) with the state-of-the-art deep models on *ISPRS Vaihingen challenge* OFFLINE VALIDATION SET, where the values in bold are the best and the values underlined are the second best.

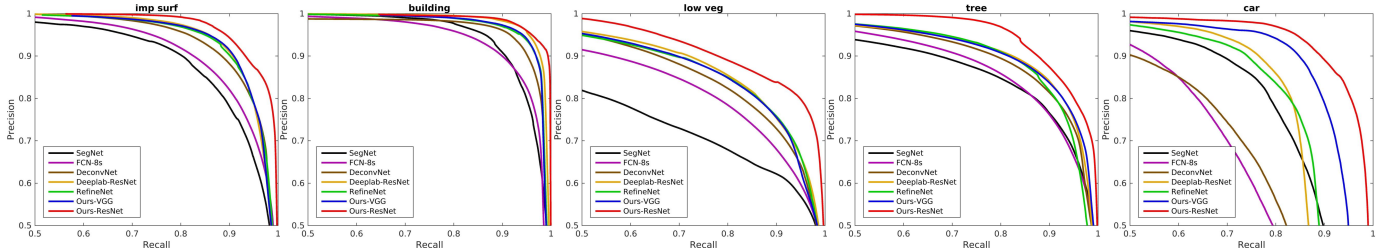| Model | imp surf | | building | | low veg | | tree | | car | | Avg. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | IoU | F1 | IoU | F1 | IoU | F1 | IoU | F1 | IoU | F1 | mean IoU | mean F1 |
| SegNet | 66.85 | 80.13 | 76.10 | 86.43 | 50.56 | 68.65 | 69.71 | 82.15 | 62.38 | 76.83 | 65.12 | 78.83 |
| FCN-8s | 75.26 | 85.89 | 80.51 | 89.20 | 65.58 | 79.21 | 70.49 | 82.69 | 45.84 | 62.87 | 67.54 | 79.97 |
| DeconvNet | 80.27 | 89.06 | 87.19 | 93.16 | 68.57 | 81.36 | 74.91 | 85.65 | 51.93 | 68.36 | 72.57 | 83.52 |
| Deeplab-ResNet | 82.20 | 90.23 | 91.22 | 95.41 | 71.12 | 83.12 | 76.93 | 86.96 | 56.78 | 72.43 | 75.65 | 85.63 |
| RefineNet | 80.08 | 88.94 | 88.62 | 93.97 | 70.69 | 82.83 | 76.00 | 86.36 | 68.35 | 81.56 | 76.75 | 86.73 |
| Ours-VGG | 82.70 | 90.53 | 89.54 | 94.48 | 69.00 | 81.66 | 76.17 | 86.47 | 76.89 | 86.93 | 78.86 | 88.02 |
| Ours-ResNet | **85.86** | **93.76** | **92.45** | **96.19** | **76.26** | **87.62** | **83.77** | **90.61** | **81.14** | **89.81** | **83.90** | **91.60** |



Figure 10: Precision-recall (PR) curves of all the comparing deep models on *ISPRS Vaihingen challenge* OFFLINE VALIDATION SET. Categories from left to right: impervious surface (imp surf), building, low vegetation (low veg), tree , car. (Best viewed in color)

4) SegNet: Badrinarayanan et al. (Badrinarayanan et al., 2015) propose SegNet for semantic segmentation of road scene, in which the decoder uses pooling indices in the encoder to perform non-linear *up-sampling*. It provides competitive performance while works faster than most of the other models.

5) DconvNet: Deconvolutional network (DconvNet) is proposed by Noh et al. (Noh et al., 2015) for semantic segmentation, which is composed of *deconvolution* and *un-pooling* layers. It achieves the state-of-the-art performance on PASCAL VOC 2012 (Everingham et al., 2015).

6) Deeplab-ResNet: Chen et al. (Chen et al., 2015) propose Deeplab-ResNet based on three 101-layer ResNet (He et al., 2016), which achieves the state-of-the-art performance on PASCAL VOC 2012 (Everingham et al., 2015). Actually, they use three-scale (0.5, 0.75 and 1 the size of input image) images as input to three 101-layer ResNet respectively, and then fuse three outputs as final prediction.

7) RefineNet: RefineNet is proposed by Lin et al. (Lin et al., 2016) for semantic segmentation, which is based on ResNet (He et al., 2016). It achieves the state-of-the-art performance on seven benchmarks, such as PASCAL VOC 2012 (Everingham et al., 2015) and NYUDv2(Silberman et al., 2012). Here, we take RefineNet based on 101-layer ResNet for comparison.

It should be noted that all the experimental settings for the above models are the same, except for two aspects. Firstly, their training hyper-parameter values used in the Caffe framework (Jia et al., 2014) are different. This is because it may need different hyper-parameter values (such as learning rate) to make them converge when training different deep models. Secondly, all the models are trained based on the widely used *transfer learning* (Yosinski et al., 2014; Penatti et al., 2015; Hu et al., 2015; Xie et al., 2015) in the field of *deep learning*. Specifically, except for our models, all the other models are trained by finetuning their corresponding best models pre-trained on PASCAL VOC 2012 (Everingham et al., 2015) on semantic segmentation task. For our models, only the parameters of the encoder part (see Fig. 1) are initialized with the pre-trained models. Furthermore, the influence of *transfer learning* on our models is analyzed in Section 4.7.

**Benchmark Comparing Methods:** By submitting the results of test set to the *ISPRS challenge* organizer, ScasNet is also compared with other competitors' methods on benchmark test. The details of these methods (including our methods) are listed as follows, where the names in brackets are the short names on the

challenge evaluation website [*]:

1) Ours-ResNet (**'CASIA2'**): The single self-cascaded network with the encoder based on a variant of 101-layer ResNet (Zhao et al., 2016). In our method, only raw image data is used for training. Specifically, 3-band IRRG images are used for Vaihingen and only 3-band IRRG images obtained from raw image data (i.e., 4-band IRRGB images) are used for Potsdam. Moreover, we do not use the elevation data (DSM and NDSM), additional hand-crafted features, model ensemble strategy or any postprocessing.

2) SVL-features + DSM + Boosting + CRF (**'SVL_*'**): The method as baseline implemented by the challenge organizer (Gerke, 2015). In addition to the standard SVL-features (Gould et al., 2011), they also use NDVI (Normalized Digital Vegetation Index), saturation and NDSM features. Then, an Adaboost-based classifier is trained. A CRF (Conditional Random Field) model is applied to obtain final prediction. For comparison, 'SVL_6' is compared for Vaihingen and 'SVL_3' (no CRF) for Potsdam.

3) CNN + NDSM + Deconvolution (**'UZ_1'**): The method proposed by (Volpi and Tuia, 2017). They use an downsample-then-upsample architecture , in which rough spatial maps are first learned by convolutions and then these maps are upsampled by *deconvolution*. NDSM data is used in their method.

4) CNN + DSM + NDSM + RF + CRF (**'ADL_3'**): The method proposed by (Paisitkriangkrai et al., 2016). They apply both CNN and hand-crafted features to dense image patches to produce per-pixel category probabilities. Random forest (RF) classifier is trained on hand-crafted features and the output probabilities are combined with those generated by the CNN. CRF is applied as a postprocessing step.

5) FCN + DSM + RF + CRF (**'DST_2'**): The method proposed by (Sherrah, 2016). They use a hybrid FCN architecture to combine image data with DSM data. Then, CRF is applied as a postprocessing step.

6) FCN + SegNet + VGG + DSM + Edge (**'DLR_8'**): The method proposed by (Marmanis et al., 2016). They use a multi-scale ensemble of FCN, SegNet and VGG, incorporating both image data and DSM data. Moreover, they combine semantic labeling with informed edge detection.

