

Embedding Domain-Invariant Building Segmentation Information On Change Detection Model

Minh-Khoa Le^{1,2}, Dinh Phong Vo, and Bac Le^{1,2}

¹ Faculty of Information Technology, University of Science, Ho Chi Minh City, Vietnam

² Vietnam National University, Ho Chi Minh City, Vietnam
18120415@student.hcmus.edu.vn; dpvo@tuta.io; lhbac@fit.hcmus.edu.vn

Abstract. The use of remote sensing data in combination with deep learning methods has resulted in great successes on tasks such as building segmentation and change detection. While some studies have demonstrated that supervised deep learning approaches can outperform traditional methods which use algebra or transformation techniques with hand-crafted features, these methods are often limited to small datasets and struggle with distribution shifts. Advances in the field of Domain Adaptation have led to the development of architectures with improved ability to learn domain-invariant features. In this paper, we propose a domain adaptation segmentation model for building segmentation and apply it to the problem of change detection. In addition, empirical evaluations show that our proposed methods can outperform some competitive baselines on popular benchmark datasets.

Keywords: deep learning, remote sensing, change detection, building segmentation.

1 Introduction

Satellite imagery is a rich source of information to solve many humanitarian issues [21], [24]. As time passes, more satellite data can be accessed easier and more freely [10]. As a result, the goal of quickly detecting or calculating insights from satellite images becomes possible, realizing the desire of quickly identifying or assessing damage from calamities.

Building segmentation and change detection are two crucial problems in remote sensing. Building segmentation is a classification task that extracts building from a satellite image [17]. Change detection is the process of identifying objects' differences in the same geographical location at different times [20]. They are widely used in many applications, such as agricultural monitoring [12].

Due to the explosive growth in the quantity and quality of satellites, deep learning models, which results in an abundant source of high-quality satellite images, have shown outstanding performance in remote sensing tasks [3], [6]. However, supervised deep learning power is heavily rely on hand-labelling datasets,

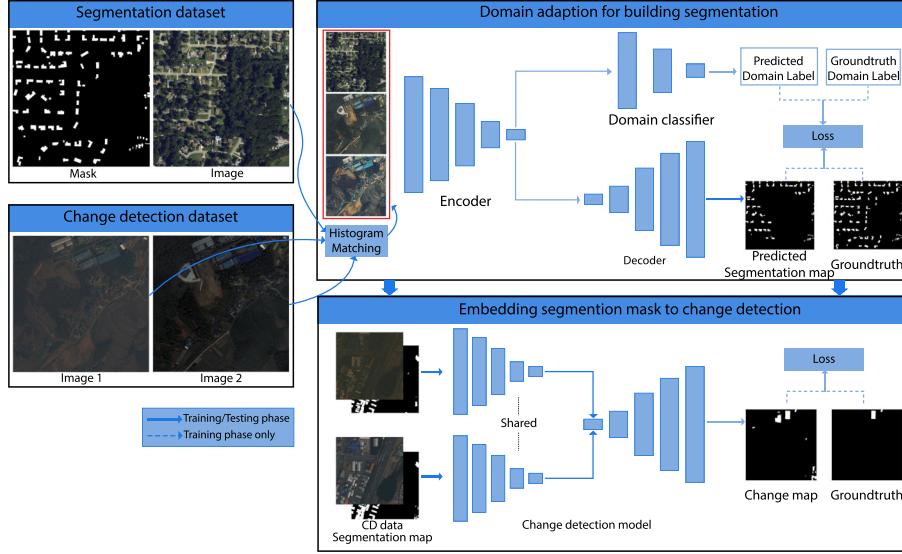


Fig. 1. Overview of our proposed method including two parts: domain adaptation for building segmentation and embedding segmentation mask to change detection. In the first part, after being applied histogram matching algorithm, images from segmentation and change detection datasets are fed into our domain adaptation segmentation model to extract segmentation masks and predict image domains. In training phase, losses are calculated to optimize the model. In the second part, segmentation informations are used further with images from change detection dataset to predict change map. In training phase, we use predicted map and groundtruth to optimize our model.

which is resources consuming [4]. Moreover, distribution shifts between remote sensing datasets cause the model’s poor performance across multiple datasets.

