# Building Instance Classification Using Street View Images

Jian Kang[a], Marco Körner[b], Yuanyuan Wang[a], Hannes Taubenböck[c], Xiao Xiang Zhu[a,d,*]

[a]*Signal Processing in Earth Observation (SiPEO), Technical University of Munich (TUM), 80333 Munich, Germany*
[b]*Chair of Remote Sensing Technology, Technical University of Munich (TUM), 80333 Munich, Germany*
[c]*German Remote Sensing Data Center (DFD) (IMF), German Aerospace Center (DLR), 82234 Wessling, Germany*
[d]*Remote Sensing Technology Institute (IMF), German Aerospace Center (DLR), 82234 Wessling, Germany*

## Abstract

*This is the pre-print version, to read the final version please go to ISPRS Journal of Photogrammetry and Remote Sensing, Elsevier. (https://doi.org/DOI: 10.1016/j.isprsjprs.2018.02.006).* Land-use classification based on spaceborne or aerial remote sensing images has been extensively studied over the past decades. Such classification is usually a patch-wise or pixel-wise labeling over the whole image. But for many applications, such as urban population density mapping or urban utility planning, a classification map based on individual buildings is much more informative. However, such semantic classification still poses some fundamental challenges, for example, how to retrieve fine boundaries of individual buildings. In this paper, we proposed a general framework for classifying the functionality of individual buildings. The proposed method is based on Convolutional Neural Networks (CNNs) which classify façade structures from street view images, such as Google StreetView, in addition to remote sensing images which usually only show roof structures. Geographic information was utilized to mask out individual buildings, and to associate the corresponding street view images. We created a benchmark dataset which was used for training and evaluating CNNs. In addition, the method was applied to generate building classification maps on both region and city scales of several cities in Canada and the US.

*Keywords:* CNN, Building instance classification, Street view images, OpenStreetMap

## 1. Introduction

The classification of land cover from Earth Observation (EO) images in complex urban environments has been a focus in remote sensing over the past decades [1, 2, 3, 4, 5, 6]. Beyond, high resolution spaceborne and aerial images are one of the handful information sources for monitoring urban development on large scales.

However, the transfer from land cover to land use in EO-data is complex and relies mostly on the geometry and the appearance of individual buildings and the patterns they group together [7, 8, 9, 10, 11, 12, 13, 14, 15, 16]. The correlation of physical indicators such as building volumes, density or alignment has been used to infer the usage of buildings, e.g. as commercial areas (e.g. Figure 1(a)), residential areas (e.g. Figure 1(b)) or industrial areas (e.g. Figure 1(c)). Nevertheless, such pattern analysis can not be directly transferable to the classification of individual buildings as we go to a finer level of urban intrinsic scale. For example, Figure 1(a) shows a commercial area comprised of multiple high-rise buildings. However, the label "commercial area" cannot be assigned to all the building instances within it. As illustrated in Figure 2, the corresponding street view images show that the commercial area is comprised of a few apartments, office buildings, and one church. This also applies to the example shown in Figure 1 (b) and (c), where both the residential and industrial areas are comprised of buildings with different functionalities. As can be seen, land-use classification at a level of individual buildings is not a trivial task. Usually, such a classification map is only obtainable through city cadastral databases, not accessible or sometimes even not existent. Updating such databases without automatic methods can be very labor intensive. Hence, automatically achieving a building instance-level classification is necessary and can be beneficial for applications related with urban planning. Towards an automatic classification of individual buildings, the challenges are twofold. Firstly, remote

---

*Corresponding author (Email:xiao.zhu@dlr.de)

(a) commercial        (b) residential        (c) industrial

Figure 1: Examples of land-use classification



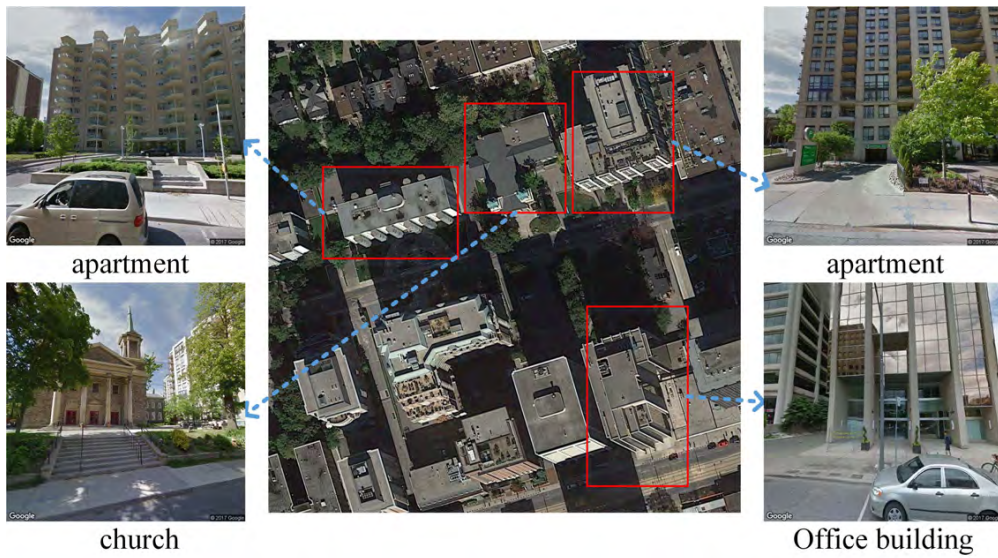apartment        apartment

church        Office building

Figure 2: The commercial land-use area as shown in Figure 1 (a), along with the street view images for some buildings selected by the red rectangles. These buildings do not belong to the same category, even though they are located in the same land-use area. Besides, compared to the roof structures, the information of façade structures displayed in street view images is richer and more sufficient to be used for building instance classification.

sensing images usually only contain roof structures due to their nadir-looking imaging geometry. The visual difference of the roofs between certain building classes, e.g. apartments and office buildings, can be subtle, as an example shown in Figure 2. Secondly, the extraction of building footprints directly from remote sensing images is still under preliminary research. A clear segmentation of building footprints usually requires height information which comes at an additional cost. In this paper, we propose a general framework to tackle the abovementioned challenges, which exploits the information extraction from freely available street view images and online geographic maps. Specifically, façade structures shown in online street view images are sufficiently rich for building functionality classification, and the online map services, such as OpenStreetMap [17] or Google Maps, can provide the building footprints which can be associated to street view images via their geographic locations. As shown in Figure 2, the façades displayed in street view images reveal much more details of different types of buildings than the corresponding roof patches. Therefore, building instances are classified based on their geo-tagged street view images in the proposed method, and the inferred labels are then linked to individual building footprints through spatial clustering. We also build a benchmark dataset of building street view images to train Convolutional Neural Networks (CNNs) for the classification over large areas, as CNN has been demonstrated its powerful ability in the tasks of this sort [18, 19, 20].

In a summary, the contributions of this paper are listed as follows:

- Proposed a general framework for land-use classification at a level of individual buildings.

- Built a street view benchmark dataset for training building instance CNN classifiers based on façade structures. The dataset utilized in this paper can be downloaded via `www.sipeo.bgu.tum.de/downloads/BIC_GSV.tar.gz`

- The obtained building classification maps demonstrated the potentials for many innovative urban analysis, e.g. very high resolution urban population density mapping, urban social structure understanding, city economy structure analysis and general urban planning.

## 2. Related work

Feature extraction from remote sensing images plays a vital role in land-use classification. Handcrafting features have been well studied for decades, such as scale-invariant feature transform (SIFT) [21] encoded by bag of visual words (BoVW) [22, 23, 24], multiple textural features [25], 3D features derived from a digital surface model [26] and features learned by sparse coding methods [27, 28, 29, 30, 31, 32, 33, 13, 34].

Recently, many approaches based on deep learning techniques have emerged [35, 36, 37]. Chen *et al.* [38] proposes a hierarchical feature extraction method via stacked autoencoders, which merges both spectral and spatial information of hyperspectral images for land-use classification. In [39], deep belief networks are employed for the feature learning in remote sensing scene classification. Both [40] and [41] investigate the possibility of transferring features learned by CNN from *ImageNet* dataset [42] to achieve remote sensing image classification by fine-tuning procedures. To improve the composition-based inference of land-use classes, multiscale CNN-based approaches are developed in [43, 44, 45]. By exploiting deep Boltzmann machine, a novel weakly supervised learning approach for object detection in remote sensing images is introduced [46]. For effectively dealing with the problem of object rotation variations, a rotation-invariant CNN model is proposed in [47]. Based on greedy layerwise unsupervised pretraining, [48] proposes a novel unsupervised deep feature extraction method. Taking advantage of geographical information from OpenStreetMap, a fully convolutional neural network is trained to achieve pixel-wise classifications in optical images on large scales [49]. Recurrent Neural Network (RNN) is also proved to be efficient for classifying sequence-based data like hyperspectral images [50]. An end-to-end fully Conv-Deconv network for unsupervised spectral-spatial feature extraction in hyperspectral images has been proposed in [51]. In order to better interpret land-uses of Synthetic Aperture Radar (SAR) images in urban areas, [52] proposes a pseudo-siamese CNN for identifying corresponding patches in very-high-resolution (VHR) optical and SAR remote sensing imagery. Surveys about the applications of deep learning techniques to land-use classification with remote sensing images are proposed in [53, 54].

Even the abovementioned literature is of course not exhaustive, none of them have explicitly addressed the land-use classification at a level of individual buildings.
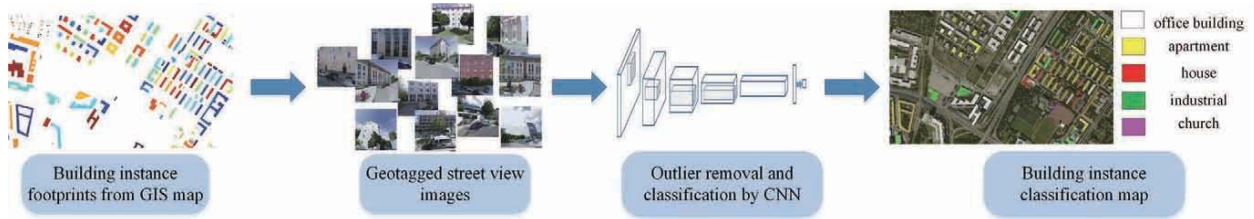
Figure 3: The proposed workflow for land-use classification at a level of individual buildings.

## 3. Overall workflow

As illustrated in Figure 3, the proposed workflow for building instance classification contains the following steps:

- Retrieval of building footprints and associated street view images.

- Outlier removal by the pretrained CNN on Places2 dataset [19].

- Building instance classification by the CNN trained on our benchmark dataset.

### 3.1. Retrieval of building footprints and street view images

The building footprints and their geographic locations, can be retrieved from online geographic information systems (GIS), such as OpenStreetMap or Google Maps. For example, the building footprints of the area shown in Figure 2 are displayed in Figure 4, along with the associated GPS coordinates (latitude, longitude). The color is randomly assigned to indicate different building instances. Given these GPS coordinates, we can download the corresponding Google StreetView images [55] which show façade structures of individual buildings, since the retrieved images can display these specific locations by the closest panoramas.

