

Recurrent Multiresolution Convolutional Networks for VHR Image Classification

John Ray Bergado, Claudio Persello, *Senior Member, IEEE*, and Alfred Stein

Abstract—Classification of very high resolution (VHR) satellite images has three major challenges: 1) inherent low intra-class and high inter-class spectral similarities, 2) mismatching resolution of available bands, and 3) the need to regularize noisy classification maps. Conventional methods have addressed these challenges by adopting separate stages of image fusion, feature extraction, and post-classification map regularization. These processing stages, however, are not jointly optimizing the classification task at hand. In this study, we propose a single-stage framework embedding the processing stages in a recurrent multiresolution convolutional network trained in an *end-to-end* manner. The feedforward version of the network, called *FuseNet*, aims to match the resolution of the panchromatic and multispectral bands in a VHR image using convolutional layers with corresponding downsampling and upsampling operations. Contextual label information is incorporated into *FuseNet* by means of a recurrent version called *ReuseNet*. We compared *FuseNet* and *ReuseNet* against the use of separate processing steps for both image fusion, e.g. pansharpening and resampling through interpolation, and map regularization such as *conditional random fields*. We carried out our experiments on a land cover classification task using a Worldview-03 image of Quezon City, Philippines and the ISPRS 2D semantic labeling benchmark dataset of Vaihingen, Germany. *FuseNet* and *ReuseNet* surpass the baseline approaches in both quantitative and qualitative results.

Index Terms—Convolutional networks, recurrent networks, land cover classification, VHR image, deep learning.

I. INTRODUCTION

CLASSIFICATION of very high resolution (VHR) remotely sensed images allows us to automatically produce maps at a level of detail comparable to conventional in-situ mapping methods. Due to the high spatial resolution of such images, automated classification comes with a set of challenges. One challenge is the inherent low intra-class and high inter-class spectral similarities, inhibiting discrimination of the classes of interest. Conventional methods address this challenge by extracting spatial-contextual features from the image such as texture-describing measures, e.g. gray level co-occurrence matrix (GLCM) [1] and local binary patterns (LBP) [2], or products of morphological operators [3], [4]. This step is crucial for obtaining discriminative features and accurate classification. However, such feature extraction methods are often disjoint from the supervised classifier, and, hence, not optimized for the task at hand. Deep learning offers a framework to build end-to-end classifiers—directly learning the predictions from the inputs with minimal or no pre-classification and post-classification steps. Features automatically extracted by deep learning based classifiers such

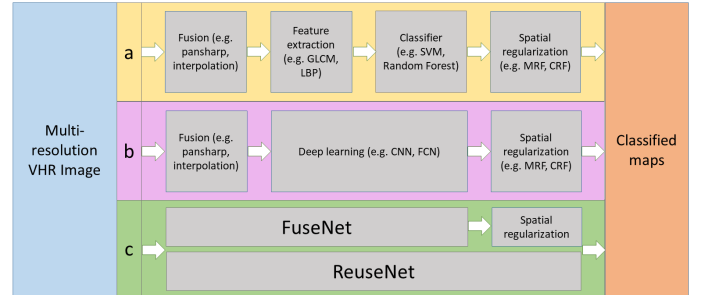


Fig. 1. Illustration comparing a standard (a), state-of-the-art (b), and proposed (c) pipelines for classifying multiresolution VHR images.

as convolutional neural networks (CNN) [5] perform better than intermediate handcrafted features [6], [7]. These networks automatically learn spatial-contextual features directly from the input VHR image—effectively integrating the feature extraction step into the training of the classifier as shown in Figure 1 (b). The design of network architecture, inspired by the model of the visual cortex [8], makes CNN suitable for image analysis and land cover classification.

Another challenge in the classification of VHR images is the multiresolution nature of the images acquired by space-borne sensors. Most VHR satellite images (e.g. Quickbird, Worldview, IKONOS, and Pleiades) capture panchromatic (PAN) images in a spatial resolution four times of the multispectral (MS) bands. Such mismatch in spatial resolution of the images requires an additional step to fuse these images before performing the semantic analysis. Pansharpening and interpolation-based resampling are common techniques for fusing a multiresolution image [9]. Similar to conventional feature extraction methods, the operations to fuse multiresolution bands of a VHR image add another separate pre-classification step that is disjoint from the training of the supervised classifier, and, hence, may not be optimal for the task at hand. CNN extracts a hierarchy of spatial features at different resolutions. We can exploit the multiresolution nature of the VHR data to design a CNN that performs fusion and feature extraction at the same time.

Literature shows that classification accuracy can be improved by using post-classification spatial regularization [10], [11], [12]. Methods employing graphical models, such as conditional random fields (CRF) and Markov random fields (MRF), provide a way to perform this spatial regularization step. Similar to the two pre-classification steps described above, a post-classification map regularization technique adds another step independent of the training of the classifier

Acknowledgement here.

Manuscript received Month DD, 2017; revised Month DD, 2017.

itself—further including a separate objective function to be optimized. For classifying a multiresolution VHR image, a typical classification pipeline would be composed of three main stages: a pre-classification step performing image fusion and feature extraction, a supervised learning algorithm performing the classification, and a post-classification step regularizing the maps obtained from the supervised classification algorithm. This conventional approach is shown in Figure 1 (a).

Convolutional networks have been recently applied to classify remotely sensed images with very high resolution [11], [13], [14], [7], [15], [16], [17]. But, aside from [18] which used a combination of patch-based CNN and stacked autoencoders to fuse PAN and MS images, the majority of the works did not address the problem of multiresolution VHR images. A patch-based CNN [7] and a fully convolutional network (FCN) [15] were used to detect informal settlements from a pansharpened VHR image. Fully convolutional networks were also used to classify urban objects in VHR images both acquired in aerial and space-borne sensors [11], [13], [14], [16]. Moreover, [11], [13], [17] also utilized a separate post-classification step for map regularization. In this paper, we design a novel single-stage network performing image fusion, classification, and map regularization of a multiresolution VHR image in an end-to-end manner.

A. Contributions

We propose a multiresolution convolutional network, called FuseNet, and its recurrent version, called ReuseNet, to perform image fusion, classification, and map regularization of a multiresolution VHR image in an end-to-end fashion. We summarize the main contributions of this paper in: image fusion, map regularization, and sensitivity analysis of network parameters.

1) *Image fusion*: We propose a convolutional network learning how to fuse a multiresolution VHR image, extract spatial features, and classify the latter into classes of interest all at the same time. We call this network FuseNet. It uses convolutional layers with corresponding downsampling and upsampling operations to learn to match and fuse the multiresolution channels of a multispectral VHR image.

2) *Contextual label dependency through network recurrence*: We incorporate recurrence in the FuseNet architecture to model contextual *label-to-label* dependencies and effectively regularize classification maps. We call this improved version ReuseNet. Contextual label dependencies are incorporated in ReuseNet by feeding classification scores of a previous FuseNet instance to a succeeding one. Moreover, we introduce and compare a novel method to initialize the parameters and initial score maps of a ReuseNet.

3) *Sensitivity analysis*: We analyze the sensitivity of the network to some of its chosen hyperparameters. We investigate the effect of varying a number of hyperparameters of our network to the classification performance. The considered hyperparameters are: the bottleneck feature map dimensions, the number of convolutional layers, the input patch sizes, the upsampling operations, and the number of FuseNet instances within a ReuseNet.

II. CONVOLUTIONAL NETWORKS

A. Background

Convolutional neural networks are a variant of artificial neural networks connected in a sequential feedforward fashion employing convolutional and pooling (aggregational) operations. Convolutional operations greatly reduce the number of learnable parameters and allows the network to use the same filter to detect the same spatial pattern over different parts of an image. Pooling with downsampling enables the network to learn some degree of translational invariance.

Recurrent neural networks are artificial neural networks employing feedback connection, i.e. connections form a directed cycle. For example, the Jordan network [19] has connections from the output units back to the hidden units. Two key concepts namely, parameter sharing and graph unfolding [20, pp. 369–372], allow these networks to accept input sequences of variable lengths while maintaining model complexity—making recurrent networks widely applied to sequential data. However, parameter sharing and graph unfolding can also be used to design networks for different purposes, e.g. application to non-sequential data, while still taking advantage of the benefits, such as model compactness, from the two concepts [21].

B. Deep Networks as Data-flow Graphs

We can generalize any variant of deep networks by seeing them as data-flow graphs—a graph representing how a set of input data are processed along a possibly branching chain of functions, in the end producing a final set of outputs. Using such a model, we define the networks by three elements: the sets of data they take as an input, the operations they perform in each function block, and the intermediate and final set of outputs they produce. Aside from these three key elements of data-flow graphs, details of a unique configuration and instance of a convolutional network are defined by its *hyperparameters* and *parameters* respectively. Hyperparameters are associated with the configuration of a network architecture and are set to fixed values before training the network. Parameters are values associated to a specific network instance and are learned during network training.

1) *Input*: A convolutional network receives as an input either the whole image itself to be classified or a subset of it, called an input patch. The dimension of this patch is defined by the patch size hyperparameter M and the number of bands B . A convolutional network accepts an $N \times B \times M \times M$ array of pixel values as an input (in the case of the image patches having equal height and width), N being the number of patches processed by the network in parallel. Aside from the input image patch, the corresponding reference image can also be considered as an input in terms of data-flow graphs since no operation precedes it.

