Weakly Supervised Silhouette-based Semantic Scene Change Detection

Ken Sakurada, Mikiya Shibuya, Weimin Wang

Abstract—This paper presents a novel semantic scene change detection scheme with only weak supervision. A straightforward approach for this task is to train a semantic change detection network directly from a large-scale dataset in an end-to-end manner. However, a specific dataset for this task, which is usually labor-intensive and time-consuming, becomes indispensable. To avoid this problem, we propose to train this kind of network from existing datasets by dividing this task into change detection and semantic extraction. On the other hand, the difference in camera viewpoints, for example, images of the same scene captured from a vehicle-mounted camera at different time points, usually brings a challenge to the change detection task. To address this challenge, we propose a new siamese network structure with the introduction of correlation layer. In addition, we create a publicly available dataset for semantic change detection to evaluate the proposed method. The experimental results verified both the robustness to viewpoint difference in change detection task and the effectiveness for semantic change detection of the proposed networks. Our code and dataset are available at https://github.com/xdspacelab/sscdnet.

I. INTRODUCTION

Semantically understanding scene changes, such as semantic scene change detection, is one of the new problems that have attracted attention in the fields of computer vision, remote sensing, and natural language processing [1]–[5]. Change detection methods have been comprehensively studied and applied to many kinds of tasks, such as detecting anomaly using surveillance and satellite cameras, inspecting infrastructure [6], managing disaster [7], [8], and automating agriculture [9]. However, the existing methods of change detection specify a few detection targets, such as pedestrians and vehicles, for each application. In cases where images contain various kinds of scene changes, more semantic information except for these targets is required for better discrimination in other advanced applications, such as updating city model for autonomous driving [10].

Semantic scene change detection is a challenging task to detect and label scene changes on *each input image* (Fig. 1). There are several types of scene changes in terms of *stuff and thing* classes [11] (Table I). Figure 2 shows examples of each change type. One of the most straightforward methods is comparing the results of (I) pixel-wise semantic segmentation or (II) object detection between input images. Both of the methods (I) and (II) perform well for thing-to-thing changes of *different class*. However, for example, in the case of thing-to-thing changes of *same class and different instance* such as Fig. 2 (a), the straightforward method fails to detect the instance changes. The same thing applies to other change types (Fig. 2 (b), (c)).

The authors are with Artificial Intelligence Research Center, National Institute of Advanced Industrial Science and Technology, Tokyo, Japan, (e-mail: {k.sakurada,mikiya-shibuya,weimin.wang}@aist.go.jp)

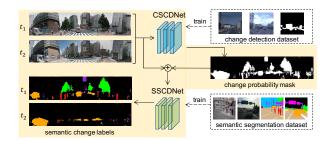


Fig. 1. Overview of the proposed method. First, the CSCDNet takes an image pair as input, which is trained using a change detection dataset, and outputs one change probability mask. Thereafter, the input image pair and the estimated change mask are fed into the SSCDNet, which is trained using a dataset synthesized from a semantic image segmentation dataset. Finally, the SSCDNet estimates the pixel-wise semantic labels of each input image.

Moreover, viewpoint changes of vehicular imagery are larger than those of images taken by surveillance and satellite cameras, which makes it complicated to detect changes between images with large variances of scene depth due to the problems of image correspondence, appearance change and occlusion. Needless to say, although large-scale training datasets make it possible to estimate semantic changes with an end-to-end learning approach directly, it is labor-intensive to create large-scale semantic change detection datasets for each class definition of applications in terms of collecting and labeling images.

In order to overcome these difficulties, we propose a novel semantic change detection scheme with only weak supervision by dividing this task into change detection and semantic extraction (Fig.1). The proposed method is composed of the two convolutional neural networks (CNNs), a correlated siamese change detection network (CSCDNet), and a silhouette-based semantic change detection network (SSCDNet). First, the CSCDNet takes an image pair as input and outputs one change probability mask. Thereafter, the input image pair and the estimated change mask are fed into the SSCDNet. Finally, it estimates the pixel-wise semantic labels of each input image.