7) SegNet + DSM + NDSM (**'ONE_7'**): The method proposed by (Audebert et al., 2016). They fuse the output of two multi-scale SegNets, which are trained with IRRG images and synthetic data (NDVI, DSM and NDSM) respectively.

8) CNN + DSM + SVM (**'GU'**): In their method, both image data and DSM data are used to train a CNN. Moreover, CNN is trained on six scales of the input data. Finally, a SVM maps the six predictions into a single-label.

9) CNN + DSM (**'AZ_1'**): In their method, a CNN with encoder-decoder architecture is used. The input to the network includes six channels of IRRGB, NDVI, and NDSM, which are concatenated together.

10) SegNet + NDSM (**'RIT_2'**): In their method, two SegNets are trained with RGB images and synthetic data (IR, NDVI and NDSM) respectively. Then, feature fusion in the early stages is performed.

### *4.4.* **Evaluation Metrics**

To assess the quantitative performance, two overall benchmark metrics are used, i.e., *F1 score* (F1) and *intersection over union* (IoU). F1 is defined as

$$\text{F1} = 2\frac{\text{Pre} \times \text{Rec}}{\text{Pre} + \text{Rec}}, \ \text{Pre} = \frac{tp}{tp + fp}, \ \text{Rec} = \frac{tp}{tp + fn}. \tag{9}$$

Here, $tp$, $fp$ and $fn$ are the number of true positives, false positives and false negatives, respectively.

IoU is defined as:

$$\text{IoU}(\mathcal{P}_m, \mathcal{P}_{gt}) = \frac{|\mathcal{P}_m \cap \mathcal{P}_{gt}|}{|\mathcal{P}_m \cup \mathcal{P}_{gt}|}, \tag{10}$$

where $\mathcal{P}_{gt}$ is the set of ground truth pixels and $\mathcal{P}_m$ is the set of prediction pixels, '$\cap$' and '$\cup$' denote *intersection* and *union* operations, respectively. $|\cdot|$ denotes calculating the number of pixels in the set.

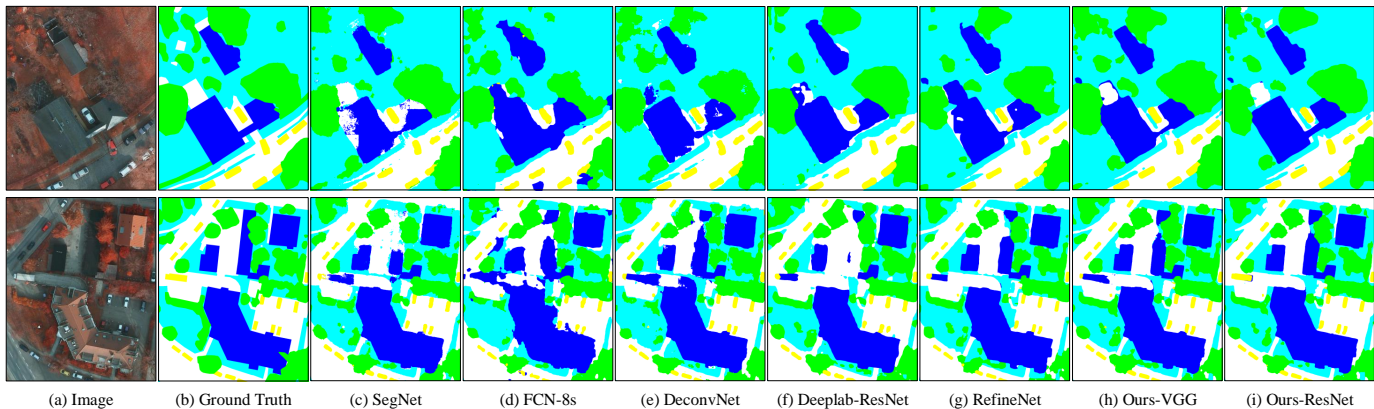| (a) Image | (b) Ground Truth | (c) SegNet | (d) FCN-8s | (e) DeconvNet | (f) Deeplab-ResNet | (g) RefineNet | (h) Ours-VGG | (i) Ours-ResNet |

Figure 11: Qualitative comparison with the state-of-the-art deep models on *ISPRS Potsdam challenge* OFFLINE VALIDATION SET. The label includes six categories: impervious surface (imp surf, white), building (blue), low vegetation (low veg, cyan), tree (green), car (yellow) and clutter/background (red).

Table 3: Quantitative comparison (%) with the state-of-the-art deep models on *ISPRS Potsdam challenge* OFFLINE VALIDATION SET, where the values in bold are the best and the values underlined are the second best.

| Model | imp surf | | building | | low veg | | tree | | car | | Avg. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | IoU | F1 | IoU | F1 | IoU | F1 | IoU | F1 | IoU | F1 | mean IoU | mean F1 |
| SegNet | 85.42 | 92.14 | 91.17 | 95.38 | 78.58 | 88.01 | 75.71 | 86.17 | 88.12 | 93.68 | 83.80 | 91.08 |
| FCN-8s | 77.55 | 87.35 | 79.94 | 88.85 | 71.95 | 83.68 | 69.53 | 82.02 | 79.68 | 88.69 | 75.73 | 86.12 |
| DeconvNet | 87.08 | 93.09 | 93.12 | 96.44 | 77.59 | 87.38 | 71.67 | 83.50 | 92.28 | 95.98 | 84.35 | 91.28 |
| Deeplab-ResNet | 88.23 | 93.75 | 94.39 | 97.11 | 78.85 | 88.18 | 74.50 | 85.39 | 87.11 | 93.11 | 84.62 | 91.51 |
| RefineNet | 86.80 | 92.93 | 91.13 | 95.36 | 78.69 | 88.07 | 73.51 | 84.74 | 92.75 | 96.24 | 84.58 | 91.47 |
| Ours-VGG | 88.68 | 94.00 | 94.12 | 96.97 | 80.67 | 89.30 | **77.86** | **87.55** | 94.07 | 96.94 | 87.08 | 92.95 |
| Ours-ResNet | **90.06** | **94.77** | **96.27** | **98.10** | **80.83** | **89.40** | 76.86 | 86.92 | **94.90** | **97.38** | **87.78** | **93.31** |

To evaluate the performance of different comparing deep models, we compare the above two metrics on each category, and the mean value of metrics to assess the average performance. Furthermore, precision-recall (PR) curve is drawn to qualify the relation between *precision* and *recall* on each category. Specifically, the predicted score maps are first binarized using different thresholds varying from 0 to 1. Then by comparing these binarized results with the ground truth, a series of precision-recall values can be obtained to plot the PR curve.

When compared with other competitors' methods on benchmark test (ISPRS, 2016), besides the F1 metric for each category, the *overall accuracy* (Overall Acc.) derived from the pixel-based confusion matrix (ISPRS, 2016) is also compared to assess the global performance.

It should be noted that all the metrics are computed using an alternative ground truth in which the boundaries of objects have been eroded by a 3-pixel radius. The eroded areas are ignored during evaluation, so as to reduce the impact of uncertain border definitions.