2) *Operations*: *Convolutions* are the main operations used by convolutional networks. A convolution applies a linear operation on an input image/feature map using a set of K' learnable kernels. Applying a kernel \mathbf{w} , composed of a $K \times K' \times G \times G$ array of learnable parameters, on a $K \times H \times W$ input feature map \mathbf{x} , where G is the kernel size, K is the

number of kernels in each set of kernels, and H and W are the height and width of the feature map, produces a $K' \times H' \times W'$ output feature map \mathbf{x}' . The output at the i' row and j' column of the k' feature map is given by:

$$\mathbf{x}'_{k'i'j'} = \sum_{k=1}^K \sum_{p=1}^G \sum_{q=1}^G \mathbf{x}_{kij} \cdot \mathbf{w}_{kk'pq} + b_{k'} \quad (1a)$$

$$i = i' + p - \lceil \frac{G}{2} \rceil \quad (1b)$$

$$j = j' + q - \lceil \frac{G}{2} \rceil \quad (1c)$$

where $b_{k'}$ is the learnable bias parameter associated with the k' feature map. The width and height of the output feature map are given by:

$$H' = \lfloor \frac{H - G + 2Z}{S} + 1 \rfloor \quad (2a)$$

$$W' = \lfloor \frac{W - G + 2Z}{S} + 1 \rfloor \quad (2b)$$

where the *zero-padding* Z is the number of rows and columns of zeros added to the border of the input feature map and the convolutional *stride* S is the number of units separating contiguous receptive fields of the kernel on the input feature map.

Nonlinearity is applied after the linear operation of a convolution. Since applying a series of linear operations can be reduced to a single linear operation, an elementwise nonlinear function applied between each convolution allows the network to learn more complex input to output mapping. A common choice is the rectifier function

$$\mathbf{x}'_{i'j'k'} = \max(0, \mathbf{x}_{ijk}) \quad (3)$$

or a variation of it [22], [23], [24].

Pooling takes an aggregate of values over local regions of the input. A common choice of a pooling function is the average or maximum function. In contrast to convolution, a basic pooling does not have any learnable parameters. Originally, pooling was used to give the network a small degree of translation invariance by summarizing values of the input on non-overlapping windows ($S = G$)—also downsampling the input by a factor of S , with proper zero-padding.

Upsampling operations are applied to increase the spatial dimensions of input feature maps. Upsampling is important specifically if the network needs to produce output predictions of the same size as the input, i.e. we want to produce a label for each pixel in the $M \times M$ input patch. One way to upsample is by employing resampling techniques such as nearest neighbor or bilinear interpolation [10]. The original fully convolutional network (FCN) [25] learns the upsampling operation using *backwards convolution* (or more technically fitting called *transposed convolution*).

Merging combines two or more sets of feature maps in a network either by addition or by concatenation. Addition is an elementwise operation performed between feature maps—adding each unit with corresponding indices—hence, all the three dimensions (K , H , W) must be the same for all inputs [25]. Concatenation stacks the input feature maps depthwise—hence, only the spatial dimensions (H , W) must be the same.

3) *Outputs*: In data-flow graph terms, the outputs of a convolutional network consist of all the intermediate feature maps, the final class score maps, and the corresponding loss and accuracy calculated using the class score maps and the reference labels. Final class score maps correspond to the units in the last layer of a neural network and its dimension depends on how the task is defined. Authors in [16] categorize the approaches to this task into three variants: 1) patch classification, 2) subpatch labeling, and 3) full patch labeling. In patch classification, we assign a single label to the patch, i.e. the label corresponds to the class of the central pixel of the patch [6], [16], [7]. In subpatch labeling, we assign labels on a smaller part of the patch corresponding to the area near the center of the patch [16]. Finally, in full patch labeling, we assign labels to all the pixels in the patch [25], [13], [26], [16], [15]. The last method, aside from being more efficient, also decouples the limit of the input patch size to the number of downsampling operations in the network.

C. Training Deep Networks

We train the network by minimizing an objective function in terms of the parameters of the network. For classification involving C classes, a cross-entropy loss function is often used given by:

$$E_N(\mathbf{w}) = - \sum_{n=1}^N \mathbf{t}_n \cdot \log(\mathbf{y}_n) \quad (4)$$

where E is the loss function value evaluated over N samples, \mathbf{t}_n is a binary vector encoding the target class labels (with the index corresponding to a class having a value of 1 and 0 otherwise), \cdot denotes the dot product, and \mathbf{y}_n is the class score maps of a sample n calculated using a *softmax* activation function:

$$\mathbf{y}_{kij} = \frac{\exp(\mathbf{x}_{kij})}{\sum_{c=1}^C \exp(\mathbf{x}_{cij})} \quad (5)$$

In this equation, \mathbf{y} is the softmax score and \mathbf{x} is the last set of feature maps containing unnormalized class scores at location ij .

A stochastic version of the backpropagation with gradient descent algorithm is often used to minimize the objective function [27]. The training is finished after a fixed number of epoch or when a certain convergence criterion is met. We can infer predictions from the final trained network instance by truncating the loss evaluation in the computational graph and taking the index of the maximum class score map value along the class score dimension by

$$y_{ij} = \operatorname{argmax}_c \mathbf{y}_{cij} \quad (6)$$

where \mathbf{y} and y are the class score and prediction for location ij respectively.

D. Regularizing Deep Networks

Deep networks are often prone to overfit the training set. Overfitting occurs when a model reports high accuracy during

training but performs poorly on unseen test data. Regularization approaches address the overfitting problem using three common methods: *data augmentation*, *weight decay*, and *early-stopping*. Data augmentation technique increases the number of training samples by permuting them with applicable rotational and/or translational transformations. Data augmentation helps the network to learn relevant invariances that may be present in the input. Weight decay modifies the loss function by

$$Q(\mathbf{w}) = E(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \quad (7)$$

adding a penalty proportional to the square of the l_2 -norm of the weight vector \mathbf{w} . The weight decay λ hyperparameter controls the contribution of this penalty to the loss function. Early stopping prematurely stops the training when a criterion measured from a validation set is met.

III. PROPOSED APPROACH

In this paper, we propose a multiresolution convolutional network, called FuseNet, and its recurrent version, called ReuseNet, to perform an end-to-end fusion, classification, and map regularization of a multi-resolution VHR image. ReuseNet is built on top of a fully convolutional network architecture learning to: 1) fuse PAN and MS bands of a VHR image, 2) perform land cover classification on the fused images, and 3) spatially regularize the resulting classification.

A. FuseNet

The architecture of FuseNet is inspired by several encoder-decoder like convolutional network structures [28], [26], [16] where the first set of layers learn deep features by a series of convolution, nonlinearity, and maximum pooling with downsampling operations, followed by a second set of layers using upsampling and nonlinearity operations to restore the resolution of the original input image. The main difference of FuseNet with these encoder-decoder architectures is the two initial separate streams of the downsampling part of the network that learns how to fuse two images of different resolution. FuseNet is specifically designed for VHR satellite images with PAN band and MS bands of ground sampling distance ratio of four (e.g. Quickbird, Worldview 2/3, Pleiades, Ikonos). But the architecture can be further generalized to fuse any number of images with different spatial resolutions.

FuseNet accepts two sets of input: an image patch \mathbf{x}_{PAN} of dimensions $N \times 1 \times 4M \times 4M$ taken from PAN image and another patch \mathbf{x}_{MS} of dimensions $N \times 4 \times M \times M$ taken from the corresponding location in the MS image. It performs two series of convolution, nonlinearity, and maximum pooling with downsampling to \mathbf{x}_{PAN} such that the spatial dimensions of the intermediate feature maps match the spatial dimensions of \mathbf{x}_{MS} . The nonlinearity operations use an exponential linear activation function [24]. The second input is linearly projected in k dimensions using 1×1 convolutions such that k matches the number of intermediate feature maps extracted from the first set of input—ensuring a balanced contribution from the two streams of feature. FuseNet then merges the linear

projection of \mathbf{x}_{MS} with intermediate feature maps extracted from \mathbf{x}_{PAN} via a concatenation operation.

Additional series of convolution, nonlinearity, and maximum pooling with downsampling operations are applied to the merged feature maps producing the set of feature maps with smallest spatial dimensions—called *bottleneck* [16]. FuseNet then upsamples the bottleneck back to the resolution of \mathbf{x}_{PAN} using transposed convolutions. The resulting set of feature maps is linearly projected again using 1×1 convolutions such that the number of feature maps matches C . The final class score maps \mathbf{y}_{scores} are obtained by applying a softmax activation function. This series of operations can be formulated as a function composition given by:

$$\mathbf{y}_{scores} = s(l_1(u(d_1(d_0(\mathbf{x}_{PAN}) \oplus l_0(\mathbf{x}_{MS})))))) \quad (8)$$

where d_i is a series of convolution, nonlinearity, and maximum pooling with downsampling operations, u is the series of upsampling operations via transposed convolution, l_i are the linear projections via 1×1 convolutions, s is the softmax function, and \oplus denotes merging via concatenation. Details of each operation, including the hyperparameter values and dimensions of some chosen intermediate feature maps, are provided in Table I. A cross-entropy function following Equation 4 is used to compute the loss in each iteration. Unlabeled pixels are assigned a loss function value of zero.