The SSCDNet can be trained with the dataset synthesized from commonly available semantic image segmentation datasets, such as the Mapillary Vistas dataset [12], to avoid creating a new dataset for semantic change detection. The estimation accuracy of the SSCDNet depends on that of change detection. However, in the case of images captured from a vehicle-mounted camera at different time points, existing change detection methods suffer from estimation errors due to differences in camera viewpoints. Hence, we propose a new siamese network architecture with the introduction of correlation layers, named as the CSCDNet, which is trained

TABLE I

Change types and the applicability of each semantic change detection method. The straightforward methods (I) and (II) are comparing estimated labels of (I) pixel-wise semantic segmentations and (II) object detections, respectively.

| change type | | | method | | | |
|----------------------------------|---------------------------------|----------|----------|----------|--|--|
| change type | | (I) | (II) | ours | | |
| thing-to-thing | different class | / | / | 1 | | |
| uning-to-uning | same class & different instance | X | Х | √ | | |
| stuff-to-stuff | different class | / | X | ✓ | | |
| stuff-to-stuff | same class | Х | X | ✓ | | |
| thing-to-stuff stuff-to-thing | different class | 1 | Х | 1 | | |

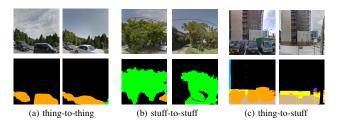


Fig. 2. Examples of change types. (a) same class and different instance (cars to other ones), (b) same class *stuff* (growth and seasonal change of vegetation), *thing-to-stuff* (car to wall and concrete ground).

using a change detection dataset. The CSCDNet can deal with differences in camera viewpoints and achieves state of the art performance on the panoramic change detection (PCD) dataset [13]. Additionally, we incorporate the data augmentation for the input change mask in the training step to improve the robustness of the SSCDNet to change detection errors. For evaluating the proposed methods, we have created the panoramic semantic change detection (PSCD) dataset in the hopes of accelerating researches in the field of dynamic scene modeling.

Our main contributions are as follows:

- We propose a novel semantic change detection network that can be trained with only weak supervision from existing datasets.
- Our siamese change detection network, which uses correlation layers that can deal with differences in camera viewpoints, achieves state of the art performance on the PCD dataset.
- We create the first publicly available dataset for semantic scene change detection.

This paper is organized as follows. In Sec.II, we summarize the related work. Section III explains the details of the proposed network and the training method. Section IV shows the experimental results. Section V presents our conclusions.

II. RELATED WORK

Many methods for temporal scene modeling have been proposed. However, most of them focused on detecting changes or estimating the length of time that each part of a scene exists for. Semantic recognition is required for advanced applications based on dynamic modelings, such as autonomous driving and augmented reality. This section explains the reason for the proposal of the semantic change

detection method using commonly available semantic image segmentation datasets.

Change Detection

Change detection methods are classified into several categories depending on types of target scene changes and available information. Change detection in 2D (image) domain is the most standard approach, especially for surveillance and satellite cameras [14]–[17], which are accurately aligned. A typical approach models the appearance of the scene from a set of images captured at different times, against which a newly captured query image is compared to detect changes [18].

Some studies formulate the problem in a 3D domain. Schindler et al. proposed the probabilistic temporal inferences model based on the visibility of each 3D point reconstructed from images taken from multiple viewpoints at different times [19]. The work by Matzen et al. [20] is classified into the same category. In terms of application, the works by Taneja et al. [21], [22], and Sakurada et al. [7] might be the closest to our research.

In recent years, significant efforts have been made to change detection using machine learning, especially for deep neural networks (DNNs) [6], [10], [13], [23], [24]. There are mainly two types of formulations, "patch similarity estimation" and "pixel-wise segmentation". Patch similarity estimation has been studied for not only change detection but also feature, stereo, and image matchings [25]–[29]. Pixel-wise change detection has been further studied in the context of anomaly detection, background subtraction, and moving object detection [24], [30], [31].

Semantic Change Detection

There are few studies on semantic change detection because most of change detection studies that specify their target domain, such as moving object, forest, and do not explicitly recognize semantic classes of change. The work by Suzuki et al. [5] proposes a method to classify a change mask using multi-scale feature maps extracted using a CNN. It does not consider the problem of detecting changes and estimating correspondences between input images and the change mask (e.g., in Fig.2, there are different change regions between two input images). Daudt et al. [1], [2] detected land surface changes between satellite images. In the case of land surface change detection of satellite images, unlike scene change detection, it is unnecessary to estimate correspondences between input images and the change mask because the change regions between the input images are common. However, for street-level scene change detection, the estimation is necessary because scene objects can appear, disappear, and move.