### 4.5. Comparison with Deep Models

To evaluate the effectiveness of the proposed ScasNet, the comparisons with five state-of-the-art deep models on the three challenging datasets are presented as follows:

1) **Massachusetts Building Test Set**: As the global visual performance (see the 1st row in Fig. 6) and local close-ups (see the last two rows in Fig. 6) show, SegNet, FCN-8s and DeconvNet have difficulty in recognizing confusing size-varied buildings. For fine-structured buildings, FCN-8s performs incomplete and inaccurate labeling while SegNet and DeconvNet do better. The results of Deeplab-ResNet, RefineNet
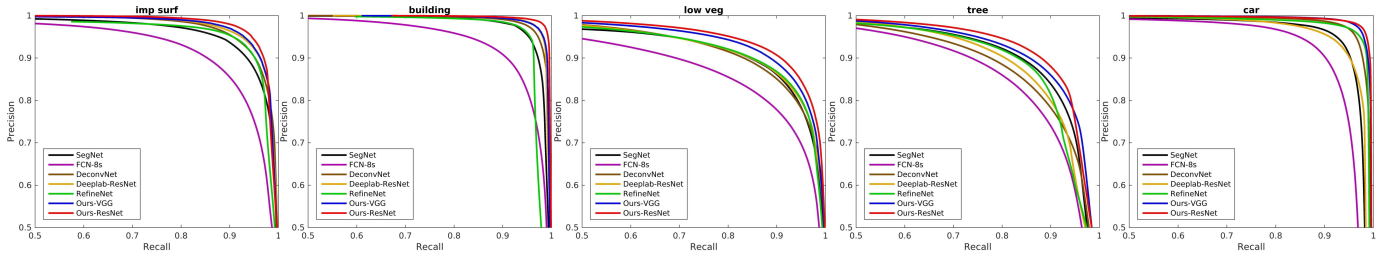
---

*http://www2.isprs.org/commissions/comm3/wg4/results.html

16

Figure 12: Precision-recall (PR) curves of all the comparing deep models on *ISPRS Potsdam challenge* OFFLINE VALIDATION SET. Categories from left to right: impervious surface (imp surf), building, low vegetation (low veg), tree , car. (Best viewed in color)

Table 4: Quantitative comparison (%) with other competitors' methods on *ISPRS Vaihingen challenge* ONLINE TEST SET, where the values in bold are the best and the values underlined are the second best. The names in brackets are the short names on the challenge evaluation website.

| Method | imp surf | building | low veg | tree | car | Overall Acc. |
|---|---|---|---|---|---|---|
| SVL-features + DSM + Boosting + CRF (**'SVL_6'**) | 86.00 | 90.20 | 75.60 | 82.10 | 45.40 | 83.20 |
| CNN + NDSM + Deconvolution (**'UZ_1'**) | 89.20 | 92.50 | 81.60 | 86.90 | 57.30 | 87.30 |
| CNN + DSM + NDSM + RF + CRF (**'ADL_3'**) | 89.50 | 93.20 | 82.30 | 88.20 | 63.30 | 88.00 |
| FCN + DSM + RF + CRF (**'DST_2'**) | 90.50 | 93.70 | 83.40 | 89.20 | 72.60 | 89.10 |
| FCN + SegNet + VGG + DSM + Edge (**'DLR_8'**) | 90.40 | 93.60 | 83.90 | <u>89.70</u> | 76.90 | 89.20 |
| SegNet + DSM + NDSM (**'ONE_7'**) | <u>91.00</u> | <u>94.50</u> | <u>84.40</u> | **89.90** | <u>77.80</u> | <u>89.80</u> |
| Ours-ResNet (**'CASIA2'**) | **93.20** | **96.00** | **84.70** | **89.90** | **86.70** | **91.10** |

and Ours-VGG are relatively good, but they tend to have more false negatives (blue). Ours-ResNet generates more coherent labeling on both confusing and fine-structured buildings. Table 7 summarizes the quantitative performance. As it shows, Ours-VGG achieves almost the same performance with Deeplab-ResNet, while Ours-ResNet achieves more decent score. Fig. 8 shows the PR curves of all the deep models, in which both Our-VGG and Our-ResNet achieve superior performances.

2) **Vaihingen Challenge Validation Set**: As shown in Fig. 9, SegNet, FCN-8s, DeconvNet and RefineNet are sensitive to the cast shadows of buildings and trees. They can not distinguish similar manmade objects well, such as buildings and roads. Meanwhile, for fine-structured objects, these methods tend to obtain inaccurate localization, especially for the car. The results of Deeplab-ResNet are relatively coherent, while they are still less accurate. Ours-VGG and Ours-ResNet show better robustness to the cast shadows. They can achieve coherent labeling for confusing manmade objects. Moreover, fine-structured objects also can be labeled with precise localization using our models. The quantitative performance is shown in Table 2. As can be seen, the performance of our best model outperforms other advanced models by a considerable margin on each category, especially for the car. Furthermore, the PR curves shown in Fig. 10 exhibit that, our best model performs better on all the given categories.

3) **Potsdam Challenge Validation Set**: As Fig. 11 shows, all the five comparing models are less effective in the recognition of confusing manmade objects. They are not robust enough to the occlusions and cast shadows. For fine-structured objects like the car, FCN-8s performs less accurate localization, while other four models do better. Although the labeling results of our models have a few flaws, they can achieve relatively more coherent labeling and more precise boundaries. Table 3 summarizes the quantitative performance. As it shows, in labeling the VHR images with such a high resolution of 5cm, all these models achieve decent results. Still, the performance of our best model exceeds other advanced models by a considerable margin, especially for the car. Moreover, as the PR curves in Fig. 12 show, our best model presents very decent performance.

In short, the above comparisons show that, on one hand, the proposed ScasNet has strong recognition ability for confusing manmade objects in VHR images. Meanwhile, ScasNet is quite robust to the occlusions and cast shadows, and it can perform coherent labeling even for very uneven regions. These results demonstrate

Table 5: Quantitative comparison (%) with other competitors' methods on *ISPRS Potsdam challenge* ONLINE TEST SET, where the values in bold are the best and the values underlined are the second best. The names in brackets are the short names on the challenge evaluation website.

| Method | imp surf | building | low veg | tree | car | Overall Acc. |
|---|---|---|---|---|---|---|
| SVL-features + DSM + Boosting (**'SVL_3'**) | 84.00 | 89.80 | 72.00 | 59.00 | 69.80 | 77.20 |
| CNN + DSM + SVM (**'GU'**) | 87.10 | 94.70 | 77.10 | 73.90 | 81.20 | 82.90 |
| CNN + NDSM + Deconvolution (**'UZ_1'**) | 89.30 | 95.40 | 81.80 | 80.50 | 86.50 | 85.80 |
| CNN + DSM (**'AZ_1'**) | 91.40 | 96.10 | 86.10 | 86.60 | 93.30 | 89.20 |
| SegNet + NDSM (**'RIT_2'**) | <u>92.00</u> | <u>96.30</u> | 85.50 | 86.50 | <u>94.50</u> | 89.40 |
| FCN + DSM + RF + CRF (**'DST_2'**) | 91.80 | 95.90 | <u>86.30</u> | <u>87.70</u> | 89.20 | <u>89.70</u> |
| Ours-ResNet (**'CASIA2'**) | **93.30** | **97.00** | **87.70** | **88.40** | **96.20** | **91.10** |

Table 6: Quantitative comparison (%) between 3-scale test (0.5, 1 and 1.5 times the size of raw image) and 1-scale test on *ISPRS Vaihingen & Potsdam challenge* ONLINE TEST SET.

| Benchmark | Method | imp surf | building | low veg | tree | car | Overall Acc. |
|---|---|---|---|---|---|---|---|
| Vaihingen | 1-scale test (**'CASIA3'**) | 92.70 | 95.50 | 83.90 | 89.40 | 86.70 | 90.60 |
| | 3-scale test (**'CASIA2'**) | 93.20 | 96.00 | 84.70 | 89.90 | 86.70 | 91.10 |
| Potsdam | 1-scale test (**'CASIA3'**) | 93.40 | 96.80 | 87.60 | 88.30 | 96.10 | 91.00 |
| | 3-scale test (**'CASIA2'**) | 93.30 | 97.00 | 87.70 | 88.40 | 96.20 | 91.10 |

the effectiveness of our multi-scale contexts aggregation approach. On the other hand, ScasNet can label size-varied objects completely, resulting in accurate and smooth results, especially for the fine-structured objects like the car. This demonstrates the validity of our refinement strategy.