We described above the default configuration of FuseNet, called FuseNet_{low}, performing the fusion at the lower (MS image) resolution. We also tested a network, called FuseNet_{skip}, adding skip connections to some lower-level feature maps of FuseNet_{low} [25]. Figure 2 shows a diagram illustrating the general architecture of FuseNet_{skip}. Additionally, we experimented with a FuseNet performing the fusion at the resolution of the PAN image, called FuseNet_{high} which is more similar to pansharpening—upsampling \mathbf{x}_{MS} first before fusing them with \mathbf{x}_{PAN} . Tables I and II show details of the operations, including dimensions of intermediate output feature maps, used by the FuseNet variants. The format is adapted from [29]. \mathbf{x}_{PAN} and \mathbf{x}_{MS} denote input patches from the PAN and MS images respectively. IFM and BFM correspond to intermediate and bottleneck feature maps. Convolutions are denoted as “conv(kernel size G)-(number of kernels K)”. Maximum pooling operations (maxpool) are fixed to have pooling size G_p and stride S_p equal to two. Upsampling operations are denoted as “ups(number of kernels K)-(upsampling factor)”. Merging operations are denoted as concat for concatenation and add for addition. Fixed upsampling can either be via pansharpening or bilinear interpolation. Consecutive Batch Normalization [30] and exponential linear activation [24] operations between convolutions and pooling are omitted for brevity. Finally, operations shared by separate streams of feature align with the columns of these streams.

FuseNet implements a full patch labeling approach since it produces labeled image patches of the same dimensions as the input PAN image patch. Inference of final classification map is given by Equation 6 and can be applied to an input image of variable spatial dimension. Application to input of variable size is made possible by the fully-convolutional nature of the network [25]—allowing it to be applied as an image filter [13]

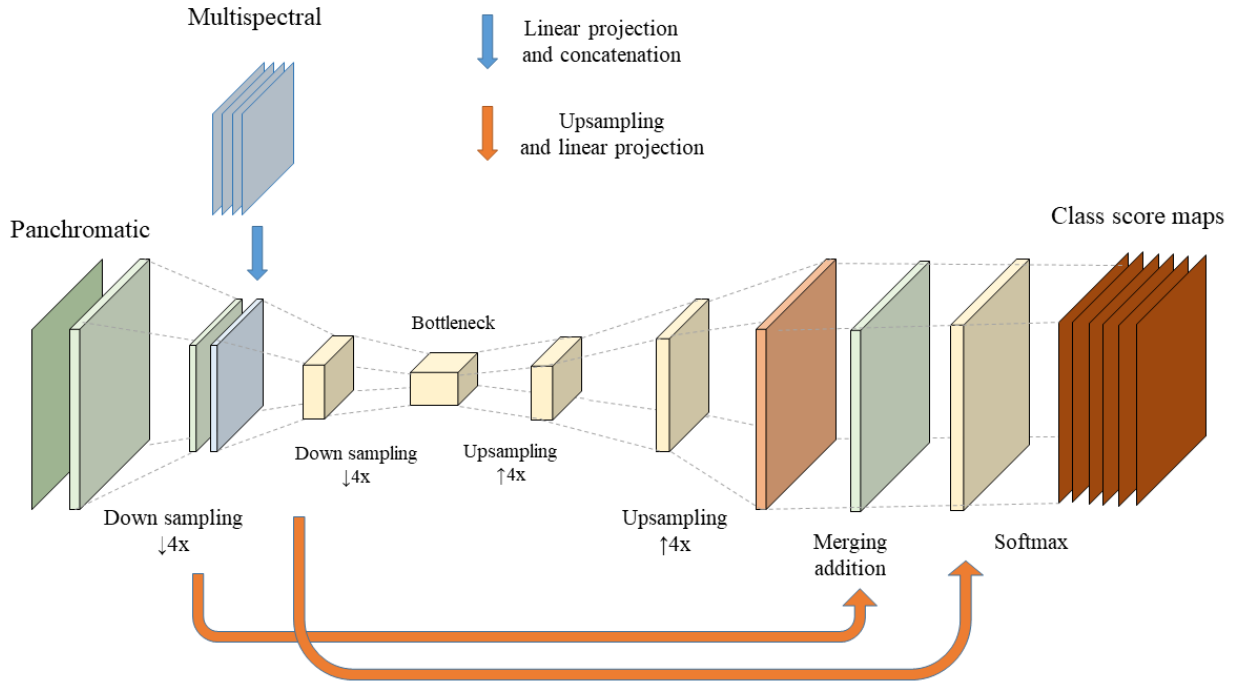


Fig. 2. The general architecture of FuseNet with skip connections (FuseNet_{skip}). FuseNet accepts two streams of input: one from PAN image patches and another from MS image patches. It applies convolutional and pooling layers with downsampling to extract spatial features and at the same time match the resolution of the two streams of input. Similar operations are performed to the output of the merged streams of input arriving at a feature map with smallest spatial dimensions (bottleneck). From there, upsampling operations using transposed convolutions are performed to restore the resolution back to the resolution of the PAN image patches. Skip connections are implemented using appropriate upsampling and linear projections.

TABLE I
DETAILED OPERATIONS OF FUSENET_{low}, FUSENET_{high}, AND FUSENET_{pansharp/bilinear}.

FuseNet _{low}		FuseNet _{high}		Net _{pansharp/bilinear}	
$\mathbf{x}_{PAN}(1 \times 4M \times 4M)$	$\mathbf{x}_{MS}(4 \times M \times M)$	$\mathbf{x}_{PAN}(1 \times 4M \times 4M)$	$\mathbf{x}_{MS}(4 \times M \times M)$	$\mathbf{x}_{PAN}(1 \times 4M \times 4M)$	$\mathbf{x}_{MS}(4 \times M \times M)$
conv13-16 maxpool conv7-32 maxpool	conv1-32		ups2-16 ups2-8 conv1-4 IFM1 ($5 \times 4M \times 4M$)		fixed upsampling
		concat		concat	
		IFM3 ($5 \times 4M \times 4M$)		IFM3 ($4 \times 4M \times 4M$)	
IFM1 ($32 \times M \times M$)	IFM2 ($32 \times M \times M$)		conv13-16 maxpool conv32-7 maxpool		
	concat				
	IFM3 ($64 \times M \times M$)				
		conv3-64 maxpool conv3-128 maxpool			
		BFM ($128 \times M/4 \times M/4$)			
		ups2-128 ups2-64 ups2-32 ups2-16 conv1-6			
		IFM4 ($6 \times 4M \times 4M$)			
		softmax			

to any input with spatial dimensions of at least equal to the FCN's effective receptive field.

B. ReuseNet

ReuseNet builds on top of the architecture of FuseNet by incorporating recurrent connections. Incorporation of this recurrent architecture in a full patch labeling approach enables

the network to learn contextual label-to-label dependencies by feeding output score maps of a FuseNet instance to another instance of itself as an input. Such dependencies are similar to what graphical model (e.g. CRF/MRF) based methods learn in a post-classification regularization inference. For instance, a fully-connected CRF [31] solves an energy function that penalizes label configurations based on a unary term, often

TABLE II
DETAILED OPERATIONS OF FUSENET_{skip}.

FuseNet _{skip}		
$\mathbf{x}_{PAN}(1 \times 4M \times 4M)$	$\mathbf{x}_{MS}(4 \times M \times M)$	Skip-connected layers
conv13-16 maxpool conv7-32 maxpool	conv1-32	
IFM1 ($32 \times M \times M$)	IFM2 ($32 \times M \times M$)	
concat		
IFM3 ($64 \times M \times M$)		
conv3-64 maxpool conv3-128 maxpool		
BFM ($128 \times M/4 \times M/4$)		
ups2-128 ups2-64 ups2-32 ups2-16 conv1-6		
IFM6 ($6 \times 4M \times 4M$)		IFM1 ups6-4
		IFM5 ups6-8
		IFM7 ($6 \times 4M \times 4M$)
		IFM8 ($6 \times 4M \times 4M$)
add		
IFM4 ($6 \times 4M \times 4M$)		
softmax		

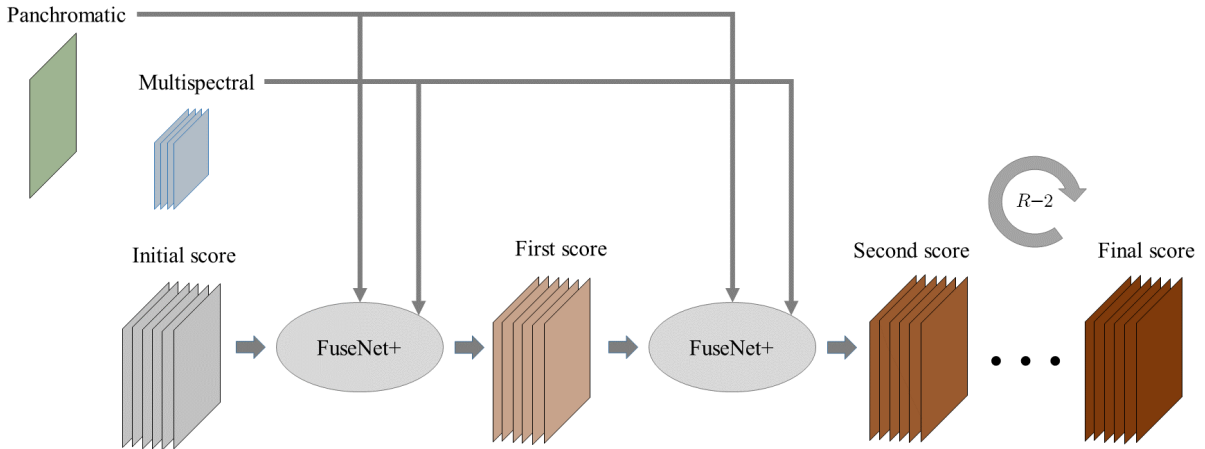


Fig. 3. The general architecture of ReuseNet with R FuseNet+ instances. FuseNet+ employs exactly the same operations as FuseNet except for the first layer applying additional sets of convolutional filters on the input score maps. ReuseNet accepts three streams of input: 1) \mathbf{x}_{PAN} , 2) \mathbf{x}_{MS} , and 3) score maps of the same resolution as \mathbf{x}_{PAN} . It applies the same operations employed by a FuseNet in R cycles, taking the output score map of the previous cycle as an input.