An important characteristic of satellite images is the off-nadir angle, which is the angle between the satellite and the target at the time of collection. Satellites typically capture images at large off-nadir angles to collect more data. However, this can result in poor geometric consistency for tall objects when multiple off-nadir angles are introduced. There are few datasets that consist of multiple-angle images, resulting in a significant amount of satellite data going unused. Shen et al. have introduced a new dataset called S2Looking [19] for change detection tasks, which includes off-nadir angles ranging from -35° to 35° .

Recently, the relationship between segmentation and change detection has been explored. Image segmentation technique provides significant advantages for remote sensing imagery analysis, pointing out many buildings in different scales [14]. Researchers have improved their model’s performance by embedding segmentation masks or network architecture in change detection [23], [27].

Our contributions. To summarize, we make the following contributions.

1. We propose Domain Adaptation Segmentation Model (DASM) for building segmentation. We also apply many pre-processing methods to achieve better evaluation scores on unseen data with distribution shifts.
2. Additionally, we improve a Siamese Unet Change Detection Model by embedding building segmentation information into input data.
3. We conduct many experiments to demonstrate our methods against baselines and state-of-the-arts on popular remote sensing datasets, where our methods show better F1-Score than some previous works.

2 Related Work

2.1 Building segmentation for satellite imagery

Building segmentation is a binary segmentation task that aims to extract buildings from satellite images. Many researchers have designed their networks based on UNet architecture because of its excellent performance on the generic image segmentation task [13]. There are also many efforts to enhance the UNet architecture for better performance. Iglovikov et al. introduce the TernaNetV2 [14] with the improvement from UNet by using WideResNet-38 network with in-place activated batch normalization as the encoder, showing superior performance in the building segmentation task. Nevertheless, those methods suffered drastically from image contrast, brightness, and off-nadir angles [12].

2.2 Change detection

Change detection became possible after the existence of satellite images. In the past, image differencing was the primary technique [20]. However, it easily suffers from changes that are not in favor [28]. For example, seasonal changes can create differences in color at different times, dramatically decreasing the image differencing method performance. Recently, deep learning methods have been developed with improved accuracy since the rise of Convolutional Neural Networks (CNNs) [3], [6]. Some of them are post-classification methods whose structures contain two classifiers to classify the bi-temporal images [2], [26]. Then, the change map can be obtained by comparing classification maps.

Many attempts have been done to train CNNs directly and generate change maps from satellite images. They focus on improving the model's feature extraction and feature discrimination abilities. Daudt et al. introduce three fully convolutional models based on simplified U-Net, and Siamese Net [6]. However, deep-learning-based methods requires a large amount of computing resources, which are not always available. Moreover, performance of those models decrease vastly when heterogeneous data are introduced.

2.3 Domain adaptation

One of the main domain adaptation methods is domain-invariant feature learning. It focuses on aligning source and target domains by producing domain-

invariant feature representations [25]. In particular, a network can produce domain-invariant feature representations if the output feature follows the same distribution, whether input comes from the source or target domain.

Adversarial-based domain adaptation method aims to minimize cross-domain discrepancy via an adversarial objective. A method called Domain Adversarial Training of Neural Networks (DANN) uses domain confusion loss to confuse high-level layers about the input’s domain [9]. The proposed architecture consists of 3 components: feature extractor, domain classifier, and label predictor. By minimizing the label prediction loss while maximizing the domain classification loss, the feature extractor is trained to generate domain-invariant feature representations and perform well on the label prediction task. On the other hand, those methods are partly invariant to input domains that we need to retrain our models when a new data domain comes. Additionally, domain adaptation methods perform not good as traditional supervised manner [9].

3 Proposed Methods

Our idea is using segmentation masks generated from our segmentation model as an additional information channel of the input images, making our change detection model more aware of building entities in the image.