### 4.6. Comparison on Benchmark Test

To further evaluate the effectiveness of the proposed ScasNet, comparisons with other competitors' methods on the two challenging benchmarks are presented as follows:

1) **Vaihingen Challenge**: On benchmark test of Vaihingen[*], Fig. 14 and Table 4 exhibit qualitative and quantitative comparisons with different methods, respectively. As shown in Fig. 14, other methods, even though the elevation data is used, are less effective for labeling confusing manmade objects and fine-structured objects simultaneously. In contrast, our method can obtain coherent and accurate labeling results. Moreover, our method can achieve labeling with smooth boundary and precise localization, especially for fine-structured objects like the car. As Table 4 shows, the quantitative performances of our method also outperform other methods by a considerable margin, especially for the car.

2) **Potsdam Challenge**: On benchmark test of Potsdam[†], qualitative and quantitative comparison with different methods are exhibited in Fig. 15 and Table 5, respectively. As shown in Fig. 15, all the comparing methods obtain good results, while more coherent and accurate results are achieved by our method. In addition, our method shows better robustness to the cast shadows. Meanwhile, as can be seen in Table 5, the quantitative performances of our method also outperform other methods by a considerable margin on all the categories.

As the above comparisons demonstrate, the proposed multi-scale contexts aggregation approach is very effective for labeling confusing manmade objects. Thus, our method can perform coherent labeling even for the regions which are very hard to distinguish. Meanwhile, our refinement strategy is much effective for accurate labeling. This results in a smooth labeling with accurate localization, especially for fine-structured objects like the car. Furthermore, both of them are collaboratively integrated into a deep model with the well-designed residual correction schemes. As a result, our method outperforms other sophisticated methods by the date of submission, even though it only uses a single network based on only raw image data. Other

---

[*]http://www2.isprs.org/vaihingen-2d-semantic-labeling-contest.html
[†]http://www2.isprs.org/potsdam-2d-semantic-labeling.html

Table 7: Ablation experiments (%) on *ISPRS Vaihingen challenge* OFFLINE VALIDATION SET. 'MSC' denotes aggregating multi-scale contexts in a parallel stack shown in Fig. 2(c). 'MSC+SC' denotes sequentially aggregating multi-scale contexts in a self-cascaded manner. 'MSC+SC+CReC' denotes sequentially aggregating multi-scale contexts in a self-cascaded manner and adding residual correction schemes in context aggregation, as shown in Fig. 2(b). 'Ref' denotes adding refinement. 'Ref+RReC' denotes adding refinement and residual correction schemes in refinement process, as the rightmost part of Fig. 1 shows.

| Model | imp surf | | building | | low veg | | tree | | car | | Avg. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | IoU | F1 | IoU | F1 | IoU | F1 | IoU | F1 | IoU | F1 | mean IoU | mean F1 |
| Baseline | 76.74 | 86.84 | 82.33 | 90.31 | 67.77 | 80.79 | 72.62 | 84.14 | 40.71 | 57.86 | 68.04 | 79.99 |
| +MSC | 75.67 | 86.15 | 86.38 | 92.70 | 66.56 | 79.92 | 73.92 | 85.00 | 40.80 | 57.95 | 68.67 | 80.34 |
| +MSC+SC | 77.13 | 87.38 | 87.01 | 93.05 | 68.80 | 81.52 | 74.98 | 85.70 | 46.35 | 62.98 | 70.90 | 82.13 |
| +MSC+SC+CReC | 80.10 | 88.95 | 87.72 | 93.26 | 68.92 | 81.68 | 75.15 | 85.81 | 56.07 | 71.48 | 73.59 | 84.24 |
| +MSC+SC+CReC+Ref | 80.61 | 89.26 | 89.06 | 94.21 | 70.57 | 82.75 | 76.39 | 86.62 | 61.34 | 76.04 | 75.59 | 85.78 |
| +MSC+SC+CReC+Ref+RReC | 82.70 | 90.53 | 89.54 | 94.48 | 69.00 | 81.66 | 76.17 | 86.47 | 76.89 | 86.93 | 78.86 | 88.02 |

Table 8: Quantitative comparison (%) between with and without using finetuning technique of the encoder part on *ISPRS Vaihingen challenge* OFFLINE VALIDATION SET. The model used to initialize the encoder part is pre-trained on PASCAL VOC 2012 (Everingham et al., 2015). 'w/o' denotes without using finetuning, 'w/' denotes using finetuning.

| Model | Finetuning | imp surf | | building | | low veg | | tree | | car | | Avg. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | IoU | F1 | IoU | F1 | IoU | F1 | IoU | F1 | IoU | F1 | mean IoU | mean F1 |
| VGG ScasNet | w/o | 80.55 | 89.23 | 87.23 | 93.18 | 67.18 | 80.06 | 73.94 | 84.88 | 71.35 | 82.88 | 76.05 | 86.05 |
| | w/ | 82.70 | 90.53 | 89.54 | 94.48 | 69.00 | 81.66 | 76.17 | 86.47 | 76.89 | 86.93 | 78.86 | 88.02 |
| ResNet ScasNet | w/o | 78.78 | 88.92 | 85.99 | 92.30 | 59.03 | 75.57 | 72.34 | 84.49 | 60.36 | 75.73 | 71.30 | 83.40 |
| | w/ | 85.86 | 93.76 | 92.45 | 96.19 | 76.26 | 87.62 | 83.77 | 90.61 | 81.14 | 89.81 | 83.90 | 91.60 |

competitors either use extra data such as DSM and model ensemble strategy, or employ structural models such as CRF.

To evaluate the performance brought by the three-scale test ( 0.5, 1 and 1.5 times the size of raw images), we submit the single scale test results to the challenge organizer. The evaluation results are listed in Table 6. As can be seen, all the categories on Vaihingen dataset achieve a considerable improvement except for the car. A possible reason is that, our refinement strategy is effective enough for labeling the car with the resolution of 9cm. Moreover, there is virtually no improvement on Potsdam dataset. Maybe for such a high resolution of 5cm, the influence of multi-scale test is negligible.

## *4.7.* **Model Analysis**

To evaluate the performance brought by each aspect we focus on in the proposed ScasNet, the ablation experiments of VGG ScasNet are conducted. Table 7 lists the results of adding different aspects progressively. The encoder (see Fig. 1) which is based on a VGG-Net variant (Chen et al., 2015) is taken as the baseline. As it shows, compared with the baseline, the overall performance of fusing multi-scale contexts in the parallel stack (see Fig. 2(c)) only improves slightly. By contrast, there is an improvement of near 3% on mean IoU when our approach of self-cascaded fusion is adopted. Moreover, when residual correction scheme is dedicatedly employed in each position behind multi-level contexts fusion, the performance improves even more. These improvements further demonstrate the effectiveness of our multi-scale contexts aggregation approach and residual correction scheme. As can be seen, the performance of each category indeed improves when successive refinement strategy is added, but it doesn't seem to work very well. However, when residual correction scheme is elaborately applied to correct the latent fitting residual in multi-level feature fusion, the performance improves once more, especially for the car.

To evaluate the effect of *transfer learning* (Yosinski et al., 2014; Penatti et al., 2015; Hu et al., 2015; Xie et al., 2015), which is used for training ScasNet, the quantitative performance brought by initializing the encoder's parameters (see Fig. 1) with pre-trained model (i.e., finetuning) are listed in Table 8. As it shows, the performance of VGG ScasNet improves slightly, while ResNet ScasNet improves significantly. These results indicate that, it is very difficult to train deep models sufficiently with so small remote sensing datasets, especially for the very deep models, e.g., the model based on 101-layer ResNet. Therefore, the ScasNet benefits from the widely used *transfer learning* in the field of *deep learning*.