taken as the negative logarithm of the class scores [10], and a pairwise term, adding a penalization for pixels with different labels based on image-space and feature-space distances. For ease of notation, let the series of operations performed by FuseNet (Equation 8) be given by the function f :

$$\mathbf{y} = f(\mathbf{x}_{PAN}, \mathbf{x}_{MS}) \quad (9)$$

where the \mathbf{x} 's are the input of FuseNet, and \mathbf{y} is the class score map resulting from this input. The operations performed by ReuseNet are given by a recurrent variant g :

$$\mathbf{y}_1 = g(\mathbf{x}_{PAN} \oplus \mathbf{y}_0, \mathbf{x}_{MS}) \quad (10a)$$

$$\mathbf{y}_r = g(\mathbf{x}_{PAN} \oplus \mathbf{y}_{r-1}, \mathbf{x}_{MS}) \quad (10b)$$

where the r score map is obtained by applying the same function to a combination of the previous $r - 1$ score map

and the original FuseNet input as a new input. The recurrent variant g (denoted as FuseNet+ in Figure 3) applies exactly the same operations as f except for the first operation that instead of only taking \mathbf{x}_{PAN} as an input, this operation takes the concatenation of \mathbf{x}_{PAN} and a class score map \mathbf{y}_r associated to the network instance r . Figure 3 shows a diagram illustrating the general architecture of ReuseNet.

We tested ReuseNet with several number of FuseNet instances (2, 3, and 4), calling each ReuseNet- R where R is the number of FuseNet instances within the ReuseNet. We also investigated various methods for initializing weights and initial score maps \mathbf{y}_0 of ReuseNet. Plain ReuseNet initializes the score maps with zeros, while ReuseNet_{map-init} initializes the score maps using scores from a pre-trained FuseNet showing the best results in the fusion comparison experiments. We

further extend $\text{ReuseNet}_{\text{map-init}}$ by initializing the weights of the FuseNet instance in the ReuseNet with the same FuseNet that provides the initial score maps. We call this extension $\text{ReuseNet}_{\text{map-weights-init}}$.

C. Perspective on Learning the Fusion Approach and Incorporating Recurrence

Conventional approaches to classify multiresolution images require a separate step to match the resolution of the images. One way is to spatially sharpen the MS images using the PAN image (also called pansharpening) [32]. Another possible way is to resample images to match a specific resolution using nearest neighbor or bilinear interpolation. However, these standard fusion techniques are performed independently from the classification problem and are suboptimal. FuseNet provides a streamlined approach including the fusion of the multiresolution images within the learning of the classifier. We expect that coupling and learning the fusion method within a supervised classifier will outperform an approach based on a separate fusion method.

The parameter sharing across FuseNet instances in a ReuseNet is consistent with the definition of a recurrent network, i.e. a recurrent network is a feedforward network that keeps on reusing the same set of weights to cycle through a sequence. The authors in [21] view such incorporation of recurrence as a way to increase the contextual window size, equivalent to the patch size M in a patch classification approach, of their patch classification based approach while controlling the capacity of the network via inter-instance weight sharing. While both increase in contextual window size and capacity control of a CNN-based image patch classifier helps to improve the latter's performance, the first benefit is lost in a full patch labeling approach. In a fully convolutional network implementing full patch labeling, the contextual window size does not change as recurrent operations are added to the network since the contextual window size is equivalent to the effective size of the receptive field of the network. The effective size of the receptive field of the network depends on kernel sizes and strides of the network's convolutional and pooling operations, which are fixed and the same across instances.

In the proposed ReuseNet, recurrence integrates *contextual label information* to our model by considering class score maps as inputs to each FuseNet instance. This allows the model to learn label-to-label dependencies in addition to the spatial contextual information learned from the pixel values, *pixel-to-label* dependencies. This is a form of structured output prediction [33] where interdependencies between outputs may be expressed in terms of constraints restricting permissible output combinations or a more flexible form such as spatial dependencies across different output variables. Graphical models such as conditional random fields [34] are commonly used for such structured prediction tasks. ReuseNet uses operations in a deep convolutional network to learn features from both the input image and class scores—integrating the learning of label-to-label dependencies from the data instead of explicit image-space and feature-space distances as represented in a

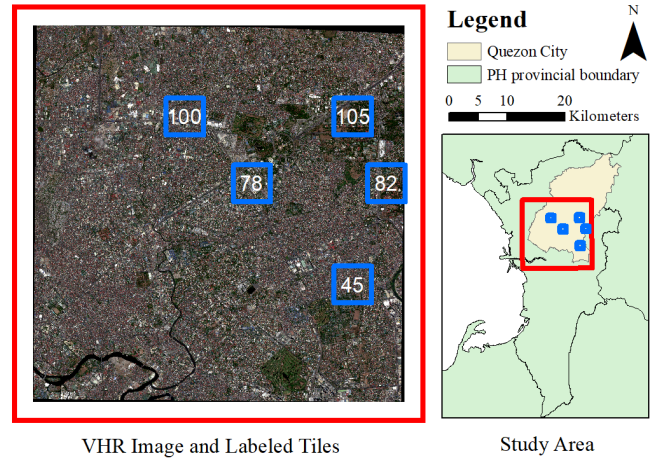


Fig. 4. Figure showing the true color VHR image together with the locations of the labeled tiles (in blue squares) and the study area: Quezon City, Philippines.

TABLE III
NUMBER OF LABELED PIXELS IN EACH TILE

Tile	Number of labeled pixels	Set
100	2178768	Training
105	2173602	Training
45	2063971	Validation
78	1977336	Test
82	1961955	Test

pairwise potential of CRF. This allows ReuseNet to be trained end-to-end as opposed to a two-stage approach applying a post-classification MRF/CRF as done in [35] and [13].

IV. DATA AND EXPERIMENTAL SETUP

A. Dataset Description

1) *Worldview-03 Quezon City dataset*: we evaluated the proposed networks in the land cover classification of a dataset covering Quezon City, Philippines. The dataset is composed of a Worldview-03 satellite image of the city acquired on 17th April 2016 and corresponding manually prepared reference images for five chosen tiles (subsets) of the satellite image. The satellite image has a PAN band of 0.3 m resolution and four MS bands (near-infrared, red, green, and blue) of 1.2 m resolution. Reference images were prepared via photointerpretation and set to have the same spatial resolution as the PAN image. The whole satellite image was first divided into regularly-sized image tiles. PAN image tiles have a dimension of 3200 pixels \times 3200 pixels, while MS image tiles have a dimension of 800 pixels \times 800 pixels. Five non-adjacent tiles were sparsely labeled—annotating a pixel with a label belonging to one of the following six classes: impervious surface, building, low vegetation, tree, car, and clutter. Two of the five labeled tiles were used for training (100 and 105), one for validation (45), and the remaining two for testing (78 and 82). Training samples are composed of pairs of image patches with dimensions $M \times M$ (taken from the MS image) and $4M \times 4M$ (taken from the PAN image tile). Figure 4 shows the VHR image and the corresponding locations of labeled tiles

in the study area while Table III shows the number of labeled pixels in each image tile. Training samples were normalized to have a value between zero and one. The reference image patches have been converted into a “one-hot” encoding—a vector having zero values except for the index corresponding to the code of the class.

2) *ISPRS Vaihingen dataset*: for the ReuseNet experiments, we utilized the ISPRS 2D semantic labeling benchmark dataset of Vaihingen as an additional dataset [36]. We adopted the experimental setup used in [13], [16], employing the same training and validation tiles, to provide comparable results. We followed the sampling done in [13], except that data augmentation was not applied—resulting in less training samples. The method discussed in [37] was employed to extract the normalized DSM.

B. Comparison of methods

For the image fusion part, we compared FuseNet against two other baseline approaches: one using pansharpening and another using bilinear interpolation to match the resolution of \mathbf{x}_{MS} to the resolution of \mathbf{x}_{PAN} . We call these two baseline approaches $\text{Net}_{\text{pansharp}}$ and $\text{Net}_{\text{bilinear}}$. $\text{Net}_{\text{pansharp}}$ applies Gram-Schmidt pansharpening technique [38]. Only the pansharpened image is fed as an input to $\text{Net}_{\text{pansharp}}$. $\text{Net}_{\text{bilinear}}$ upsamples the resolution of the MS image to match the resolution of the PAN image using bilinear interpolation. The upsampled MS images are then merged to the PAN image using concatenation. The architecture of the network after the fusion is kept the same to have a fair comparison among the different approaches (see details of the FuseNet variants in Table I). Additionally, we compared a SegNet [26] trained on the first three principal components of the pansharpened image, since SegNet only accepts three inputs. We found that discarding one band (NIR) considerably degrades the results.