III. WEAKLY SUPERVISED SILHOUETTE-BASED SEMANTIC SCENE CHANGE DETECTION

There are many types of label definitions for semantic image segmentation depending on the applications; for example, ground-level images of indoor and outdoor scenes [12], [32], aerial and satellite images [33], [34]. Additionally,

TABLE II

Details of the datasets used in the experiments. *(The CSCDNet is trained with only image pairs of a scene and their change masks of the PSCD dataset.)

| Dataset | PCD [13] TSUNAMI GSV | | Vistas [12] | PSCD (This work) | |
|------------------|-------------------------|--------|------------------|---------------------|--|
| Number of images | 100 | 100 | 20,000 | 500 | |
| Original size | 1024 × 224 | | various | 1024 × 224 | |
| Crop size | 224 × 224 | | - | 224 × 224 | |
| Size in training | 256 × 256 | | 256×256 | 256×256 | |
| Paired | √ | | - | √ | |
| Change mask | √ | | - | √ | |
| Semantic label | - | | ✓ | √ | |
| Alignment | medium | coarse | - | coarse | |
| Training target | CSCDNet | | SSCDNet | CSSCDNet | |
| | | | | *(CSCDNet) | |

the definition of *change* (e.g., whether changes of moving objects, display of digital screens, the light of a lamp, transparent barriers, growth of plants, a pool of water, and seasonal changes of vegetation are ignored or not) depends on the application. Thus, there is a large number of combinations of change and semantic definitions. Clearly, it is time-consuming to create semantic change detection datasets for each application. Furthermore, as mentioned above (Sec.II), it is necessary to estimate correspondences between input images and the change mask because the existing change detection datasets do not explicitly contain that information.

To solve these problems, the proposed method includes two CNNs, namely, the CSCDNet and the SSCDNet. This separated architecture enables the method to train the semantic change detection system with change detection datasets and commonly available semantic image segmentation datasets. The rest of this section explains the details regarding the weakly supervised method.

A. Overview

Figure 1 shows an overview of the proposed semantic change detection method. First, the CSCDNet takes an image pair as input, which is trained using a change detection dataset and outputs the change probability of each pixel as one change mask image. Subsequently, the input image pair and the estimated change mask are fed into the SSCDNet, which is trained using a dataset synthesized from a semantic image segmentation dataset. Finally, the SSCDNet estimates the pixel-wise semantic labels of each input image. It should be noted here that the SSCDNet can estimate semantic change labels and correspondences between input images and the change mask simultaneously.

We conjecture that these semantic label estimations and splitting the change mask into the input images can be trained using a commonly available semantic image segmentation dataset, such as the Mapillary Vistas dataset [12], and that semantic information can improve the accuracy of the change mask estimation. Table II shows the details of the datasets used in this paper. The experimental results show the effectiveness of this strategy (Sec.IV). The details of the training dataset synthesis and the network architectures are explained in the following subsections.

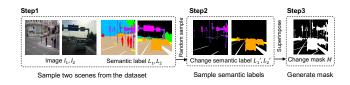


Fig. 3. Synthesis of training dataset for the SSCDNet from semantic image segmentation dataset.

B. Dataset Synthesis from Semantic Segmentation Dataset

Here, we consider the problem of estimating pixel-wise semantic change labels of each input image from an image pair and the change mask. There are several possible methods for generating training datasets to solve this problem. A simulator using a photorealistic rendering, such as Virtual KITTI [35], SYNTHIA [36] and SceneNet RGB-D [37] datasets, is one solution. Although photorealistic images might be effective for pre-training, fine-tuning is necessary to address the domain gaps between synthetic and real images. To bridge the gap, Shrivastava et al. proposed the method to learn a model to improve the realism of a simulator's output using unlabeled real data [38]. However, it is difficult to directly apply this method to natural scene images, which are more complicated than their target domains. Alternatively, synthesis using real images can be applied. Dwibedi et al. proposed the synthetic method to generate large annotated instance datasets in a cut and paste manner [39]. Their study might be the closest to our method.

Figure 3 shows an overview of the proposed training dataset synthesis for the SSCDNet from a semantic image segmentation dataset. First, two RGB images I_1, I_2 , and their semantic label images L_1, L_2 are randomly sampled from the semantic image segmentation dataset. Thereafter, the change semantic label images L_1', L_2' are generated by sampling n semantic labels randomly and removing the others from each semantic label image $(1 \le n \le \min(n_{\max}, N_i - 1))$. N_i represents the number of the classes that the semantic label image L_i contains. The maximum number of class samplings n_{\max} should be decided depending on the number of classes of the semantic segmentation dataset. Finally, the change mask is generated by superimposing the randomly sampled semantic labels as binary silhouettes M.