### 3.2. Outlier removal by pretrained CNN on Places2 dataset

Due to the uncontrolled quality of street view images, many of them cannot be directly utilized for the building classification. For example, as shown in Figure 5, one retrieved image is taken from the building interior and the other two buildings are occluded by a vehicle and trees on the side-walks. Therefore, the corresponding façade structures are not available for classifying these buildings. These outliers can severely influence the classification results. For removing them, we employ the released VGG16 model [56] trained on Places2 dataset [19] to preliminarily screen the street view images, as this architecture has achieved the highest top-1 accuracy[1]. The dataset contains almost 10 million scene photos, labeled with 476 scene categories and attributes, which include the building-related categories, i.e. [*apartment, church, house, industrial area, museum, building facade, embassy, hospital, parking garage, hotel*]. Only the images belonging to the abovementioned categories are preserved for the follow-up classification.

### 3.3. Building instance classification

To train a building instance classifier, we first build a corresponding street view benchmark dataset, which contains totally 19658 images from eight classes, i.e. *apartment, church, garage, house, industrial, office building, retail and roof*, and there are around 2500 images for each building class, as shown in Figure 6 and 7. The geo-tagged images are downloaded through Google StreetView API[2], with the associated metadata[3], i.e. the image size and pitch value are set to be $512 \times 512$ pixels and 10 degrees, respectively. As illustrated in Figure 8, all the street view images are located over several cities of the US and Canada, e.g. Montreal, New York and Denver, and their associated ground truth building labels are extracted from OpenStreetMap[4]. The descriptions for the building classes are demonstrated in Table 1.

---

[1] https://github.com/metalbubble/places365
[2] https://developers.google.com/maps/documentation/streetview/
[3] https://developers.google.com/maps/documentation/streetview/intro
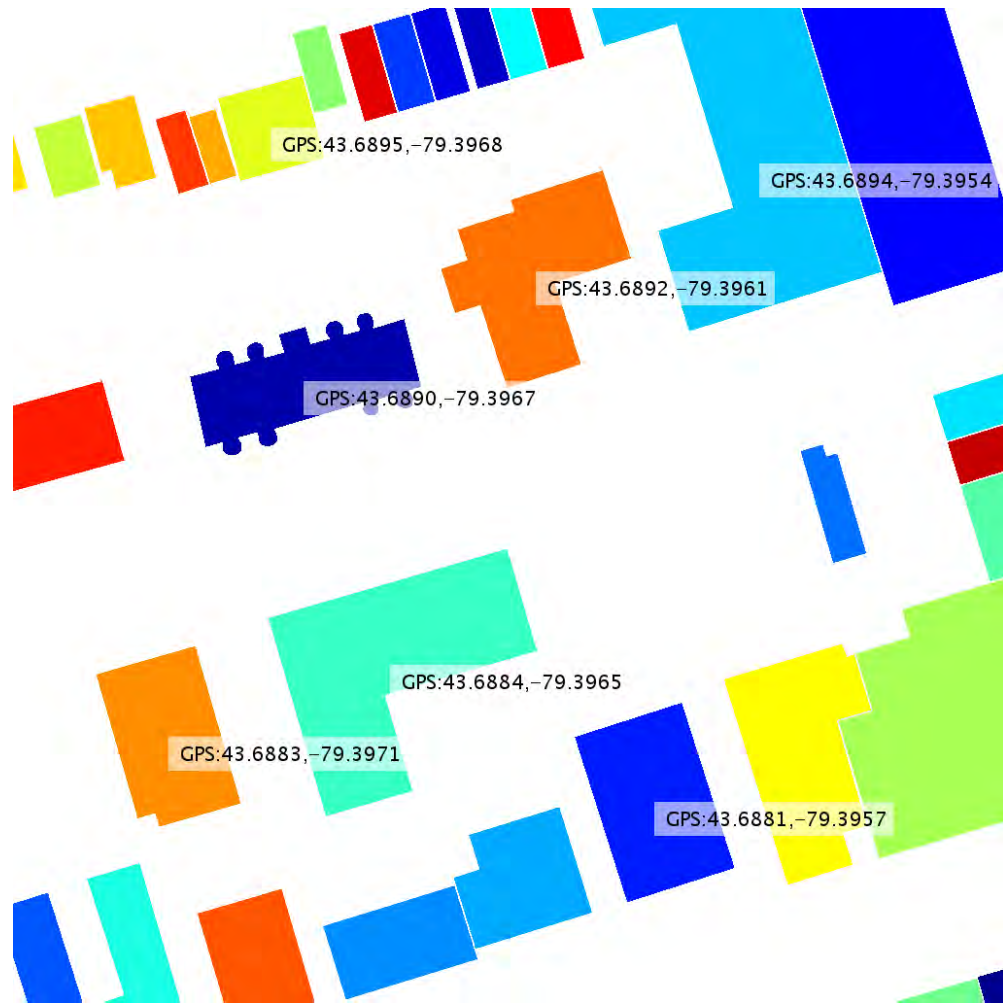[4] http://wiki.openstreetmap.org/wiki/Map_Features#Building

Figure 4: Geographic information (GPS) retrieved from Google Maps of the remote sensing image in Figure 2, with the color randomly assigned to each building mask.



Indoor        Vehicle occlusion        Tree occlusion

Figure 5: Outlier examples of the retrieved street view images. We can see that there is no available information of building façades for the classification.

Table 1: Building class descriptions from OpenStreetMap

| | |
|---|---|
| apartment | A building arranged into individual dwellings, often on separate floors. May also have retail outlets on the ground floor. |
| church | A building that was built as a church. |
| garage | A building suitable for the storage of one or possibly more motor vehicle or similar. |
| house | A dwelling unit inhabited by a single household (a family or small group sharing facilities such as a kitchen). |
| industrial | A building where some industrial process takes place. |
| office building | A building where non-specific commercial activities take place. |
| retail | A building primarily used for selling goods that are sold to the public. |
| roof | A structure that consists of a roof with open sides, such as a rain shelter, and also gas stations. |

Since the dataset is not sufficiently large to train a CNN with millions of parameters from the scratch, we choose to fine-tune a pretrained CNN with our dataset. It is common that a pretrained CNN on a large dataset such as *ImageNet* [18] can be well adapted to other new tasks with small scale datasets, since low-level features such as corners and edges generated by prior layers of CNN are general in different images. The high-level image representations extracted by posterior layers are dependent on different tasks. Therefore, fine-tuning the layers of the pretrained CNN with the new dataset has been proven to be an efficient way for the adaptation of the CNN to a new training task.

To further improve the classification robustness, the street view images for each building instance are classified, and the building class can be obtained in a decision level. Assuming there are $M$ street view images retrieved of the study building instance, the final building class $y$ can be determined by

$$y = \operatorname*{argmax}_{i} \frac{1}{M} \sum_{j=0}^{M-1} f_i^{(j)}, \tag{1}$$

where $f_i^{(j)}$ is the $i$th element of the CNN *softmax* layer output $\mathbf{f}^{(j)}$, which denotes the probability distribution over the whole building classes, and $j$ is the index of the classified street view image. For example, Figure 9 shows a building to be classified and its corresponding street view images from four different positions. After the right two images filtered by the outlier removal step, we can calculate the final probability distribution vector by averaging those of the left two images and obtain the building class accordingly.

## 4. Experiments

We train several state-of-the-art CNN architectures, e.g. AlexNet [57], VGG [56] and ResNet [58] by fine-tuning all the convolutional layers with our benchmark dataset, and demonstrate the corresponding training and testing performances. Among those networks, we choose the best one for generating building classification maps both on region and city scales.

### 4.1. Training

As illustrated in Figure 8, we split the whole dataset into two parts: 17600 images for training (2200 images for each building class) and 2058 images for testing. Note that all the testing images are retrieved from different cities with those utilized for training. In order to monitor the training status of networks, we randomly select 3200 images from the training samples to be the validation data. We train four different networks i.e. AlexNet, VGG16, ResNet18 and ResNet34 following the same procedure: Convolutional layers of all these networks are initialized by those pretrained with *ImageNet*, and fully connected layers are randomly initialized following a uniform distribution.
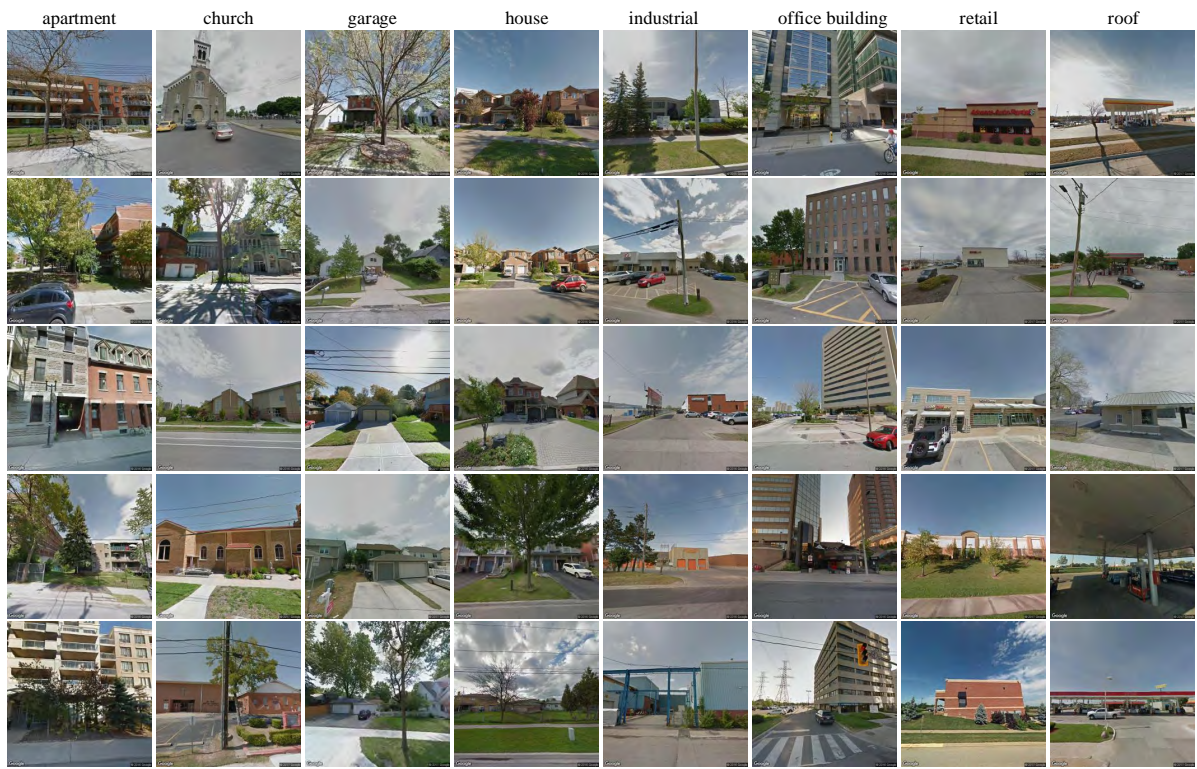
Figure 6: Examples of the benchmark dataset. It totally contains 19658 street view images of buildings with eight classes, i.e. *apartment, church, garage, house, industrial, office building, retail and roof*. The images are downloaded from Google StreetView [55], and the associated labels are jointly retrieved from OpenStreetMap based on the geographic information.
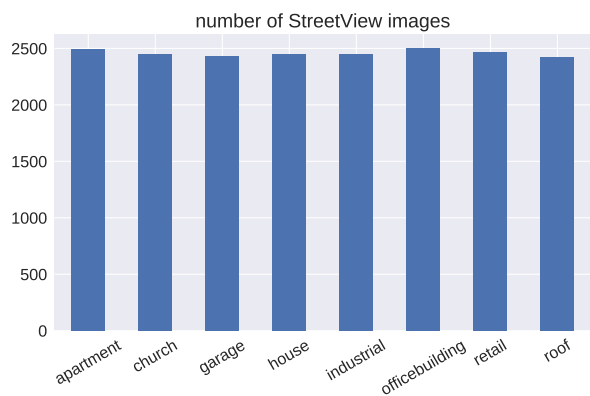