Table 9: Complexity comparison (%) with the state-of-the-art deep models.

| | SegNet | FCN-8s | DeconvNet | Deeplab-ResNet | RefineNet | Ours-VGG | Ours-ResNet |
|---|---|---|---|---|---|---|---|
| Model size | 112M | 512M | 961M | 503M | 234M | 151M | 481M |
| Time | 17s | 11s | 18s | 47s | 21s | 11s | 33s |

To further verify the validity of each aspect of our ScasNet, features of some key layers in VGG ScasNet are visualized in Fig. 13. For clarity, we only visualize part of features in the last layers before the pooling layers, more detailed visualization can be referred in the Appendix B of supplementary material. As shown in Fig. 13(a) and (b), the 1st-layer convolutional filters tend to learn more meaningful features after funetuning, which indicates the validity of *transfer learning*. As Fig. 13(c) and (d) indicate, the layers of the first two stages tend to contain a lot of noise (e.g., too much littery texture), which could weaken the robustness of ScasNet. That is a reason why they are not incorporated into the refinement process.

As can be seen in Fig. 13(e), the responses of feature maps outputted by the encoder tend to be quite messy and coarse. However, as shown in Fig. 13(f), coherent and intact semantic responses can be obtained when our multi-scale contexts aggregation approach is used. Moreover, as Fig. 13(g) shows, much low-level details are recovered when our refinement strategy is used. The boundary responses of cars and trees can be clearly seen.

Fig. 13(h), (i) and (j) visualize the fused feature maps before residual correction, the feature maps learned by inverse residual mapping $\mathcal{H}[\cdot]$ (see Fig. 4) and the fused feature maps after residual correction, respectively. As Fig. 13(h) shows, there is much information lost when two feature maps with semantics of different levels are fused. Nevertheless, as shown in Fig. 13(j), these deficiencies are mitigated significantly when our residual correction scheme is employed. That is, as Fig. 13(i) shows, the inverse residual mapping $\mathcal{H}[\cdot]$ could compensate for the lack of information, thus counteracting the adverse effect of the latent fitting residual in multi-level feature fusion.

Table 9 compares the complexity of ScasNet with the state-of-the-art deep models. The time complexity is obtained by averaging the time to perform single scale test on 5 images (average size of $2392 \times 2191$ pixels) with a GTX Titan X GPU. As it shows, ScasNet produces competitive results on both space and time complexity.

## 5. Conclusion

In this work, a novel end-to-end self-cascaded convolutional neural network (ScasNet) has been proposed to perform semantic labeling in VHR images. The proposed ScasNet achieves excellent performance by focusing on three key aspects: 1) A self-cascaded architecture is proposed to sequentially aggregate global-to-local contexts, which are very effective for confusing manmade objects recognition. Technically, multi-scale contexts are first captured on the output of a CNN encoder, and then they are successively aggregated in a self-cascaded manner; 2) With the acquired contextual information, a coarse-to-fine refinement strategy is proposed to progressively refine the target objects using the low-level features learned by CNN's shallow layers. Therefore, the coarse labeling map is gradually refined, especially for intricate fine-structured objects; 3) A residual correction scheme is proposed for multi-feature fusion inside ScasNet. It greatly corrects the latent fitting residual caused by the semantic gaps in features of different levels, thus further improves the performance of ScasNet. As a result of these specific designs, ScasNet can perform semantic labeling effectively in a manner of global-to-local and coarse-to-fine.

Extensive experiments verify the advantages of ScasNet: 1) On both quantitative and visual performances, ScasNet achieves extraordinarily more coherent, complete and accurate labeling results while remaining better robustness to the occlusions and cast shadows than all the comparing advanced deep models; 2) ScasNet outperforms the state-of-the-art methods on two challenging benchmarks by the date of submission: *ISPRS 2D Semantic Labeling Challenge* for Vaihingen and Potsdam, even not using the available elevation data, model ensemble strategy or any postprocessing; 3) ScasNet also shows extra advantages on both space and
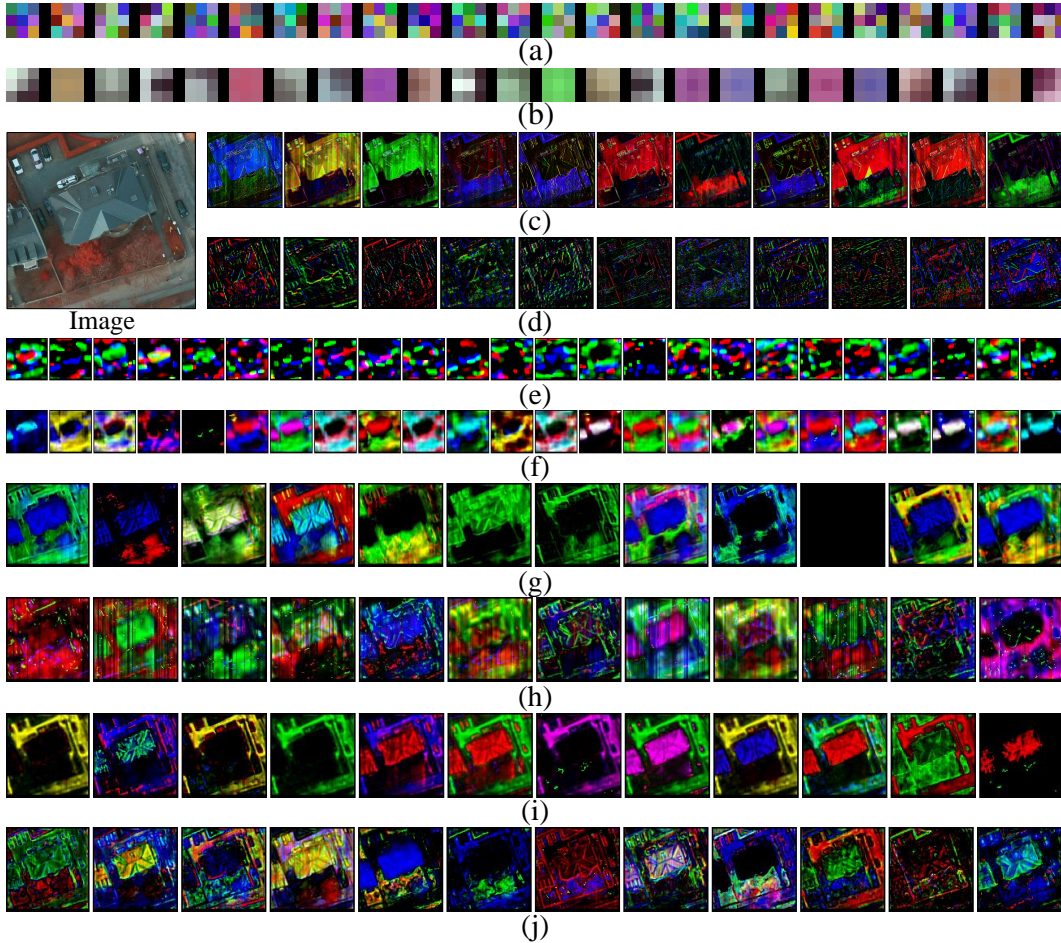
Figure 13: Filters and feature maps learned by VGG ScasNet. For better visualization, the visuals are mapped to full channel range and combined in some cases to occupy all RGB channels. In these features, colored regions denote strong responses and deep black regions otherwise. (a) The 1st-layer convolutional filters ($3 \times 3 \times 3 \times 64$) without finetuning. (b) The 1st-layer convolutional filters ($3 \times 3 \times 3 \times 64$) with finetuning. (c) The last convolutional feature maps in the 1st stage ($400 \times 400 \times 64$). (d) The last convolutional feature maps in the 2nd stage ($201 \times 201 \times 128$). (e) Feature maps outputted by the encoder ($51 \times 51 \times 512$). (f) Feature maps after our multi-scale contexts aggregation approach ($51 \times 51 \times 512$). (g) Feature maps after our refinement strategy ($101 \times 101 \times 256$). (h) Fused feature maps before residual correction ($101 \times 101 \times 256$). (i) Feature maps learned by inverse residual mapping $\mathcal{H}[\cdot]$ (see Fig. 4). (j) Fused feature maps after residual correction ($101 \times 101 \times 256$).

time complexity compared with some complex deep models.

