

Continuous Cross-resolution Remote Sensing Image Change Detection

Hao Chen, Haotian Zhang, Keyan Chen, Chenyao Zhou, Song Chen, Zhengxia Zou and Zhenwei Shi*, *Member, IEEE*

Abstract—Most contemporary supervised Remote Sensing (RS) image Change Detection (CD) approaches are customized for equal-resolution bitemporal images. Real-world applications raise the need for cross-resolution change detection, aka, CD based on bitemporal images with different spatial resolutions. Given training samples of a fixed bitemporal resolution difference (ratio) between the high-resolution (HR) image and the low-resolution (LR) one, current cross-resolution methods may fit a certain ratio but lack adaptation to other resolution differences. Toward continuous cross-resolution CD, we propose scale-invariant learning to enforce the model consistently predicting HR results given synthesized samples of varying resolution differences. Concretely, we synthesize blurred versions of the HR image by random downsampled reconstructions to reduce the gap between HR and LR images. We introduce coordinate-based representations to decode per-pixel predictions by feeding the coordinate query and corresponding multi-level embedding features into an MLP that implicitly learns the shape of land cover changes, therefore benefiting recognizing blurred objects in the LR image. Moreover, considering that spatial resolution mainly affects the local textures, we apply local-window self-attention to align bitemporal features during the early stages of the encoder. Extensive experiments on two synthesized and one real-world different-resolution CD datasets verify the effectiveness of the proposed method. Our method significantly outperforms several vanilla CD methods and two cross-resolution CD methods on the three datasets both in in-distribution and out-of-distribution settings. The empirical results suggest that our method could yield relatively consistent HR change predictions regardless of varying bitemporal resolution ratios. Our code is available at https://github.com/justchenhao/SILI_CD.

Index Terms—High-resolution remote sensing image, Cross-resolution change detection, Scale-invariant learning, Implicit neural representation, Attention mechanism.

I. INTRODUCTION

The work was supported by the National Key Research and Development Program of China (Grant No. 2022ZD0160401), the National Natural Science Foundation of China under Grant 62125102, the Beijing Natural Science Foundation under Grant JL23005, and the Fundamental Research Funds for the Central Universities. (Corresponding Author: Zhenwei Shi (shizhenwei@buaa.edu.cn)).

Hao Chen, Haotian Zhang, Keyan Chen, Chenyao Zhou, and Zhenwei Shi are with the Image Processing Center, School of Astronautics, Beihang University, Beijing 100191, China, and with the Beijing Key Laboratory of Digital Media, Beihang University, Beijing 100191, China, and with the State Key Laboratory of Virtual Reality Technology and Systems, Beihang University, Beijing 100191, China, and also with the Shanghai Artificial Intelligence Laboratory, Shanghai 200232, China.

Song Chen is with the Department of Journalism and Communications, Jeonuk National University, Jeonju-si 54896, South Korea.

Zhengxia Zou is with the Department of Guidance, Navigation and Control, School of Astronautics, Beihang University, Beijing 100191, China, and also with the Shanghai Artificial Intelligence Laboratory, Shanghai 200232, China.

REMOTE sensing (RS) image change detection (CD) refers to identifying changes of interest objects or phenomena in the scene by comparing co-registered multi-temporal RS images taken at different times in the same geographical area [1]. Change detection could be applied to various applications, e.g., urban planning [2], disaster management [3], agricultural surveys [4], and environmental monitoring [5].

The availability of high-resolution (HR) remote sensing (RS) images enables monitoring changes on Earth's surface at a fine scale. Existing deep learning-based techniques, such as convolutional neural networks (CNNs) [6] and vision transformers [7], are widely applied in RS CD [8]. Despite promising results, most existing supervised CD approaches are customized for handling equal-resolution bitemporal images and are insufficient to adapt to cross-resolution conditions, where bitemporal images have different resolutions.

Real-world applications raise the need for change recognition based on multi-temporal images across resolutions. We identify roughly two scenarios: 1) the long-term CD task, with a relatively low-resolution (LR) pre-event image and an HR post-event one considering earlier satellite observations (e.g., decades before) have relatively lower spatial resolution than those obtained by current satellite sensors; 2) the event/disaster rapid response task, with an archived HR pre-event image of a certain area and a relatively LR post-event image, considering the lack of real-time availability of HR satellite data, due to its smaller spatial coverage and longer revisit period, compared to LR data.

To handle the cross-resolution RS CD, aka., change detection based on bitemporal images with different spatial resolutions, most current methods [9–15] align the two inputs in the image space, either by downsampling the HR image [12, 13] or upsampling the LR image in a fixed (e.g., bilinear/cubic interpolation) [11] or a learnable manner [9, 10]. Recent attempts [14–16] align the bitemporal resolution differences in the feature space, e.g., upsampling the feature map of the LR image by considering that of the HR one [14].

Despite current progress in cross-resolution CD, a model [9–15] trained with a fixed resolution difference (i.e., difference ratio, e.g., 4 or 8) may work well for a certain condition, but may not be suitable for situations of other resolution differences, which limits its real-world applications. To fill this gap, different from existing approaches that are specifically designed for a fixed bitemporal resolution difference, we explore a continuous cross-resolution CD method that enables a single model to adapt arbitrary difference ratios between bitemporal

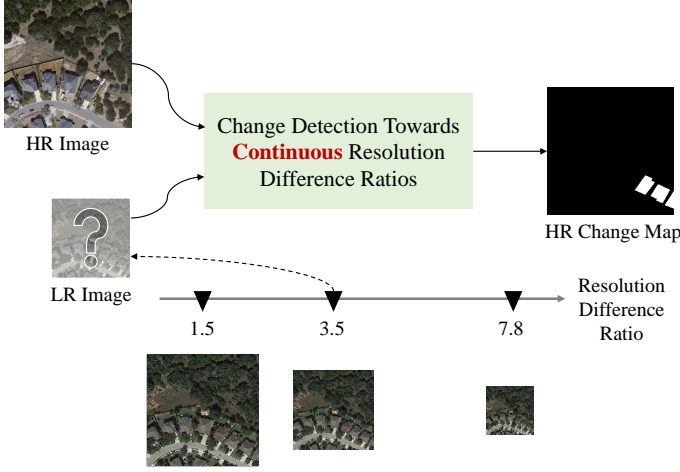


Fig. 1. Illustration of continuous cross-resolution change detection, i.e., CD towards varying resolution difference ratios between the HR image and the relatively LR image.

images. Different from the traditional cross-resolution CD task that applies a fixed resolution difference for assessment, the continuous cross-resolution CD task evaluates the CD model on validation/testing samples with varying bitemporal resolution difference ratios that may be different from that of the training samples. As shown in Fig. 1, given an HR image and a relatively LR image, our goal is to obtain the HR change map regardless of resolution difference ratios.

To achieve this, we propose a scale-invariant training pipeline to learn an embedding space that is insensitive to scale changes of input images. Given original training samples with a fixed resolution ratio, we synthesize samples with random resolution differences by downsampling HR images and swapping bitemporal regions. We enforce the model yielding HR predictions for these synthesized samples, thus improving the ability to adapt different resolution ratios. We then incorporate the coordinate-based method, namely implicit neural representation (INR) [17], to decode pixel-wise change information from the embedding space and corresponding pixel positions. Specifically, a multi-level feature embedding space is learned for a trade-off between semantic accuracy and spatial details [18]. Different from existing CD methods that employ sophisticated multi-level feature fusion (e.g., UNet [5, 19–27] or FPN[28–30]) to yield HR predictions, our coordinate-based approach implicitly learns land-cover shapes that may benefit handling LR images with blurry low-quality objects. Furthermore, we propose bitemporal interaction on the early-level features to further fill the resolution gap by applying transformers [31] to model correlations between bitemporal pixels within the local windows on the feature maps. Motivated by the fact that spatial resolution differences directly affect the local textures and image details are locally correlated without long-range dependency [32], only local information of the LR and HR patches may be sufficient to model correlations between cross-resolution pixels.

The contribution of our work can be summarised as follows:

- We propose a scale-invariant methodology whereby an embedding space insensitive to scale changes is learned

for cross-resolution RS image CD. Unlike extant approaches that are tailored to specific difference ratios between bitemporal resolutions, our method is capable of adapting to continuous resolution difference ratios.

- We introduce coordinate-based representations to decode the HR change mask from the embedding space by implicitly learning the shape of objects of interest, therefore benefiting recognizing blurred objects in the LR image. Moreover, we incorporate local in-window interactions between bitemporal features to equip the model to better adapt to resolution disparities across bitemporal images.
- Extensive experiments on two synthetic and one real-world cross-resolution CD datasets validate the effectiveness of the proposed method. Our approach outperforms several extant methods for cross-resolution CD as well as vanilla CD methods in both in-distribution and out-of-distribution settings.

The rest of this paper is organized as follows. Sec. II introduces related work of existing methods of vanilla CD and those handling bitemporal resolution differences. Sec. III presents the proposed scale-invariant learning with implicit neural networks. Experimental results are given in Sec. IV, and the Conclusion is drawn in Sec. V.

II. RELATED WORK

A. Deep Learning-based optical RS Image CD

The past several years have witnessed remarkable progress in supervised change detection for optical remote sensing imagery using deep learning (DL) techniques. Advanced DL techniques, e.g., CNNs [6], fully convolutional neural networks (FCN) [18], and transformers [31] have been widely applied in the field of RS CD [8].

The predominant recent attempts have aimed to enhance the discriminative capacity of CD models by incorporating advanced network backbones (e.g., HRNet[33, 34], vision transformers[35–38]) and network structures (e.g., dilated convolution [39, 40], deformable convolution [41, 42], attention mechanism [2, 14, 22, 23, 26, 40, 43–55], and flow field-based model [51, 56]), devising multi-level bitemporal feature fusion strategies (e.g., UNet [5, 19–27] or FPN[28–30]), employing multi-task learning (e.g., additional supervision of land-cover maps for each temporal [23, 57–59], boundary supervision of the change edge map [38, 60, 61]), combining generative adversarial network (GAN) loss [21, 62], training with more diverse synthetic data [27, 63, 64], and fine-tuning from a pre-trained model (e.g., self-supervised pre-training [65, 66] and supervised pre-training [67]). Note that the paper mainly focuses on binary change detection. The additional supervision of land-cover maps for each temporal [59] could also improve the binary change detection performance apart from the purpose of identifying the semantic change categories.

Context modeling, encompassing both spatial context and spatial-temporal context, is crucial for discerning changes of relevance and filtering out extraneous changes across bitemporal images. Among the aforementioned attempts, attention mechanisms, including channel attention[22, 23, 40, 43–46, 68], spatial attention [22, 23, 40, 43, 44], self-attention

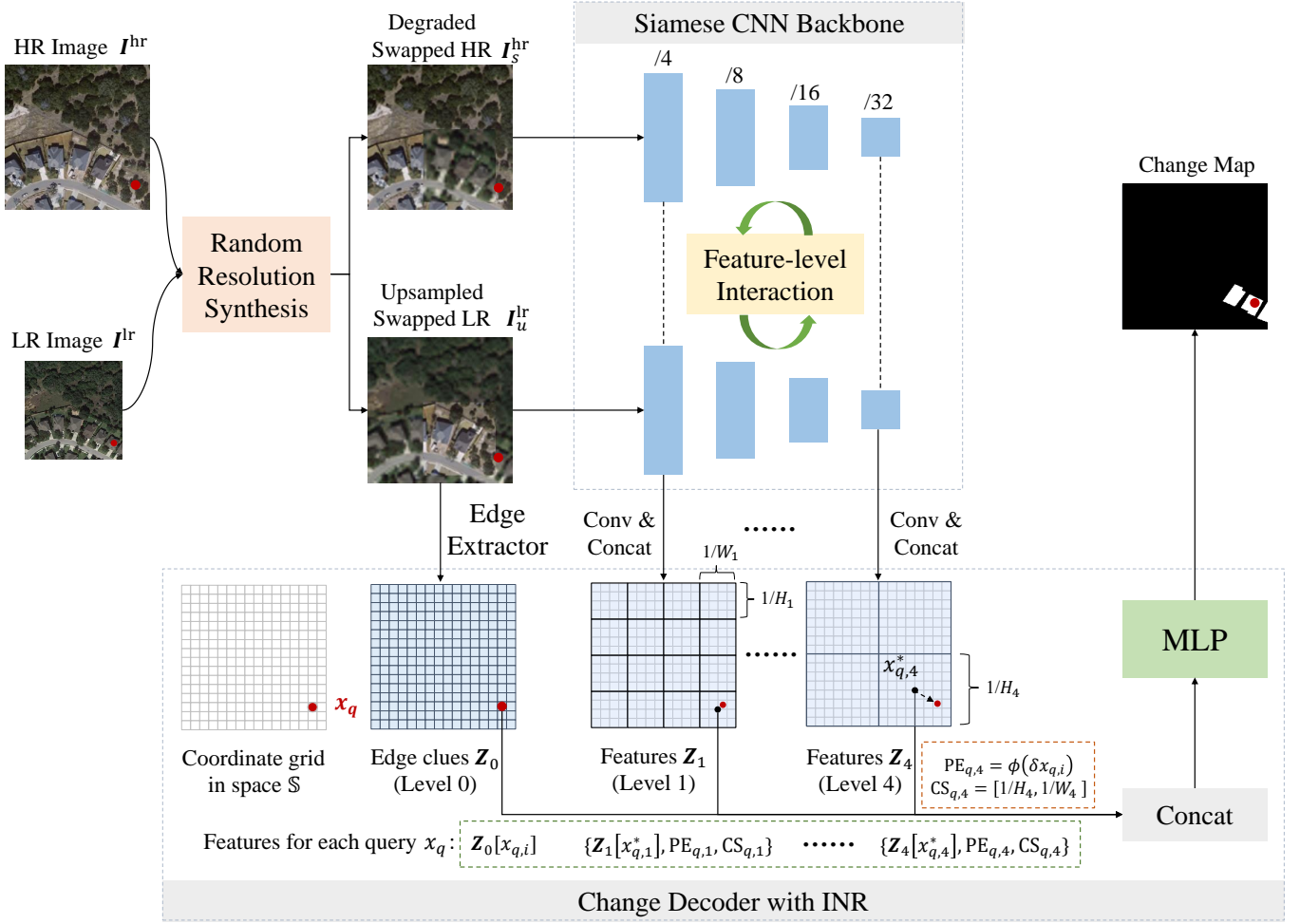


Fig. 2. Illustration of the overall pipeline of the proposed Scale-Invariant Learning with Implicit neural networks (SILI). We aim to learn a multi-level embedding space invariant to resolution differences across bitemporal images by enforcing the model generating HR change masks for synthesized samples with random resolution difference ratios. We leverage implicit neural representations that encapsulate the local mask shape to decode the HR mask from dense coordinate queries and corresponding multi-level features, including learnable edge clues. Note that we calculate positional encoding (PE) of the relative coordinate between the query and corresponding feature cell. The cell scale (CS) of each level is also fed into the decoder.

[2, 46, 52–54], and cross-attention [30, 55, 69–71], have been extensively leveraged as conventional context modeling techniques for the CD task. Early works that incorporate spatial context have primarily focused on employing attention mechanisms as feature enhancement modules, applying them either separately to each temporal image [23, 72] or to fused bitemporal features [22, 40, 43–45], lacking the exploit of the temporal-related correlations. More recent works have explicitly modeled spatial-temporal relations by employing cross-attention [30, 70, 71] or self-attention/transformers on bitemporal features [42, 47, 48, 50, 51, 53, 73–76]. For instance, the Bitemporal Image Transformer (BIT) [53] applies multi-head self-attention to sparse visual tokens extracted from the per-pixel feature space of bitemporal images to efficiently model their global spatial-temporal relations.