We use a Unet [1] architecture as the segmentation model in the first step. However, we observe that the training and test data come from two different domains, which may make the model perform poorly. Therefore, we introduce the domain adaptation technique in the building segmentation task to address this problem. Moreover, pre-processing methods are applied to improve the model including per-image normalization, histogram matching.

3.1 Domain adaption for building segmentation

Baseline segmentation model A baseline model inspired by TerausNetV2 [14] is chosen to solve domain adaption building segmentation. Concretely, the model has an encoder-decoder architecture, and the backbone is ImageNet pre-trained ResNet-50 [11]. ResNet has been proposed and is well-known for the ease of training networks substantially deeper than those used previously [11]. Moreover, skip connections are added at each convolution block, allowing features from the input image and encoders to flow directly to high-level feature maps, improving localization accuracy and speeding up the training process [18].

Domain adaptation segmentation model On baseline segmentation model, we apply the Domain-Adversarial Neural Networks (DANN) technique [9] to construct a segmentation model called Domain-Adaptation Segmentation Model (DASM). The expectation is to improve the baseline’s accuracy on the test set and minimize the effects of distribution shifts. We aim to confuse the decoder module about the image domains and generate the features invariant to domains. After the encoder module, a domain adaptation component - a simple neural

network with dropout layers - domain classifier is added to classify images from different domains (the domain classifier in Fig. 1).

Denote that $G_e(\cdot; \theta_e)$ is the encoder with parameters θ_e , $G_s(\cdot; \theta_s)$ is the decoder with parameters θ_s , and $G_c(\cdot; \theta_c)$ is domain classifier with parameter θ_c . For an input image x_i with the domain label d_i and the building mask y_i , segmentation loss and domain classification loss is defined as:

$$\begin{aligned} L_s^i(\theta_e, \theta_s) &= L_s(G_s(G_e(x_i; \theta_e); \theta_s), y_i) \\ L_c^i(\theta_e, \theta_c) &= L_c(G_c(G_e(x_i; \theta_e); \theta_c), d_i) \end{aligned}$$

For N sample images with the first n samples with labels from source domain - \mathcal{D}_S and the last $N - n$ samples without labels from target domain - \mathcal{D}_T , total loss of the model are defined as:

$$\begin{aligned} E(\theta_e, \theta_s, \theta_c) &= L_s(\theta_e, \theta_s) - L_c(\theta_e, \theta_c) \\ &= \frac{1}{n} \sum_{i=1}^n L_s^i(\theta_e, \theta_s) - \lambda \left(\frac{1}{n} \sum_{i=1}^n L_c^i(\theta_e, \theta_c) + \frac{1}{N-n} \sum_{i=n+1}^N L_c^i(\theta_e, \theta_c) \right) \end{aligned} \quad (1)$$

with λ is a hyperparameter that controls the importance of the domain loss. We can minimize equation 1 by finding the saddle point $\hat{\theta}_e, \hat{\theta}_s, \hat{\theta}_c$ with respect to

$$(\hat{\theta}_e, \hat{\theta}_s) = \underset{\theta_e, \theta_s}{\operatorname{argmin}} E(\theta_e, \theta_s, \hat{\theta}_c) \quad (2)$$

$$\hat{\theta}_c = \underset{\theta_c}{\operatorname{argmax}} E(\hat{\theta}_e, \hat{\theta}_s, \theta_c) \quad (3)$$

Equation 1 calculates the loss function by subtracting domain loss from segmentation loss. We expect the subtraction makes the loss converge to a saddle point, where domain loss is maximized while minimizing the segmentation loss. Such subtraction can be implemented easily by a gradient reversal layer [9]. This layer has no parameters and reverses the gradient's sign by multiplying a negative constant during the backpropagation phase. By applying the gradient reversal layer, a saddle point can be found using the gradient descent algorithm:

$$\theta_e \leftarrow \theta_e - \mu \left(\frac{\partial L_s^i}{\partial \theta_e} - \lambda \frac{\partial L_c^i}{\partial \theta_e} \right) \quad (4)$$

$$\theta_s \leftarrow \theta_s - \mu \frac{\partial L_s^i}{\partial \theta_s} \quad (5)$$

$$\theta_c \leftarrow \theta_c - \mu \lambda \frac{\partial L_c^i}{\partial \theta_c} \quad (6)$$

with μ is learning rate. Converging to the saddle point, encoder fails to distinguish domains. As a result, the encoder generates domain-invariant features.