C. Network Architecture

Correlated Siamese Change Detection Network (CSCD-

Net): We propose the CSCDNet to overcome the limitation of the camera viewpoints of the previous methods. Figure 4 shows an overview of the network architecture of the proposed method. As mentioned in Sec.II, Sakurada et al. [13] found that the comparison between feature maps extracted from input images using a CNN trained with large-scale image recognition datasets [40] is effective for scene change detection task. To incorporate this advantage, we chose the siamese network architecture based on the ResNet-18 [41] which was pretrained on the ImageNet [42] dataset as the encoder of the CSCDNet. Each feature map extracted from two input images in the encoder is concatenated with each

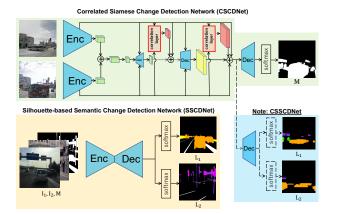


Fig. 4. Network architectures of the CSCDNet, the SSCDNet, and the CSSCDNet. The architecture of the CSSCDNet is based on the CSCDNet and its output layer is replaced with that of the SSCDNet.

decoder's output and fed into the next layer of the decoder whose architecture is based on the network by [43].

Furthermore, for the situation of an image pair with a large viewpoint difference, this difference has to be considered in the design of the network structure to improve the detection accuracy. Exploiting the dense optical flow estimated by the other methods [40] is not efficient in terms of optimization. Therefore, we inserted correlation layers [44], which are utilized for the estimation of optical flow and stereo matching, into the siamese network.

The CSCDNet takes images I_1 and I_2 captured at times t_1 and t_2 as an input. Each pixel value is normalized in [-1,1]. The change mask, as the ground-truth, M_g , is provided to the output of the network as training data. After the final convolution layer, the feature maps are evaluated by the following pixel-wise binary cross-entropy loss:

$$\mathcal{L}_c = -\sum_{\mathbf{x}} t(\mathbf{x}) \ln(p_c(\mathbf{x})) + (1 - t(\mathbf{x})) \ln(1 - p_c(\mathbf{x})), (1)$$

where \mathbf{x} , $t(\mathbf{x})$ and $p_c(\mathbf{x})$ represent the pixel coordinates of the output change mask, the ground-truth, and predictions computed using each output feature maps by a pixel-wise softmax, respectively.

Silhouette-based Semantic Change Detection Network (SSCDNet): The architecture of the SSCDNet is based on the combination of U-Net based on ResNet-18 [41], [43]. Their main differences are the input and output parts. The SSCDNet takes images I_1 , I_2 and M, which are concatenated in the channel dimension as a seven-channel image, for the input. Moreover, after the final convolution layer, the output feature maps are split in half (the bottom of Fig.4), and each of the feature maps is evaluated by the following pixel-wise cross-entropy loss:

$$\mathcal{L}_s = -\sum_{\mathbf{x}} \sum_{k} t_1(\mathbf{x}, k) \ln(p_1(\mathbf{x}, k)) + t_2(\mathbf{x}, k) \ln(p_2(\mathbf{x}, k)),$$
(2)

where k is an index of classes $(1 \le k \le K, K)$: the number of classes), $t(\mathbf{x}, k)$ represents the ground-truth with 1-of-K coding scheme, $p(\mathbf{x}, k)$ represents predictions computed

TABLE III

 F_1 score and mIoU of change detection for TSUNAMI and GSV datasets. Siamese-CDResNet represents the CSCDNet without correlation layers. The CSCDNet consistently outperforms the other methods.

| | | F ₁ score (mIoU) | |
|-----------------------------|---------------|-----------------------------|---------------|
| | TSUNAMI | GSV | Average |
| DenseSIFT [13] | 0.649 (-) | 0.528 (-) | 0.589 (-) |
| CNN-feat [13] | 0.723 (-) | 0.639 (-) | 0.681 (-) |
| DeconvNet [10] | 0.774 (-) | 0.614 (-) | 0.694 (-) |
| WS-Net [24] | - (-) | - (-) | - (0.477) |
| FS-Net [24] | - (-) | - (-) | - (0.588) |
| CDNet [40] | 0.848 (0.811) | 0.695 (0.672) | 0.772 (0.741) |
| CosimNet- 3layer-12 [45] | 0.806 (-) | 0.692 (-) | 0.749 (-) |
| Siamese- CDResNet (Ours) | 0.850 (0.815) | 0.718 (0.691) | 0.784 (0.753) |
| CSCDNet (Ours) | 0.859 (0.824) | 0.738 (0.706) | 0.799 (0.765) |

from each output feature maps by a pixel-wise softmax.

Correlated Siamese Semantic Change Detection Network (CSSCDNet): For a comparative study, we proposed the CSSCDNet as a naive method in the case that the semantic change detection dataset is available. The architecture is based on the CSCDNet. After the final convolution layer, the output feature maps are split in half, and each of the feature maps is evaluated by the pixel-wise cross-entropy loss in the same manner as the SSCDNet (in the dash line box of Fig.4).