Different from existing context modeling approaches in CD, we introduce local-window self-attention over bitemporal pixels belonging to each small non-overlapping image window. Our motivation stems from the notion that disparities in spatial resolution reflect differences in local textural detail

within images. Comparing local regions between bitemporal images of varying resolutions may therefore suffice to align their features. Although Swin Transformer [77] whose core is local-window self-attention has been applied in the CD task [36, 37, 78], it is treated as the mere network backbone, therefore processing bitemporal images independently without modeling their temporal correlations.

Furthermore, most current CD methods have been principally formulated under the assumption of consistent spatial resolution across the bitemporal images. They are thus inadequate for application to cross-resolution paradigms. Our proposed model, in contradistinction, is specifically designed to adapt to resolution differences across bitemporal images.

B. Handling Bitemporal Resolution Differences

In light of their real-world applicability, cross-resolution change detection (also termed different-resolution change detection), operable on remote sensing imagery of heterogeneous resolution obtained through different sensors, has claimed burgeoning interest [9–15]. This article mainly focuses on the

supervised cross-resolution CD of optical RS images. Works addressing heterogeneous image change detection [13, 16] for synthetic aperture radar (SAR) and optical data fall outside the purview of this article.

There exist two predominant categories of cross-resolution CD techniques: those operating in image space and those operating in feature space. 1) Image-space alignment [9–12]: first calibrate the spatial scales of bitemporal images and then apply conventional CD methods to the aligned images. The simplest way is to upsample LR images to HR resolution via bilinear/cubic interpolation [11] or down-sample HR images to LR ones [12]. More recently, super-resolution techniques have been deployed for low-to-high-resolution transformation in a learnable fashion [9, 10]. 2) Feature-space alignment: align feature representations via interpolation [14, 15]. One Recent work [14] applies transformers to learn correlations between the upsampled LR features and original HR ones, achieving semantic alignment across resolutions.

Most existing methods have been formulated solely for scenarios exhibiting a fixed resolution difference, thus inadequate when the resolution discrepancy between bitemporal images varies. Towards more practical real-world applications, we propose a method adaptable to variable resolution differences. Specifically, we learn a scale-invariant embedding space insensitive to changes in resolution via enforcing the model outputs HR CD results regardless of the downsampling factor applied to input HR images. The synthetic reconstructions of randomly downsampled HR images narrow the resolution gap between HR and LR images, thereby achieving adaptability to continuous resolution differences.

C. Implicit Neural Representation

Implicit Neural Representation (INR), also known as coordinate-based neural representations, is essentially a continuously differentiable function that facilitates transformations from coordinates to signals [79]. Originally stemming from the field of 3D reconstruction, INR is used to represent the object shape [80] and 3D scenes [81] as a replacement for explicit representations such as point clouds, meshes, or voxels. Thanks to the design of coordinate-based representation, INR exhibits an ability to model images of variant resolutions, thus being employed in image processing tasks such as super-resolution [82, 83], semantic segmentation [84, 85], and instance segmentation [86].

Recently, INR has been applied in the field of RS [87–93], including 3D RS scene reconstruction [87] and segmentation [88, 89], 2D RS image synthesis [90], and super-resolution [91, 92]. However, the employment of INR for RS 2D image understanding remains limited, particularly for the task of CD which has received scarce exploration. For the task of cross-resolution CD, we incorporate INR to enhance model adaptability to cross-temporal resolution differences. Our motivation is that the INR may enable the implicit encoding of the shape of change objects and extraction of the corresponding HR change mask from latent space based on coordinate queries, regardless of the resolution difference across bitemporal images.

III. SCALE-INVARIANT LEARNING WITH IMPLICIT NEURAL REPRESENTATIONS FOR CROSS-RESOLUTION CHANGE DETECTION

In this section, we first give an overview of the proposed scale-invariant cross-resolution method and then introduce its three main components. Finally, implementation details are given.

A. Overview

Cross-resolution change detection aims to obtain an HR change mask based on bitemporal images with different resolutions (i.e., an LR image and an HR image). Towards real-world applications, here we propose Scale-Invariant Learning with Implicit neural networks (SILI) for handling varying resolution difference ratios across bitemporal images.

The essence of the proposed method is to learn a scale-invariant embedding space regardless of the resolution discrepancies between bitemporal images, and decode the high-resolution change mask with dense coordinate queries and corresponding multi-level features by leveraging the implicit neural representations encapsulating the shape of local changes.

Fig. 2 illustrates the proposed SILI. It has three main components, including random resolution image synthesis, image encoder with feature-level bitemporal interaction, and change decoder based on implicit neural representations.

1) **Random resolution image synthesis.** Rather than manipulating the original bitemporal images with a fixed resolution difference ratio, we perform random downscaling reconstruction on the HR image (i.e., downsampling succeeded by upsampling) to narrow the resolution gap between HR and LR images. Moreover, we propose random bitemporal region swapping to further improve the model adaptability to scale variance. For more details see Sec. III-B.

2) **Image encoder.** Given synthesized bitemporal images, a normal Siamese CNN backbone (e.g., ResNet-18 [6]) is employed to obtain multi-level image features for each temporal. Bitemporal features of certain levels are interacted with each other by leveraging local-window self-attention to reduce the semantic gap caused by resolution differences. Details of the bitemporal feature interaction see Sec. III-C

3) **Change decoder.** Instead of upsampling multi-level bitemporal features to the target size with traditional interpolation, we incorporate coordinate-based representations to decode the label for each position by feeding corresponding multi-level features and position embeddings to a multilayer perceptron (MLP) that implicitly learns the shape of local changes. See Sec. III-D for more details.

B. Random Resolution Image Synthesis

In the training phase, we introduce scale-invariant learning by compelling the change detection model to generate HR change masks for synthesized bitemporal images subject to random scale manipulation, thus enhancing the adaptation capacity for handling continuous resolution difference ratios across bitemporal images. Specifically, given bitemporal images ($\mathbf{I}^{\text{lr}} \in \mathbb{R}^{H_{\text{lr}} \times W_{\text{lr}} \times 3}$, $\mathbf{I}^{\text{hr}} \in \mathbb{R}^{H_{\text{hr}} \times W_{\text{hr}} \times 3}$) with resolution

differences ratio ($r_d = H_{hr}/H_{lr}$), we design three main steps: 1) upsampling the LR image to the same size as that of the HR image, 2) perform random downsampling reconstruction on the HR image, 3) random region swap between bitemporal images. Note that the resolution differences ratio defines as the ratio of the ground resolution of the LR image and that of the HR image. For example, $r_d = 4$ for bitemporal images with ground resolution 0.5m/pixel and 2m/pixel. Fig. 3 illustrates the overall process of random resolution image synthesis.

1) **Upsampling LR image.** We first upsample the LR image I^l_r to the size of I^{hr} via bicubic interpolation. Instead of learning upsampled LR images via training a super-resolution reconstruction model, we aim to learn a scale-invariant CD model that is able to handle degraded constructions with essentially lower resolutions than the HR image. Formally, The upsampled LR image $I^l_u \in \mathbb{R}^{H_{hr} \times W_{hr} \times 3}$ is given by

$$I^l_u = \text{upsampling}(I^l_r, r_d). \quad (1)$$

2) **Random downsampling reconstruction.** To acclimatize the model to various resolution differences, we synthesize degraded variants of the HR image through downsampling by a random ratio, thereafter rescaling to the initial size. Formally, we randomly sample a ratio r from the Uniform distribution $r \sim U[1, r_d]$. The downsampling reconstruction version $I^{hr}_d \in \mathbb{R}^{H_{hr} \times W_{hr} \times 3}$ can be given by

$$I^{hr}_d = \text{upsampling}(\text{downsampling}(I^{hr}, r), r), \quad (2)$$

where bicubic interpolation is applied to implement upsample and downsample.

3) **Random bitemporal region swap.** Considering that simply downsampling the HR image may not wholly fill the gap between the HR and the LR image captured by different sensors, we further propose to swap a randomly selected region between bitemporal images. Such operation can be viewed as a form of image-level bitemporal interaction, allowing the CD model to process LR and HR data concurrently, which may benefit learning more scale-invariant representations. Formally, the swapped bitemporal images $I^{hr}_s, I^{lr}_s \in \mathbb{R}^{H_{hr} \times W_{hr} \times 3}$ are given by

$$I^{hr}_s, I^{lr}_s = \text{swap}(I^{hr}_d, I^{lr}_u, u, v, \text{crop_size}), \quad (3)$$

where $(u, v), u \sim U[1, W_{hr} - \text{crop_size}], v \sim U[1, H_{hr} - \text{crop_size}]$ is the coordinate of the upper-left point of the cropped region and crop_size is the size of the swapped region. crop_size is default set to half of W_{hr} .

Note that in the inference/testing phase, we only perform the first step, i.e., rescale the LR image to the size of the HR image. In other words, we do not apply the random downsampling reconstruction and random bitemporal region swap in the testing phase.

C. Image Encoder with Bitemporal Local Interaction

Given synthesized bitemporal images, we employ an off-the-shell Siamese CNN backbone (i.e., ResNet-18) for generating multi-level features $X^i_j \in \mathbb{R}^{H_j \times W_j \times C_j}$ for each temporal image $i \in \{1, 2\}$. Note that $j \in \{1, 2, 3, 4\}$ denotes the level of the generated features with the size of $H_j \times W_j, H_j =$

$H_{hr}/2^{(j+1)}, W_j = W_{hr}/2^{(j+1)}$. Instead of encoding bitemporal images independently without interaction, we supplement feature interaction between bitemporal image features of a certain level to refine them via modeling local spatial-temporal correlations thus benefiting feature extraction at the next level.

Concretely, we introduce local-window self-attention [77] over bitemporal pixels within each non-overlapping window on the feature map of a certain level. Our incentive resides in that discrepancies in spatial resolution between images may predominantly mirror local texture variances in land cover and thus leveraging local correlations may in turn benefit aligning features of bitemporal images with different resolutions.

Fig. 4 illustrates the proposed bitemporal interaction based on local-window self-attention. For bitemporal features X^1_j, X^2_j of a certain level j , we evenly partition them into non-overlapped windows. Let $X^{1*}_{j,n}, X^{2*}_{j,n} \in \mathbb{R}^{WS \times WS \times C_j}, n \in \{1, \dots, N_w\}$ be bitemporal features within each window, where WS is the window size, n denotes the window index in N_w partitioned windows. We apply multi-head self-attention (MSA) on bitemporal patches within each local window. Formally, the refined bitemporal features $X^{1*}_{j,n}, X^{2*}_{j,n}$ of level j are given by

$$X^{1*}_{j,n}, X^{2*}_{j,n} = \text{Transformer_Encoder}(X^1_{j,n}, X^2_{j,n}), \quad (4)$$

where a vanilla transformer encoder [31] is employed to implement multi-head self-attention. Note that we apply shared learnable local positional embeddings (PE) [31] for each window. The PE could encode temporal and local spatial position information, thus helping model spatial-temporal relations. The transformer encoder consists of N_{te} transformer layers. Our empirical evidence (see Sec. IV-E) suggests local context modeling at early layers is sufficient. Please refer to [31] for more details on the transformer layer. Note that the calculation of local-window self-attention for each window could be processed in parallel in GPU.

After gleaned bitemporal multi-level features, we further transform the output features to a uniform dimension C by applying one 1×1 convolution layer to each level. The transformed bitemporal image features $Z^1_j, Z^2_j \in \mathbb{R}^{H_{hr} \times W_{hr} \times 2C}$ of each level j are then fused via channel-wise concatenation. The resulting multi-level features $\{Z_j | j \in \{1, 2, 3, 4\}\}$ are given by

$$Z_j = \text{Concat}(Z^1_j, Z^2_j). \quad (5)$$

Apart from the multi-level features from the vanilla backbone, we extract handcrafted low-level edges from bitemporal images as spatial clues to obtain high-resolution change masks in the subsequent change decoder. It is motivated by the evidence [94–96] that the incorporation of handcrafted edge features (e.g., Canny [97], Sobel, or Prewitt operator) within the deep neural networks benefits the change detection performance. As the Canny operator could obtain more clean and accurate edges than the Sobel operator, we chose the Canny operator to extract low-level edge clues. Here, we simply utilize the Canny operator on each dimension of bitemporal images to obtain handcrafted edge features which are then fed into the change decoder. Formally, the edge features

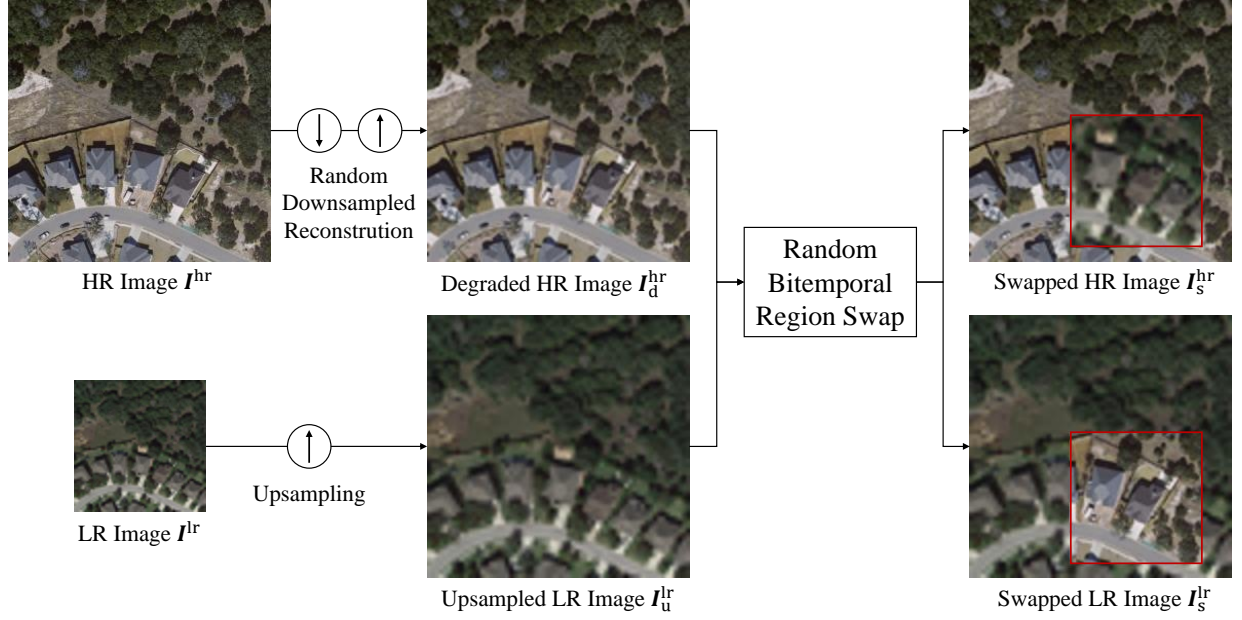


Fig. 3. Illustration of random resolution image synthesis, including 1) upsampling the LR image, 2) random downsampling and reconstructing the HR image, and 3) random region swap between the upsampled LR image and the degraded HR image.

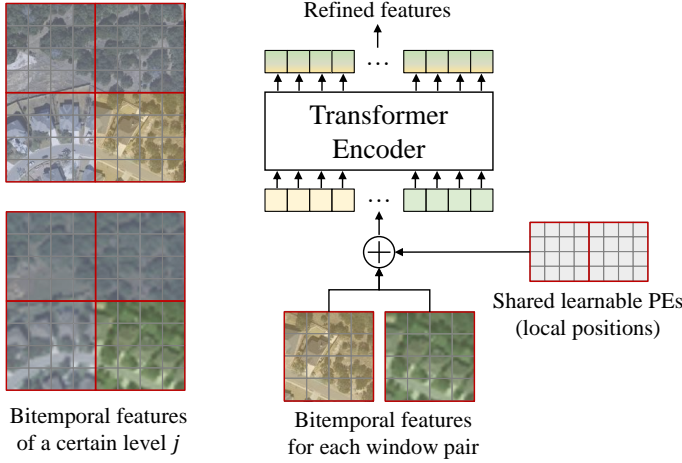


Fig. 4. Illustration of bitemporal interaction based on local-window self-attention. For bitemporal features of a certain level j , we perform multi-head self-attention within each window pair to better semantically align features from different resolution images.