Pre-processing methods We call combination of DASM and pre-processing methods as Improved Domain-Adaptation Segmentation Model (I-DASM).

Per-image normalization Image normalization is simple and the first process to solve the domain shift problem. It brings the images with different ranges of pixel intensity into the same range. As a result, images from different domains share the same contrast.

Histogram matching In different lighting conditions, satellite sensors, the observed colors may differ from the actual colors [22]. By adjusting gamma, saturation via histogram matching, we can correct colors to produce an output image that shares an equivalent color distribution with the input image.

Given a pair of images, a grayscale source image S that needed to be applied histogram matching and a grayscale reference image R . Each image has a probability density function $p^s(a)$ and $p^r(a)$ where a is a grayscale value and $p^s(a)$ and $p^r(a)$ is the probability of that value for source and reference image. The probability can be computed from its color histogram:

$$p^s(a_i) = \frac{n_i^s}{n^s}, \quad p^r(a_i) = \frac{n_i^r}{n^r} \quad (7)$$

where n_i is the frequency of the pixel value a_i , and n is the total number of pixels in the image. To apply histogram matching, we need to compute its cumulative distribution function for each image.

$$F^s(a_j) = \sum_{i=0}^j p^s(a_i), \quad F^r(a_j) = \sum_{i=0}^j p^r(a_i), \quad j = 0, \dots, G-1 \quad (8)$$

where G is the total number of pixel values. For each source pixel value $G_i^s \in [0, 255]$, we need to find a reference pixel value G_j^r for which $F^s(G_i^s) = F^r(G_j^r)$, which is the result of a histogram matching function $M(G_i^s) = G_j^r$. Finally, to find the output image, we apply $M()$ for each pixel of the reference image.

We choose a collection of various images from Alabama dataset as reference images. Each source image is compared with the collection, and select the most similar to it. To compare two images, we use Pearson's Correlation Coefficient [8] as the similarity metric. The higher the metric is, the more similar they are.

3.2 Embedding Segmentation to Change Detection

Baseline change detection model We use the baseline Siamese UNet [18] architecture with ResNet-50 [11] as the backbone. The network employs an encoder-decoder design based on FC-Siam-Diff [6]. The encoding layers of the network are separated into two streams of similar structure with shared weights as in a traditional Siamese network [7] to extract features from two input images. Each has 5 convolution layers and 5 max-pooling layers to continuously extract features from input images by increasing the number of channels and decreasing the size of images twice in each layer.

Skip connections are added as an improvement to the architecture. Moreover, each skip connection takes input features from both encoders, gets the absolute value of their subtraction before feeding them to the decoders. It makes the network simpler and helps it learn to compare the differences between images [6].

Improved Change Detection Model (ICDM) We decide to embed the segmentation mask into the input images. It is the most straightforward way to help our change detection model emphasize the differences between images, especially the object differences. The input image is extended by adding one more channel of probability building segmentation mask (Fig. 1). Concretely, the segmentation mask is the output of the I-DASM.

4 Experiments

4.1 Experiments details

We collected Alabama satellite images from Bing Maps, combined with US Building Footprint [16] to create a building segmentation dataset called Alabama dataset. All 10,200 samples from the Alabama dataset are used for the segmentation task as the training set and reference images for the histogram matching algorithm. For the test set, we use every S2Looking image patch pair.

There is no ground-truth segmentation labels in the S2Looking dataset, so we use change detection mask as the label to compare segmentation models. With a pair of input RGB images from the change detection dataset, segmentation models result in 2 masks. We extract buildings from segmentation masks into a set of objects. And then using the Hungarian algorithm [15] to map building masks between two sets. Every mask that is not mapped is classified as a changed building and presents in the change detection mask.