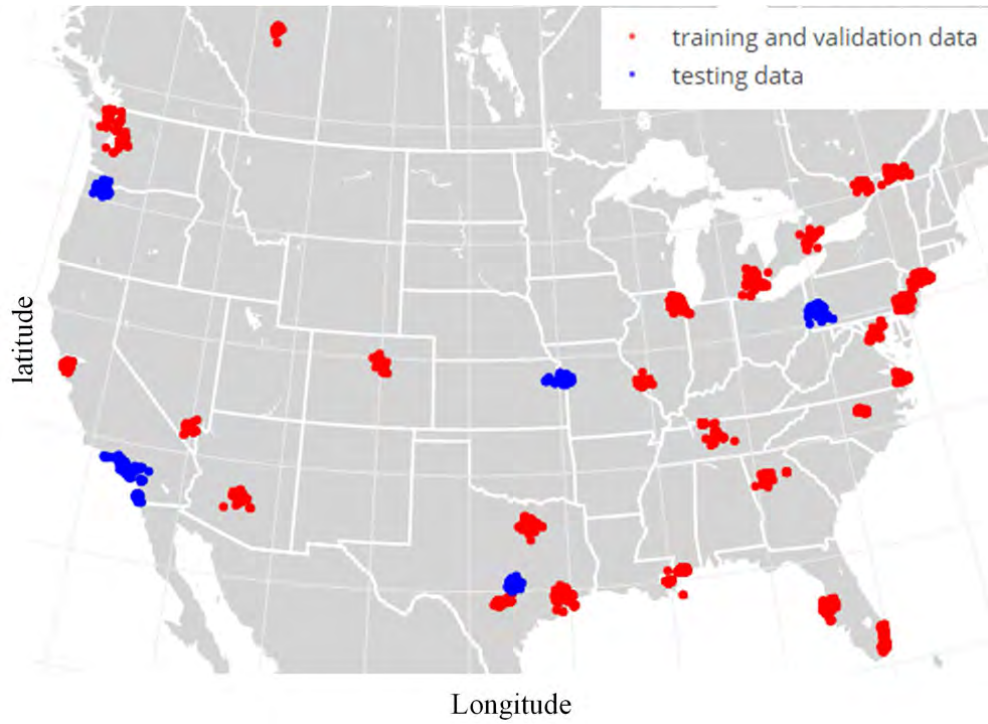Figure 7: Number of street view images of each building class.

Figure 8: GPS locations of our benchmark dataset. We split all the images into two parts: one for training (17600 images) and the others for testing (2058 images). Note that all the testing images are located in different cities with the training ones.



(a)                                     (b)

Figure 9: (a) Illustration of different looking-angles for the same building (red rectangular). (b) The corresponding retrieved street view images: Left column shows the obvious building façades, while the right two images are outliers. In order to improve the robustness, we classify several street view images for one building, and fuse their classification labels in a decision level.
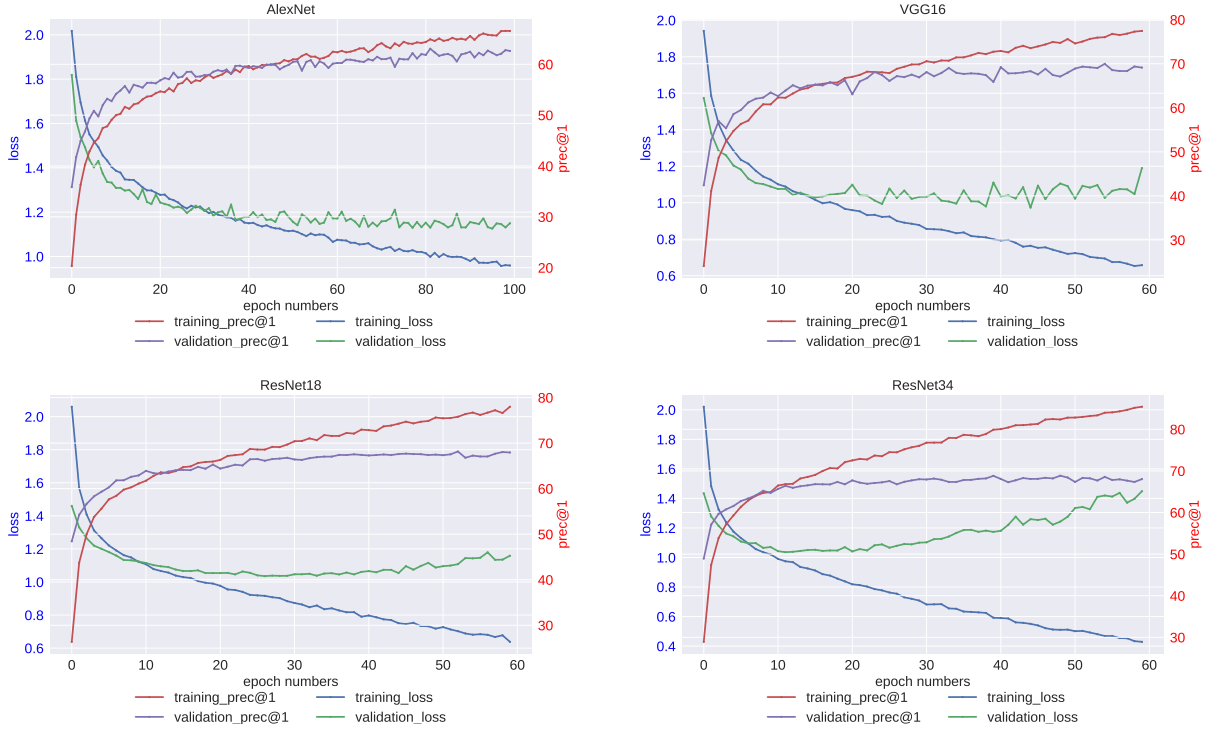
8

Figure 10: The learning and top 1-precision curves of the four networks, i.e. AlexNet (Top-left), VGG16 (Top-right), ResNet18 (Bottom-left) and ResNet34 (Bottom-right). It can be seen that training losses of the four networks reduce as the epochs increase. Besides, the validation learning curve of AlexNet converges until 80 epochs, and those of the other three networks can converge within 60 epochs. Overfitting behaviors are found in ResNet18 and ResNet34, and it is more severe in ResNet34. One plausible reason is that the total parameter number of ResNet34 (21 million) is more than that of ResNet18 (11 million). As shown by top-1 precisions, AlexNet can achieve about 65%, while the other networks can obtain about 70%.

Each training batch contained in total 64 images. The stochastic gradient descent algorithm with a learning rate of $\eta = 5 \cdot 10^{-4}$ and a momentum value of $p = 0.9$ was employed for training. To adjust the learning rate, we decayed its value by a factor of 0.1 in every 30 epochs. Cross-entropy loss was utilized for training with the weight decay parameter of $w = 10^{-5}$. The neurons of fully connected layers were dropped out by a probability of 25%. To augment the training data, we randomly cropped $224 \times 224$ pixels from the original $256 \times 256$ pixels and randomly flipped the cropped images horizontally. All the experiments were implemented with Pytorch[5] and carried out by one NVIDIA TITAN X (Pascal) 12GB GPU.

As shown in Figure 10, we plot the learning curves of both training and validation data, and calculate the corresponding top 1-precision values during training. It can be seen that training losses of the four networks reduce as the number of epochs increases. Besides, the validation learning curve of AlexNet converges until 80 epochs, and those of the other three networks can converge within 60 epochs. Overfitting behaviors are found in ResNet18 and ResNet34, and it is more severe in ResNet34. One plausible reason is that the total parameter number of ResNet34 (21 million) is more than that of ResNet18 (11 million). As shown by the top-1 precisions, AlexNet can achieve about 65%, while the other networks can obtain about 70%. For the follow-up evaluations, we choose ResNet18 trained until 40 epochs and 25 epochs of ResNet34 and compare the performances of those four networks.

### 4.2. Testing

As illustrated in Figure 11 and 12, we demonstrate the normalized confusion matrices of all the trained networks evaluated by our test data, and the associated F1 scores of the eight building classes, respectively. F1 score ($F_1$), also
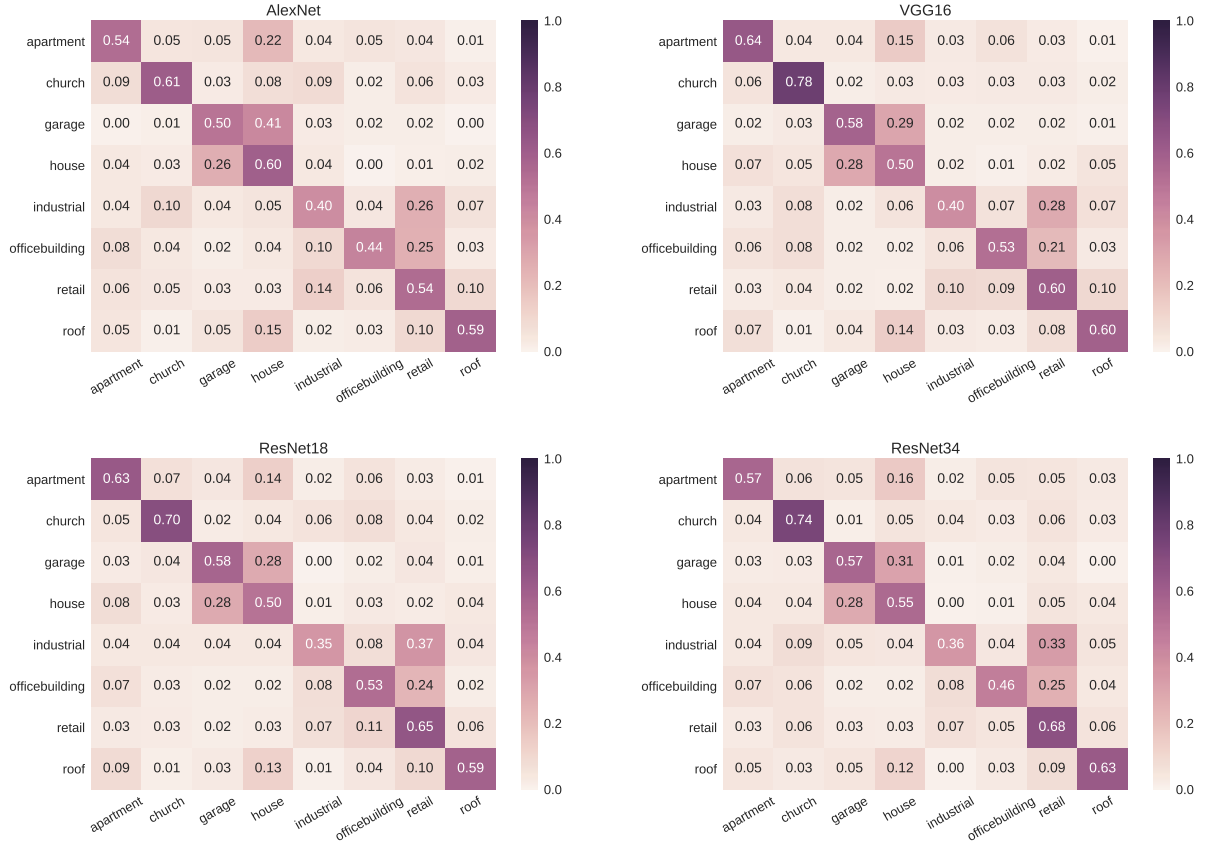
---