$X_0 \in \mathbb{R}^{H_{hr} \times W_{hr} \times 3}$ are given by the channel-wise summation of that from each temporal image as follows:

$$X_0 = \text{Canny}(I_s^{hr}) + \text{Canny}(I_u^{lr}). \quad (6)$$

Note that for simplicity, we directly perform pixel-wise addition of bitemporal Canny features to obtain the edge clues.

Similarly, the handcrafted edge X_0 goes through a relatively large kernel (7×7) convolution layer to obtain the learned edge clues $Z_0 \in \mathbb{R}^{H_{hr} \times W_{hr} \times 3}$ that are then fed into the subsequent change decoder.

D. Change Decoder with Implicit Neural Representation

Given multi-level bitemporal image features and edge clues, we aim to decode the HR change mask $CM \in \mathbb{R}^{H_{hr} \times W_{hr}}$ by

leveraging implicit neural representation (INR), viz. feeding dense coordinates alongside corresponding image features to a learnable MLP that implicitly represents the shape of local changes. Our motivation is that the INR may assist in reconstructing the detailed shape of the degraded land cover of change from the LR image by leveraging fine features from HR images. The key is to learn implicit neural networks f_θ (typically an MLP) over coarse resolution feature maps to define continuous representations that could yield the HR change mask according to the coordinate queries of the HR grid.

Now, we define a continuous normalized 2D space $\mathbb{S} = \{x = (u, v) | u, v \in [0, 1]\}$. Images or feature maps of different levels can be evenly distributed in the space \mathbb{S} where each cell in the grid is assigned a 2D coordinate according to its center position. For instance, given a position indexed h -th, w -th ($h \in \{0, 1, \dots, H-1\}, w \in \{0, 1, \dots, W-1\}$) in an grid of size $H \times W$, its coordinate in space \mathbb{S} is $(u, v) = (\frac{1}{2H} + \frac{h}{H}, \frac{1}{2W} + \frac{w}{W})$.

Fig. 5 illustrates the coordinate relations between the HR grid and a relatively LR grid with respect to the space \mathbb{S} . We only show one dimension (width direction) for a better view. Here, the HR grid denotes the dense coordinate queries while the LR grid denotes the feature map from a certain level.

Query features from relatively LR feature maps. Let x_q be the coordinate of point q in an HR grid with respect to \mathbb{S} . Given one query x_q , we need first to collect corresponding features on the coarse feature map Z_j by calculating the nearest cell to the query for each level $j \in \{1, 2, 3, 4\}$. The matched coordinate $x_{q,j}^*$ for Z_j can be given by

$$x_{q,j}^* = (\frac{1}{2H_j} + \frac{h^*}{H_j}, \frac{1}{2W_j} + \frac{w^*}{W_j}), \quad (7)$$

where (h^*, w^*) is the corresponding coordinate of the matched

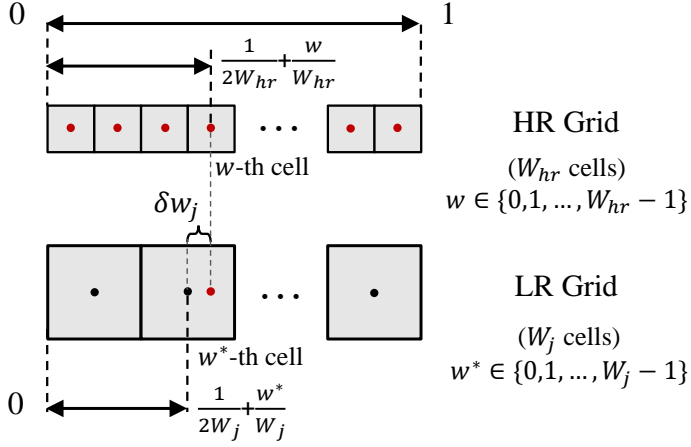


Fig. 5. Illustration of grid coordinates in the normalized space \mathbb{S} . Here, the HR grid denotes the dense coordinate queries and the LR grid denotes the feature map from a certain level j . Only the width direction is shown for a better view. We also demonstrate the coordinate matching between the query (w -th) cell in the HR grid and its nearest (w^* -th) cell in the LR grid.

point in the discrete space of \mathbf{Z}_j . h^*, w^* can be calculated as follows:

$$h^* = \text{round}\left(\frac{H_j}{H_{hr}}\left(\frac{1}{2} + h\right) - \frac{1}{2}\right) \quad (8)$$

$$w^* = \text{round}\left(\frac{W_j}{W_{hr}}\left(\frac{1}{2} + w\right) - \frac{1}{2}\right), \quad (9)$$

where h, w is the height/width index of q in the HR image, respectively.

Because the feature map of level 0 has the same resolution as the HR coordinate grid, the corresponding feature vector $\mathbf{Z}_0[x_q]$ can be directly obtained by the query coordinate x_q . Note that for a tradeoff between accuracy and efficiency (see IV-E), the input coordinate grid is downsampled by a factor of 2 compared to the original HR image.

Relative positional encoding (PE). We further calculate the relative coordinate encoding $\text{PE}_{q,j} \in \mathbb{R}^{C_{pe}}$ for level j between x_q and $x_{q,j}^*$:

$$\text{PE}_{q,j} = \phi(\delta x_{q,j}) = \phi(x_q - x_{q,j}^*), \quad (10)$$

where $\phi(\cdot)$ denotes the position encoding function [31] to transform the 2D coordinate to high dimensional vectors that are more capable of representing high-frequency signals. $\delta x_{q,j}$ is the relative coordinate between the query and its nearest grid cell center.

Encode cell scale. Considering that the grid cell of different resolutions occupies different spatial scopes, to distinguish features from different levels, we also combine cell scale, i.e., the absolute height and width of a cell with respect to the continuous space \mathbb{S} as follows:

$$\text{CS}_{q,j} = \left[\frac{1}{H_j}, \frac{1}{W_j}\right], \quad (11)$$

where $\text{CS}_{q,j} \in \mathbb{R}^2$ is the cell scale for feature level $j, j \in \{1, 2, 3, 4\}$.

Decode change probability. As shown in Fig. 2, after obtaining multi-level features and corresponding PEs and cell

scales, an MLP is employed to decode the change probability for each query as follows:

$$\mathbf{P}[x_q] = f_\theta(\text{Concat}(\mathbf{Z}_0[x_q], \{\mathbf{Z}_j[x_{q,j}^*], \text{PE}_{q,j}, \text{CS}_{q,j}\}_{j=1}^4)), \quad (12)$$

where $\text{Concat}(\cdot)$ denotes to channel-wise concatenate the input items. $\mathbf{P} \in \mathbb{R}^{H_{hr} \times W_{hr} \times 2}$ is the change score maps where the 2D vector for each position indicates the probability of change or not.

E. Implementation Details

CNN backbone. We employ the off-the-shell ResNet-18 as the CNN backbone. Its intermediate multi-level features (channel dimensions C_j are 64, 128, 256, 512, respectively for level $j \in \{1, 2, 3, 4\}$) are transformed to the same dimension $C = 64$ via one convolution layer. We apply the bitemporal interaction with a local window size $\text{WS} = 8$ at level $j = 1, 2, 3$.

Change decoder. The channel dimension C_{pe} of the relative positional encoding is set to 24. The implicit neural network f_θ is implemented by a three-layer MLP with BatchNorm and ReLU in between. The output channel dimension of each MLP is "64, 64, 2", respectively.

Loss function. In the training phase, we minimize the cross-entropy loss between the predicted change probability map \mathbf{P} and ground truth to optimize the network parameters. In the inference phase, the change mask can be obtained by per-pixel Argmax operation on the channel dimension of \mathbf{P} .

IV. EXPERIMENTAL RESULTS

A. Experimental Setup

To evaluate the proposed cross-resolution CD model, we conduct experiments on two synthesized cross-resolution CD datasets (construct the LR image by downsampling), and one real-world cross-resolution CD dataset (the LR and HR image pair captured by different-resolution satellite sensors).

LEVIR-CD(4 \times). LEVIR-CD [2] is a widely used building CD dataset, which contains 637 pairs of bitemporal VHR (0.5m/pixel) images, each size of 1024×1024 . We follow the default dataset split [2], including 445/64/128 samples for training/validation/testing, respectively. We further crop each sample into small patches of size 256×256 with no overlap. To synthesize cross-resolution scenarios, we downsample the post-event (t2) image for each sample by a ratio of 4. In this way, we obtain the simulated LEVIR-CD(4 \times), where the post-event image has a 4 times spatial resolution lower than that of the pre-event (t1) image.

SV-CD(8 \times). Season-varying change detection (SV-CD) [98] is another widely used binary CD dataset. It contains 11 pairs of VHR (0.031m/pixel) RGB images with sizes ranging from 1900×1000 to 4725×2700 pixels. The changes in buildings, cars, and roads are taken into consideration. We follow the default dataset split, which contains 10000/3000/3000 cropped samples each size of 256×256 for training/validating/testing, respectively. For each sample, we downsample the post-event image by a ratio of 8, thus obtaining the synthesized SV-CD(8 \times).

DE-CD(3.3 \times). DynamicEarthNet [99] is a multi-class land use and land cover (LULC) segmentation and change detection dataset for daily monitoring of Earth’s surface. It covers 75 areas of interest (AOIs) around the world and consists of samples captured in the range from 2018-01-01 to 2019-12-31. Each AOI provides high-time-frequency (daily) Planet imagery (3m/pixel) and monthly LULC per-pixel annotations, as well as monthly Sentinel-2 imagery (upsampled to match the size of the Planet image) whose original spatial resolution is 10m/pixel. Each sample has a size of 1024×1024 . We reorganize the original dataset for cross-resolution change detection. We collect the time-aligned (monthly) image and label data, where the Sentinel-2 data in each month in 2018 is for pre-event and the Planet data captured 1 year later than the Sentinel-2 is for post-event. For simplicity, we only focus on the changes in the land cover belonging to impervious surfaces. We exclude those without changes of interest and therefore obtain 506 samples, which are then randomly split into 354/51/101 samples for training/validating/testing. In this way, we have aggregated DE-CD(3.3 \times) where the bitemporal resolution difference ratio is around 3.3. Similarly, we cropped each sample into 256×256 small patches with no overlap.

To evaluate the proposed method, we set the following models for comparison:

1) **Base**. Our baseline consists of a CNN backbone (ResNet-18) and a change decoder with the channel-wise concatenated input of bitemporal transformed features (channel dimension of 64) at each level ($j \in \{1, 2, 3, 4\}$) from the encoder. Similar to our INR decoder, the baseline decoder has three-layer convolutions (output dimensions of 64,64,2) with BatchNorm and ReLU in between.

2) **SILI**. Our proposed SILI model with the random resolution image synthesis, a ResNet-18-based encoder with bitemporal feature interactions, and a change decoder with INR.

Training details. Data augmentation techniques, including random flip, and Gaussian blur are applied. We employ SGD with a batch size of 8, a momentum of 0.9, and a weight decay of 0.0005. The initial learning rate is 0.01 and linearly decays to 0 until 200 epochs. We evaluate the model using the validating set at the end of each training epoch. The best validating model is evaluated on the test set.

Evaluation Metrics. We use the F1-score regarding the change category as the evaluation metrics. Precision, recall, and Intersection over Union (IoU) belonging to the change category are also reported. These indices can be defined by:

$$\begin{aligned} F1 &= \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} \\ \text{precision} &= \frac{TP}{TP + FP} \\ \text{recall} &= \frac{TP}{TP + FN} \\ \text{IoU} &= \frac{TP}{TP + FN + FP} \end{aligned} \quad (13)$$

where TP, FP, FN are the number of true positives, false positives, and false negatives, respectively.

B. Overall Comparison

We make a comparison with several state-of-the-art conventional change detection methods, including three pure convolutional-based methods (FC-EF [20], FC-Siam-Diff [20], FC-Siam-Conc [20]), and six attention-based methods (STANet [2], IFNet [22], IFNet [22], SNUNet [68], BIT [53], ICIFNet [47], DMINet [55]). We also compare two CD methods (SUNet [15], SRCDNet [10]) specifically for the scenario of different resolutions across bitemporal images.

- FC-EF [20]. Image-level fusion method where bitemporal images are channel-wise concatenated to be fed into an FCN.
- FC-Siam-Diff [20]. Feature-level fusion method where a Siamese FCN is employed to obtain multi-level features for each temporal image, then bitemporal feature differencing is calculated for fusing temporal information.
- FC-Siam-Conc [20]. Feature-level fusion method where channel-wise concatenation is used for fusing temporal information.
- STANet [2]. Metric-based method, which incorporates multi-scale self-attention to enhance the discriminative capacity for bitemporal features.
- IFNet [22]. Feature-level concatenation method, which employs channel/spatial attention on the concatenated bitemporal features at each level of the decoder. Deep supervision is applied on each level for better training of the intermediate layers.
- SNUNet [68]. Feature-level concatenation method, which employs NestedUNet [100] to extract multi-level bitemporal features. Channel attention and deep supervision are applied on each level of the decoder.
- BIT [53]. Feature-level differencing method, which expresses the input images into a few visual words (tokens), and models spatiotemporal context in the token-based space by transformers to efficiently benefit per-pixel representations.
- ICIFNet [47]. Feature-level differencing method, which integrates CNN and Transformer to parallelly extract multi-level bitemporal features. Cross-attention is applied to fuse parallel features at each level.
- DMINet [55]. Feature-level fusion method, which combines self-attention and cross-attention on bitemporal features of each level to perform temporal interactions, and uses both feature differencing and concatenation parallelly to obtain the change information. Deep supervision is also applied for better performance.
- SUNet [15]. Feature-space alignment method, which designs an asymmetric convolutional network in the early stage of the encoder to spatially align HR/LR images. Handcrafted edge maps for each bitemporal image are also fed into the model as auxiliary information. For a fair comparison, we implement it by upsampling the LR image to the size of the HR image to eliminate the loss of small targets.
- SRCDNet [10]. Image-space alignment method, which jointly optimizes a GAN-based image super-resolution model and a change detection model. For a fair com-

parison, due to the inaccessibility of the ground truth HR version for the LR image, we use the pair of HR images and their downsampled version to train the super-resolution model and apply it to the LR image to obtain the upsampled LR image in the inference phase.

We implement the above change detection models using their public codes with default hyperparameters. Note that for adapting these conventional CD methods to the cross-resolution CD task, we resize the LR image to the size of the HR image by cubic interpolation before feeding them into the CD model.

Table I reports the overall comparison results on the LEVIR-CD(4 \times), SV-CD(8 \times), and DE-CD(3.3 \times) test sets. In this setting, each compared model is tested by the bitemporal samples with fixed resolution difference ratios the same as in the training phase. Quantitative results show that our proposed method consistently outperforms the compared conventional CD methods as well as cross-resolution CD methods in terms of F1/IoU/OA scores across the three datasets. Note that as the pure convolutional-based methods (FC-EF, FC-Siam-Conc, and FC-Siam-Diff) fail to fit the DE-CD(3.3 \times) training set, therefore their performance scores are omitted.

Comparison with conventional CD methods. We can observe from Table I that the conventional change detection models with feeding image-space aligned bitemporal inputs by interpolating LR images to the size of HR images can be viewed as strong counterparts in the cross-resolution setting. For example, the recent transformer-based methods (e.g., BIT and ICIFNet) could yield competitive even superior performance over specially designed cross-resolution CD models (SUNet and SRCNet). It indicates that state-of-the-art conventional CD models can be effectively adapted to the cross-resolution change detection task via simple interpolation-based image-level alignment. Despite the common design of the model structure without sophisticated multi-scale feature fusion strategies (e.g., UNet-based incremental aggregation [22, 55, 68]), or transformer structures [47, 53], our proposed method with the MLP-based change decoder could surpass extant methods.