IV. EXPERIMENTS

To evaluate the effectiveness of our approach, we performed three experiments. The first experiment is the accuracy evaluation of the change detection with the CSCDNet on the PCD dataset [13]. The proposed siamese change detection networks with and without correlation layers and other existing methods are compared. The second experiment is an accuracy evaluation of the semantic change detection with the SSCDNet using datasets synthesized from the Mapillary Vistas dataset [12]. The data augmentation of the change mask is also evaluated, which improves the robustness of the SSCDNet against change detection errors of the CSCDNet. In the final experiment, we applied our semantic change detection method to the PSCD dataset, which is different from the training dataset of the SSCDNet, and show the effectiveness of our approach.

A. Panoramic Semantic Change Detection (PSCD) dataset

For the quantitative evaluation of the proposed approach, we have created a new dataset named the PSCD dataset, which opens up new vistas for semantic change detection. The PSCD dataset comprises 500 panoramic image pairs. Each pair consists of images I_1, I_2 taken at two different time points t_1 , and t_2 . These panoramic images, which are taken in urban and tsunami-damaged areas, are downloaded from Google Street View.

The PSCD dataset contains the change binary masks C_1, C_2 , the semantic labels S_1, S_2 , the instance labels D_1, D_2 , the attributes A_1, A_2 (3D object, 2D texture, (digital) display) and the visibilities V_1, V_2 (glass, mirror, and wire

fence). The image annotation was performed by a team of 37 well-trained image annotators, and the average annotation time was approximately 156 minutes per image pair. We defined the 67 semantic classes based on those of the Mapillary Vistas dataset [12], and integrated the original classes into the N=11 classes based on the map updating applications as shown in Fig.7. The annotation data and the metadata for downloading the Google Street View images is made publicly available at https://github.com/xdspacelab/sscdnet.

B. Experimental Settings

1) Training dataset generation: We generated training datasets for the CSCDNet, the SSCDNet, and the CSSCDNet from the PCD, the Mapillary Vistas, and the PSCD datasets, respectively. Table II shows the details of the dataset. The PCD dataset is composed of panoramic image pairs I_1 , I_2 taken at two different time points t_1 , and t_2 , and the change mask M_g . From the image set $[I_1, I_2, M_g]$, patch images are cropped by sliding and resized. Furthermore, data augmentation is performed by rotating the patches. Thus, 12,000 sets of image patches were generated. The PSCD dataset is resized and cropped, and data augmentation is performed in the same way as the PCD dataset.

We also generated training datasets for the SSCDNet from the Mapillary Vistas dataset [12]. The Mapillary Vistas dataset for research use contains 20,000 scene images and the pixel-wise semantic labels with 66 semantic classes (including an unlabeled class). We integrate them into the following 11 classes: animal, vehicle, barrier, area, structure, lane marking, vegetation, traffic, others, debris, and no change. Furthermore, the PSCD dataset and the subset TSUNAMI of the PCD dataset used in the final experiment contain much debris. Hence, we added 150 debris images (100 and 50 images for training and validation, respectively) into the dataset used for the dataset synthesis mentioned in Sec.III-B. We selected the value of $n_{\rm max}$ as N-1=10 based on the ablation study. Figure 3 shows an example of the dataset synthesized by the proposed method.

- 2) Data augmentation for robustness to change detection error: If the change masks that are synthesized from semantic segmentation datasets are directly used in the training of the SSCDNet, the trained SSCDNet can be vulnerable to errors in change detection. To improve the robustness of the SSCDNet to change detection error, we perform the data augmentation for change mask in training. Specifically, the change mask is randomly applied to one of the four morphological transformations (erosion, dilation, opening, closing) with a random kernel size k ($1 \le k \le 20$). We expect that the semantic label information can reduce the error of semantic change detection due to the error of change detection by simulating the change mask.
- 3) Training details: The CSCDNet, the SSCDNet, and the CSSCDNet are trained using eight Nvidia Tesla P100 GPUs using the PyTorch framework. We used the batch size of 32. The numbers of iteration for the CSCDNet, the SSCDNet, and the CSSCDNet are 3×10^4 , 1×10^5 , and 1×10^5 , respectively. The Adam algorithm, with a learning rate of 2×10^{-4} , is used. The evaluations of the estimation

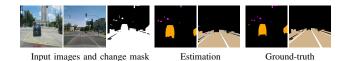


Fig. 5. Example of results estimated by the SSCDNet.