We compared ReuseNet against FuseNet using fully-connected CRF [31] (FuseNet+CRF) to assess the capability of our classifiers to spatially regularize the classification results. The FuseNet+CRF baseline is similar to the approach adopted in [10], [13] but applied to PAN and MS images with different spatial resolutions. Spatial and feature space distances in the pairwise potentials of the fully-connected CRF are computed from the PAN image. We performed a grid-search of the CRF parameters, i.e. the weights and standard deviations of the appearance and smoothness kernels, and used the set of the parameters with the highest accuracy on the validation tile. We fixed the number of iterations to 10 for the mean field approximation algorithm used to perform inference in a fully-connected CRF.

We also performed a sensitivity analysis of a few chosen hyperparameters of $\text{FuseNet}_{\text{low}}$. We varied the bottleneck feature map dimensions, number of convolutional layers (in the downsampling part of the network), input patch sizes, and upsampling methods—performing the experiments in this order. We took the hyperparameter value that maximizes the overall accuracy on the validation tile and fix it for the succeeding sets of experiments. We experimented using bottleneck feature map dimensions: 16×16 , 8×8 , 4×4 , 2×2 ,

and 1×1 . After fixing the bottleneck feature map dimension, we increased the number of convolutional layers preceding the last downsampling operation—effectively increasing the number of convolutional layers from 8 to 14 in steps of two. We investigated varying patch sizes of $(4M, M)$: (32, 8), (64, 16), (96, 24), (128, 32). For the upsampling operations, we explored two additional methods using nearest neighbor and bilinear interpolation to upsample the feature map and then performing 3×3 convolutions after each upsampling operation.

We trained all the networks using a set of 17409 image patches taken from the training tiles and used 8255 image patches taken from the validation tile for early-stopping. We performed a random sampling with the constraint that the pixel near the center of the image patch is labeled. This may produce overlapping patches unlike the systematic gridwise sampling approach used in [13]. Gridwise sampling reduces the number of training patches since the reference images is sparsely labeled, only around five percent of the pixels are labeled. The total loss value computed over a mini-batch is the total loss of all pixels divided by the number of labeled pixels within the mini-batch.

The FuseNets are trained using backpropagation with stochastic gradient descent setting the initial learning rate $\eta = 0.01$, momentum $\alpha = 0.9$, mini-batch size $N = 32$, and maximum number of epochs $\mathcal{T} = 240$. We decrease the learning rate in a stepwise manner as done in [39]—multiplying it by a factor of 0.1 after 60 and 180 epochs. The weights were initialized as in [40]. We did not find dropout to be helpful; hence, we only used an l_2 -weight decay penalty—setting $\lambda = 0.001$ —and a variant of early-stopping to regularize FuseNet. For early stopping, the classification accuracy on the validation set is calculated every epoch and the last model with the best validation accuracy is fixed to be the final instance of the model.

The FuseNet instances within a ReuseNet are identical, sharing the same network configuration and parameters. Each instance also couples a cross-entropy loss function with each of their score map. The total objective loss of a ReuseNet is the average of the cross-entropy loss values from all the FuseNet instances. We also used the same backpropagation with stochastic gradient descent setting as training a FuseNet with the initial learning rate $\eta = 0.01$, momentum $\alpha = 0.9$, mini-batch size $N = 32$, and maximum number of epochs $\mathcal{T} = 240$. Likewise, we decreased the learning rate in a stepwise manner—multiplying it by a factor of 0.1 after 60 and 180 epochs. For regularization, we only used an l_2 -weight decay penalty—setting $\lambda = 0.001$. We can infer classification map from a ReuseNet in the same manner of inference as a FuseNet, with one additional option: to extract different predictions from each FuseNet instance.

For applying ReuseNet on the ISPRS Vaihingen dataset, we employed a feedforward network similar to the No-downsampling FCN proposed by [13] truncating the last two layers (fc5 and fc6) before softmax activation and entirely removing all maximum pooling without downsampling operations. With only convolutional layers (without pooling), we call this network AllConvNet. The network was trained on 12717 training patches as opposed to the 123330 train-

ing patches in [13]. Although having less parameters and having trained with a smaller number of training samples, AllConvNet provided comparable results with the original No-downsampling FCN while requiring less operations. We trained AllConvNet for 150000 iterations as reported in [13]. ReuseNet versions of AllConvNet were applied to the ISPRS Vaihingen dataset and were compared to the best results of both [13] and [16]. All the networks in this additional set of experiments were trained using a variant of SGD proposed in [41].

C. Accuracy Assessment

We compared the results of the different approaches using global measures: 1) overall classification accuracy (OA), 2) the Kappa coefficient (κ), 3) average class accuracy (AA), 4) and average class-F1 scores (F1). OA is given by:

$$OA = \frac{\sum_{i=1}^C n_{ii}}{n} \quad (11)$$

where n_{ii} is the number of samples classified as class i in both the the predictions and reference images, n is the total number of labeled samples in the reference images, and C is the number of classes, whereas κ is given by:

$$\kappa = \frac{n \sum_{i=1}^C n_{ii} - \sum_{i=1}^C n_{i+n_i}}{n^2 - \sum_{i=1}^C n_{i+n_i}} \quad (12)$$

where n_{i+} and n_{+i} are the number of samples classified as class i in the predictions and reference images respectively. Both OA and κ provides the rate of correctly classified pixels with the latter compensating for random agreement in classification. These global measures, however, are biased toward frequently occurring classes—meaning, classes with less frequencies have relatively little impact to the two measures. Unlike OA and κ , AA and F1 provides average of measures independent of class distribution. AA is given by:

$$AA = \frac{1}{C} \sum_{i=1}^C \frac{n_{ii}}{n_{i+}} \quad (13)$$

while F1 is given by:

$$F1 = \frac{1}{C} \sum_{i=1}^C \frac{2 \frac{n_{ii}}{n_{i+}} \frac{n_{ii}}{n_{+i}}}{\frac{n_{ii}}{n_{i+}} + \frac{n_{ii}}{n_{+i}}} \quad (14)$$

AA computes the average within-class rate of correctly classified pixels, while F1 calculates the harmonic mean of the precision (user's accuracy) and recall (producer's accuracy). We also observe and comment on the quality of the resulting classified maps.

V. RESULTS AND DISCUSSION

A. FuseNet

Table IV shows the results of accuracies comparing different fusion approaches. The numerical results are evaluated using all the labeled pixels in the two test tiles (see Table III for the

TABLE IV
COMPARISON OF FUSION APPROACHES

Network	OA (%)	κ (%)	AA (%)	F1 (%)
Net _{bilinear}	84.76	78.70	81.99	77.48
Net _{pansharp}	86.87	81.53	82.76	77.86
SegNet [26]	88.11	83.17	83.96	77.01
FuseNet _{high}	88.03	83.18	89.79	79.06
FuseNet _{low}	91.63	88.03	92.91	82.90
FuseNet _{skip}	91.90	88.43	93.46	81.74

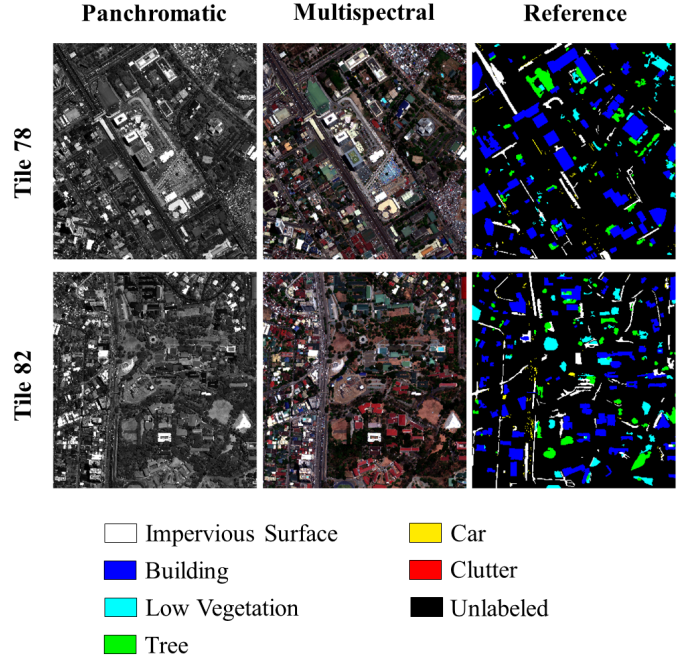


Fig. 5. PAN, MS, and reference images in the tiles used for testing. Corresponding legend is shown.

total number test samples). FuseNet_{skip} scores the highest in all the four numerical metrics, except for F1 where FuseNet_{low} scores the highest. FuseNet_{low} outperforms both the variants using fixed upsampling (Net_{pansharp}, SegNet, and Net_{bilinear}) and the variants learning the upsampling but fusing at the scale of the image with higher resolution (FuseNet_{high}). Observing each metric: FuseNet_{low} gains about 3–6% in OA, 4–9% in κ , 3–10% in AA, and 1–5% in F1 against the other baselines (with the exemption of FuseNet_{skip}). FuseNet_{skip} further increases the numerical results of FuseNet_{low} in the first three metrics by about 0.2–0.5% but degrades the F1 by about 1.2%.