Comparison with cross-resolution CD methods. Quantitative results show that our proposed method significantly precedes the compared cross-resolution methods on the three datasets. Worth noting that our baseline is comparable or even superior to our counterparts. It indicates the effectiveness yet simpleness of our image-level alignment design that turns a naive CD model adapting cross-resolution scenarios.

Visual comparison. Fig. 6 illustrates the visual results of the compared change detection model on the LEVIR-CD(4 \times), SV-CD(8 \times), and DE-CD(3.3 \times) test sets under the fixed cross-resolution setting. We use different colors to denote TP (white), TN (black), FP (red), and FN (green). Note that results of some early pure-convolutional CD models (FC-EF, FC-Siam-Conc, and FC-Siam-Diff) are not included for a better view. We can observe that the proposed model could achieve better predictions across the three datasets. For instance, as shown in Fig. 6 (b) where three new build-ups appear on the left of the region, conventional CD methods are struggling to recognize these changes of interest due to their weak textures

caused by downsampling the post-event image. SUNet tends to overestimate the change areas resulting in relatively lower precision. Our method could yield relatively accurate results despite blurred regions that occurred changes. It may be due to our designed change decoder that learns implicitly the shape of changes by using dense coordinate querying an INR MLP, therefore recovering HR changes of interest even if given LR degraded inputs.

Comparison of model efficiency. To a fair comparison, all the models are trained and tested on a computing server equipped with a single NVIDIA RTX 3090 GPU. Table II reports the number of model parameters (Params.), floating-point operations per second (FLOPs), and GPU training time of each method. The input to the model has a size of $256 \times 256 \times 3$. The reported time corresponds to the duration required to complete one epoch of training on the LEVIR-CD dataset using a batch size of 8. The results show that the proposed method outperforms the recent DMINet and ICIFNet with smaller model parameters and less computational cost. From Table II and Table I, we can observe that the proposed method achieves significant accuracy improvement compared to our baseline while utilizing a modest increase in model parameters and maintaining acceptable computational consumption.

C. Handling Continuous Resolution Difference Ratios

To further verify the model adaptation ability for continuous cross-resolution conditions, we feed samples of varying resolution difference ratios (r_d) across bitemporal images into the CD model that are trained on a fixed difference ratio setting. For a fair comparison, we apply image-space alignment by interpolating the LR image to an HR reconstruction before feeding it to each CD model.

Let r_{d0} be the original resolution difference ratio of the training samples. r_{d0} equals 4, 8, and 3.3 for LEVIR-CD, SV-CD, and DE-CD datasets, respectively. Based on the resolution difference ratio in the validation/testing phase compared to that during training, we primarily have two settings: in-distribution and out-of-distribution settings. For simplicity, we denote values between 1 to r_{d0} as in-distribution ratios, and those larger than r_{d0} as out-of-distribution ratios. Given one HR bitemporal sample from LEVIR-CD and SV-CD datasets, we downsample the post-event HR image with different scales to obtain samples with varying ratios. For the real-world DE-CD dataset, because of the lack of real HR pre-event images, we downsample the post-event HR image for in-distribution conditions and further downsample the pre-event LR image for out-of-distribution conditions.

Table III, Table IV, and Table V report the cross-resolution performance of different CD models on the LEVIR-CD, SV-CD, and DE-CD test sets, respectively. Quantitative results show our proposed method not only consistently outperforms other methods in terms of F1/IoU scores across the three datasets in the in-distribution settings, but also exhibits significant advantages in the out-of-distribution settings.

We can observe that most of the methods achieve optimal results under a certain in-distribution ratio, while in the out-of-distribution setting, as the resolution difference ratio increases,

TABLE I

COMPARISON RESULTS ON THE THREE CD TEST SETS. THE BEST RESULTS ARE MARKED IN **BOLD**. ALL THE SCORES ARE DESCRIBED AS PERCENTAGES (%).

	LEVIR-CD(4×)					SV-CD(8×)					DE-CD(3.3×)				
	Pre. / Rec. / F1 / IoU / OA					Pre. / Rec. / F1 / IoU / OA					Pre. / Rec. / F1 / IoU / OA				
FC-EF [20]	79.57 / 71.48 / 75.31 / 60.40 / 97.64					74.25 / 45.03 / 56.06 / 38.95 / 91.67					-				
FC-Siam-Conc [20]	84.23 / 69.90 / 76.40 / 61.81 / 97.82					73.11 / 50.87 / 59.99 / 42.85 / 92.00					-				
FC-Siam-Diff [20]	86.12 / 60.15 / 70.83 / 54.83 / 97.50					76.10 / 56.68 / 64.97 / 48.12 / 92.79					-				
STANet [2]	57.87 / 45.47 / 50.93 / 34.16 / 95.58					83.06 / 70.73 / 76.41 / 61.82 / 94.85					11.80 / 46.89 / 18.85 / 10.41 / 97.61				
IFNet [22]	86.81 / 80.85 / 83.73 / 72.01 / 98.41					94.94 / 79.62 / 86.61 / 76.38 / 97.09					26.20 / 52.64 / 34.98 / 21.20 / 98.84				
SNUNet [68]	89.67 / 81.00 / 85.11 / 74.09 / 98.57					92.61 / 83.80 / 87.98 / 78.55 / 97.30					38.21 / 37.16 / 37.68 / 23.21 / 99.27				
BIT [53]	89.57 / 82.11 / 85.68 / 74.94 / 98.61					97.09 / 84.80 / 90.53 / 82.69 / 97.91					62.05 / 33.38 / 43.41 / 27.72 / 99.48				
ICIFNet [47]	87.84 / 84.62 / 86.20 / 75.75 / 98.63					95.68 / 90.56 / 93.05 / 87.00 / 98.40					63.50 / 25.04 / 35.92 / 21.89 / 99.47				
DMINet [55]	89.66 / 84.28 / 86.89 / 76.82 / 98.72					97.77 / 89.76 / 93.60 / 87.96 / 98.55					71.47 / 33.84 / 45.93 / 29.81 / 99.53				
SUNet [15]	64.12 / 93.54 / 76.08 / 61.40 / 97.03					63.55 / 97.98 / 77.10 / 62.73 / 93.13					32.60 / 71.00 / 44.68 / 28.77 / 98.96				
SRCDNet [10]	66.29 / 84.18 / 74.17 / 58.94 / 97.04					91.30 / 91.89 / 91.59 / 84.49 / 98.01					39.62 / 33.22 / 36.13 / 22.05 / 99.30				
Base	89.56 / 84.24 / 86.81 / 76.70 / 98.71					96.11 / 89.00 / 92.42 / 85.90 / 98.28					58.50 / 27.38 / 37.30 / 22.93 / 99.45				
Ours	90.55 / 86.30 / 88.38 / 79.18 / 98.86					95.29 / 93.36 / 94.32 / 89.24 / 98.67					61.35 / 42.32 / 50.10 / 33.42 / 99.50				

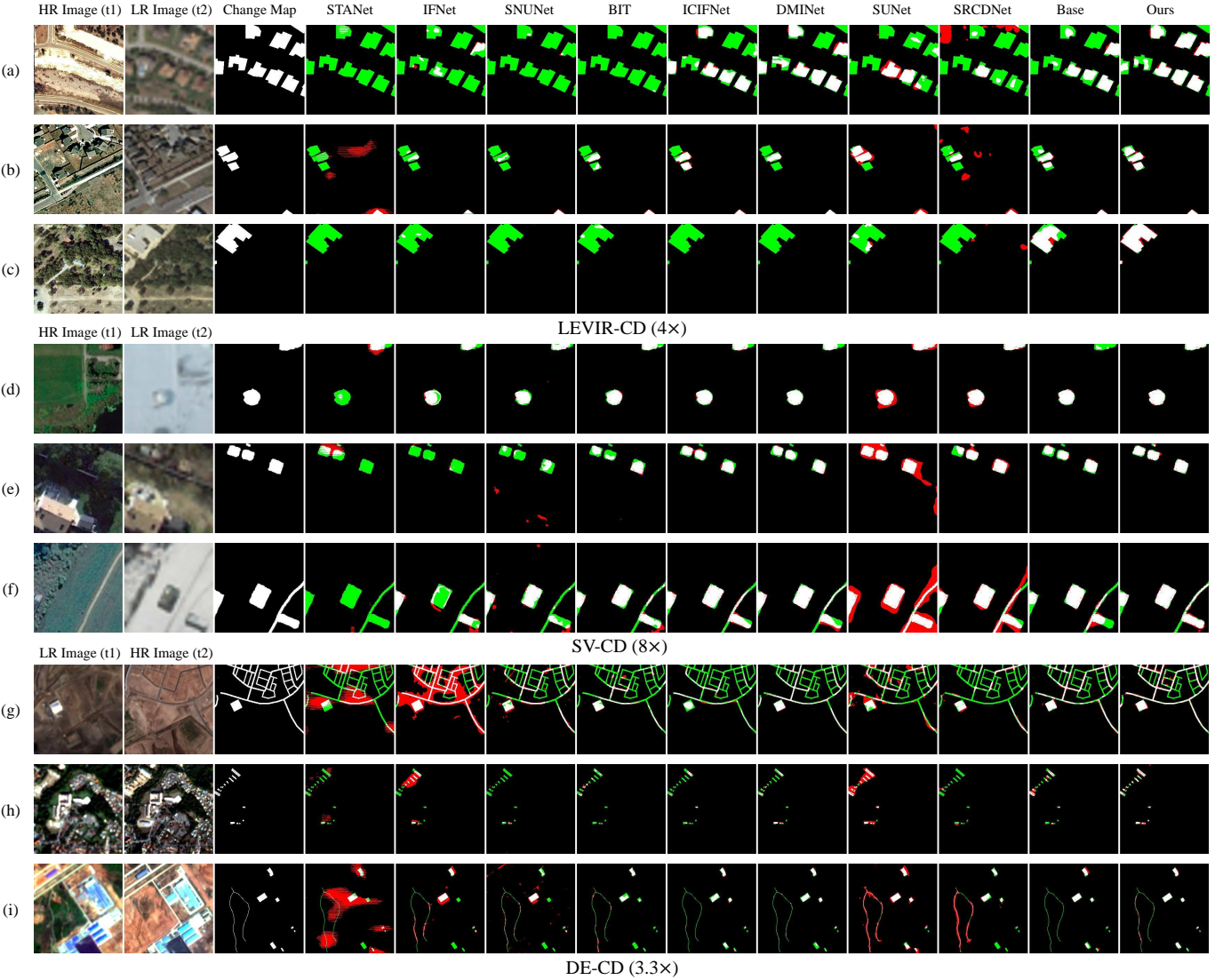


Fig. 6. Visual results of the compared methods on the three datasets. For a better view, we use white for true positive, black for true negative, red for false positive, and green for false negative.

TABLE II

COMPARISON OF MODEL EFFICIENCY. WE REPORT THE NUMBER OF MODEL PARAMETERS (PARAMS.), FLOATING-POINT OPERATIONS PER SECOND (FLOPS), AND TRAINING TIME FOR ONE EPOCH ON THE LEVIR-CD TRAINING SET. THE INPUT IMAGE TO THE MODEL HAS A SIZE OF $256 \times 256 \times 3$. THE BATCH SIZE IS SET TO 8.

Model	Params.(M)	FLOPs (G)	Training Time (s)
FC-EF [20]	1.35	7.14	52.33
FC-Siam-conc [20]	1.55	10.64	64.99
FC-Siam-diff [20]	1.35	9.44	63.21
STANet [2]	16.93	26.32	270.74
IFNet [22]	50.71	164.70	372.66
SNUNet [68]	12.03	109.76	305.82
BIT [53]	3.55	17.42	140.46
ICIFNet [47]	25.83	50.50	555.37
DMINet [55]	6.76	28.38	172.29
SUNet [15]	15.56	79.78	321.31
SRCDNet [10]	12.77	30.98	81.86
Base	11.97	11.42	53.18
Ours	13.06	17.5	103.81

the performance decreases. It is not surprising that the optimal ratio of most methods is less than the difference ratio of the training data. It is because these models train the Siamese encoder to adapt to both HR and LR data, which may result in a compromise of an in-between resolution. It is worth noting that the proposed method exhibits nearly consistent performance for each in-distribution setting, ranging from $1\times$ to $8\times$, on the SV-CD dataset, while our base model on the $1\times$ setting is much inferior to (i.e., 1.3 points of the F1 score drops) that on the $8\times$ setting. It may be attributed to our design of the scale-invariant learning framework as well as the change decoder which implicitly represents the detailed shape of land covers of interest. We can also observe that the proposed method achieves larger performance boosts compared to other methods in the case of out-of-distribution compared to in-distribution settings. For example, in the LEVIR-CD test set, our method significantly outperforms the counterpart (e.g., DMINet) by 16 points in terms of F1 score in the out-of-distribution ($8\times$) setting, compared to 1.3 points in the in-distribution ($4\times$) setting. Moreover, we can observe that some early approaches, e.g., FC-EF, FC-Siam-Conc, FC-Siam-Diff, and SUNet somehow exhibit relatively insufficient yet stable performance across different resolution differences. Some recent advanced CD methods such as DMINet and ICIFNet deliver promising performance in scenarios with small resolution differences but their performance declines significantly in cases of the large resolution difference settings (e.g., over 20 percent drops in terms of F1 score on the LEVIR-CD dataset of $8\times$ setting). It may be because these methods tend to overfit the known patterns and struggle to adapt to unseen ones. Overall, the proposed method demonstrates a balanced performance, consistently outperforming others across all cross-resolution settings.

To better illustrate the cross-resolution adaptability of our method, we display the F1-score curve of different models under varying resolution difference ratios on the LEVIR-CD, SV-CD, and DE-CD test sets in Fig. 7. We can observe that our method substantially shows more stability and better accuracy

than other methods.

Fig. 8, Fig. 9, and Fig. 10 also illustrate the visual results of compared models on these datasets with varying bitemporal resolution difference ratios. The visual comparison also verifies the cross-resolution adaptability of the proposed method. For instance, Fig. 9 shows some newly built ground facilities on the left side of the region. Our method can obtain consistent accurate predictions across varying difference ratios while most other compared methods fail to recognize the change of interest under the out-of-distribution ratios (e.g., $12\times$).

Apart from the setting of cross-resolution training/testing, i.e., the model is trained on samples with fixed resolution difference ratios and then validated on samples with different cross-resolution conditions, we also perform the setting of original-resolution training and cross-resolution testing, i.e., the model is trained on equal-resolution samples from the original CD training set and then validated on those with varying cross-resolution conditions.

Table VI reports the cross-resolution performance of different models on the LEVIR-CD dataset set. Each compared model is trained on the HR training samples with equal bitemporal resolution from the original LEVIR-CD dataset. In the training phase, we perform random downsampled reconstruction on the pre-event image by a ratio from Uniform distribution $r \sim U[1, 8]$. Similarly, we downsample the post-event HR image using different scales to obtain cross-resolution samples in the testing phase. Quantitative results show that the proposed method consistently outperforms other methods in terms of F1/IoU scores on testing samples with different cross-resolution ratios. We can observe from the results that most methods achieve the best results when the ratio is equal to 1, while the performance decreases when the ratio increases. For instance, DMINet exhibits comparable performance to our method when the ratio equals 1, but when the ratio increases to 8, its performance is dramatically dropped by nearly 90 percent, while our method could maintain acceptable performance. The results further indicate the cross-resolution adaptability of the proposed method.

D. Ablation Studies

We perform ablation experiments on the three critical components of the proposed methods, i.e., Random Resolution Synthesis (RRS), Implicit Neural Decoder (IND), and Bitemporal Local Interaction (BLI). We start from the baseline (Base) and incrementally supplement the above three components to evaluate their respective gains to the CD performance.