We use Precision, Recall, and F1-Score to evaluate models on both tasks.

$$\text{Precision} = \frac{TP}{TP + FP} \quad \text{Recall} = \frac{TP}{TP + FN} \quad (9)$$

$$\text{F1-Score} = \frac{2 \times TP}{(TP + FP) + (TP + FN)} \quad (10)$$

with TP, FP, FN correspond to the number of true-positive, false-positive, and false-negative predicted pixels for class 1 and 0.

In training phase on both tasks, we use Adam optimizer with a learning rate of 0.01. After every ten steps, the learning rate will be decreased ten times. A combination of cross-entropy loss and Dice loss is the objective function. For change detection task, we randomly select six 512×512 parts for each pair of images and the change detection mask using the random crop regarding the mask to avoid overfitting because of the class-imbalance. A augmentation pipeline, including random shift, scale and rotation, random RGB shift, brightness and contrast, is embedded in pre-processing phase to improve the robustness of the model. Every method was trained for 20 epochs with this setting.

4.2 Experiments on segmentation models

We conduct many experiments to compare methods for the segmentation task. Training set is the Alabama dataset, and test sets are both Alabama and S2Looking [19] datasets. The training setting was described in the section 4.1.

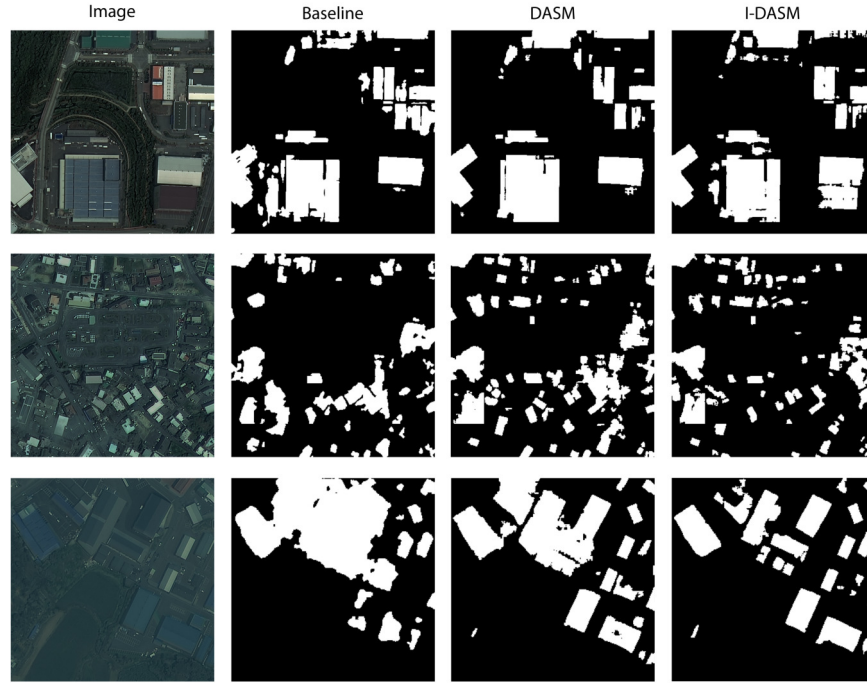


Fig. 2. Images from S2Looking dataset and segmentation results of models. Each row contains a sample from the dataset and outputs from models.

In table 1, we show the F1-Score of our models on two different tasks: segmentation in the Alabama dataset and change detection in the S2Looking [19] dataset. Although the baseline shows excellent performance on the segmentation task - at 0.8421, it is not good at the change detection in S2Looking [19] with a deficient F1-Score. Our DASM improves the change detection F1-Score while preserving a high segmentation F1-Score. I-DASM has performed better on both two tasks than DASM with pre-processing methods. It significantly improves the change detection F1-Score with a small segmentation performance trade-off.