TABLE IV

mIoU of SSCDNet for synthetic data from Mapillary Vistas dataset.

| DA for test | | - | | \checkmark |
|------------------------|----------|---------------|-------|--------------|
| DA for training | - | ✓ | - | ✓ |
| mIoU | 0.570 | 0.544 | 0.428 | 0.494 |
| | | | | |
| | • . | | | |
| I_1 I_2 | Traine | ed without DA | L_1 | L_2 |
| Change mask with noise | M' Train | ned with DA | Cha | nge mask M |

Input images Predictions by SSCDNet Ground-truth

Fig. 6. Example of results estimated by the SSCDNet trained with data augmentation of change mask. The left images I_1 , I_2 and M' show inputs of the SSCDNet. The top and bottom images in the middle column show the prediction results by the SSCDNet trained without and with the data augmentation of change mask. The right images L_1 , L_2 and M show the ground-truth of the semantic labels and the change mask.

accuracies of the CSCDNet using the PCD dataset and the CSSCDNet using the PSCD dataset are performed using the five-fold cross-validation.

C. Evaluation

- 1) Change detection for the PCD dataset: Table III shows F_1 scores and mean intersection-over-union (mIoU) of each method for TSUNAMI and GSV datasets. The CSCDNet outperforms the other methods in terms of both F_1 scores and mIoU. Furthermore, the improvements in the scores for GSV are more significant than those of TSUNAMI. The main reason is that GSV contains more precise changes and the camera viewpoint differences are relatively larger than TSUNAMI because of the differences in their scene depths. The CSCDNet can accurately detect the precise scene changes dealing with the differences in camera viewpoints.
- 2) Accuracy of the SSCDNet for synthetic data: Figure 5 shows an example of the results estimated using the SSCDNet. The SSCDNet can accurately estimate semantic changes on each input image even if there are overlapping areas of change between input images. Table IV shows the mIoU of the SSCDNet for the synthetic validation data from the Mapillary Vistas dataset. There are four combinations of training and test datasets with or without the data augmentation of the aforementioned change mask. In the case of test data without data augmentation, namely, the input change mask is quite accurate, the SSCDNet trained using the dataset without the data augmentation performs better than one trained with the augmentation. However, in the

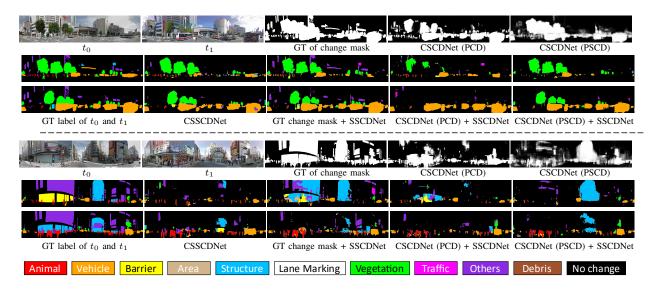


Fig. 7. Examples of semantic scene change detection for the PSCD dataset. One failure case is shown in the lower part.

TABLE V

IoU of the semantic change detection for the PSCD dataset.

| | CSCDNet + SSCDNet | | | GT mask - | + SSCDNet | CSSCDNet | |
|--------------------------|-------------------|--------|----------------------|-----------|------------|----------|-------------|
| Training data (CD / SCD) | PCD / | Vistas | PSCD (mask) / Vistas | | - / Vistas | | PSCD (full) |
| DA for training | - | ✓ | - | √ | - | √ | n/a |
| mIoU | 0.187 | 0.181 | 0.220 | 0.208 | 0.348 | 0.329 | 0.288 |

case of test data with the augmentation, namely, the input change mask has some errors, the SSCDNet trained using the augmentation outperforms the other. Figure 6 shows an example of results estimated by the SSCDNet trained with the data augmentation of the change mask. The estimation results obtained using the SSCDNet trained without the data augmentation of the change mask have errors due to the errors from the input change mask. However, the SSCDNet trained with the augmentation can accurately predict the semantic change labels while being more robust to the effects of errors of the input change mask.

3) Semantic change detection for the PSCD dataset: Figure 7 shows examples of the semantic change detection results for the entire process of our proposed method. Table V shows the mIoU of each method for the PSCD dataset. In the top rows of Fig.7, the CSCDNet can accurately detect scene changes, although some detection errors are owing to reflections from window-glasses and advertisement boards on the buildings because of the lack of training data. (The CSCDNet trained with the PCD dataset can detect some changes of small advertisement boards but not those of vegetations, and vice versa for the CSCDNet trained using the PSCD dataset.) Certainly, if the semantic change detection dataset, of which the creation is labor-intensive, is available, the strategy of the end-to-end learning for semantic change detection can be applied, and the performance is almost the best (CSSCDNet). However, even if the dedicated dataset is unavailable, the SSCDNet can estimate semantic scene changes for each input image successfully depending on the change detection accuracy. The lower part of Fig.7 shows failure cases due to a lack of training data.