[5] http://pytorch.org/

Figure 11: The associated normalized confusion matrices of the four networks evaluated on the test images, i.e. AlexNet (Top-left), VGG16 (Top-right), ResNet18 (Bottom-left) and ResNet34 (Bottom-right).

known as F-measure, is a criteria to measure classification accuracy, which considers both the precision $p$ and the recall $r$. It is defined as

$$F_1 = 2 \cdot \frac{p \cdot r}{p + r}. \tag{2}$$

Moreover, the overall precisions, recalls and F1 scores of the four networks are demonstrated in Table 2. From the results, we can see that the classification performance of AlexNet is worse than the other three networks. For the classes of apartment, church, garage, industrial and office building, VGG16 achieves the highest F1 score, and for the other classes, ResNet34 is the best among them. According to the overall accuracies shown in Table 2, we choose the trained VGG16 model for the upcoming generation of building classification maps of the study areas.

Table 2: Overall precisions, recalls and F1 scores

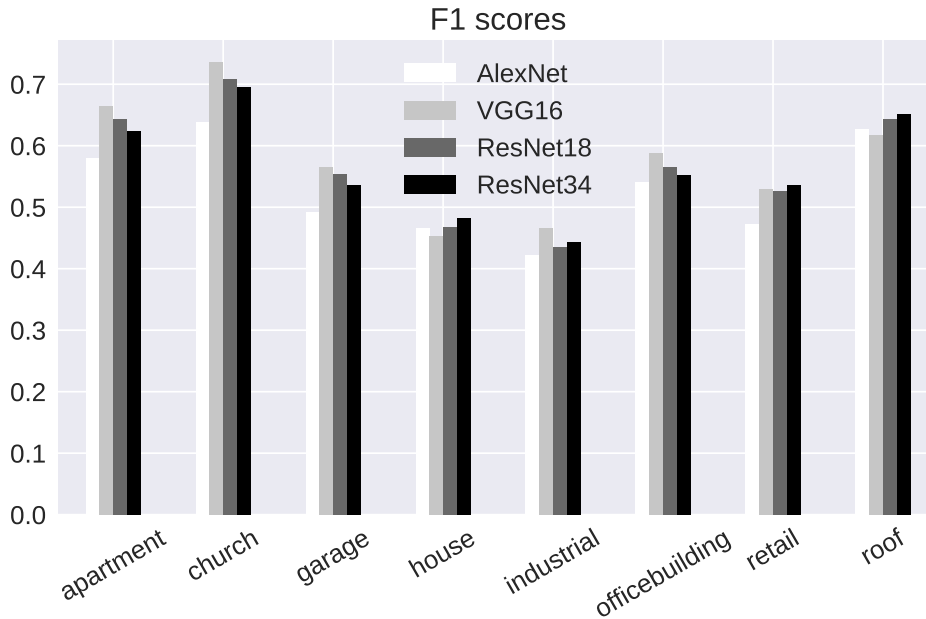| network | precision | recall | F1 score |
|---------|-----------|--------|----------|
| AlexNet | 0.55 | 0.53 | 0.53 |
| VGG16 | **0.59** | **0.58** | **0.58** |
| ResNet18 | 0.58 | 0.57 | 0.57 |
| ResNet34 | **0.59** | 0.57 | 0.56 |

10

Figure 12: F1 score performances of the four trained networks on the eight building classes. For the classes of apartment, church, garage, industrial and office building, VGG16 achieves the highest F1 score, and for the other classes, ResNet34 is the best among them.

### 4.3. Building classification maps of study areas

#### 4.3.1. Maps of study areas in Vancouver and Fort Worth

One testing area in Vancouver (image is from Google Earth) can be seen in Figure 13. The associated ground truth and predicted building classification maps are present in Figure 14, where different colors represent different building classes. We also draw the corresponding confusion matrix of the inferred result in Figure 15. The total number of building instances in this area is 196. Our result predicts 93 apartments, 10 churches, 13 garages, 24 houses, 1 industrial building, 21 office buildings, 26 retails and 1 roof. 7 buildings are not classified, since no corresponding street view images are found. Moreover, the confidence score for the class of each building is shown by the opacity of the associated color mask, i.e. the higher the opacity, the larger the confidence score and vice versa. From the results, we can see that this study area is mainly composed of apartments, which indicates that it is a residential district and there may be a high population density of this area. Our analysis is confirmed by the 2011 census population density of Vancouver downloaded from the website[6], as shown in Figure 15 (Right). The white rectangle in the figure marks the study area which has the highest population density (over $10000/km^2$). Such classification map gives an insight into the social structure of an residential area. For example, the houses and retails are both grouped together at the right corner of this district.

Another testing area is located in Fort Worth, shown by the red rectangle area in Figure 16. The ground truth and predicted building classification maps are present in Figure 17. The associated confusion matrix is demonstrated in Figure 18. The total number of buildings in this area is 316. Our result predicts 34 apartments, 30 churches, 2 houses, 19 industrial buildings, 152 office buildings, 28 retails and 6 roofs. There are no street view images for the remaining 45 buildings. According to the predicted result, we can see that this area is mainly composed of office buildings, which indicates that it is a business district and may locate in the center of Fort Worth.
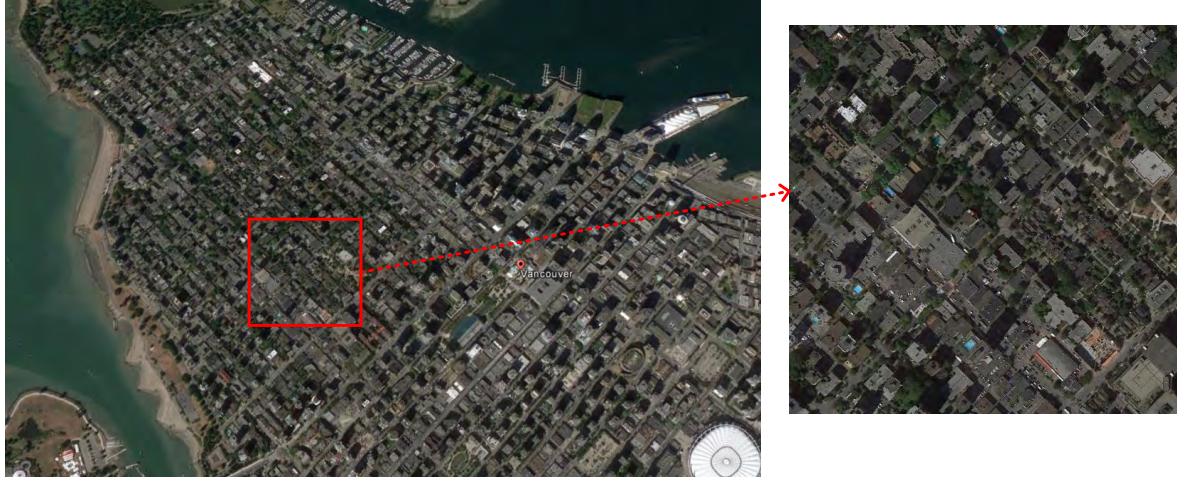
---

[6]https://blogs.ubc.ca/maps/2013/07/03/vancouverpopulationdensity/

Figure 13: Illustration of one study area in Vancouver (image is from Google Earth).



Classification map (ground truth)          Predicted classification map (VGG16)

- apartment
- church
- garage
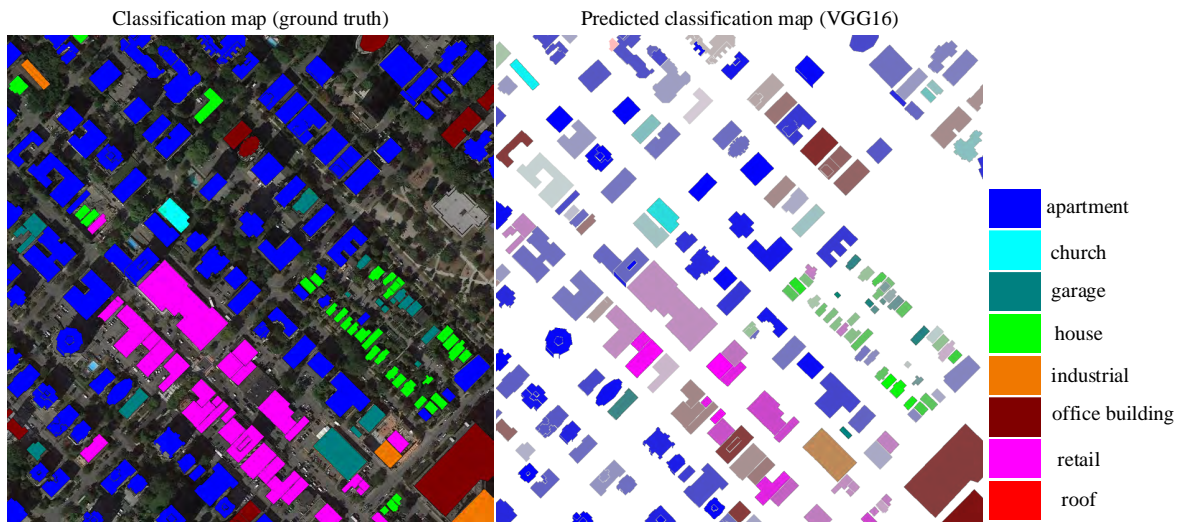- house
- industrial
- office building
- retail
- roof