Table 1. Result for the evaluated segmentation methods

Method	Alabama F1-Score	S2Looking F1-Score
Baseline (Section 3.1)	0.8412	0.1520
DASM (Section 3.1)	0.8020	0.1950
I-DASM (Section 3.1)	0.8336	0.2084

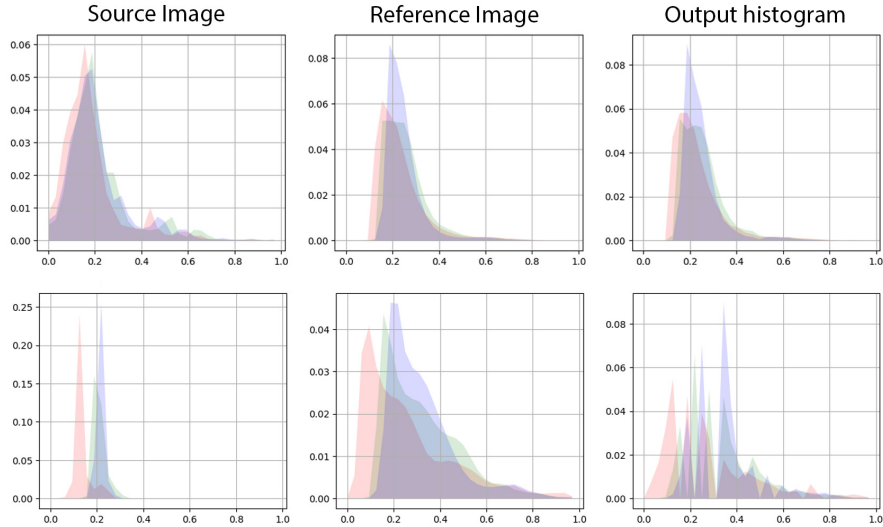


Fig. 3. Histograms of source, reference, and output images. Each row presents a sample from dataset. Each graph has red, green, and blue histograms corresponding to the red, green, and blue channels.

Fig. 2 shows the baseline model easily suffers from the colors. When some buildings and backgrounds overlap, it fails to distinguish and classifies background pixels as building pixels. DASM no more suffers from the landscape colors and recognizes the border between buildings and backgrounds well without the need to be provided any segmentation label. The improvement proves that our encoder has produced features partly invariant to input domains.

With low-contrast source images from S2Looking and reference images from the Abalama, histogram matching results in images with the same color distribution as reference images, better brightness and contrast. First sample in Fig. 3 shows that the output histogram resembles the reference one. In second row of Fig. 3, the algorithm results in significant deviations between some neighboring bins in the histogram when the difference between source and reference images is vast. Additionally, histogram matching does not give attention to the image’s context and even tries to match the building colors with landscape colors.

4.3 Experiments on change detection models

In table 2, we observed that our baseline achieved F1-Score, at 0.5621, higher than some previous methods, such as FC-Siam-Diff [6]. By providing the baseline building information, our method is improved by over 1%, getting a higher F1-Score at 0.5737. Our methods archive higher F1-Scores than other methods, such as FC-Siam-Diff [6], and STANet-Base [5]. For S2Looking [19] dataset, CDNet [3] performs better than our model, at an F1-Score of 0.6056, around 3% higher than our object-level change detection model.