Better performance was achieved when the CSCDNet was trained using the PSCD dataset rather than being trained on the PCD dataset, which indicates that only the change detection dataset of the same domain as the target data should be used if it is available. Furthermore, the SSCDNet using the ground-truth change mask performs better than the CSSCDNet, which is trained using the full set of the semantic change detection dataset. Hence, the SSCDNet will exhibit higher performance when accurate change mask information is available by other methods [7], [21], [22] and sensors.

V. CONCLUSIONS

We proposed a novel semantic change detection scheme with only weak supervision. The proposed method is composed of the two CNNs, the CSCDNet and the SSCDNet. The CSCDNet can deal with the difference of camera viewpoints and achieves state of the art change detection performance for the PCD dataset. The SSCDNet can be trained with dataset synthesized from semantic image segmentation datasets to avoid creating a new dataset for semantic change detection. To evaluate the effectiveness of the proposed method, we created the first publicly available dataset for semantic scene change detection, named as the PSCD dataset. Experimental results with this dataset verified the effectiveness of the proposed scheme in the semantic change detection task.

ACKNOWLEDGMENT

This work is partially supported by KAKENHI 18K18071 and the New Energy and Industrial Technology Development Organization (NEDO).

REFERENCES

- R. C. Daudt, B. Le Saux, and A. Boulch, "Fully convolutional siamese networks for change detection," in 2018 25th IEEE International Conference on Image Processing (ICIP). IEEE, 2018, pp. 4063– 4067
- [2] "Multitask learning for large-scale semantic change detection," Computer Vision and Image Understanding, vol. 187, p. 102783, 2019.
- [3] H. Jhamtani and T. Berg-Kirkpatrick, "Learning to Describe Differences Between Pairs of Similar Images," in EMNLP, 2018.
- [4] D. H. Park, T. Darrell, and A. Rohrbach, "Robust Change Captioning," in *ICCV*, 2019, pp. 4624–4633.
- [5] T. Suzuki et al., "Semantic Change Detection," in ICARCV, 2018.
- [6] S. Stent, R. Gherardi, B. Stenger, and R. Cipolla, "Detecting Change for Multi-View, Long-Term Surface Inspection," in BMVC, 2015.
- [7] K. Sakurada, T. Okatani, and K. Deguchi, "Detecting Changes in 3D Structure of a Scene from Multi-view Images Captured by a Vehicle-Mounted Camera," in CVPR, 2013.
- [8] K. Sakurada, T. Okatani, and K. M. Kitani, "Massive City-scale Surface Condition Analysis using Ground and Aerial Imagery," in ACCV, 2014.
- [9] J. Dong, J. G. Burnham, B. Boots, G. Rains, and F. Dellaert, "4D Crop Monitoring: Spatio-Temporal Reconstruction for Agriculture," in *ICRA*, 2017.
- [10] P. F. Alcantarilla, S. Stent, G. Ros, R. Arroyo, and R. Gherardi, "Street-View Change Detection with Deconvolutional Networks," in *Robotics: Science and Systems*, 2016.
- [11] A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollár, "Panoptic Segmentation," in CVPR, 2019, pp. 9404–9413.
- [12] G. Neuhold, T. Ollmann, S. Rota Bulo, and P. Kontschieder, "The Mapillary Vistas Dataset for Semantic Understanding of Street Scenes," in *ICCV*, 2017.
- [13] K. Sakurada and T. Okatani, "Change Detection from a Street Image Pair using CNN Features and Superpixel Segmentation," in BMVC, 2015.
- [14] D. Crispell, J. Mundy, and G. Taubin, "A Variable-Resolution Probabilistic Three-Dimensional Model for Change Detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 50, no. 2, 2012.
- [15] A. Huertas and R. Nevatia, "Detecting Changes in Aerial Views of Man-Made Structures," in *ICCV*, 1998.
- [16] T. Pollard and J. L. Mundy, "Change Detection in a 3-d World," in CVPR, 2007.
- [17] R. J. Radke, S. Andra, O. Al-Kofahi, and B. Roysam, "Image Change Detection Algorithms: A Systematic Survey," TIP, vol. 14, no. 3, 2005.
- [18] K. Wang, C. Gou, and F.-Y. Wang, "M4CD: A robust change detection method for intelligent visual surveillance," *IEEE Access*, vol. 6, pp. 15505–15520, 2018.
- [19] G. Schindler and F. Dellaert, "Probabilistic temporal inference on reconstructed 3D scenes," in CVPR, 2010.
- [20] K. Matzen and N. Snavely, "Scene Chronology," in ECCV, 2014.
- [21] A. Taneja, L. Ballan, and M. Pollefeys, "Image based detection of geometric changes in urban environments," in ICCV, 2011.
- [22] —, "City-Scale Change Detection in Cadastral 3D Models Using Images," in CVPR, 2013.
- [23] A. Fujita, K. Sakurada, T. Imaizumi, R. Ito, S. Hikosaka, and R. Nakamura, "Damage Detection from Aerial Umages via Convolutional Neural Networks," in MVA, 2017.
- [24] S. H. Khan, X. He, F. Porikli, M. Bennamoun, F. Sohel, and R. Togneri, "Learning Deep Structured Network for Weakly Supervised Change Detection," in *IJCAI*, 2017.
- [25] Z. Chen, X. Sun, L. Wang, Y. Yu, and C. Huang, "A deep visual correspondence embedding model for stereo matching costs," in *ICCV*, 2015.
- [26] T.-Y. Lin, Y. Cui, S. Belongie, and J. Hays, "Learning deep representations for ground-to-aerial geolocalization," in CVPR, 2015.
- [27] E. Simo-Serra, E. Trulls, L. Ferraz, I. Kokkinos, P. Fua, and F. Moreno-Noguer, "Discriminative learning of deep convolutional feature point descriptors," in *ICCV*, 2015.
- [28] S. Zagoruyko and N. Komodakis, "Learning to compare image patches via convolutional neural networks," in CVPR, 2015.
- [29] J. Zbontar and Y. LeCun, "Stereo matching by training a convolutional neural network to compare image patches," *JMLR*, vol. 17, no. 1, 2016.
- [30] M. Camplani, L. Maddalena, G. M. Alcover, A. Petrosino, and L. Salgado, "A benchmarking framework for background subtraction in rgbd videos," in *International Conference on Image Analysis and Processing*, 2017.