Figure 14: The predicted building classification map (Right), along with the ground truth (Left), where different colors represent different building classes. The total number of building instances in this area is 196. Our result predicts 93 apartments, 10 churches, 13 garages, 24 houses, 1 industrial building, 21 office buildings, 26 retails and 1 roof. 7 buildings are not classified, since no corresponding street view images are found. Moreover, the confidence score for the class of each building is shown by the opacity of the associated color mask, i.e. the higher the opacity, the larger the confidence score and vice versa.
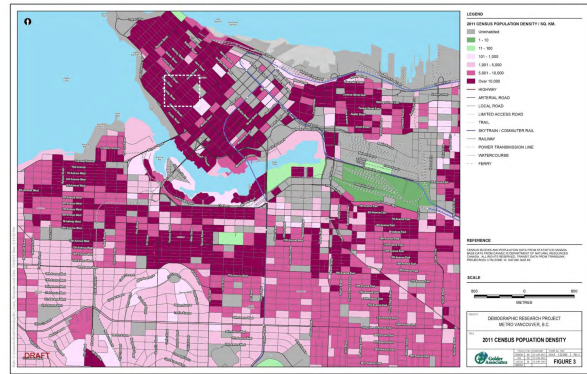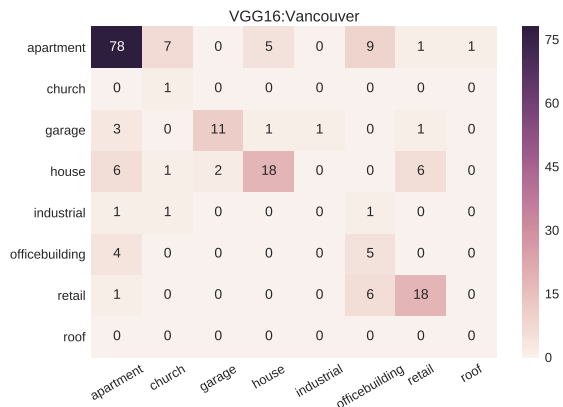
Figure 15: (Left) The confusion matrix of the classification result of the area in Vancouver. We can see that this area is mainly composed by apartments. (Right) 2011 census population density of Vancouver. The white rectangle indicates the study area, which has a high population density of over $10000/km^2$
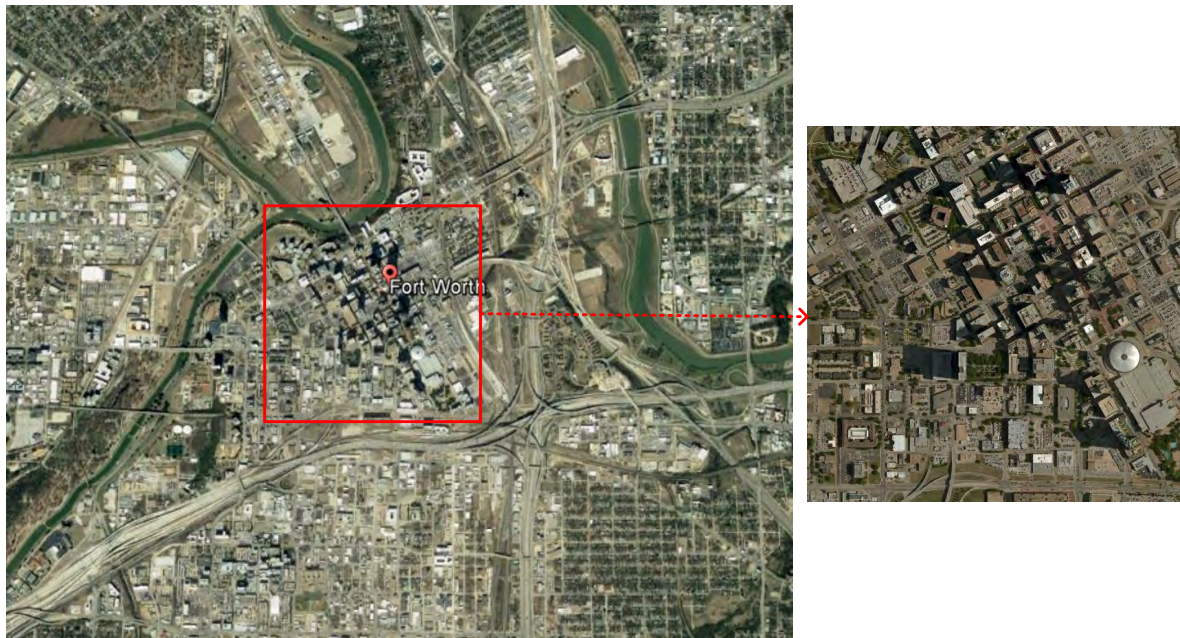


Figure 16: Illustration of one study area in Fort Worth (image is from Google Earth).

Classification map (ground truth)      Predicted classification map (VGG16)

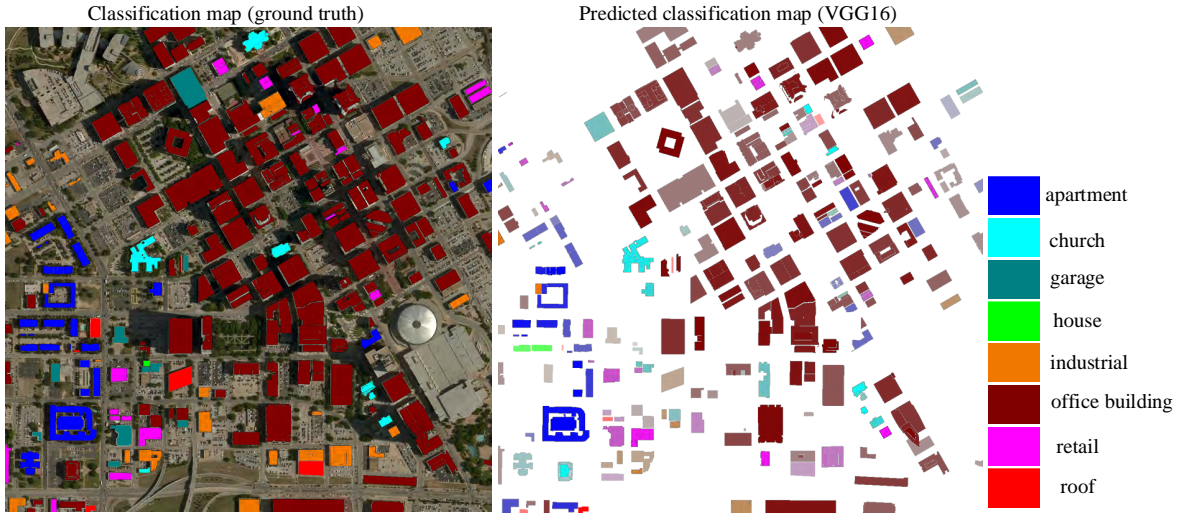| | |
|---|---|
| ■ | apartment |
| ■ | church |
| ■ | garage |
| ■ | house |
| ■ | industrial |
| ■ | office building |
| ■ | retail |
| ■ | roof |

Figure 17: The predicted building classification map (Right), along with the ground truth (Left). The total number of buildings in this area is 316. Our result predicts 34 apartments, 30 churches, 2 houses, 19 industrial buildings, 152 office buildings, 28 retails and 6 roofs. There are no street view images for the remaining 45 buildings.

**VGG16:Fort Worth**

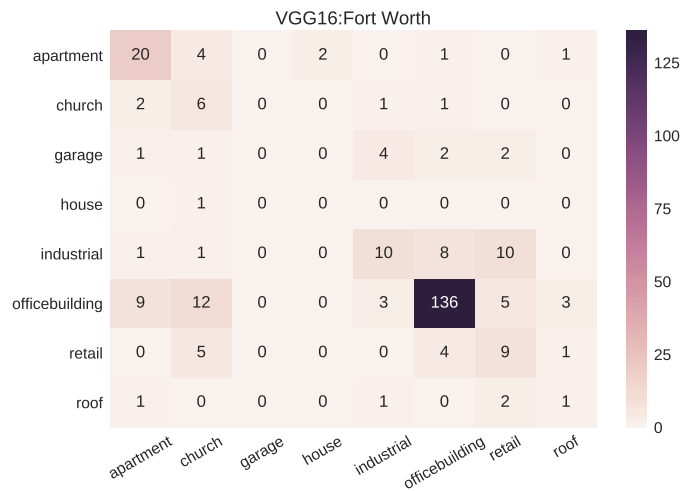| | apartment | church | garage | house | industrial | officebuilding | retail | roof |
|---|---|---|---|---|---|---|---|---|
| apartment | 20 | 4 | 0 | 2 | 0 | 1 | 0 | 1 |
| church | 2 | 6 | 0 | 0 | 1 | 1 | 0 | 0 |
| garage | 1 | 1 | 0 | 0 | 4 | 2 | 2 | 0 |
| house | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| industrial | 1 | 1 | 0 | 0 | 10 | 8 | 10 | 0 |
| officebuilding | 9 | 12 | 0 | 0 | 3 | 136 | 5 | 3 |
| retail | 0 | 5 | 0 | 0 | 0 | 4 | 9 | 1 |
| roof | 1 | 0 | 0 | 0 | 1 | 0 | 2 | 1 |

Figure 18: The confusion matrix of the classification result of the area in Fort Worth. We can see that this area is mainly composed by office buildings, which indicates that it is a business district and may locate in the center of Fort Worth.
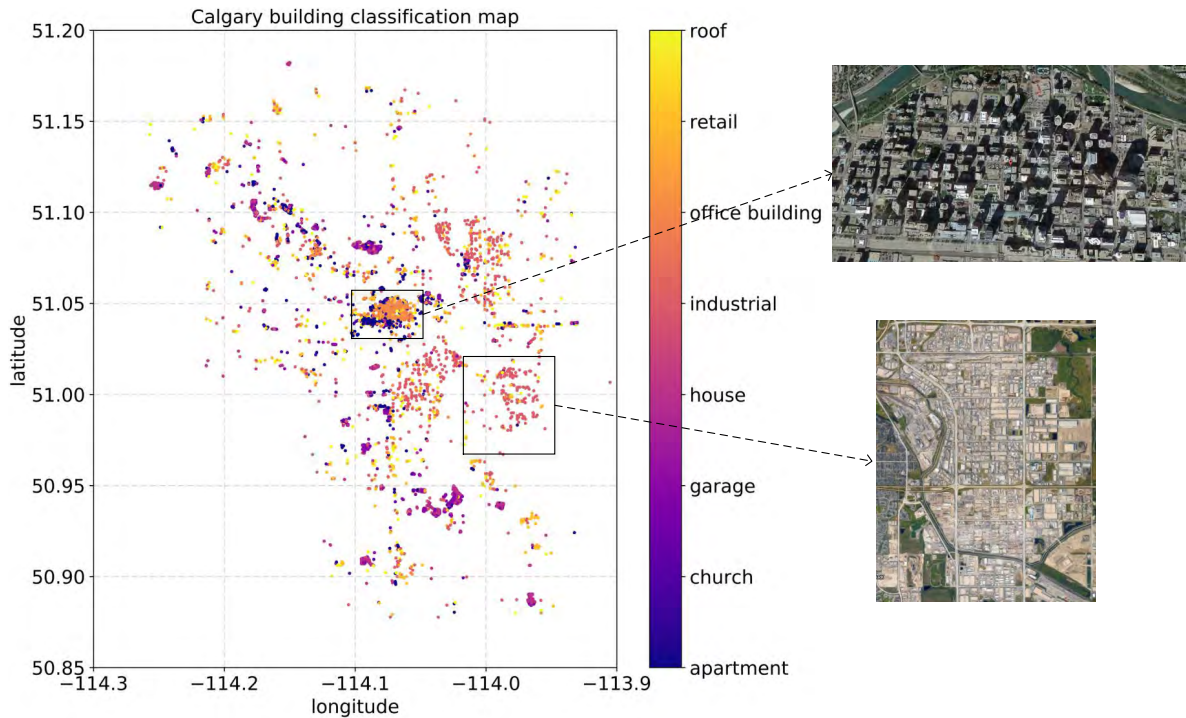
14

Figure 19: The city-scale building classification map of Calgary, where each classified building instance is displayed as a colored point with GPS coordinates. It is obvious that there are three main industrial districts, and the downtown area is crowded by office buildings. Correspondingly, we also present the remote sensing images of one industrial and the downtown areas (black rectangles). Such classification map can infer that Calgary is an industry city with single central business district to which the three main industrial blocks are located.