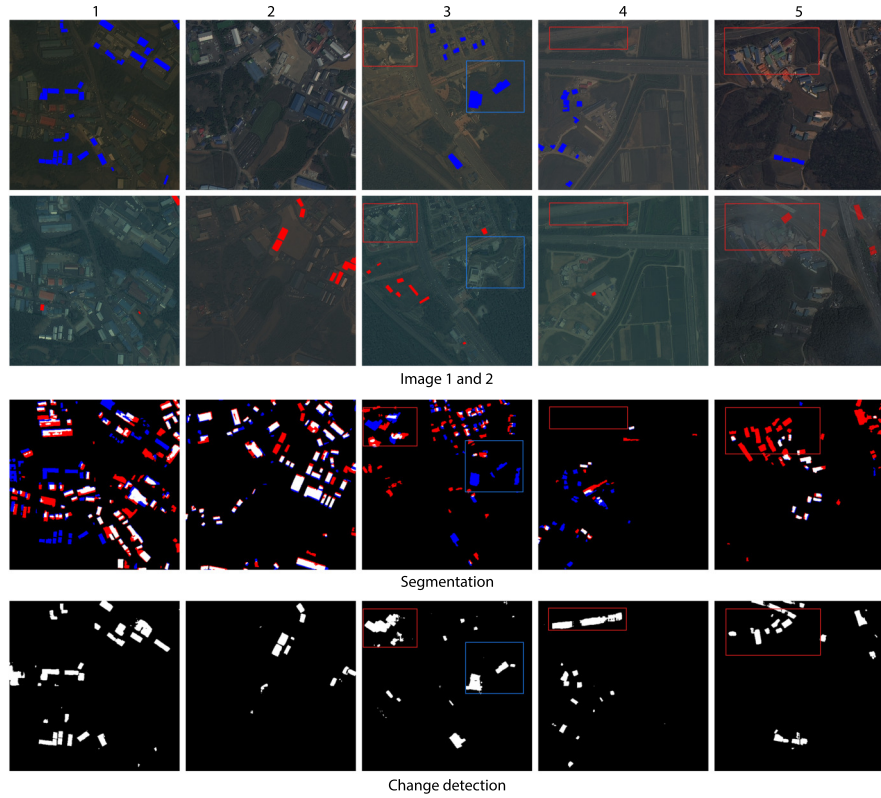


Fig. 4. Result of S2Looking dataset on segmentation and change detection tasks. First two rows present pairs of before and after images. Each image contains red and blue objects representing buildings that only exist in the other image at the same place. The segmentation row shows our I-DASM model results on both before and after images: White pixels are buildings appear in both images, blue and red pixels only appear in images 1 and 2, respectively. The last row is the outputs of our change detection model.

Fig. 4 shows our ICDM predictions on the S2Looking dataset. We observe that our model performs well and points out many change pixels. It is robust to the difference between two images: color distribution and brightness. Moreover, our model can discriminate between the change in construction and the change in non-construction. Although there are changes in buildings and background in the blue rectangle in sample 3, the CD model can distinguish them pretty well.

Our method still suffered from the off-nadir angle. Although there is no change in buildings, the off-nadir angles are various in the two images. Our model also fails to recognize buildings when clouds cover buildings in an image. In sample 5 in Fig. 4, both segmentation and change detection models perform well on clear places but fail completely in the opaque area.

Table 2. Result of our change detection methods and other methods on S2Looking.

Method	Precision	Recall	F1-Score
FC-EF [6]	0.8136	0.0895	0.1613
FC-Siam-Diff [6]	0.8329	0.1576	0.2650
STANet-Base [5]	0.2575	0.5629	0.3534
CDNet [3]	0.6748	0.5493	0.6056
I-DASM (Section 3.1)	0.1393	0.4133	0.2084
Baseline (Section 3.2)	0.5731	0.5514	0.5621
ICDM (Section 3.2)	0.5454	0.6051	0.5737

5 Conclusion

In this paper, we address the issue of distribution shifts in building segmentation tasks. We propose a framework for reducing the gap between different domains and improving the performance of the model on the test set without sacrificing its performance on the training set. Our method also offers a potential solution for improving change detection models using available data without the need for extensive manual labeling.

As a potential future improvement, we could combine histogram matching and DANN into a single model. We could also explore better ways to use the segmentation masks to improve the change detection model.

Acknowledgment

We would like to express our heartfelt gratitude to our mentor and professor for their invaluable support and guidance throughout our research journey. Their expertise and wisdom have been instrumental in shaping our research and providing us with valuable insights and feedback.

References

1. Alsabhan, W., Alotaiby, T.: Automatic building extraction on satellite images using unet and resnet50. *Computational Intelligence and Neuroscience* **2022**, 1–12
2. Cao, C., Dragičević, S., Li, S.: Land-use change detection with convolutional neural network methods. *Environments* **6**(2) (2019)
3. Chen, H., Li, W., Shi, Z.: Adversarial instance augmentation for building change detection in remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing* **60**, 1–16 (2022)
4. Chen, H., Qi, Z., Shi, Z.: Remote sensing image change detection with transformers. *IEEE Transactions on Geoscience and Remote Sensing* **60**, 1–14 (2022)
5. Chen, H., Shi, Z.: A spatial-temporal attention-based method and a new dataset for remote sensing image change detection. *Remote Sensing* **12**(10) (2020)
6. Daudt, R.C., Saux, B.L., Boulch, A.: Fully convolutional siamese networks for change detection (2018)