We have two relevant observations: 1) learning the fusion can improve the classification of PAN and MS VHR images with different resolutions; 2) fusing at the scale of the image with lower resolution results in better classification than performing the fusion at the scale of the image with higher resolution. The first point demonstrates our expected effectiveness of coupling and learning the fusion operation within a supervised classifier. One explanation for the second point could be the placement of upsampling layers. Introducing upsampling layers early in the network—as done in FuseNet_{high}—may produce artifacts that can degrade its performance.

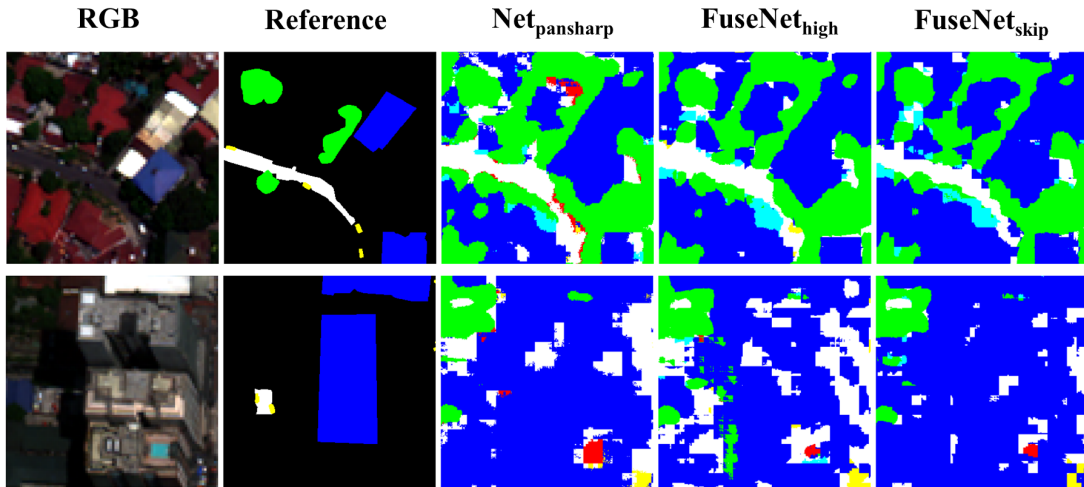


Fig. 6. Two subsets from the test tiles showing, from right to left, the satellite image (natural color), reference image, and classification maps from selected two FuseNet variants (FuseNet_{high} and FuseNet_{skip}) and one baseline method (Net_{pansharp}).

Figure 5 shows the PAN, MS, and reference images of the tiles used for testing. Figure 6 shows the classification results of two FuseNet variants (FuseNet_{high} and FuseNet_{skip}) and one baseline method (Net_{pansharp}) on two selected areas of the test tiles. The most noticeable misclassifications are found in large and high-rise buildings, in both test tiles, and an overpassing road in tile 78. The facades and rooftops of the buildings are often mistaken to be impervious surfaces by the classifiers; while the overpassing road is mistaken to be a building. These regions can appear to have similar spectral characteristics and can only be distinguished by presence of other cues such as appearing to be elevated. However, with the absence of elevation information, such cues are not directly incorporated in the input data. Manually distinguishing arguably vaguely-defined classes such as low-vegetation and impervious surface can also be problematic, especially in the PAN image, with the lack of ancillary information such as elevation. Adding a digital elevation model or a digital surface model can help address the misclassification of these regions. The cars are also generally misclassified by all the classifiers which is, aside from being underrepresented in terms of the number of labeled pixels, due to the lack of spatial resolution of the MS bands and the cars' spectral similarity with other classes (such as impervious surface and buildings) in the PAN band. Overall, FuseNet_{skip} generally has less errors in the facade of large buildings, lessen the artifacts noticeably present in the other techniques, and has better delineation of classes with irregular boundaries such as trees and low-vegetation—providing the best classification results among all the FuseNet variants. We, therefore, apply recurrence to FuseNet_{skip} architecture to build the ReuseNet instances.

B. ReuseNet

1) *Worldview-03 Quezon City dataset*: Table V shows the accuracies obtained by comparing different classification techniques on the Worldview-03 Quezon City dataset. We found that both the ReuseNet instances and the baseline method FuseNet+CRF improves the numerical results of the plain

TABLE V
COMPARISON OF MAP REGULARIZATION APPROACHES ON
WORLDVIEW-03 QUEZON CITY DATASET

Network	OA (%)	κ (%)	AA (%)	F1 (%)
FuseNet	91.90	88.43	93.46	81.74
FuseNet+CRF	93.07	90.08	94.71	81.72
ReuseNet-2	92.82	89.69	94.09	82.64
ReuseNet-3	92.98	89.88	94.54	85.42
ReuseNet-4	93.49	90.58	94.53	86.67
ReuseNet-5	92.74	89.53	92.78	87.29

FuseNet_{skip} gaining around: 0.9–1.5% in OA, 1.2–2.1% in κ , and 0.6–1.2% in AA. For the F1, however, FuseNet+CRF method performs worse than the plain FuseNet losing 0.02%; while all the other ReuseNet instances improves the F1 by around 0.9–5.5%. ReuseNet-4 outperforms all the other classifiers in all the metrics except for AA and F1—where both ReuseNet-3 and FuseNet+CRF outperform it by some margin in AA (0.01% and 0.18% respectively) and ReuseNet-5 considerably outperforms it in F1 by 0.62%. In particular, all the ReuseNets consistently show better F1 compared to both FuseNet and FuseNet+CRF—gaining almost 6%. These expected relatively smaller gains in numerical accuracy is consistent with what the author in [13] found—applying a post-classification CRF to an FCN to classify extremely high resolution aerial imagery increases the overall classification accuracy by around 0.1–1.0%. More noticeable changes are expected in the resulting improved regularity of the classified maps.

The numerical results above supports our assertion that introducing contextual label information through recurrence in an FCN applying a full-patch labeling approach can improve the classification of a VHR image. Such incorporation of label information allows our classifier to learn both pixel-to-label and label-to-label contextual dependencies. We can develop an intuition of these two dependencies by using an analogy to photointerpretation. We can easily imagine that it is easier to label a pixel when viewed with its neighboring pixels. This

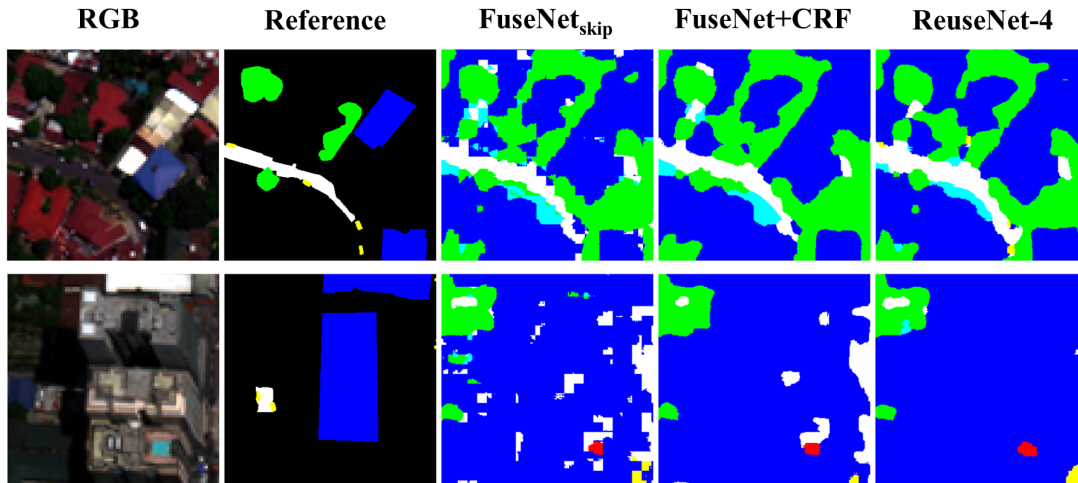


Fig. 7. Two subsets from the test tiles showing, from right to left, the satellite image (natural color), reference image, and classification maps from FuseNet_{skip}, FuseNet_{skip}+CRF, and ReuseNet-4. All ReuseNets reported are “plain” meaning initial score maps are filled with zeros. Reusenet- R denotes a ReuseNet composed of R number of FuseNet instances.

setup is analogous to the improvements a spatial-contextual classifier, like a CNN applying a patch classification, approach bring over a simple pixel-based classifier. But we can also see that it is easier to label a pixel when, aside from viewing its neighboring pixels, its surrounding pixels’ labels are given. With contextual label information, the classifier can learn and leverage class spatial co-occurrences. Additionally, we observe that adding more FuseNet instances to the ReuseNet until $R = 4$ increases the score of all metrics, except for the average class accuracy where ReuseNet-3 marginally outperforms ReuseNet-4. Adding one more instance only improves the F1 score and degrades the other three metrics. We can interpret this addition of FuseNet instances as a way to increase ReuseNet’s capacity to refine contextual label information fed to it as latter FuseNet instances receive more refined labels.