Table VII reports the ablation results of our method on the LEVIR-CD($4\times$), SV-CD($8\times$), and DE-CD($3.3\times$) test sets. The F1-score of each model is listed for comparison. Quantitative results show that the three components of SILI bring consistent performance improvements across different datasets.

Ablation on RRS. As shown in Table VII, compared to baseline, random resolution synthesis brings in significant improvements across the three datasets. It is not surprising because such a design can be viewed as a data augmentation approach, that synthesizes degraded reconstructions with various intrinsic resolutions. For the cross-resolution CD task, our

TABLE III

CROSS-RESOLUTION COMPARISON ON THE LEVIR-CD TEST SET WITH VARYING BITEMPORAL RESOLUTION DIFFERENCE RATIOS. WE SYNTHESIZE LR IMAGES BY DOWNSAMPLING POST-EVENT IMAGES. THE BEST RESULTS FOR EACH CROSS-RESOLUTION SETTING ARE MARKED IN **BOLD**. ALL THE SCORES ARE DESCRIBED AS PERCENTAGES (%). THESE MODELS ARE TRAINED ON THE SAMPLES FROM THE LEVIR-CD(4 \times) TRAINING SET.

	In-distribution testing (F1 / IoU)					Out-of-distribution testing (F1 / IoU)		
	1 \times	1.3 \times	2 \times	3 \times	4 \times	5 \times	6 \times	8 \times
FC-EF [20]	73.67 / 58.31	74.49 / 59.35	75.11 / 60.14	75.49 / 60.62	75.31 / 60.40	74.41 / 59.24	72.52 / 56.88	66.83 / 50.18
FC-Siam-Conc [20]	78.63 / 64.78	79.33 / 65.74	79.32 / 65.73	78.34 / 64.39	76.40 / 61.81	73.55 / 58.16	69.91 / 53.73	61.57 / 44.48
FC-Siam-Diff [20]	75.39 / 60.51	76.29 / 61.66	76.06 / 61.36	74.17 / 58.94	70.83 / 54.83	66.12 / 49.39	60.35 / 43.22	46.57 / 30.35
STANet [2]	42.03 / 26.61	47.96 / 31.54	55.75 / 38.64	57.77 / 40.62	50.93 / 34.16	36.49 / 22.32	17.22 / 9.42	3.97 / 2.02
IFNet [22]	74.24 / 59.03	77.77 / 63.62	81.02 / 68.10	83.48 / 71.65	83.73 / 72.01	82.10 / 69.64	78.39 / 64.47	66.24 / 49.52
SNUNet [68]	85.13 / 74.11	86.77 / 76.62	87.52 / 77.81	87.24 / 77.36	85.11 / 74.09	76.63 / 62.12	59.56 / 42.41	17.66 / 9.68
BIT [53]	86.28 / 75.86	86.42 / 76.09	86.66 / 76.47	86.67 / 76.48	85.68 / 74.94	81.06 / 68.16	70.40 / 54.33	28.73 / 16.78
ICIFNet [47]	86.24 / 75.80	86.44 / 76.11	86.78 / 76.65	86.84 / 76.75	86.20 / 75.75	83.63 / 71.87	78.95 / 65.22	59.26 / 42.10
DMINet [55]	86.28 / 75.87	86.49 / 76.20	86.85 / 76.75	86.96 / 76.93	86.89 / 76.82	83.78 / 72.08	79.10 / 65.43	57.40 / 40.26
SUNet [15]	75.51 / 60.65	75.53 / 60.69	75.67 / 60.86	75.98 / 61.27	76.08 / 61.40	75.70 / 60.90	75.01 / 60.02	69.96 / 53.80
SRCDNet [10]	75.87 / 61.12	76.46 / 61.89	76.77 / 62.30	76.30 / 61.68	74.17 / 58.94	69.38 / 53.12	60.13 / 42.99	29.23 / 17.12
Base	86.63 / 76.42	86.88 / 76.81	87.27 / 77.41	87.49 / 77.76	86.81 / 76.70	84.16 / 72.65	77.95 / 63.87	44.83 / 28.89
Ours	87.01 / 77.01	87.65 / 78.02	88.21 / 78.90	88.55 / 79.44	88.38 / 79.18	86.73 / 76.57	84.31 / 72.87	73.13 / 57.64

TABLE IV

CROSS-RESOLUTION COMPARISON ON THE SV-CD TEST SET WITH VARYING BITEMPORAL RESOLUTION DIFFERENCE RATIOS. WE SYNTHESIZE LR IMAGES BY DOWNSAMPLING POST-EVENT IMAGES. THE BEST RESULTS FOR EACH CROSS-RESOLUTION SETTING ARE MARKED IN **BOLD**. ALL THE SCORES ARE DESCRIBED AS PERCENTAGES (%). THESE MODELS ARE TRAINED ON THE SAMPLES FROM THE SV-CD(8 \times) TRAINING SET.

	In-distribution testing (F1 / IoU)					Out-of-distribution testing (F1 / IoU)		
	1 \times	2 \times	4 \times	5 \times	8 \times	9 \times	10 \times	12 \times
FC-EF [20]	55.88 / 38.77	55.88 / 38.77	55.96 / 38.85	55.99 / 38.88	56.06 / 38.95	56.06 / 38.94	56.06 / 38.95	56.03 / 38.92
FC-Siam-Conc [20]	64.31 / 47.39	64.52 / 47.62	63.97 / 47.03	63.35 / 46.36	59.99 / 42.85	58.39 / 41.23	56.84 / 39.70	53.88 / 36.87
FC-Siam-Diff [20]	67.35 / 50.77	67.41 / 50.84	67.09 / 50.48	66.86 / 50.22	64.97 / 48.12	63.54 / 46.56	61.84 / 44.76	58.51 / 41.35
STANet [2]	72.22 / 56.52	73.39 / 57.97	76.97 / 62.56	77.72 / 63.55	76.41 / 61.82	73.93 / 58.64	71.35 / 55.47	67.05 / 50.43
IFNet [22]	81.54 / 68.83	81.97 / 69.45	84.17 / 72.66	85.34 / 74.42	86.61 / 76.38	86.34 / 75.97	85.69 / 74.97	83.77 / 72.07
SNUNet [68]	75.08 / 60.10	79.29 / 65.69	87.77 / 78.21	89.62 / 81.20	87.98 / 78.55	85.54 / 74.74	83.45 / 71.60	80.04 / 66.73
BIT [53]	85.59 / 74.81	86.82 / 76.70	90.07 / 81.94	90.98 / 83.46	90.53 / 82.69	88.11 / 78.75	85.16 / 74.15	80.07 / 66.77
ICIFNet [47]	91.20 / 83.83	91.56 / 84.44	92.83 / 86.63	93.25 / 87.35	93.05 / 87.00	91.95 / 85.10	90.65 / 82.90	88.02 / 78.60
DMINet [55]	92.14 / 85.42	92.52 / 86.09	93.66 / 88.07	93.92 / 88.54	93.60 / 87.96	93.00 / 86.91	92.05 / 85.27	89.59 / 81.14
SUNet [15]	67.12 / 50.51	67.33 / 50.76	69.88 / 53.70	72.44 / 56.80	77.10 / 62.73	77.55 / 63.33	77.69 / 63.52	76.90 / 62.46
SRCDNet [10]	78.14 / 64.13	82.07 / 69.59	89.19 / 80.49	90.67 / 82.93	91.59 / 84.49	90.76 / 83.08	89.29 / 80.66	85.19 / 74.19
Base	91.12 / 83.69	91.47 / 84.28	92.51 / 86.06	92.87 / 86.68	92.42 / 85.90	90.64 / 82.88	88.09 / 78.71	82.32 / 69.96
Ours	94.07 / 88.80	94.11 / 88.88	94.26 / 89.14	94.30 / 89.22	94.32 / 89.24	93.55 / 87.87	92.80 / 86.57	90.50 / 82.65

TABLE V

CROSS-RESOLUTION COMPARISON ON THE DE-CD TEST SET WITH VARYING BITEMPORAL RESOLUTION DIFFERENCE RATIOS. FOR IN-DISTRIBUTION TESTING, WE SYNTHESIZE RELATIVELY HR IMAGES COMPARED TO REAL PRE-EVENT LR IMAGES BY DOWNSAMPLING POST-EVENT IMAGES. FOR OUT-OF-DISTRIBUTION TESTING, WE FURTHER DOWNSAMPLE PRE-EVENT IMAGES TO SYNTHESIZE LR IMAGES. THE BEST RESULTS FOR EACH CROSS-RESOLUTION SETTING ARE MARKED IN **BOLD**. ALL THE SCORES ARE DESCRIBED AS PERCENTAGES (%). THESE MODELS ARE TRAINED ON THE SAMPLES FROM THE DE-CD(3.3 \times) TRAINING SET.

	In-distribution testing (F1 / IoU)					Out-of-distribution testing (F1 / IoU)		
	1 \times	1.3 \times	2 \times	3 \times	3.3 \times	4 \times	5 \times	6 \times
STANet [2]	18.81 / 10.38	18.72 / 10.33	18.78 / 10.37	18.82 / 10.39	18.85 / 10.41	19.05 / 10.53	18.48 / 10.18	18.18 / 10.00
IFNet [22]	32.35 / 19.30	34.18 / 20.62	34.66 / 20.96	34.84 / 21.09	34.98 / 21.20	29.05 / 17.00	23.77 / 13.49	20.21 / 11.24
SNUNet [68]	32.49 / 19.40	35.76 / 21.77	37.43 / 23.02	37.69 / 23.22	37.68 / 23.21	27.88 / 16.20	22.68 / 12.79	18.45 / 10.16
BIT [53]	39.53 / 24.63	42.24 / 26.78	43.18 / 27.53	43.35 / 27.68	43.41 / 27.72	38.94 / 24.18	34.30 / 20.70	30.20 / 17.78
ICIFNet [47]	31.84 / 18.93	34.79 / 21.06	35.80 / 21.80	35.91 / 21.89	35.92 / 21.89	32.21 / 19.20	30.15 / 17.75	29.31 / 17.17
DMINet [55]	31.60 / 18.77	30.89 / 18.27	30.26 / 17.83	30.09 / 17.71	30.02 / 17.66	29.31 / 17.17	27.30 / 15.81	23.74 / 13.47
SUNet [15]	43.69 / 27.95	44.48 / 28.60	44.66 / 28.75	44.68 / 28.77	44.68 / 28.77	42.84 / 27.25	40.56 / 25.44	37.58 / 23.14
SRCDNet [10]	34.44 / 20.80	35.46 / 21.55	35.96 / 21.92	36.10 / 22.02	36.13 / 22.05	32.14 / 19.14	29.35 / 17.20	28.28 / 16.47
Base	32.96 / 19.73	35.97 / 21.93	37.09 / 22.76	37.28 / 22.91	37.30 / 22.93	36.02 / 21.97	33.56 / 20.16	30.65 / 18.10
Ours	47.88 / 31.48	49.86 / 33.21	50.23 / 33.54	50.15 / 33.54	50.09 / 33.42	48.45 / 31.97	45.45 / 29.41	41.71 / 26.35

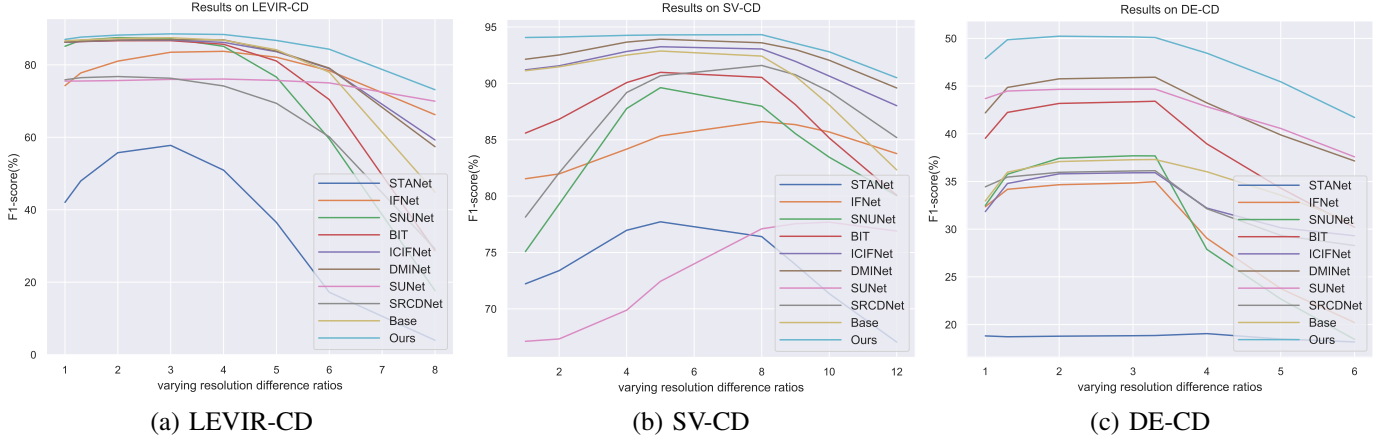


Fig. 7. F1-score comparison using varying bitemporal resolution difference ratios on the LEVIR-CD, SV-CD, and DE-CD test sets, respectively. The F1-score is reported.

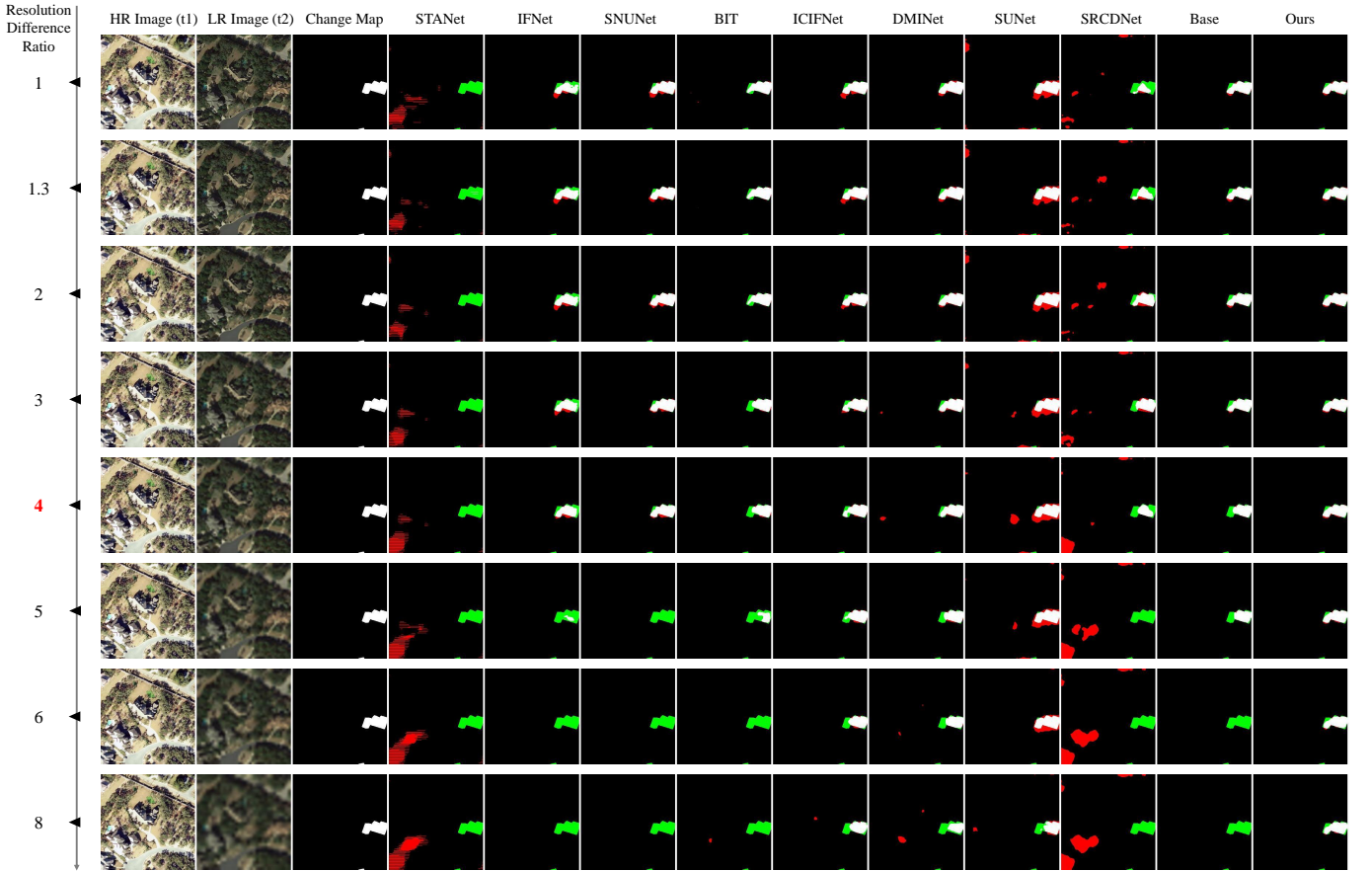


Fig. 8. Visual comparison of different methods on the LEVIR-CD test set with varying bitemporal resolution difference ratios. We synthesize LR images by downsampling with different scales the post-event image. For a better view, we use white for true positive, black for true negative, red for false positive, and green for false negative.