### 4.3.2. City-scale Maps of Calgary, Boston and Toronto

As shown in Figure 19, 21 and 24, we provide the city-scale building classification maps of Calgary, Boston and Toronto based on classifying the retrieved 6124, 64389 and 45978 building street view images, respectively, where each classified building instance is displayed as a colored point with its GPS coordinate. Besides, Figure 20, 23 and 26 demonstrate the associated numerical proportions of building classes based on the classification results. In order to quantitatively analyze the performance, 1000 buildings in each city are randomly selected and their associated building tags from OSM are retrieved according to their GPS locations. The classification performances of the three cities are demonstrated in Table 3, 4 and 5, respectively. Table 3 demonstrates that the overall accuracy of the classification result in Calgary is around 0.7, given the retrieved 1000 building tags from OSM. As illustrated in Table 4, by comparing with the building tags from OSM, the overall accuracy of the classification map in Boston can reach around 0.55. According to Table 5, more than 75% buildings in Toronto can be accurately classified.

As shown by the classification map of Calgary, it is obvious that there are three main industrial districts, and the downtown area is crowded by office buildings. Correspondingly, we also present the remote sensing images of one industrial and the downtown areas (black rectangles). Such classification map can infer that Calgary is an industry city with single central business district to which the three main industrial blocks are located.

According to the map of Boston, houses obviously dominate the buildings of the city, and they are located around the city center. Besides, Boston is also with one single central business district (noted by the black dashed rectangle), since most office buildings and apartments locate in this area. The associated remote sensing image is demonstrated at the bottom of Figure 21. As shown by the distribution maps of office buildings, houses and apartments plotted in Figure 22, we can see that the densities of office buildings and apartments decrease from the center to its outside, while it is contrary of the house density. In addition, Boston is not an industry city, since no large block of industrial
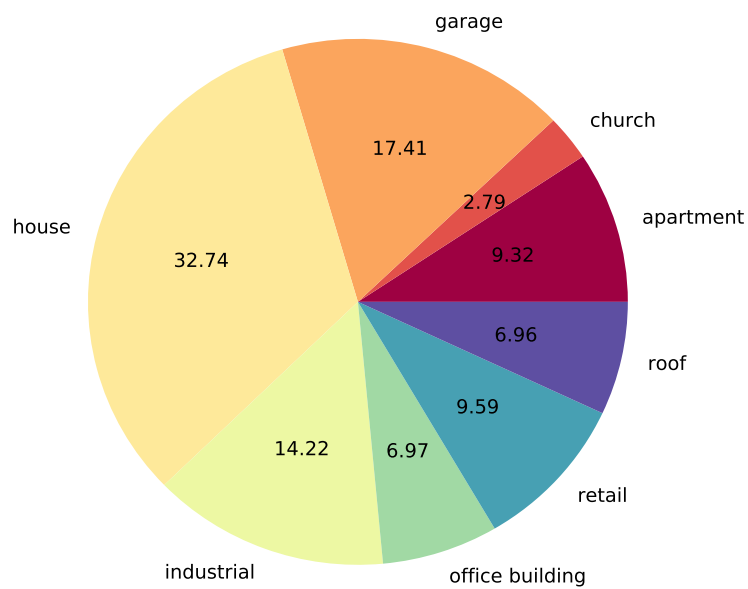
15

Figure 20: Pie chart of the building class proportions of the predicted buildings of Calgary.

Table 3: Classification performance of randomly selected 1000 buildings of Calgary

|  | precision | recall | F1 score | support |
|---|---|---|---|---|
| apartment | 0.54 | 0.77 | 0.64 | 56 |
| church | 0.00 | 0.00 | 0.00 | 1 |
| garage | 0.41 | 0.90 | 0.57 | 124 |
| house | 0.97 | 0.62 | 0.75 | 630 |
| industrial | 0.51 | 0.80 | 0.63 | 82 |
| office building | 0.65 | 0.19 | 0.29 | 58 |
| retail | 0.33 | 0.37 | 0.35 | 43 |
| roof | 0.15 | 0.83 | 0.25 | 6 |
| **overall** | 0.78 | 0.64 | 0.66 | 1000 |

Table 4: Classification performance of randomly selected 1000 buildings of Boston

|  | precision | recall | F1 score | support |
|---|---|---|---|---|
| apartment | 0.35 | 0.42 | 0.38 | 137 |
| church | 0.06 | 0.80 | 0.11 | 5 |
| garage | 0.51 | 0.38 | 0.43 | 221 |
| house | 0.69 | 0.61 | 0.65 | 546 |
| industrial | 0.07 | 0.25 | 0.11 | 4 |
| office building | 0.58 | 0.62 | 0.60 | 60 |
| retail | 0.20 | 0.42 | 0.27 | 19 |
| roof | 0.62 | 0.62 | 0.62 | 8 |
| **overall** | 0.58 | 0.53 | 0.55 | 1000 |

Table 5: Classification performance of randomly selected 1000 buildings of Toronto

|  | precision | recall | F1 score | support |
|---|---|---|---|---|
| apartment | 0.73 | 0.83 | 0.78 | 212 |
| church | 0.29 | 0.59 | 0.39 | 22 |
| garage | 0.18 | 0.42 | 0.25 | 33 |
| house | 0.94 | 0.73 | 0.82 | 575 |
| industrial | 0.36 | 0.79 | 0.49 | 24 |
| office building | 0.04 | 0.25 | 0.06 | 4 |
| retail | 0.84 | 0.50 | 0.63 | 117 |
| roof | 0.33 | 0.92 | 0.49 | 13 |
| **overall** | 0.82 | 0.71 | 0.75 | 1000 |

districts is observed in the classification map and the proportion of industrial buildings is very low.

From the maps of Figure 24 and 25, most apartments and office buildings are located in the center of Toronto, and most industrial buildings are distributed in the regions around it. As shown in Figure 26, the second largest proportion of building classes is apartment, which indicates that the population density is high in Toronto, especially in the city center. Besides, around 10% buildings are industrial, which indicates industry is one of the fields which contribute most to the economy of Toronto.

## 5. Discussion

In our training experiments, all the fully connected layers are initialized randomly. As for ResNet, there is only one fully connected layer (*softmax* layer) in the architecture and we utilize the network pre-trained on *ImageNet* dataset which contains totally 1000 classes. The parameters of the fully connected layer cannot be directly transferred to our task, since there are 8 classes in our dataset. While, for AlexNet and VGG16, besides the last fully connected layer (*softmax* layer), there are two more fully connected (fc) layers to be initialized. Taking VGG16 as an example, we took two experiments with the same hyperparameters for training the network on our benchmark dataset, while those two fully connected layers were initialized in two ways, i.e. initialized randomly and by the parameters pretrained with *ImageNet*. As shown in 27, VGG16 where the two fully connected layers were initialized by the pretrained parameters did accelerate the training of the network, since it can achieve higher classification accuracy than the one where the fc layers were initialized randomly during the first several epochs. However, both of them can converge to comparable classification accuracies at last.

According to the classification accuracies of eight building classes, churches are relatively easier to recognize than any other classes, since their structures are more unique, while some classes are not easily identified, e.g. retails and industrials. There are the following reasons that may influence the classification results. Firstly, since the ground truth labels come from the OSM users, manually labeling errors among some building classes exist in the benchmark dataset, especially for those with similar façade structures, e.g. some industrial and office buildings. As shown in Figure 28 (Left), the building displayed by the street view image tends to be an office building, while the building tag
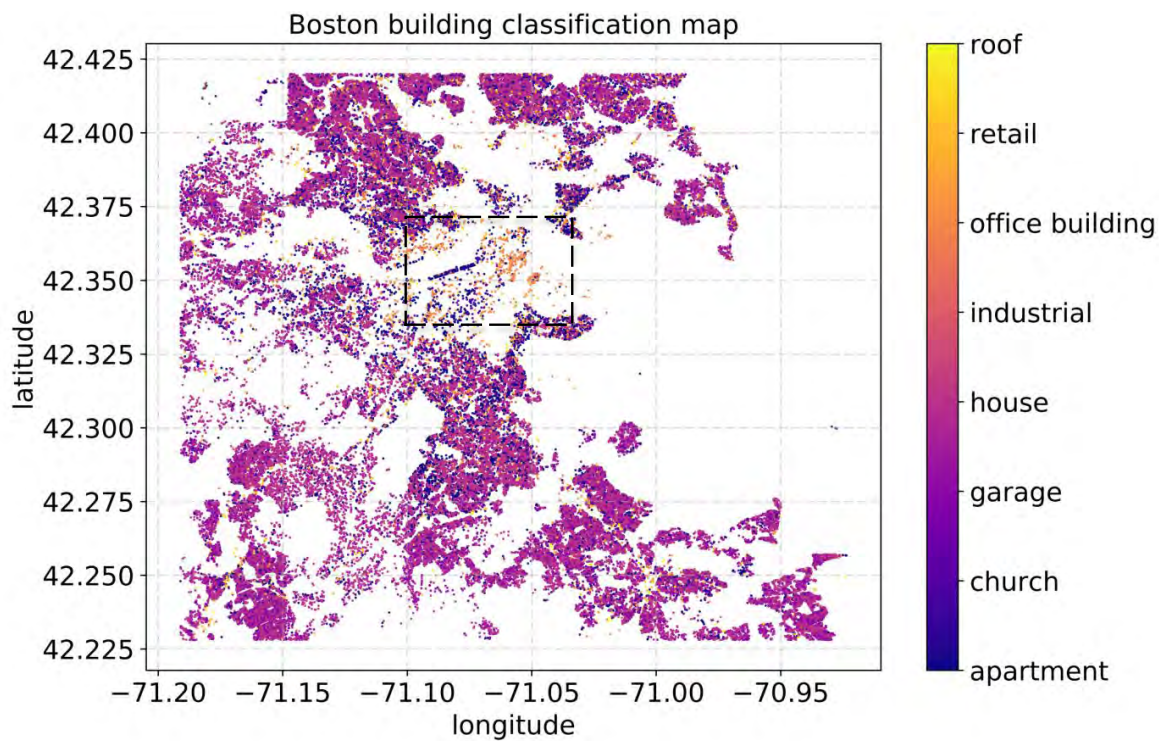
Figure 21: The city-scale building classification map of Boston, where each classified building instance is displayed as a colored point with GPS coordinates. Since most office buildings and apartments locate in the cropped area, it can be observed that Boston is with one single central business district. It is not an industry city, as no large block of industrial districts is found.
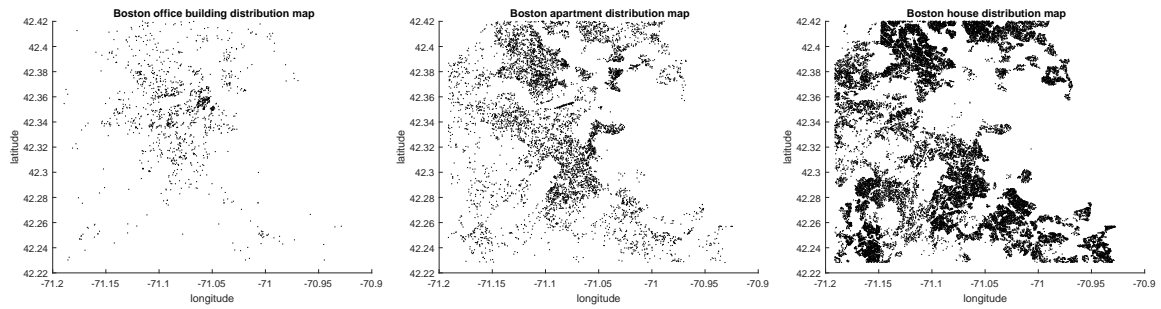
Figure 22: Office building, apartment and house distribution maps of Boston. We can see that both the densities of office buildings and apartments decrease from the center to its outside, while it is contrary of the house density.