# References

Alshehhi, R., Marpu, P. R., Woon, W. L., Mura, M. D., 2017. Simultaneous extraction of roads and buildings in remote sensing imagery with convolutional neural networks. ISPRS Journal of Photogrammetry and Remote Sensing. 130, 139–149.

Audebert, N., Saux, B. L., Lefèvre, S., 2016. Semantic segmentation of earth observation data using multimodal and multi-scale deep networks. arXiv preprint arXiv:1609.06846.

Badrinarayanan, V., Kendall, A., Cipolla, R., 2015. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. arXiv preprint arXiv:1511.00561.

Bell, S., Lawrence Zitnick, C., Bala, K., Girshick, R., 2016. Inside-outside net: Detecting objects in context with skip pooling and recurrent neural networks. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 2874–2883.

Bertasius, G., Shi, J., Torresani, L., 2016. Semantic segmentation with boundary neural fields. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 3602–3610.

Bridle, J. S., 1989. Probabilistic interpretation of feedforward classification network outputs, with relationships to statistical pattern recognition. In: Neurocomputing: Algorithms, Architectures and Applications. Springer.

Chen, L., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A. L., 2015. Semantic image segmentation with deep convolutional nets and fully connected crfs. In: International Conference on Learning Representations.

Chen, L.-C., Yang, Y., Wang, J., Xu, W., Yuille, A. L., 2016a. Attention to scale: Scale-aware semantic image segmentation. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 3640–3649.

Chen, S., Wang, H., Xu, F., Jin, Y.-Q., 2016b. Target classification using the deep convolutional networks for sar images. IEEE Transactions on Geoscience and Remote Sensing. 54 (8), 4806–4817.

Cheng, G., Han, J., 2016. A survey on object detection in optical remote sensing images. ISPRS Journal of Photogrammetry and Remote Sensing. 117, 11–28.

Cheng, G., Han, J., Lu, X., 2017a. Remote sensing image scene classification: Benchmark and state of the art. arXiv preprint arXiv:1703.00121.

Cheng, G., Wang, Y., Xu, S., Wang, H., Xiang, S., Pan, C., 2017b. Automatic road detection and centerline extraction via cascaded end-to-end convolutional neural network. IEEE Transactions on Geoscience and Remote Sensing. 55 (6), 3322–3337.

Cheng, G., Zhu, F., Xiang, S., Wang, Y., Pan, C., 2016. Accurate urban road centerline extraction from vhr imagery via multiscale segmentation and tensor voting. Neurocomputing. 205, 407–420.

Dalal, N., Triggs, B., 2005. Histograms of oriented gradients for human detection. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 886–893.

Everingham, M., Eslami, S. M. A., Gool, L. J. V., Williams, C. K. I., Winn, J. M., Zisserman, A., 2015. The pascal visual object classes challenge: A retrospective. International Journal of Computer Vision. 111 (1), 98–136.

Farabet, C., Couprie, C., Najman, L., LeCun, Y., 2013. Learning hierarchical features for scene labeling. IEEE Transactions on Pattern Analysis and Machine Intelligence. 35 (8), 1915–1929.

Gerke, M., 2015. Use of the stair vision library within the isprs 2d semantic labeling benchmark (vaihingen). Technical Report.

Gidaris, S., Komodakis, N., 2015. Object detection via a multi-region and semantic segmentation-aware cnn model. In: IEEE International Conference on Computer Vision. pp. 1134–1142.

Glorot, X., Bordes, A., Bengio, Y., 2011. Deep sparse rectifier neural networks. In: International Conference on Artificial Intelligence and Statistics. Vol. 15. p. 275.

Gong, M., Yang, H., Zhang, P., 2017. Feature learning and change feature classification based on deep learning for ternary change detection in sar images. ISPRS Journal of Photogrammetry and Remote Sensing. 129, 212–225.

Gould, S., Russakovsky, O., Goodfellow, I., , Baumstarck, P., 2011. The stair vision library (v2.5). Stanford University.

Hariharan, B., Arbeláez, P., Girshick, R., Malik, J., 2015. Hypercolumns for object segmentation and fine-grained localization. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 447–456.

He, K., Zhang, X., Ren, S., Sun, J., 2015a. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: IEEE International Conference on Computer Vision. pp. 1026–1034.

He, K., Zhang, X., Ren, S., Sun, J., 2015b. Spatial pyramid pooling in deep convolutional networks for visual recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence. 37 (9), 1904–1916.

He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 770–778.

Hu, F., Xia, G.-S., Hu, J., Zhang, L., 2015. Transferring deep convolutional neural networks for the scene classification of high-resolution remote sensing imagery. Remote Sensing. 7 (11), 14680–14707.

Ioffe, S., Szegedy, C., 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: International Conference on Machine Learning. pp. 448–456.

ISPRS, 2016. International society for photogrammetry and remote sensing. 2D Semantic Labeling Challenge.
URL http://www2.isprs.org/commissions/comm3/wg4/semantic-labeling.html

Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R. B., Guadarrama, S., Darrell, T., 2014. Caffe: Convolutional architecture for fast feature embedding. In: ACM International Conference on Multimedia. pp. 675–678.

Krizhevsky, A., Sutskever, I., Hinton, G. E., 2012. Imagenet classification with deep convolutional neural networks. In: Neural Information Processing Systems. pp. 1106–1114.

Lecun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., Jackel, L. D., 1990. Handwritten digit recognition with a back-propagation network. In: Neural Information Processing Systems. pp. 396–404.

Lecun, Y., Bottou, L., Bengio, Y., Haffner, P., 1998. Gradient-based learning applied to document recognition. Proceedings of the IEEE. 86 (11), 2278–2324.

Li, E., Femiani, J., Xu, S., Zhang, X., Wonka, P., 2015a. Robust rooftop extraction from visible band images using higher order crf. IEEE Transactions on Geoscience and Remote Sensing. 53 (8), 4483–4495.

Li, J., Huang, X., Gamba, P., Bioucas-Dias, J. M., Zhang, L., Benediktsson, J. A., Plaza, A., 2015b. Multiple feature learning for hyperspectral image classification. IEEE Transactions on Geoscience and Remote Sensing. 53 (3), 1592–1606.

Lin, G., Milan, A., Shen, C., Reid, I. D., 2016. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. arXiv preprint arXiv:1611.06612.

Liu, C., Yuen, J., Torralba, A., Sivic, J., Freeman, W., 2008. Sift flow: Dense correspondence across different scenes. European Conference on Computer Vision., 28–42.

Liu, W., Rabinovich, A., Berg, A. C., 2016a. Parsenet: Looking wider to see better. In: International Conference on Learning Representations Workshop.

Liu, Y., Fan, B., Wang, L., Bai, J., Xiang, S., Pan, C., 2017. Context-aware cascade network for semantic labeling in vhr image. In: IEEE International Conference on Image Processing.

Liu, Y., Zhong, Y., Fei, F., Zhang, L., 2016b. Scene semantic classification based on random-scale stretched convolutional neural network for high-spatial resolution remote sensing imagery. In: IEEE International Geoscience and Remote Sensing Symposium (IGARSS). pp. 763–766.

Long, J., Shelhamer, E., Darrell, T., 2015. Fully convolutional networks for semantic segmentation. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 3431–3440.