Figure 7 shows classification results of the best performing ReuseNet, the baseline method FuseNet+CRF, and the plain FuseNet. Both FuseNet+CRF and ReuseNet instances improves the quality of the resulting classified map by producing more regularized classification. We also observe that locations of the errors are carried over from the results of the FuseNet classifier from which both FuseNet+CRF and ReuseNet are based from. However, the occurrences of the errors are diminished especially on the facades of the large buildings. Detection of isolated cars in roads were also improved. Overall, results of ReuseNet-4 show better-quality classified maps by reducing noise in the classification (such as island of impossibly small buildings), further improving delineation of classes with irregular boundaries, and reducing misclassification in regions with ambiguous spectral characteristics such as facades and rooftops of high rise buildings.

2) *ISPRS Vaihingen dataset*: Table VI shows the accuracies obtained by comparing different classification techniques on the ISPRS dataset. These results are in agreement with the results from the previous dataset. All the ReuseNet versions of AllConvNet improve the results on all the four metrics except for AA and F1 of ReuseNet-2 (2.08% in AA and

TABLE VI
COMPARISON OF MAP REGULARIZATION APPROACHES ON ISPRS
VAIHINGEN DATASET

Network	OA (%)	κ (%)	AA (%)	F1 (%)
No-downsampling FCN [13]	87.17	--	--	--
CNN-FPL [16]	87.83	83.83	81.35	83.58
AllConvNet	86.98	82.71	87.17	85.46
FCN in [13] with CRF	87.90	--	--	--
ReuseNet-2	87.11	82.89	85.09	85.38
ReuseNet-3	88.08	84.18	87.29	87.24
ReuseNet-4	87.64	83.59	87.18	86.81

Unreported values in the reference are denoted by “--”

0.06% in F1 respectively). ReuseNet-3, the best performing network, considerably improves all the numerical results of the plain AllConvNet by 1.1% in OA and is comparable and even greater than the 0.73% gain after a post-classification CRF in [13], 1.47% in κ , 0.12% in AA, and 1.78% in F1. ReuseNet-3 also outperforms best results reported in both [13] and [16].

These results reconfirm that introducing contextual label information through recurrence in an FCN applying a full-patch labeling approach can improve the classification of a VHR image. Similarly, qualitative improvements—such as holes in building being filled, better delineation of all classes in general, lesser artifacts—in the resulting classified maps are observed when ReuseNet is applied as shown in Figure 8.

3) *Different initializations*: Figure 9 shows results of quantitative metrics on the three different ReuseNet initializations. There is low variation in the OA and κ . The trend of the two global scores is also inconsistent across the ReuseNet instances. For ReuseNet-2 and ReuseNet-3, the scores increases marginally (around 0.5% for OA and 0.8% for κ) when initialized with both the scores and weights from a previously-trained FuseNet. But for ReuseNet-4, there is a minor drop in both the scores (around 0.2% for both scores) when the two initialization methods are introduced. This could mean that increasing the FuseNet instances to a certain amount already provides enough room to a ReuseNet for “label refinement”

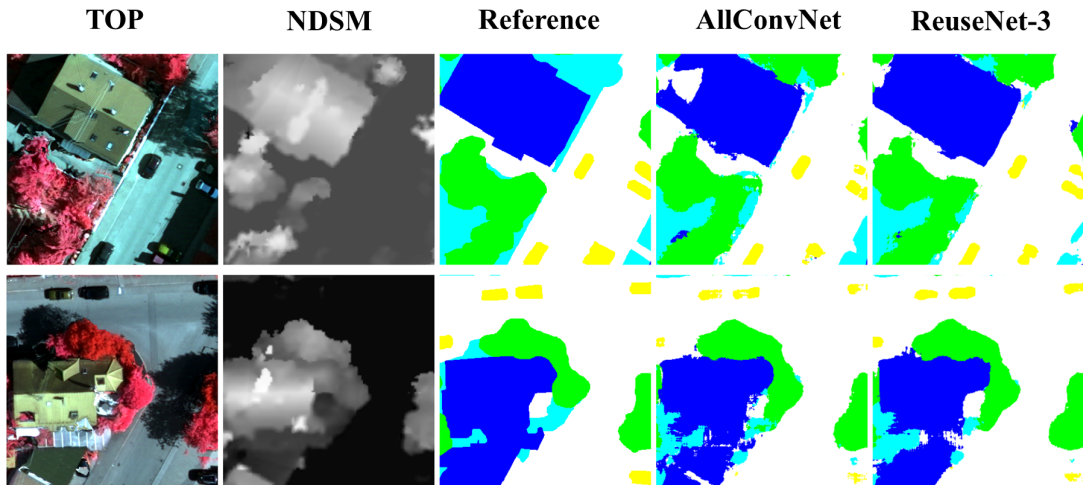


Fig. 8. Two subsets from the validation tiles of ISPRS Vaihingen dataset showing, from right to left, the true orthophoto, normalized dsm, reference image, and classification maps from AllConvNet and ReuseNet-3 (best performing ReuseNet in this dataset). All ReuseNets reported are “plain” meaning initial score maps are filled with zeros. Reusenet- R denotes a ReuseNet composed of R number of FuseNet instances.

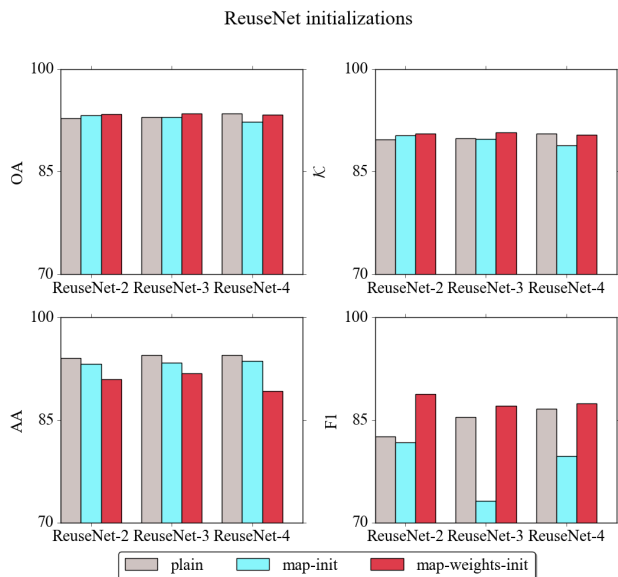


Fig. 9. Plots showing results of quantitative metrics comparing different ReuseNet initializations; plain, map-init, and map-weights-init correspond to initializing the ReuseNet with zero-score maps, scores from a previously-trained FuseNet, and scores and weights from a previously-trained FuseNet respectively.

such that gains from the initialization methods are compensated.

Introducing both initialization methods to a ReuseNet degrades the AA by around 0.9–5.2%. Applying only the initialization using scores from a FuseNet instance (map-init) degrades the F1 by around 0.9–12.2%. Interestingly, the F1 improve by around 0.8–6.1% when both initialization methods are introduced (map-weights-init). Decrease in AA can only imply an increase in false positive predictions in most of the classes; while increase in F1 could either mean decrease in false positive predictions or decrease in false negative predictions or both in most of the classes. The results therefore

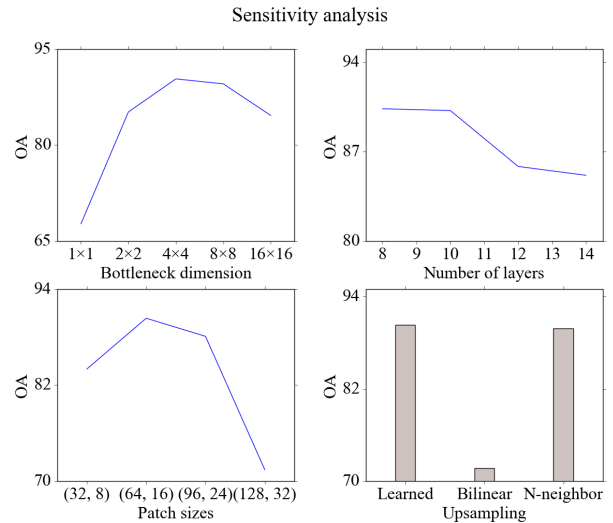


Fig. 10. Plots showing the results of sensitivity analysis. Patch sizes are written as “ $\langle(4M, M)\rangle$ ”. N-neighbor denotes nearest neighbor interpolation.

show that the initialization methods promote higher recall rate (decrease in false negatives) in underrepresented classes such as cars.

C. Sensitivity Analysis

Fig. 10 shows the results of the sensitivity analysis performed on four chosen hyperparameters of FuseNet: 1) bottleneck feature map dimensions, 2) number of convolutional layers (in the downsampling part of the network), 3) input patch sizes, and 4) upsampling methods. We got the highest validation accuracy of 90.35% using a bottleneck feature map dimension of 4×4 pixels. Decreasing the dimension more than the optimal we found severely degrades the classification resulting to large uniform areas producing stamp-like patterns (especially for 1×1). Increasing the dimensions produces much noisier classification. Fixing the bottleneck

size dimension to 4×4 and further increasing the number of convolutional layers (without downsampling) did not produce any improvements in the validation accuracy. Increasing the number of these convolutional layers within the bottleneck feature maps effectively increases the receptive field (footprint size in the input layer containing the PAN image patch) of the succeeding units by at least half of the size of kernels used in the convolutional layers. Hence, the results show that: with only eight convolutional layers (with downsampling), we can learn enough contextual information for accurate classification.