Figure 23: Pie chart of the building class proportions of the predicted buildings of Boston.
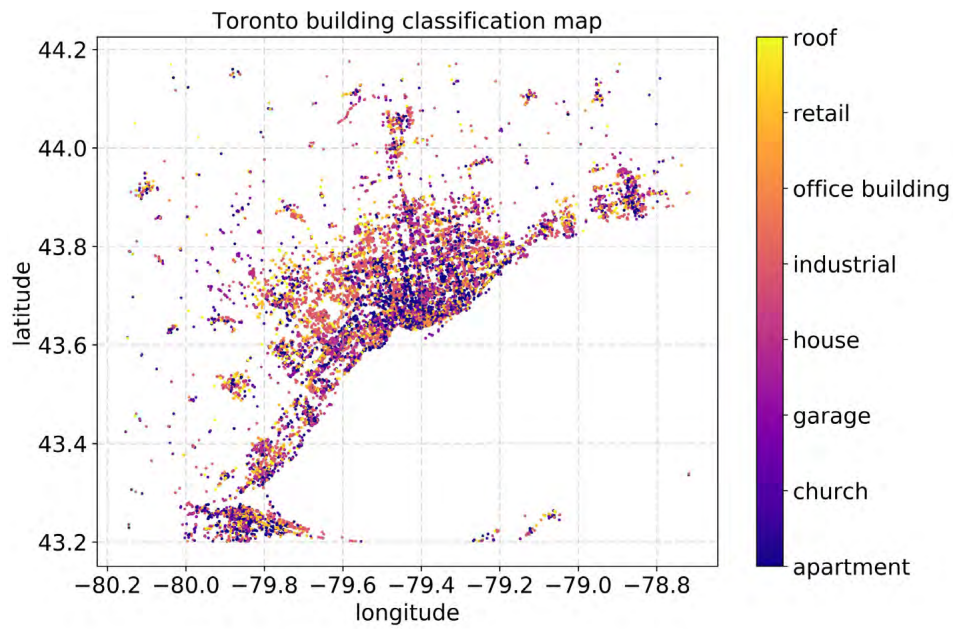
Figure 24: The city-scale building classification map of Toronto and the associated remote sensing image of the central city. Most apartments and office buildings are located in the center of Toronto, and most industrial buildings are distributed in the regions around it.
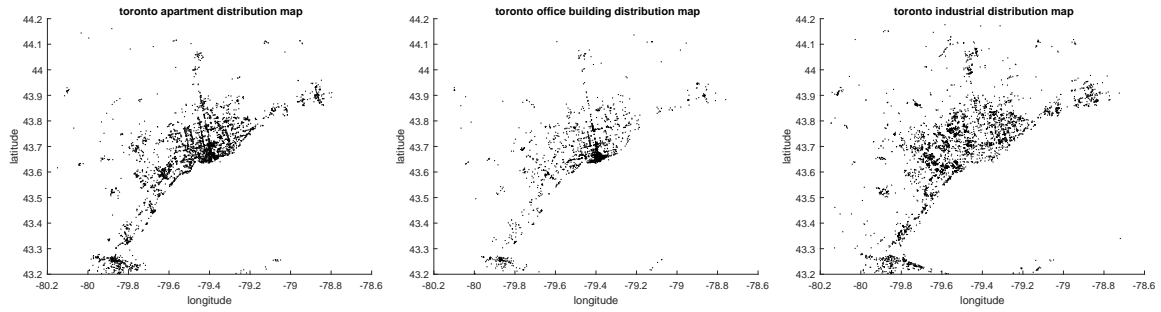
Figure 25: Apartment, office and industrial building distribution maps of Toronto.
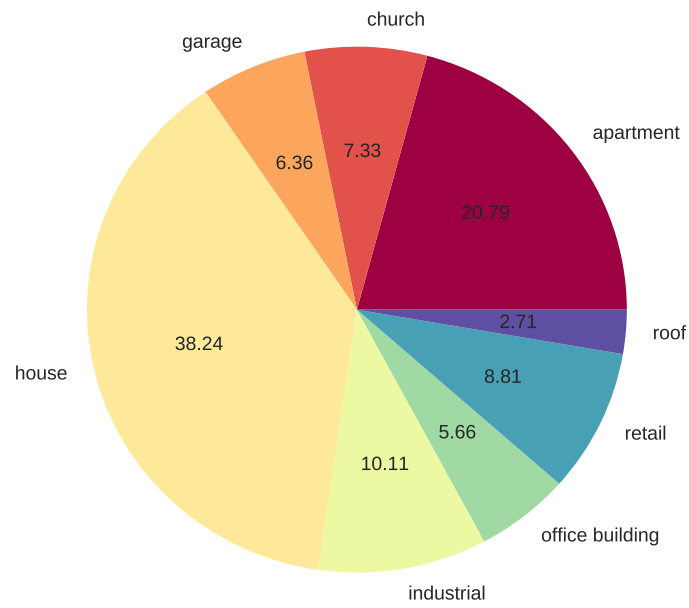


Figure 26: Pie chart of the building class proportions of the predicted buildings of Toronto.
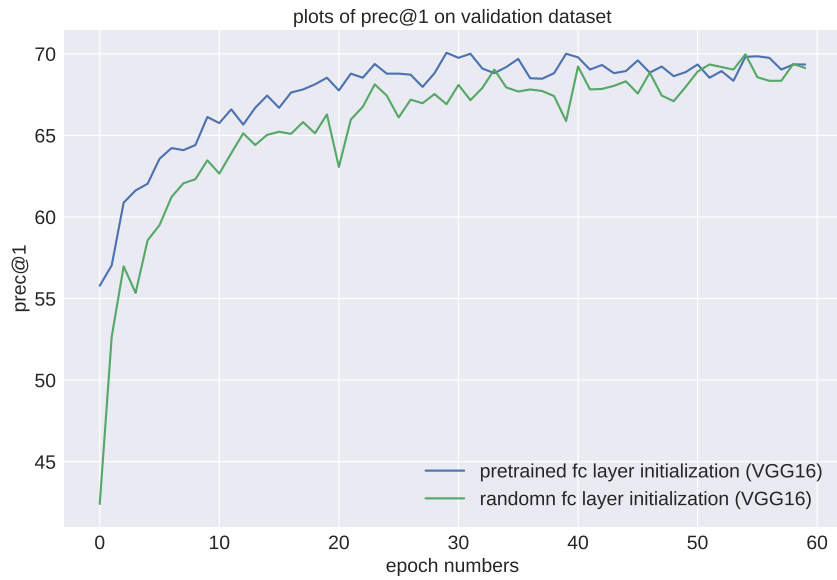
Figure 27: Top-1 precision curves of VGG16 on the validation dataset with different initializations of fully connected layers.
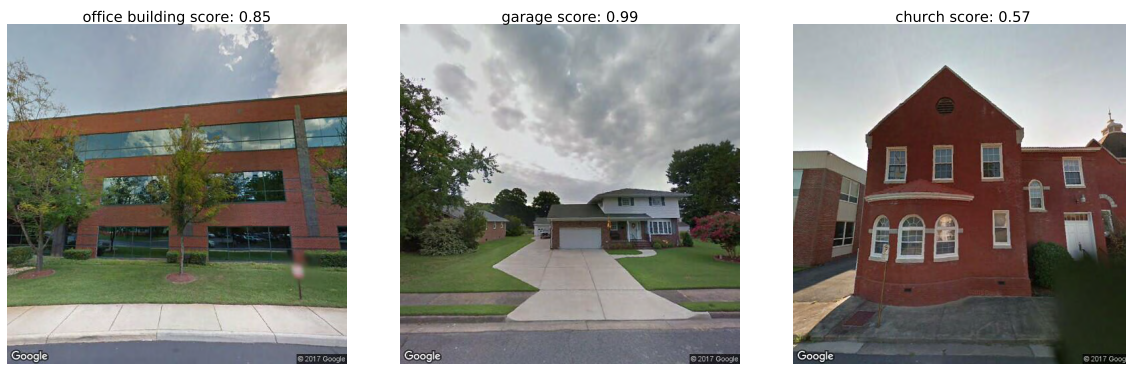


Figure 28: Some results show the reasons that may induce classification errors. (Left) The building displayed by the street view image tends to be an office building, while the ground truth building tag retrieved from OSM is industrial. (Middle) The building demonstrated in the image is a house, while it is misclassified as garage, since there are both garage and house structures demonstrated in the image. (Right) Although the building is correctly recognized, the confidence score is not so high, since the typical façade structure of church is not displayed in the retrieved image.

retrieved from OSM is industrial. Secondly, some street view images include multiple buildings of different classes, e.g. a house with a garage by its side. From Figure 28 (Middle), the building demonstrated in the image is a house, while it is misclassified as a garage. Lastly, side faces of buildings are displayed in some retrieved street view images, thus the corresponding façade features are not rich for the classification. As illustrated by Figure 28 (Right), although the building is correctly recognized, the confidence score is not so high, since the typical façade structure of church is not displayed in the retrieved image.

It is worth noting that as an alternative, a building rejection class can be added to replace the outlier removal procedure, depending on the quality of input data.

## 6. Conclusion and future work

In this paper, we presented a framework for building instance classification, which tended to provide more informative classification maps. With this approach, relatively high accuracies could be achieved for land-use classification of individual buildings. For this task, we built a street view benchmark dataset with eight building categories for training and testing. By investigating four different CNN architectures, we chose VGG16 to predict building instance classification maps on region and city scales. Such maps help us to get insight of urban areas, and have the potential for many innovative urban analysis, e.g. very high resolution urban population density mapping, urban social structure understanding, city economy structure analysis and general urban planning.