7. Dey, S., Dutta, A., Toledo, J.I., Ghosh, S.K., Lladós, J., Pal, U.: Signet: Convolutional siamese network for writer independent offline signature verification (2017)
8. Freedman, D., Pisani, R., Purves, R.: Statistics (international student edition). Pisani, R. Purves, 4th edn. WW Norton & Company, New York (2007)
9. Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., Lempitsky, V.: Domain-adversarial training of neural networks
10. Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., Moore, R.: Google earth engine: Planetary-scale geospatial analysis for everyone. *Remote Sensing of Environment* **202**, 18–27 (2017), big Remotely Sensed Data: tools, applications and experiences
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition
12. Hoeser, T., Kuenzer, C.: Object detection and image segmentation with deep learning on earth observation data: A review-part i: Evolution and recent trends. *Remote Sensing* **12** (05 2020)
13. Iglovikov, V., Shvets, A.: Terausnet: U-net with VGG11 encoder pre-trained on imagenet for image segmentation. *CoRR* (2018)
14. Iglovikov, V.I., Seferbekov, S.S., Buslaev, A.V., Shvets, A.: Terausnetv2: Fully convolutional network for instance segmentation. *CoRR* (2018)
15. Kuhn, H.: The hungarian method for the assignment problem. *Naval Research Logistic Quarterly* **2** (05 2012)
16. Microsofts: United states microsoft building footprints (2018)
17. Neupane, B., Horanont, T., Aryal, J.: Deep learning-based semantic segmentation of urban features in satellite images: A review and meta-analysis. *Remote Sensing* **13**(4) (2021)
18. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: *MICCAI* (2015)
19. Shen, L., Lu, Y., Chen, H., Wei, H., Xie, D., Yue, J., Chen, R., Lv, S., Jiang, B.: S2looking: A satellite side-looking dataset for building change detection. *Remote Sensing* **13**(24), 5094 (dec 2021)
20. SINGH, A.: Review article digital change detection techniques using remotely-sensed data. *International Journal of Remote Sensing* **10**(6), 989–1003 (1989)
21. Szpakowski, D.M., Jensen, J.L.R.: A review of the applications of remote sensing in fire ecology. *Remote Sensing* (22) (2019)
22. Tasar, O., Happy, S.L., Tarabalka, Y., Alliez, P.: ColorMapGAN: Unsupervised domain adaptation for semantic segmentation using color mapping generative adversarial networks. *IEEE Transactions on Geoscience and Remote Sensing* (10)
23. Wang, Y., Gao, L., Hong, D., Sha, J., Liu, L., Zhang, B., Rong, X., Zhang, Y.: Mask deeplab: End-to-end image segmentation for change detection in high-resolution remote sensing images. *International Journal of Applied Earth Observation and Geoinformation* **104**, 102582 (2021)
24. Weiss, M., Jacob, F., Duveiller, G.: Remote sensing for agricultural applications: A meta-review. *Remote Sensing of Environment* p. 111402 (2020)
25. Wilson, G., Cook, D.J.: A survey of unsupervised deep domain adaptation (2018)
26. Wu, C., Zhang, L., Du, B.: Kernel slow feature analysis for scene change detection. *IEEE Transactions on Geoscience and Remote Sensing* (4), 2367–2384 (2017)
27. Wu, L., Zhang, Z., Wang, Y., Liu, Q.: A segmentation based change detection method for high resolution remote sensing image. In: Li, S., Liu, C., Wang, Y. (eds.) *Pattern Recognition*. pp. 314–324. Springer Berlin Heidelberg, Berlin, Heidelberg
28. Xu, L., Zhang, S., He, Z., Guo, Y.: The comparative study of three methods of remote sensing image change detection. In: 2009 17th International Conference on Geoinformatics. pp. 1–4 (2009)