Lowe, D. G., 2004. Distinctive image features from scale-invariant keypoints. International Journal of Computer Vision. 60 (2), 91–110.

Lu, X., Yuan, Y., Zheng, X., 2017a. Joint dictionary learning for multispectral change detection. IEEE Transactions on Cybernetics. 47 (4), 884–897.

Lu, X., Zheng, X., Yuan, Y., 2017b. Remote sensing scene classification by unsupervised representation learning. IEEE Transactions on Geoscience and Remote Sensing. PP (99), 1–10.

Luus, F. P., Salmon, B. P., van den Bergh, F., Maharaj, B., 2015. Multiview deep learning for land-use classification. IEEE Geoscience Remote Sensing Letters 12 (12), 2448–2452.

Maggiori, E., Tarabalka, Y., Charpiat, G., Alliez, P., 2017. Convolutional neural networks for large-scale remote-sensing image classification. IEEE Transactions on Geoscience and Remote Sensing. 55 (2), 645–657.

Marmanis, D., Schindler, K., Wegner, J. D., Galliani, S., Datcu, M., Stilla, U., 2016. Classification with an edge: improving semantic image segmentation with boundary detection. arXiv preprint arXiv:1612.01337.

Mas, J. F., Flores, J. J., 2008. The application of artificial neural networks to the analysis of remotely sensed data. International Journal of Remote Sensing. 29 (3), 617–663.

Matikainen, L., Karila, K., 2011. Segment-based land cover mapping of a suburban area-comparison of high-resolution remotely sensed datasets using classification trees and test field points. Remote Sensing. 3 (8), 1777–1804.

Mnih, V., 2013. Machine learning for aerial image labeling. Ph.D. thesis, University of Toronto.

Mostajabi, M., Yadollahpour, P., Shakhnarovich, G., 2015. Feedforward semantic segmentation with zoom-out features. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 3376–3385.

Nair, V., Hinton, G. E., 2010. Rectified linear units improve restricted boltzmann machines. In: International Conference on Machine Learning. pp. 807–814.

Nogueira, K., Mura, M. D., Chanussot, J., Schwartz, W. R., dos Santos, J. A., 2016. Learning to semantically segment high-resolution remote sensing images. In: IEEE International Conference on Pattern Recognition. pp. 3566–3571.

Noh, H., Hong, S., Han, B., 2015. Learning deconvolution network for semantic segmentation. In: IEEE International Conference on Computer Vision. pp. 1520–1528.

Paisitkriangkrai, S., Sherrah, J., Janney, P., van den Hengel, A., 2016. Semantic labeling of aerial and satellite imagery. IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing. 9 (7), 2868–2881.

Penatti, O. A., Nogueira, K., dos Santos, J. A., 2015. Do deep features generalize from everyday objects to remote sensing and aerial scenes domains. In: IEEE Conference on Computer Vision and Pattern Recognition Workshop. pp. 44–51.

Pinheiro, P. O., Lin, T.-Y., Collobert, R., Dollár, P., 2016. Learning to refine object segments. arXiv preprint arXiv:1603.08695.

Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention - MICCAI. pp. 234–241.

Rumelhart, D. E., Hinton, G. E., Williams, R. J., 1986. Learning representations by back-propagating errors. Nature. 323 (6088), 533–536.

Sherrah, J., 2016. Fully convolutional networks for dense semantic labelling of high-resolution aerial imagery. arXiv preprint arXiv:1606.02585.

Silberman, N., Hoiem, D., Kohli, P., Fergus, R., 2012. Indoor segmentation and support inference from rgbd images. In: European Conference on Computer Vision. pp. 746–760.

Simonyan, K., Zisserman, A., 2015. Very deep convolutional networks for large-scale image recognition. In: International Conference on Learning Representations. pp. 1–14.

Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., Salakhutdinov, R., 2014. Dropout: a simple way to prevent neural networks from overfitting. Journal of Machine Learning Research. 15 (1), 1929–1958.

Volpi, M., Tuia, D., 2017. Dense semantic labeling of subdecimeter resolution images with convolutional neural networks. IEEE Transactions on Geoscience and Remote Sensing. 55 (2), 881–893.

Wen, D., Huang, X., Liu, H., Liao, W., Zhang, L., 2017. Semantic classification of urban trees using very high resolution satellite imagery. IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing. 10 (4), 1413–1424.

Xie, M., Jean, N., Burke, M., Lobell, D., Ermon, S., 2015. Transfer learning from deep features for remote sensing and poverty mapping. arXiv preprint arXiv:1510.00098.

Xu, X., Li, J., Huang, X., Mura, M. D., Plaza, A., 2016. Multiple morphological component analysis based decomposition for remote sensing image classification. IEEE Transactions on Geoscience and Remote Sensing. 54 (5), 3083–3102.

Xue, Z., Li, J., Cheng, L., Du, P., 2015. Spectralspatial classification of hyperspectral data via morphological component analysis-based image separation. IEEE Transactions on Geoscience and Remote Sensing. 53 (1), 70–84.

Yosinski, J., Clune, J., Bengio, Y., Lipson, H., 2014. How transferable are features in deep neural networks. In: Neural Information Processing Systems. pp. 3320–3328.

Yu, F., Koltun, V., 2016. Multi-scale context aggregation by dilated convolutions. In: International Conference on Learning Representations.

Yuan, Y., Mou, L., Lu, X., 2015. Scene recognition by manifold regularized deep learning architecture. IEEE Transactions on Neural Networks and Learning Systems. 26 (10), 2222–2233.

Zeiler, M. D., Fergus, R., 2014. Visualizing and understanding convolutional networks. In: European Conference on Computer Vision. pp. 818–833.

Zeiler, M. D., Krishnan, D., Taylor, G. W., Fergus, R., 2010. Deconvolutional networks. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 2528–2535.

Zhang, C., Pan, X., Li, H., Gardiner, A., Sargent, I., Hare, J., Atkinson, P. M., 2017. A hybrid mlp-cnn classifier for very fine resolution remotely sensed image classification. ISPRS Journal of Photogrammetry and Remote Sensing. pp, 1–12.

Zhang, L., Zhang, L., Tao, D., Huang, X., 2012. On combining multiple features for hyperspectral remote sensing image classification. IEEE Transactions on Geoscience and Remote Sensing. 50 (3), 879–893.

Zhang, P., Gong, M., Su, L., Liu, J., Li, Z., 2016. Change detection based on deep feature representation and mapping transformation for multi-spatial-resolution remote sensing images. ISPRS Journal of Photogrammetry and Remote Sensing. 116, 24–41.

Zhang, Q., Seto, K. C., 2011. Mapping urbanization dynamics at regional and global scales using multi-temporal dmsp/ols nighttime light data. Remote Sensing of Environment 115 (9), 2320–2329.

Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J., 2016. Pyramid scene parsing network. arXiv preprint arXiv:1612.01105.

Zhao, W., Du, S., 2016. Learning multiscale and deep representations for classifying remotely sensed imagery. ISPRS Journal of Photogrammetry and Remote Sensing. 113, 155–165.

Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., , Torralba, A., 2015. Object detectors emerge in deep scene cnns. In: International Conference on Learning Representations.
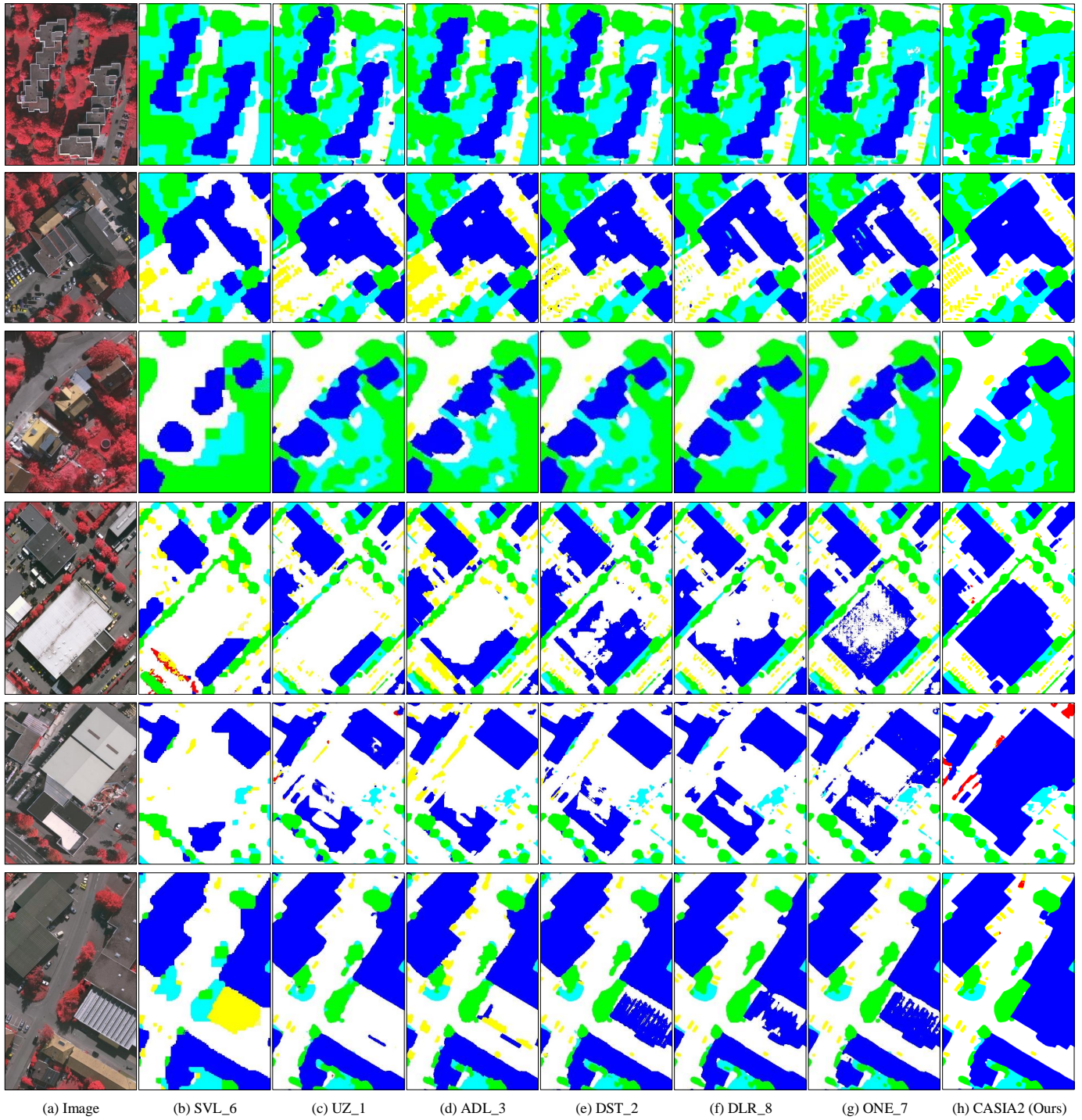
Figure 14: Qualitative comparison with other competitors' methods on *ISPRS Vaihingen challenge* ONLINE TEST SET. The label includes six categories: impervious surface (imp surf, white), building (blue), low vegetation (low veg, cyan), tree (green), car (yellow) and clutter/background (red).
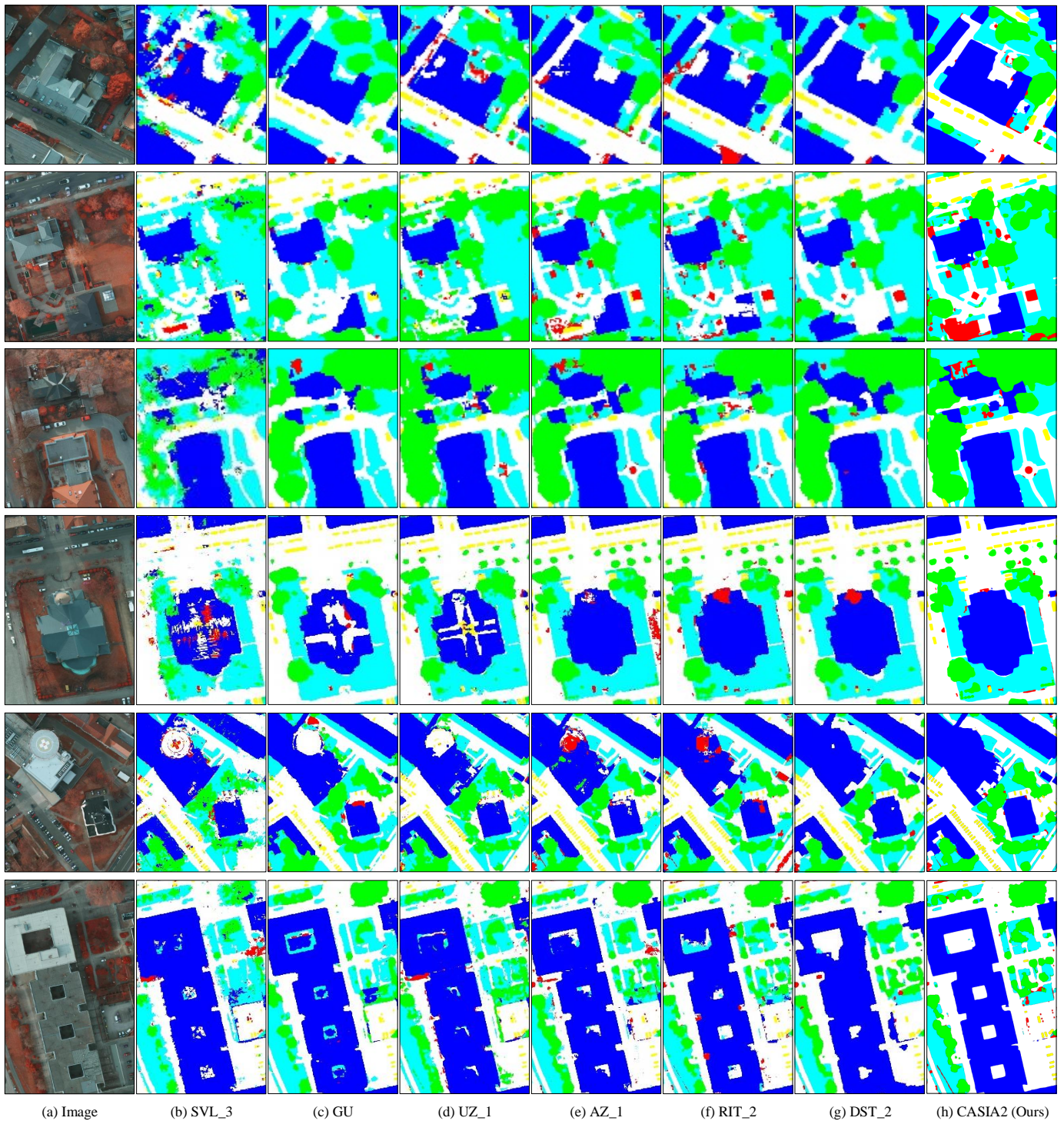
Figure 15: Qualitative comparison with other competitors' methods on *ISPRS Potsdam challenge* ONLINE TEST SET. The label includes six categories: impervious surface (imp surf, white), building (blue), low vegetation (low veg, cyan), tree (green), car (yellow) and clutter/background (red).