For the future work, to improve the classification performance, other information can be fused, e.g. text descriptions associated with social media images and text information displayed in images, e.g. brand names. Also, in order to obtain denser building classification maps, information from remote sensing images and GIS maps can be exploited for those buildings without street view images. In case that building footprints cannot be retrieved from GIS maps, a method of individual building detection in remote sensing images should be also developed.

## 7. Acknowledgment

## References

[1] J. R. Anderson, A land use and land cover classification system for use with remote sensor data, volume 964, US Government Printing Office, 1976.

[2] M. Pal, P. M. Mather, An assessment of the effectiveness of decision tree methods for land cover classification, Remote sensing of environment 86 (2003) 554–565.

[3] F. Yuan, K. E. Sawaya, B. C. Loeffelholz, M. E. Bauer, Land cover classification and change analysis of the twin cities (minnesota) metropolitan area by multitemporal landsat remote sensing, Remote sensing of Environment 98 (2005) 317–328.

[4] W. L. Stefanov, M. S. Ramsey, P. R. Christensen, Monitoring urban land cover change: An expert system approach to land cover classification of semiarid to arid urban centers, Remote sensing of Environment 77 (2001) 173–185.

[5] V. F. Rodriguez-Galiano, B. Ghimire, J. Rogan, M. Chica-Olmo, J. P. Rigol-Sanchez, An assessment of the effectiveness of a random forest classifier for land-cover classification, ISPRS Journal of Photogrammetry and Remote Sensing 67 (2012) 93–104.

[6] A. Albert, J. Kaur, M. Gonzalez, Using convolutional networks and satellite imagery to identify patterns in urban environments at a large scale, arXiv preprint arXiv:1704.02965 (2017).

[7] D. Lu, Q. Weng, Use of impervious surface in urban land-use classification, Remote Sensing of Environment 102 (2006) 146–160.

[8] P. Gong, D. J. Marceau, P. J. Howarth, A comparison of spatial feature extraction algorithms for land-use classification with spot hrv data, Remote sensing of environment 40 (1992) 137–151.

[9] J. D. Paola, R. A. Schowengerdt, A detailed comparison of backpropagation neural network and maximum-likelihood classifiers for urban land use classification, IEEE Transactions on Geoscience and remote sensing 33 (1995) 981–996.

[10] F. Pacifici, M. Chini, W. J. Emery, A neural network approach using multi-scale textural metrics from very high-resolution panchromatic imagery for urban land-use classification, Remote Sensing of Environment 113 (2009) 1276–1292.

[11] S. Khorram, J. A. Brockhaus, H. M. Cheshire, Comparson of landsat mss and tm data for urban land-use classification, IEEE transactions on geoscience and remote sensing (1987) 238–243.

[12] K. Di, D. Li, D. Li, Land use classification of remote sensing image with gis data based on spatial data mining techniques, International Archives of Photogrammetry and Remote Sensing 33 (2000) 238–245.

[13] G. Cheng, J. Han, L. Guo, Z. Liu, S. Bu, J. Ren, Effective and efficient midlevel visual elements-oriented land-use classification using vhr remote sensing images, IEEE Transactions on Geoscience and Remote Sensing 53 (2015) 4238–4249.

[14] X. Huang, D. Wen, J. Li, R. Qin, Multi-level monitoring of subtle urban changes for the megacities of china using high-resolution multi-view satellite imagery, Remote Sensing of Environment 196 (2017) 56–75.

[15] X. Huang, Q. Lu, L. Zhang, A multi-index learning approach for classification of high-resolution remotely sensed images over urban areas, ISPRS Journal of Photogrammetry and Remote Sensing 90 (2014) 36–48.

[16] X. Huang, H. Liu, L. Zhang, Spatiotemporal detection and analysis of urban villages in mega city regions of china using high-resolution remotely sensed imagery, IEEE Transactions on Geoscience and Remote Sensing 53 (2015) 3639–3657.

[17] OpenStreetMap contributors, Planet dump retrieved from https://planet.osm.org , `https://www.openstreetmap.org`, 2017.

[18] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., Imagenet large scale visual recognition challenge, International Journal of Computer Vision 115 (2015) 211–252.

[19] B. Zhou, A. Khosla, A. Lapedriza, A. Torralba, A. Oliva, Places: An image database for deep scene understanding, arXiv preprint arXiv:1610.02055 (2016).

[20] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick, Microsoft coco: Common objects in context, in: European Conference on Computer Vision, Springer, pp. 740–755.

[21] D. G. Lowe, Object recognition from local scale-invariant features, in: Computer vision, 1999. The proceedings of the seventh IEEE international conference on, volume 2, Ieee, pp. 1150–1157.

[22] Y. Yang, S. Newsam, Bag-of-visual-words and spatial extensions for land-use classification, in: Proceedings of the 18th SIGSPATIAL international conference on advances in geographic information systems, ACM, pp. 270–279.

[23] A. M. Cheriyadat, Unsupervised feature learning for aerial scene classification, IEEE Transactions on Geoscience and Remote Sensing 52 (2014) 439–451.

[24] Q. Zhu, Y. Zhong, B. Zhao, G.-S. Xia, L. Zhang, Bag-of-visual-words scene classifier with local and global features for high spatial resolution remote sensing imagery, IEEE Geoscience and Remote Sensing Letters 13 (2016) 747–751.

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

[26] H. Taubenböck, M. Klotz, M. Wurm, J. Schmieder, B. Wagner, M. Wooster, T. Esch, S. Dech, Delineation of central business districts in mega city regions using remotely sensed data, Remote sensing of Environment 136 (2013) 386–401.

[27] Z. Wang, N. M. Nasrabadi, T. S. Huang, Spatial–spectral classification of hyperspectral images using discriminative dictionary designed by learning vector quantization, IEEE Transactions on Geoscience and Remote Sensing 52 (2014) 4808–4822.

[28] S. Yang, H. Jin, M. Wang, Y. Ren, L. Jiao, Data-driven compressive sampling and learning sparse coding for hyperspectral image classification, IEEE Geoscience and Remote Sensing Letters 11 (2014) 479–483.

[29] X. Sun, N. M. Nasrabadi, T. D. Tran, Task-driven dictionary learning for hyperspectral image classification with structured sparsity constraints, IEEE Transactions on Geoscience and Remote Sensing 53 (2015) 4457–4471.

[30] I. Rigas, G. Economou, S. Fotopoulos, Low-level visual saliency with application on aerial imagery, IEEE Geoscience and Remote Sensing Letters 10 (2013) 1389–1393.

[31] F. Zhang, B. Du, L. Zhang, Saliency-guided unsupervised feature learning for scene classification, IEEE Transactions on Geoscience and Remote Sensing 53 (2015) 2175–2184.

[32] D. Tuia, R. Flamary, M. Barlaud, Nonconvex regularization in remote sensing, IEEE Transactions on Geoscience and Remote Sensing 54 (2016) 6470–6480.

[33] X. Yao, J. Han, G. Cheng, X. Qian, L. Guo, Semantic annotation of high-resolution satellite images via weakly supervised learning, IEEE Transactions on Geoscience and Remote Sensing 54 (2016) 3660–3671.

[34] D. Tuia, R. Flamary, N. Courty, Multiclass feature learning for hyperspectral image classification: Sparse and hierarchical solutions, ISPRS Journal of Photogrammetry and Remote Sensing 105 (2015) 272–285.

[35] G. Cheng, J. Han, X. Lu, Remote sensing image scene classification: Benchmark and state of the art, Proceedings of the IEEE (2017).

[36] X. Ma, H. Wang, J. Wang, Semisupervised classification for hyperspectral image based on multi-decision labeling and deep feature learning, ISPRS Journal of Photogrammetry and Remote Sensing 120 (2016) 99–107.

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

[38] Y. Chen, Z. Lin, X. Zhao, G. Wang, Y. Gu, Deep learning-based classification of hyperspectral data, IEEE Journal of Selected topics in applied earth observations and remote sensing 7 (2014) 2094–2107.

[39] Q. Zou, L. Ni, T. Zhang, Q. Wang, Deep learning based feature selection for remote sensing scene classification, IEEE Geoscience and Remote Sensing Letters 12 (2015) 2321–2325.

[40] O. A. Penatti, K. Nogueira, J. A. dos Santos, Do deep features generalize from everyday objects to remote sensing and aerial scenes domains?, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, pp. 44–51.

[41] D. Marmanis, M. Datcu, T. Esch, U. Stilla, Deep learning earth observation classification using imagenet pretrained networks, IEEE Geoscience and Remote Sensing Letters 13 (2016) 105–109.

[42] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A large-scale hierarchical image database, in: Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on, IEEE, pp. 248–255.

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

[44] F. P. Luus, B. P. Salmon, F. van den Bergh, B. T. J. Maharaj, Multiview deep learning for land-use classification, IEEE Geoscience and

Remote Sensing Letters 12 (2015) 2448–2452.

[45] Q. Liu, R. Hang, H. Song, F. Zhu, J. Plaza, A. Plaza, Adaptive deep pyramid matching for remote sensing scene classification, arXiv preprint arXiv:1611.03589 (2016).

[46] J. Han, D. Zhang, G. Cheng, L. Guo, J. Ren, Object detection in optical remote sensing images based on weakly supervised learning and high-level feature learning, IEEE Transactions on Geoscience and Remote Sensing 53 (2015) 3325–3337.

[47] G. Cheng, P. Zhou, J. Han, Learning rotation-invariant convolutional neural networks for object detection in vhr optical remote sensing images, IEEE Transactions on Geoscience and Remote Sensing 54 (2016) 7405–7415.

[48] A. Romero, C. Gatta, G. Camps-Valls, Unsupervised deep feature extraction for remote sensing image classification, IEEE Transactions on Geoscience and Remote Sensing 54 (2016) 1349–1362.

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

[50] L. Mou, P. Ghamisi, X. X. Zhu, Deep recurrent neural networks for hyperspectral image classification, IEEE Transactions on Geoscience and Remote Sensing 55 (2017) 3639–3655.

[51] L. Mou, P. Ghamisi, X. X. Zhu, Unsupervised spectral-spatial feature learning via deep residual conv-deconv network for hyperspectral image classification, IEEE Transactions on Geoscience and Remote Sensing 56 (2018) 391–406.

[52] L. H. Hughes, M. Schmitt, L. Mou, Y. Wang, X. X. Zhu, Identifying corresponding patches in sar and optical images with a pseudo-siamese cnn, arXiv preprint arXiv:1801.08467 (2018).

[53] L. Zhang, L. Zhang, B. Du, Deep learning for remote sensing data: A technical tutorial on the state of the art, IEEE Geoscience and Remote Sensing Magazine 4 (2016) 22–40.

[54] X. X. Zhu, D. Tuia, L. Mou, G. S. Xia, L. Zhang, F. Xu, F. Fraundorfer, Deep learning in remote sensing: A comprehensive review and list of resources, IEEE Geoscience and Remote Sensing Magazine 5 (2017) 8–36.

[55] D. Anguelov, C. Dulong, D. Filip, C. Frueh, S. Lafon, R. Lyon, A. Ogale, L. Vincent, J. Weaver, Google street view: Capturing the world at street level, Computer 43 (2010) 32–38.

[56] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556 (2014).

[57] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: Advances in neural information processing systems, pp. 1097–1105.

[58] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778.