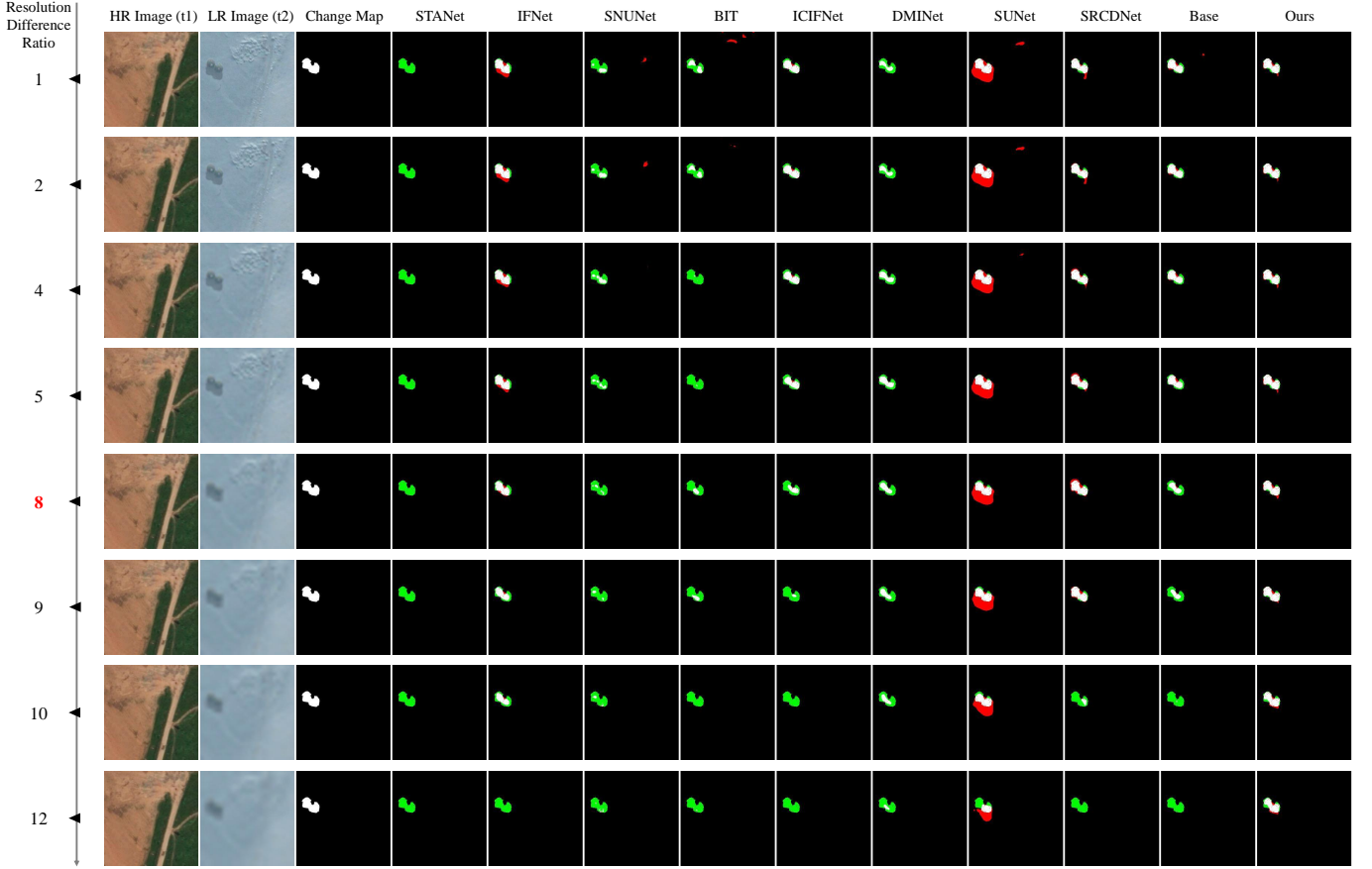


Fig. 9. Visual comparison of different methods on the SV-CD test set with varying bitemporal resolution difference ratios. We synthesize LR images by downsampling with different scales the post-event image. For a better view, we use white for true positive, black for true negative, red for false positive, and green for false negative.

TABLE VI
CROSS-RESOLUTION COMPARISON ON THE LEVIR-CD TEST SET. THESE MODELS ARE TRAINED ON THE SAMPLES FROM THE ORIGINAL ($1\times$) LEVIR-CD TRAINING SET AND ARE TESTED WITH SAMPLES WITH VARYING BITEMPORAL RESOLUTION DIFFERENCE RATIOS. WE SYNTHESIZE LR IMAGES BY DOWNSAMPLING POST-EVENT IMAGES. THE F1-SCORE AND IoU ARE REPORTED.

	$1\times$	$1.3\times$	$2\times$	$3\times$	$4\times$	$5\times$	$6\times$	$8\times$
FC-EF [20]	75.79 / 61.01	75.78 / 61.01	75.43 / 60.56	74.46 / 59.32	72.79 / 57.21	70.48 / 54.42	67.41 / 50.84	59.74 / 42.59
FC-Siam-Conc [20]	82.28 / 69.90	82.12 / 69.66	81.52 / 68.80	79.93 / 66.57	77.60 / 63.39	74.31 / 59.13	70.08 / 53.95	58.51 / 41.35
FC-Siam-Diff [20]	79.17 / 65.52	79.20 / 65.56	78.30 / 64.33	76.00 / 61.30	72.77 / 57.19	67.70 / 51.18	61.42 / 44.32	45.36 / 29.33
STANet [2]	87.27 / 77.41	86.70 / 76.53	85.14 / 74.12	72.79 / 57.22	40.43 / 25.34	16.45 / 8.96	9.10 / 4.77	7.00 / 3.62
IFNet [22]	88.11 / 78.75	88.06 / 78.67	87.42 / 77.65	85.31 / 74.39	80.95 / 68.00	73.35 / 57.91	59.29 / 42.13	21.01 / 11.74
SNUNet [68]	89.37 / 80.78	89.25 / 80.58	88.28 / 79.01	77.61 / 63.41	44.43 / 28.56	15.77 / 8.56	9.87 / 5.19	8.40 / 4.38
BIT [53]	88.54 / 79.44	88.57 / 79.48	88.23 / 78.94	86.03 / 75.49	78.66 / 64.83	61.00 / 43.88	33.69 / 20.26	7.84 / 4.08
ICIFNet [47]	88.20 / 78.89	88.18 / 78.86	87.92 / 78.44	86.48 / 76.19	82.01 / 69.50	72.62 / 57.01	54.95 / 37.88	15.92 / 8.65
DMINet [55]	89.56 / 81.09	89.45 / 80.91	89.01 / 80.20	87.29 / 77.45	82.64 / 70.42	70.60 / 54.55	48.92 / 32.38	11.01 / 5.83
SUNet [15]	78.32 / 64.37	78.21 / 64.22	77.94 / 63.85	77.41 / 63.14	76.69 / 62.19	75.42 / 60.53	73.40 / 57.97	61.83 / 44.75
SRCDDNet [10]	76.66 / 62.15	76.38 / 61.78	75.91 / 61.17	73.31 / 57.87	64.78 / 47.91	51.14 / 34.35	30.15 / 17.75	7.25 / 3.76
Base	88.63 / 79.59	88.60 / 79.53	88.23 / 78.93	86.29 / 75.88	79.25 / 65.63	57.95 / 40.79	22.08 / 12.41	0.64 / 0.32
Ours	89.70 / 81.33	89.67 / 81.27	89.24 / 80.58	88.14 / 78.80	85.79 / 75.11	81.65 / 68.99	75.11 / 60.14	65.17 / 48.33

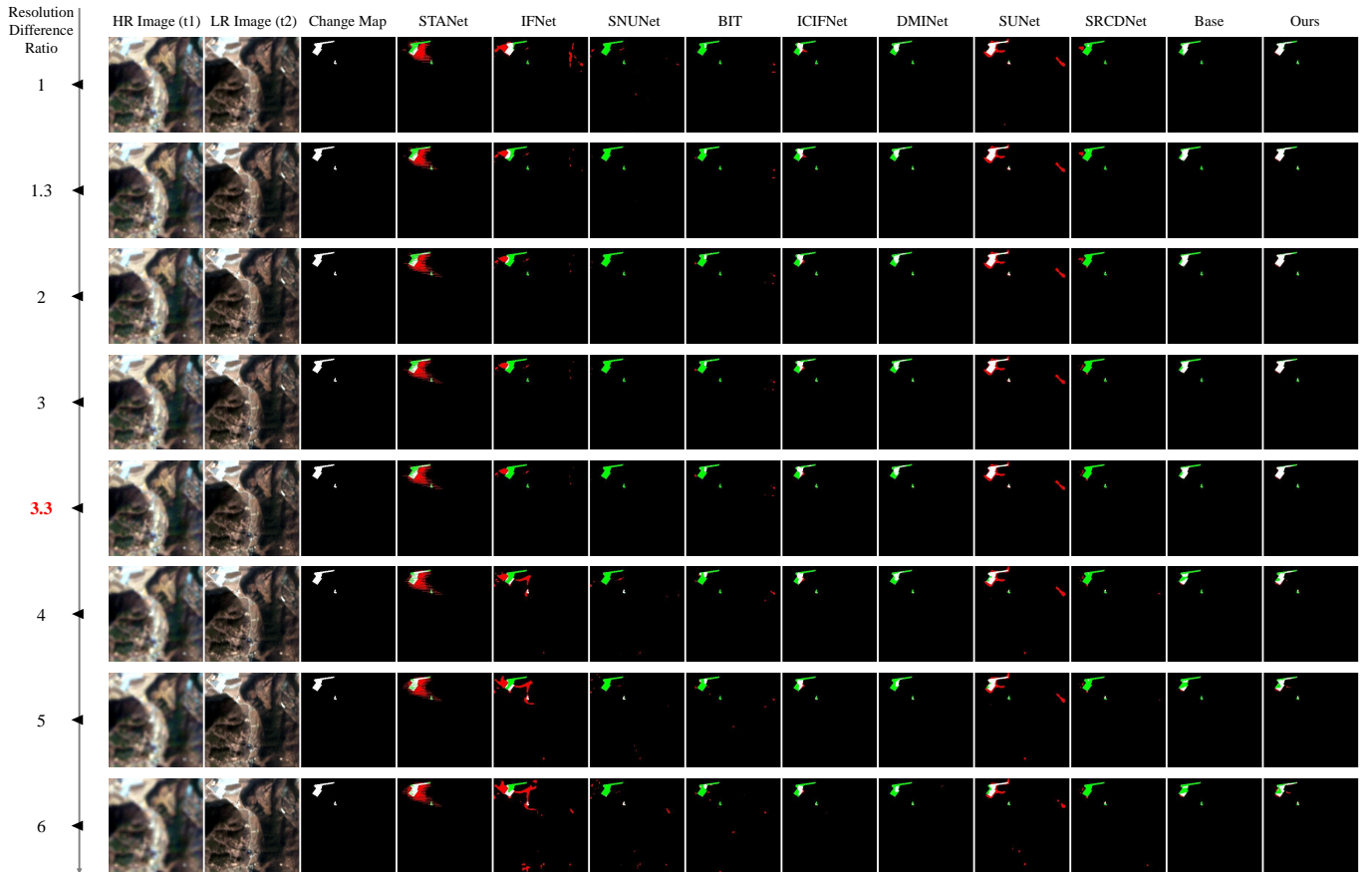


Fig. 10. Visual comparison of different methods on the DE-CD test set with varying bitemporal resolution difference ratios. Due to lacking real HR pre-event images, we synthesize relatively HR images by downsampling the post-event image, for in-distribution conditions. The real LR pre-event image is further downsampled for out-of-distribution testing. For a better view, we use white for true positive, black for true negative, red for false positive, and green for false negative.

data-level design synthesizing images with intermediate resolutions between low- and high-resolution inputs may benefit model learning by providing a progression of resolutions to reduce the resolution gap across bitemporal images.

Ablation on IND. We can observe from Table VII that our proposed INR-based change decoder could further consistently improve the baseline on the three datasets, especially on the DE(3.3 \times) dataset with relatively lower spatial resolution and with smaller pixel numbers per change area. It indicates the effectiveness of our IND in yielding the HR change mask from multi-level features, especially for recovering small objects of change. We further make a comparison to several conventional multi-level feature fusion approaches, including FPN [101] and UNet [102]. Those structures perform incremental aggregation from coarse to fine for the multi-level features (level 1 to level 4) from each temporal image. The concatenated bitemporal HR semantic features are then fed into three-layer convolutions for change classification, similar to our baseline. Quantitative results in Table VIII suggest the effectiveness of our IND for the cross-resolution CD task, compared with counterparts. Note that each model in Table VIII is trained with RRS for a fair comparison.

Ablation on BLI. Table VII demonstrates that our bitemporal local interaction produces consistent performance gains

across the three datasets. To further demonstrate the effectiveness of BLI, we also compare the commonly used global self-attention. For a fair comparison, we only replace the local-window self-attention in BLI with the global self-attention. Table X shows the comparison results on the three datasets. Note that here we only add bitemporal interaction on the features of level 1 from the encoder. Quantitative results show that our method consistently outperforms the self-attention counterpart, suggesting that local bitemporal interactions are more effective for the cross-resolution CD task. It indicates that modeling spatial-temporal correlations in the local regions between cross-resolution bitemporal images may be sufficient to align their semantic features.

E. Parametric Analysis

Effect of Random Bitemporal Region Swap. We propose to swap a random region between bitemporal images with different intrinsic spatial resolutions as a form of patch-level data augmentation to benefit the learning of scale-invariant features. The size of the swapped region, i.e., crop size, is an important hyperparameter. To explore the impact of crop size on CD performance, we perform ablation on different crop sizes for bitemporal region swapping. Our Base model is used as the baseline. Table IX reports the F1/IoU scores of

TABLE VII

ABLATION STUDY OF OUR SILI ON THREE CD DATASETS. ABLATIONS ARE PERFORMED ON THE RANDOM RESOLUTION SYNTHESIS (RRS), IMPLICIT NEURAL DECODER (IND), AND BITEMPORAL LOCAL INTERACTION (BLI). THE F1-SCORE IS REPORTED.