- [31] K. Toyama, J. Krumm, B. Brumitt, and B. Meyers, "Wallflower: Principles and practice of background maintenance," in ICCV, 1999.
- [32] I. Armeni, S. Sax, A. R. Zamir, and S. Savarese, "Joint 2D-3D-Semantic Data for Indoor Scene Understanding," arXiv preprint arXiv:1702.01105, 2017.
- [33] ISPRS, "2D Semantic Labeling Contest," http://www2.isprs.org/ commissions/comm3/wg4/semantic-labeling.html.
- [34] V. Mnih, "Machine Learning for Aerial Image Labeling," Ph.D. dissertation, University of Toronto, 2015.
- [35] A. Gaidon, Q. Wang, Y. Cabon, and E. Vig, "Virtual Worlds as Proxy for Multi-Object Tracking Analysis," in CVPR, 2016.
- [36] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and A. M. Lopez, "The SYNTHIA Dataset: A Large Collection of Synthetic Images for Semantic Segmentation of Urban Scenes," in CVPR, 2016.
- [37] J. McCormac, A. Handa, S. Leutenegger, and A. J.Davison, "SceneNet RGB-D: Can 5M Synthetic Images Beat Generic ImageNet Pretraining on Indoor Segmentation?" 2017.
- [38] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb, "Learning from simulated and unsupervised images through adversarial training," in CVPR, 2017.
- [39] D. Dwibedi, I. Misra, and M. Hebert, "Cut, Paste and Learn: Surprisingly Easy Synthesis for Instance Detection," in *ICCV*, 2017.
- [40] K. Sakurada, W. Wang, N. Kawaguchi, and R. Nakamura, "Dense Optical Flow based Change Detection Network Robust to Difference of Camera Viewpoints," arXiv preprint arXiv:1712.02941, 2017.
- [41] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in CVPR, 2016, pp. 770–778.
- [42] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in CVPR, 2009.
- [43] R. Hamaguchi and S. Hikosaka, "Building detection from satellite imagery using ensemble of size-specific detectors," in CVPR, June 2018
- [44] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. van der Smagt, D. Cremers, and T. Brox, "Flownet: Learning optical flow with convolutional networks," in *ICCV*, December 2015.
- [45] E. Guo, X. Fu, J. Zhu, M. Deng, Y. Liu, Q. Zhu, and H. Li, "Learning to measure change: Fully convolutional siamese metric networks for scene change detection," arXiv preprint arXiv:1810.09111, 2018.