We found the optimal patch sizes of 64×64 for the \mathbf{x}_{PAN} and 16×16 for \mathbf{x}_{MS} . Further increasing the patch sizes results in overclassification of a single class (impervious surface). Increasing the patch size also increases the proportion of frequently occurring classes in the training sample, possibly resulting into overclassification. Whereas, decreasing the patch size limits the contextual information incorporated in the input, and, hence, can degrade the classification results. Lastly, we find using transposed convolution for learned upsampling to perform better than using interpolation for fixed upsampling (bilinear and nearest neighbor). This result supports the expected flexibility of empirically learning the upsampling operation directly from data.

VI. CONCLUSION AND FUTURE WORKS

In this paper, we presented a recurrent multiresolution convolutional network named ReuseNet to classify VHR satellite images. The operations for fusing the bands with different resolutions are learned within convolutional layers with corresponding downsampling and upsampling operations to match the resolution of the images. Regularization of the resulting classified maps is achieved by incorporating contextual label information through the recurrent architecture of ReuseNet. Additionally, we investigated various ways to initialize ReuseNet. The effect of varying a set of chosen network hyperparameters to the classification accuracy of the network was explored. Both numerical and qualitative results show the advantages of incorporating image resolution matching and contextual label learning within the training of the classifier. To this end, we provided a single-stage classification pipeline incorporating image fusion, feature extraction, and map regularization, all combined in a convolutional network trained in an end-to-end manner.

We designed the presented network architecture such that it can easily be adapted to other multiresolution image datasets. Inclusion and leverage of contextual label information is also separate from the design of the fusion network in the sense that it can be implemented on network classifying single-resolution images. For future work, we plan to fuse images from different sensors (e.g. Sentinel-2) and classify classes of higher abstraction such as land use instead of land cover.

REFERENCES

- [1] R. Haralick, K. Shanmugan, and I. Dinstein, "Textural features for image classification," pp. 610–621, 1973. [Online]. Available: <http://dceanalysis.bigr.nl/Haralick73-Texturalfeaturesforimageclassification.pdf>
- [2] T. Ojala, M. Pietikainen, and T. Maenpaa, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 7, pp. 971–987, Jul 2002.
- [3] M. D. Mura, J. A. Benediktsson, B. Waske, and L. Bruzzone, "Morphological attribute profiles for the analysis of very high resolution images," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 48, no. 10, pp. 3747–3762, Oct 2010.
- [4] M. Fauvel, Y. Tarabalka, J. A. Benediktsson, J. Chanussot, and J. C. Tilton, "Advances in spectral-spatial classification of hyperspectral images," *Proceedings of the IEEE*, vol. 101, no. 3, pp. 652–675, March 2013.
- [5] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," in *Proceedings of the IEEE*, 1998, pp. 2278–2324.
- [6] J. R. Bergado, C. Persello, and C. Gevaert, "A deep learning approach to the classification of sub-decimeter resolution aerial images," in *Proceedings of the International Geoscience and Remote Sensing Symposium (IGARSS) 2016*, vol. 2016-November, 2016, pp. 1516–1519.
- [7] N. Mboga, C. Persello, J. Bergado, and A. Stein, "Detection of informal settlements from vhr images using convolutional neural networks," *Remote Sensing*, vol. 9, no. 11, p. 1106, 2017.
- [8] D. H. Hubel and T. N. Wiesel, "Receptive fields, binocular interaction, and functional architecture in the cat's visual cortex," *Journal of Physiology (London)*, vol. 160, pp. 106–154, 1962.
- [9] W. Ha, P. H. Gowda, and T. A. Howell, "A review of potential image fusion methods for remote sensing-based irrigation management: part ii," *Irrigation Science*, vol. 31, no. 4, pp. 851–869, Jul 2013. [Online]. Available: <https://doi.org/10.1007/s00271-012-0340-6>
- [10] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *CoRR*, vol. abs/1606.00915, 2016. [Online]. Available: <http://arxiv.org/abs/1606.00915>
- [11] S. Paisitkriangkrai, J. Sherrah, P. Janney, and A. van den Hengel, "Semantic labeling of aerial and satellite imagery," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 9, no. 7, pp. 2868–2881, 2016.
- [12] L. Wang, X. Huang, C. Zheng, and Y. Zhang, "A markov random field integrating spectral dissimilarity and class co-occurrence dependency for remote sensing image classification optimization," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 128, no. Supplement C, pp. 223–239, 2017. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0924271616302787>
- [13] J. Sherrah, "Fully convolutional networks for dense semantic labelling of high-resolution aerial imagery," *arXiv preprint arXiv:1606.02585*, 2016.
- [14] E. Maggiori, Y. Tarabalka, G. Charpiat, and P. Alliez, "Convolutional neural networks for large-scale remote-sensing image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 2, pp. 645–657, Feb 2017.
- [15] C. Persello and A. Stein, "Deep fully convolutional networks for the detection of informal settlements in vhr images," *IEEE Geoscience and Remote Sensing Letters*, vol. PP, no. 99, pp. 1–5, 2017.
- [16] M. Volpi and D. Tuia, "Dense semantic labeling of subdecimeter resolution images with convolutional neural networks," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 2, pp. 881–893, 2017.
- [17] W. Zhao, L. Jiao, W. Ma, J. Zhao, J. Zhao, H. Liu, X. Cao, and S. Yang, "Superpixel-based multiple local cnn for panchromatic and multispectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 7, pp. 4141–4156, 2017.
- [18] X. Liu, L. Jiao, J. Zhao, J. Zhao, D. Zhang, F. Liu, S. Yang, and X. Tang, "Deep multiple instance learning-based spatial-spectral classification for pan and ms imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. PP, no. 99, pp. 1–13, 2017.
- [19] M. I. Jordan, "Chapter 25 - serial order: A parallel distributed processing approach," in *Neural-Network Models of Cognition Biobehavioral Foundations*, ser. Advances in Psychology, J. W. Donahoe and V. P. Dorsel, Eds. North-Holland, 1997, vol. 121, pp. 471–495.
- [20] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org>.
- [21] P. H. O. Pinheiro and R. Collobert, "Recurrent convolutional neural networks for scene labeling," in *Proceedings of the 31st International Conference on Machine Learning (ICML)*, 2014.
- [22] A. L. Maas, A. Y. Hannun, and A. Y. Ng, "Rectifier nonlinearities improve neural network acoustic models," in *ICML 2013 Workshop on Deep Learning for Audio, Speech, and Language Processing*, 2013.

- [23] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on Imagenet classification," in *IEEE International Conference on Computer Vision (ICCV) 2015*, 2015.
- [24] D. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (ELUs)," in *International Conference on Learning Representations*, 2016.
- [25] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, 2015, pp. 3431–3440.
- [26] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [27] L. Bottou, "Stochastic gradient descent tricks," in *Neural Networks: Tricks of the Trade*. Springer Berlin Heidelberg, 2012, pp. 421–436.
- [28] H. Noh, S. Hong, and B. Han, "Learning deconvolution network for semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1520–1528.
- [29] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [30] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.
- [31] P. Krähenbühl and V. Koltun, "Efficient inference in fully connected crfs with gaussian edge potentials," in *Advances in Neural Information Processing Systems 24*, J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2011, pp. 109–117.
- [32] D. Hester, H. Cakir, S. Nelson, and S. Khorram, "Per-pixel classification of high spatial resolution satellite imagery for urban land-cover mapping," *Photogrammetric Engineering and Remote Sensing*, vol. 74, no. 4, pp. 463–471, 2008.
- [33] G. H. Bakir, T. Hofmann, B. Schölkopf, A. J. Smola, B. Taskar, and S. V. N. Vishwanathan, *Predicting Structured Data (Neural Information Processing)*. The MIT Press, 2007.
- [34] J. D. Lafferty, A. McCallum, and F. C. N. Pereira, "Conditional random fields: Probabilistic models for segmenting and labeling sequence data," in *Proceedings of the Eighteenth International Conference on Machine Learning*, ser. ICML '01. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2001, pp. 282–289.
- [35] A. D. Giorgi, G. Moser, and S. B. Serpico, "Contextual remote-sensing image classification through support vector machines, markov random fields and graph cuts," in *2014 IEEE Geoscience and Remote Sensing Symposium*, July 2014, pp. 3722–3725.
- [36] M. Cramer, "The DGPF-test on digital airborne camera evaluation — overview and test design," *Photogrammetrie - Fernerkundung - Geoinformation*, vol. 2, pp. 73–82, 2010.
- [37] Y.-K. Mousa, P. Helmholz, and D. Belton, "New dtm extraction approach from airborne images derived dsm," vol. 42, no. 1W1, 2017, pp. 75–82.
- [38] C. Laben and B. Brower, "Process for enhancing the spatial resolution of multispectral imagery using pan-sharpening," Jan. 4 2000, uS Patent 6,011,875. [Online]. Available: <https://www.google.com/patents/US6011875>
- [39] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on*, 2016.
- [40] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks." in *Aistats*, vol. 9, 2010, pp. 249–256.
- [41] M. D. Zeiler, "Adadelata: An adaptive learning rate method," *CoRR*, vol. abs/1212.5701, 2012.

Claudio Persello Biography text here.

Alfred Stein Biography text here.

John Ray Bergado Biography text here.

PLACE
PHOTO
HERE