Model	RRS	IND	BLI	LEVIR(4×)	SV(8×)	DE(3.3×)
Base	×	×	×	86.81	92.42	37.30
SILI	✓	×	×	87.48	93.49	40.73
SILI	✓	✓	×	88.04	94.16	48.86
SILI	✓	✓	✓	88.38	94.32	50.17

TABLE VIII

EFFECT OF OUR INR-BASED CHANGE DECODER. WE REPLACE INR WITH SEVERAL OFF-THE-SHELL MULTI-LEVEL FEATURE FUSION APPROACHES FOR COMPARISON. THE F1/IOU SCORE OF EACH MODEL ON THREE CD DATASETS IS REPORTED.

Decoder	LEVIR(4×)	SV(8×)	DE(3.3×)
	F1 / IoU	F1 / IoU	F1 / IoU
FPN	85.96 / 75.37	93.31 / 87.45	39.48 / 24.60
UNet	87.78 / 78.22	93.41 / 87.63	43.99 / 28.20
MLP	87.48 / 77.74	93.49 / 87.77	40.73 / 25.58
INR	88.04 / 78.64	94.16 / 89.87	48.86 / 32.33

TABLE IX

EFFECT OF RANDOM BITEMPORAL REGION SWAP ON THREE CD DATASETS. WE ALSO PERFORM ABLATIONS ON THE SIZE OF THE SWAPPED REGION. THE F1/IOU SCORES OF EACH MODEL ARE REPORTED. NOTE THAT A CROP SIZE OF 0 DENOTES NOT PERFORMING REGION SWAP. A CROP SIZE OF 256 MEANS TO SWAP THE BITEMPORAL IMAGE, I.E., THE WHOLE REGION OF THE IMAGE. WE USE OUR BASE MODEL AS THE BASELINE.

Crop size	LEVIR(4×)	SV(8×)	DE(3.3×)
	F1 / IoU	F1 / IoU	F1 / IoU
0	86.81 / 76.70	92.42 / 85.90	37.30 / 22.93
64	87.13 / 77.19	92.81 / 86.58	37.96 / 23.43
128	87.15 / 77.23	92.94 / 86.81	38.15 / 23.57
192	87.04 / 77.06	92.86 / 86.66	37.49 / 23.43
256	86.89 / 76.82	92.69 / 86.37	36.94 / 22.65

TABLE X

EFFECT OF THE LOCAL-WINDOW ATTENTION IN THE BITEMPORAL FEATURE INTERACTION. WE REPLACE LOCAL ATTENTION WITH NON-LOCAL SELF-ATTENTION FOR COMPARISON. NOTE THAT WE ONLY APPLY INTERACTION ON THE BITEMPORAL FEATURES OF LEVEL 1 FROM THE ENCODER. THE F1/IOU SCORE OF EACH MODEL ON THREE CD DATASETS IS REPORTED.

interaction	LEVIR(4×)	SV(8×)	DE(3.3×)
	F1 / IoU	F1 / IoU	F1 / IoU
×	88.04 / 78.64	94.16 / 89.87	48.86 / 32.33
non-local	88.12 / 78.76	93.87 / 88.44	48.90 / 32.36
local	88.24 / 78.96	94.23 / 89.09	49.33 / 32.74

TABLE XI

EFFECT OF THE RESOLUTION OF DENSE COORDINATE QUERIES IN THE CHANGE DECODER ON THREE CD DATASETS. WE ALSO PERFORM ABLATIONS ON WHETHER TO INTRODUCE EDGE FEATURES. THE FLOPS AND F1/IOU SCORES OF EACH MODEL ARE REPORTED. NOTE THAT ds DENOTES THE DOWNSAMPLING RATE OF THE COORDINATE QUERY MAP RELATED TO THE ORIGINAL HR IMAGE.

Edge (/ds)	FLOPs (G)	LEVIR(4×)	SV(8×)	DE(3.3×)
		F1 / IoU	F1 / IoU	F1 / IoU
×(/4)	11.52	87.38 / 77.56	93.65 / 88.06	45.72 / 29.63
learn (/4)	11.54	87.60 / 77.93	93.81 / 88.33	45.25 / 29.24
✓(/4)	11.54	87.67 / 78.04	93.94 / 88.57	48.40 / 31.93
✓(/2)	17.20	88.04 / 78.64	94.16 / 89.87	48.86 / 32.33
✓(/1)	39.84	88.05 / 78.65	93.97 / 88.62	49.36 / 32.77

TABLE XII

EFFECT OF INTRODUCING BITEMPORAL INTERACTION AT DIFFERENT STAGES (FROM LEVEL 1 TO LEVEL 4) OF THE ENCODER. THE F1-SCORE OF EACH MODEL ON THREE CD DATASETS IS REPORTED.

1	2	3	4	LEVIR(4×)	SV(8×)	DE(3.3×)
×	×	×	×	88.04	94.16	48.86
✓	×	×	×	88.24	94.23	49.33
✓	✓	×	×	88.34	94.18	49.41
✓	✓	✓	×	88.38	94.32	50.10
✓	✓	✓	✓	88.34	94.35	49.75

compared models with different crop sizes. Note that the crop size of 0 denotes not applying the bitemporal region swap. The crop size of 256 means to swap the entire image in the temporal dimension, which is equivalent to not using region swap because bitemporal images do not interact with each other at the image level. Quantitative results show that the model with random region swap significantly outperforms the baseline. It indicates the effectiveness of the proposed random bitemporal region swap. This approach can be regarded as a form of patch-level data augmentation through the interaction of bitemporal information. Notably, the optimal results are attained with a crop size of 128, with a slight performance decrease observed as the crop size increases to 192. This reduction in performance with larger crop sizes is attributed to the increased likelihood of foreground land covers appearing at the swap area's edges, introducing truncated and incomplete land cover instances that can impede the model's learning process. Therefore, we set the crop size to 128.

Effect of the resolution of coordinate query map. Our INR-based change decoder uses dense coordinate queries alongside corresponding multi-level features to obtain the HR change mask. The spatial resolution of the coordinate query map is an important hyperparameter. Let ds be the down-sampling factor of the coordinate query map relative to the original HR image. Note that we directly bilinearly interpolate the relatively LR change prediction from the decoder to match the size of the HR ground truth when applying LR coordinate queries. Table XI reports the floating-point operations per second (FLOPs), and F1/IoU scores of compared models with different ds . Note that here we use our SILI model without BIL for experiments. From the last three rows of the table, we can observe that when the resolution of queries increases,

model performance on the three datasets improves overall, yet with higher computational complexity. For a trade-off between accuracy and efficiency, we set $ds = 2$. Additionally, we also verified the effectiveness of introducing edge clues. Quantitative results in Table XI manifest adding edge clues can consistently improve the model performance on the three datasets. To further validate the efficacy of incorporating handicraft edge features, we conduct a comparison between the models with and without these features. Note that we set up two baselines, the first baseline model (i.e., $\times/(4)$) does not receive any image edge features. The second baseline (i.e., learn $(/4)$) utilizes a learnable convolution layer to extract edge features from each temporal image and subsequently aggregate them to derive edge clues. For a fair comparison, the second baseline has the same amount of additional convolution parameters as our model does. Quantitative results in Table XI show that introducing additional edge features could consistently improve the CD performance in the three datasets. It indicates the effectiveness of the incorporation of handicraft edge features and learnable features, which has also been witnessed in some recent works [94–96]. It may be because the introduction of handicraft edge features could offer additional high-frequency information that may benefit network optimization.

Which stages to introducing BLI. We introduce BLI on bitemporal image features from a certain stage of the encoder. Here, we explore which stages to introduce bitemporal interactions. We choose our SILI model without any BLI as the baseline and incrementally add bitemporal interactions from level 1 to level 4. As shown in the table XII, as the number of bitemporal interactions increases, the performance of the model in terms of F1 score broadly progressively improves. Concretely, BLI brings in significant performance gains across the three datasets in the early stages of the encoder, while in the last stage (level 4), introducing BLI achieves relatively limited improvement, or even degrades the performance. It may be because the feature discrepancy caused by the difference in radiation and intrinsic resolution between bitemporal images could be better aligned by BLI during the early stages. Therefore, our SILI introduces interactions in stages of level 1/2/3.

F. Feature Visualization

Here, we provide an example to visualize multi-level features from our model to further demonstrate the effectiveness of introducing BLI. We use a popular feature visualization technique, class activation map (CAM) [103], to show what our model learns in each stage of the encoder. CAM is basically the channel-wise weighted sum of activation maps from a certain layer in the model. We visualize the last layer of each stage in the encoder.

Fig. 11 shows the CAM visualization of our models with or without BLI. Red denotes high values while blue denotes low values. The input sample is from LEVIR-CD ($4\times$) test set. We can observe from the CAM of each level that our model can concentrate on land covers on interest (building). Features from level 1 contain more spatial details, and those

from level 4 are more semantic information but lack location precision while the intermediate levels (2/3) provide a balanced representation that well localizes semantic elements. We can also observe that our method with BLI has similar intensities between bitemporal features of no-change regions. We further show feature difference maps, i.e. absolute subtraction between bitemporal unnormalized CAMs. We can observe that positions with high bitemporal difference values of our model are mainly distributed within the red box, while the model without BLI may exhibit large difference values (e.g., level 2/4) outside the red box where contains no changes. It suggests the effectiveness of BLI in aligning bitemporal semantic features and yielding relatively lower feature differences in regions of no change.

V. CONCLUSION

In this paper, we propose a scale-invariant method with implicit neural networks to achieve continuous cross-resolution RS image CD. The scale-invariant embedding space is learned by enforcing our model predicting the HR change mask given synthesized bitemporal images with random downsampling and region swapping. Dense coordinate queries and corresponding multi-level features are used for change recognition by an MLP that implicitly represents the shape of changes. Bitemporal local interaction is further introduced at early levels of the encoder to align bitemporal feature intensities regardless of resolution differences. Extensive experiments on two synthesized and one real-world cross-resolution CD datasets verify the effectiveness of the proposed method. Our SILI significantly outperforms several conventional CD methods and two specifically designed cross-resolution CD methods on the three datasets in both in-distribution and out-of-distribution settings. Our method could yield relatively consistent HR change predictions regardless of the resolution difference between bitemporal images. The empirical results suggest that our method could well handle varying bitemporal resolution difference ratios, towards real-world applications. Future works include, 1) exploring more effective scale-invariant change detection methods from the perspective of model architecture by incorporating scale-invariant network structures, rather than indirectly enhancing scale invariance through multiscale data augmentation, 2) investigating more advanced implicit neural representation techniques and their integration into the change detection task to achieve resolution-invariant change detection, 3) exploring the combination of various handcrafted features such as LBP, HOG, with deep learning models to evaluate their potential for improving CD performance.

REFERENCES

- [1] A. SINGH, “Review article digital change detection techniques using remotely-sensed data,” *International Journal of Remote Sensing*, vol. 10, no. 6, pp. 989–1003, 1989. [Online]. Available: <https://doi.org/10.1080/01431168908903939>
- [2] H. Chen and Z. Shi, “A spatial-temporal attention-based method and a new dataset for remote sensing image change detection,” *Remote Sensing*, vol. 12, no. 10, p. 1662, 2020.

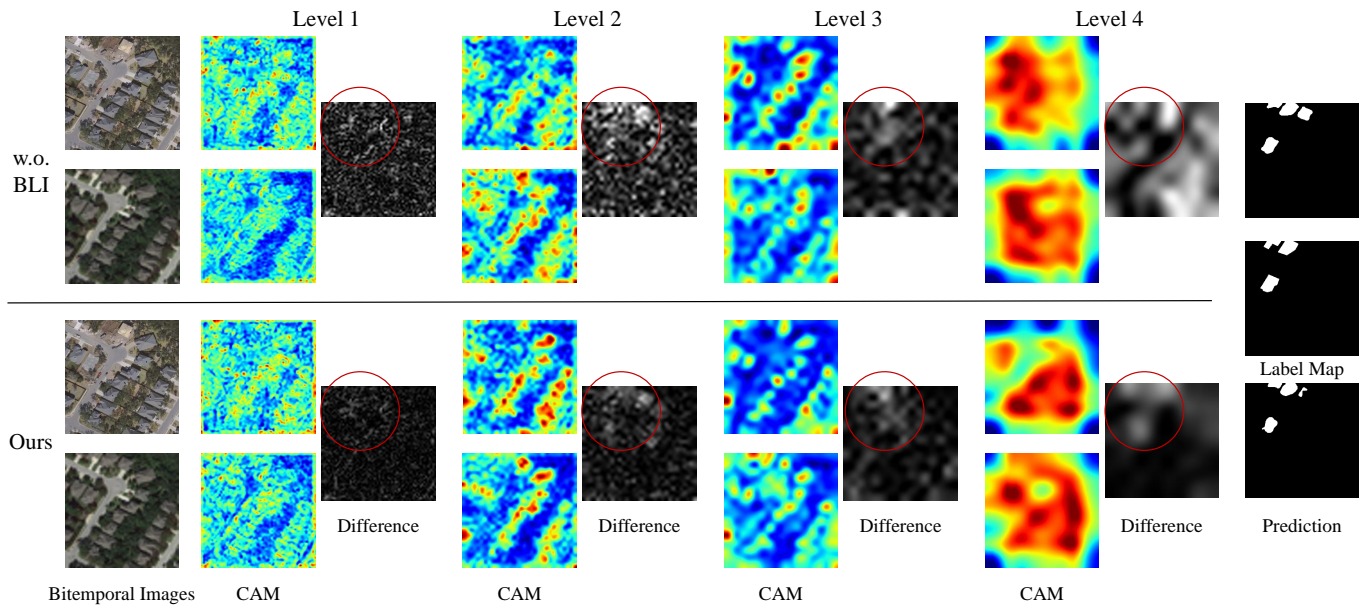


Fig. 11. Feature visualization of our models with or without bitemporal interactions. We show the class activation map (CAM) for each temporal image from level 1 to level 4. Bitemporal feature difference is also displayed to better show the effectiveness of introducing BLI. The input sample is from LEVIR-CD (4 \times) test set.

- [3] J. Z. Xu, W. Lu, Z. Li, P. Khaitan, and V. Zaytseva, "Building damage detection in satellite imagery using convolutional neural networks," 2019.
- [4] L. Bruzzone and D. Prieto, "Automatic analysis of the difference image for unsupervised change detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 38, no. 3, pp. 1171–1182, 2000.
- [5] P. P. de Bem, O. A. de Carvalho Junior, R. F. Guimarães, and R. A. T. Gomes, "Change detection of deforestation in the brazilian amazon using landsat data and convolutional neural networks," *Remote Sensing*, vol. 12, no. 6, p. 901, 2020.
- [6] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*. IEEE Computer Society, 2016, pp. 770–778.
- [7] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale."
- [8] W. Shi, M. Zhang, R. Zhang, S. Chen, and Z. Zhan, "Change detection based on artificial intelligence: State-of-the-art and challenges," *Remote Sensing*, vol. 12, no. 10, p. 1688, 5 2020.
- [9] J. Tian, D. Peng, H. Guan, and H. Ding, "RACDNet: Resolution- and alignment-aware change detection network for optical remote sensing imagery," *Remote Sensing*, vol. 14, no. 18, p. 4527, Jan. 2022, number: 18 Publisher: Multidisciplinary Digital Publishing Institute. [Online]. Available: <https://www.mdpi.com/2072-4292/14/18/4527>
- [10] M. Liu, Q. Shi, A. Marinoni, D. He, X. Liu, and L. Zhang, "Super-resolution-based change detection network with stacked attention module for images with different resolutions," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–18, 2022, conference Name: IEEE Transactions on Geoscience and Remote Sensing.
- [11] L. Wang, L. Wang, H. Wang, X. Wang, and L. Bruzzone, "SPCNet: A subpixel convolution-based change detection network for hyperspectral images with different spatial resolutions," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–14, 2022, conference Name: IEEE Transactions on Geoscience and Remote Sensing.
- [12] J. Tu, D. Li, W. Feng, Q. Han, and H. Sui, "Detecting damaged building regions based on semantic scene change from multi-temporal high-resolution remote sensing images," *ISPRS International Journal of Geo-Information*, vol. 6, no. 5, p. 131, May 2017, number: 5 Publisher: Multidisciplinary Digital Publishing Institute. [Online]. Available: <https://www.mdpi.com/2220-9964/6/5/131>
- [13] P. Zhang, M. Gong, L. Su, J. Liu, and Z. Li, "Change detection based on deep feature representation and mapping transformation for multi-spatial-resolution remote sensing images," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 116, pp. 24–41, Jun. 2016. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0924271616000563>
- [14] M. Liu, Q. Shi, J. Li, and Z. Chai, "Learning token-aligned representations with multimodal transformers for different-resolution change detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, 2022, conference Name: IEEE Transactions on Geoscience and Remote Sensing.
- [15] R. Shao, C. Du, H. Chen, and J. Li, "SUNet: Change detection for heterogeneous remote sensing images from satellite and UAV using a dual-channel fully convolution network," *Remote Sensing*, vol. 13, no. 18, p. 3750, Jan. 2021, number: 18 Publisher: Multidisciplinary Digital Publishing Institute. [Online]. Available: <https://www.mdpi.com/2072-4292/13/18/3750>
- [16] X. Zheng, X. Chen, X. Lu, and B. Sun, "Unsupervised change detection by cross-resolution difference learning," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–16, 2022, conference Name: IEEE Transactions on Geoscience and Remote Sensing.
- [17] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: representing scenes as neural radiance fields for view synthesis," *Commun. ACM*, vol. 65, no. 1, pp. 99–106, 2022.

- [18] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*. IEEE Computer Society, 2015, pp. 3431–3440.
- [19] D. Peng, Y. Zhang, and H. Guan, "End-to-end change detection for high resolution satellite images using improved unet++," *Remote Sensing*, vol. 11, no. 11, p. 1382, 2019.
- [20] R. Caye 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, oct 2018, pp. 4063–4067. [Online]. Available: https://github.com/rcdaudt/fully_convolutional_change_detection
- [21] B. Hou, Q. Liu, H. Wang, and Y. Wang, "From w-net to cdgan: Bitemporal change detection via deep learning techniques," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 58, no. 3, pp. 1790–1802, 2020.
- [22] C. Zhang, P. Yue, D. Tapete, L. Jiang, B. Shangguan, L. Huang, and G. Liu, "A deeply supervised image fusion network for change detection in high resolution bi-temporal remote sensing images," *ISPRS*, vol. 166, pp. 183–200, 2020.
- [23] Y. Liu, C. Pang, Z. Zhan, X. Zhang, and X. Yang, "Building change detection for remote sensing images using a dual-task constrained deep siamese convolutional network model," *IEEE Geoscience and Remote Sensing Letters*, pp. 1–5, 2020.
- [24] M. Papadomanolaki, S. Verma, M. Vakalopoulou, S. Gupta, and K. Karantzalos, "Detecting urban changes with recurrent neural networks from multitemporal sentinel-2 data," in *2019 IEEE International Geoscience and Remote Sensing Symposium, IGARSS 2019, Yokohama, Japan, July 28 - August 2, 2019*. IEEE, 2019, pp. 214–217.
- [25] G. Pei and L. Zhang, "Feature hierarchical differentiation for remote sensing image change detection," *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2022.
- [26] Q. Li, R. Zhong, X. Du, and Y. Du, "TransUNetCD: A Hybrid Transformer Network for Change Detection in Optical Remote Sensing Images," *IEEE Trans. Geosci. Remote Sens.*, p. 1, Apr. 2022.
- [27] H. Chen, W. Li, and Z. Shi, "Adversarial instance augmentation for building change detection in remote sensing images," *IEEE Transactions on Geoscience and Remote Sensing*, pp. 1–16, 2021.
- [28] M. Zhang and W. Shi, "A feature difference convolutional neural network-based change detection method," *TGRS*, vol. 58, no. 10, pp. 1–15, oct 2020.
- [29] T. Bao, C. Fu, T. Fang, and H. Huo, "Ppcnet: A combined patch-level and pixel-level end-to-end deep network for high-resolution remote sensing image change detection," vol. PP, pp. 1–5, 2020.
- [30] H. Jiang, X. Hu, K. Li, J. Zhang, J. Gong, and M. Zhang, "Pga-siamnet: Pyramid feature-based attention-guided siamese network for remote sensing orthoimagery building change detection," *Remote Sensing*, vol. 12, no. 3, p. 484, 2020.
- [31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*. I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, Eds., 2017, pp. 5998–6008. [Online]. Available: <http://papers.nips.cc/paper/7181-attention-is-all-you-need>
- [32] B. Guo, X. Zhang, H. Wu, Y. Wang, Y. Zhang, and Y.-F. Wang, "Lar-sr: A local autoregressive model for image super-resolution," in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 1899–1908.
- [33] Z. Liang, B. Zhu, and Y. Zhu, "High resolution representation-based Siamese network for remote sensing image change detection," *IET Image Proc.*, vol. n/a, no. n/a, Apr. 2022.
- [34] Z. Cao, M. Wu, R. Yan, F. Zhang, and X. Wan, "Detection of small changed regions in remote sensing imagery using convolutional neural network," vol. 502. IOP Publishing, jun 2020, p. 012017. [Online]. Available: <https://doi.org/10.1088%2F1755-1315%2F502%2F1%2F012017>
- [35] W. G. C. Bandara and V. M. Patel, "A transformer-based siamese network for change detection," in *IGARSS 2022 - 2022 IEEE International Geoscience and Remote Sensing Symposium*. IEEE, jul 2022, pp. 207–210.
- [36] C. Zhang, L. Wang, S. Cheng, and Y. Li, "Swinsunet: Pure transformer network for remote sensing image change detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, 2022.
- [37] J. Jiang, J. Xiang, E. Yan, Y. Song, and D. Mo, "Forest-cd: Forest change detection network based on VHR images," *IEEE Geoscience Remote Sensing Letter*, vol. 19, pp. 1–5, 2022.
- [38] Z. Chen, Y. Zhou, B. Wang, X. Xu, N. He, S. Jin, and S. Jin, "Edge-net: A building change detection method for high-resolution remote sensing imagery based on edge guidance and differential enhancement," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 191, pp. 203–222, sep 2022.
- [39] M. Zhang, G. Xu, K. Chen, M. Yan, and X. Sun, "Triplet-based semantic relation learning for aerial remote sensing image change detection," *IEEE Geosci. Remote. Sens. Lett.*, vol. 16, no. 2, pp. 266–270, 2019.
- [40] Z. Lv, F. Wang, G. Cui, J. A. Benediktsson, T. Lei, and W. Sun, "Spatial-spectral attention network guided with change magnitude image for land cover change detection using remote sensing images," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–12, 2022.
- [41] G. Cheng, G. Wang, and J. Han, "Isnet: Towards improving separability for remote sensing image change detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–1, 2022.
- [42] X. Song, Z. Hua, and J. Li, "Remote sensing image change detection transformer network based on dual-feature mixed attention," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–16, 2022.
- [43] X. Peng, R. Zhong, Z. Li, and Q. Li, "Optical remote sensing image change detection based on attention mechanism and image difference," *IEEE Transactions on Geoscience and Remote Sensing*, pp. 1–12, 2020.
- [44] A. Raza, H. Huo, and T. Fang, "EUNet-CD: Efficient UNet++ for Change Detection of Very High-Resolution Remote Sensing Images," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, Jan. 2022.
- [45] Z. Li, C. Tang, L. Wang, and A. Y. Zomaya, "Remote sensing change detection via temporal feature interaction and guided refinement," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–11, 2022.
- [46] J. Chen, Z. Yuan, J. Peng, L. Chen, H. Huang, J. Zhu, Y. Liu, and H. Li, "Dasnet: Dual attentive fully convolutional siamese networks for change detection in high-resolution satellite images," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 14, pp. 1194–1206, 2021.
- [47] Y. Feng, H. Xu, J. Jiang, H. Liu, and J. Zheng, "ICIF-Net: Intra-Scale Cross-Interaction and Inter-Scale Feature Fusion Network for Bitemporal Remote Sensing Images Change Detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–13, Apr. 2022.
- [48] F. Song, S. Zhang, T. Lei, Y. Song, and Z. Peng, "Mstdsnet-cd: Multiscale swin transformer and deeply supervised network for change detection of the fast-growing urban regions," *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2022.
- [49] J. Yuan, L. Wang, and S. Cheng, "Stransunet: A siamese transunet-based remote sensing image change detection net-

- work,” *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 9241–9253, 2022.
- [50] L. Wan, Y. Tian, W. Kang, and L. Ma, “D-tnet: Category-awareness based difference-threshold alternative learning network for remote sensing image change detection,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. PP, pp. 1–16, 2022.
- [51] M. Liu, Q. Shi, Z. Chai, and J. Li, “Pa-former: Learning prior-aware transformer for remote sensing building change detection,” *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2022.
- [52] Y. Zhou, F. Wang, J. Zhao, R. Yao, S. Chen, and H. Ma, “Spatial-temporal based multihead self-attention for remote sensing image change detection,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 10, pp. 6615–6626, oct 2022.
- [53] H. Chen, Z. Qi, and Z. Shi, “Remote sensing image change detection with transformers,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–14, 2022.
- [54] Z. Li, C. Yan, Y. Sun, and Q. Xin, “A densely attentive refinement network for change detection based on very-high-resolution bitemporal remote sensing images,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–18, 2022.
- [55] Y. Feng, J. Jiang, H. Xu, and J. Zheng, “Change detection on remote sensing images using dual-branch multilevel intertemporal network,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–15, 2023.
- [56] S. Fang, K. Li, and Z. Li, “Changer: Feature interaction is what you need for change detection,” *CoRR*, vol. abs/2209.08290, 2022.
- [57] Y. Sun, X. Zhang, J. Huang, H. Wang, and Q. Xin, “Fine-grained building change detection from very high-spatial-resolution remote sensing images based on deep multitask learning,” *IEEE Geoscience and Remote Sensing Letters*, pp. 1–5, 2020.
- [58] M. Papadomanolaki, M. Vakalopoulou, and K. Karantzalos, “A deep multitask learning framework coupling semantic segmentation and fully convolutional lstm networks for urban change detection,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 9, pp. 1–18, sep 2021.
- [59] R. Caye Daudt, B. Le Saux, A. Boulch, and Y. Gousseau, “Multitask learning for large-scale semantic change detection,” *Computer Vision and Image Understanding*, vol. 187, p. 102783, 2019. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1077314219300992>
- [60] J. Lei, Y. Gu, W. Xie, Y. Li, and Q. Du, “Boundary extraction constrained siamese network for remote sensing image change detection,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–1, 2022.
- [61] B. Bai, W. Fu, T. Lu, and S. Li, “Edge-Guided Recurrent Convolutional Neural Network for Multitemporal Remote Sensing Image Building Change Detection,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, Aug. 2021.
- [62] W. Zhao, L. Mou, J. Chen, Y. Bo, and W. J. Emery, “Incorporating metric learning and adversarial network for seasonal invariant change detection,” *IEEE Trans. Geosci. Remote. Sens.*, vol. 58, no. 4, pp. 2720–2731, 2020. [Online]. Available: <https://ieeexplore.ieee.org/document/8937747>
- [63] J. Liu, W. Xuan, Y. Gan, Y. Zhan, J. Liu, and B. Du, “An End-to-end Supervised Domain Adaptation Framework for Cross-Domain Change Detection,” *Pattern Recognition*, vol. 132, p. 108960, Dec. 2022.
- [64] F. U. Rahman, B. Vasu, J. V. Cor, J. Kerekes, and A. E. Savakis, “Siamese network with multi-level features for patch-based change detection in satellite imagery,” in *2018 IEEE Global Conference on Signal and Information Processing, GlobalSIP 2018, Anaheim, CA, USA, November 26-29, 2018*. IEEE, 2018, pp. 958–962.
- [65] O. Manas, A. Lacoste, X. G. i Nieto, D. Vazquez, and P. Rodriguez, “Seasonal contrast: Unsupervised pre-training from uncurated remote sensing data,” in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. Los Alamitos, CA, USA: IEEE Computer Society, oct 2021, pp. 9394–9403. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1109/ICCV48922.2021.00928>
- [66] Y. Wang, N. A. A. Braham, Z. Xiong, C. Liu, C. M. Albrecht, and X. X. Zhu, “SSL4EO-S12: A large-scale multi-modal, multi-temporal dataset for self-supervised learning in earth observation,” *CoRR*, vol. abs/2211.07044, 2022.
- [67] H. Chen, W. Li, S. Chen, and Z. Shi, “Semantic-aware dense representation learning for remote sensing image change detection,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–18, 2022.
- [68] S. Fang, K. Li, J. Shao, and Z. Li, “Snunet-cd: A densely connected siamese network for change detection of vhr images,” *IEEE Geoscience and Remote Sensing Letters*, pp. 1–5, 2021.
- [69] W. Zhao, X. Chen, X. Ge, and J. Chen, “Using adversarial network for multiple change detection in bitemporal remote sensing imagery,” *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2020.
- [70] H. Zhou, Y. Ren, Q. Li, J. Yin, and Y. Lin, “Masnet: Improve performance of siamese networks with mutual-attention for remote sensing change detection tasks,” *CoRR*, vol. abs/2206.02331, 2022.
- [71] F. I. Diakogiannis, F. Waldner, and P. Caccetta, “Looking for change? roll the dice and demand attention.”
- [72] J. Pan, W. Cui, X. An, X. Huang, H. Zhang, S. Zhang, R. Zhang, X. Li, W. Cheng, and Y. Hu, “MapsNet: Multi-level feature constraint and fusion network for change detection,” *Int. J. Appl. Earth Obs. Geoinf.*, vol. 108, p. 102676, Apr. 2022.
- [73] W. Wang, X. Tan, P. Zhang, and X. Wang, “A cbam based multiscale transformer fusion approach for remote sensing image change detection,” *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 6817–6825, 2022.
- [74] N. Shi, K. Chen, and G. Zhou, “A divided spatial and temporal context network for remote sensing change detection,” *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 4897–4908, 2022.
- [75] X. Song, Z. Hua, and J. Li, “Pstnet: Progressive sampling transformer network for remote sensing image change detection,” *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 8442–8455, 2022.
- [76] Q. Ke and P. Zhang, “Hybrid-TransCD: A Hybrid Transformer Remote Sensing Image Change Detection Network via Token Aggregation,” *ISPRS International Journal Geo Information*, vol. 11, no. 4, p. 263, Apr. 2022.
- [77] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows.” Montreal, QC, Canada: IEEE, 2021, pp. 9992–10002.
- [78] T. Yan, Z. Wan, and P. Zhang, “Fully transformer network for change detection of remote sensing images,” *CoRR*, vol. abs/2210.00757, 2022.
- [79] V. Sitzmann, J. N. P. Martel, A. W. Bergman, D. B. Lindell, and G. Wetzstein, “Implicit neural representations with periodic activation functions,” in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., 2020. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/hash/53c04118df112c13a8c34b38343b9c10-Abstract.html>
- [80] K. Genova, F. Cole, D. Vlasic, A. Sarna, W. T. Freeman, and

- T. A. Funkhouser, "Learning shape templates with structured implicit functions," in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, 2019, pp. 7153–7163.
- [81] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," in *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part I*, ser. Lecture Notes in Computer Science, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12346. Springer, 2020, pp. 405–421. [Online]. Available: https://doi.org/10.1007/978-3-030-58452-8_24
- [82] Y. Chen, S. Liu, and X. Wang, "Learning continuous image representation with local implicit image function," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*. Computer Vision Foundation / IEEE, 2021, pp. 8628–8638. [Online]. Available: https://openaccess.thecvf.com/content/CVPR2021/html/Chen_Learning_Continuous_Image_Representation_With_Local_Implicit_Image_Function_CVPR_2021_paper.html
- [83] X. Xu, Z. Wang, and H. Shi, "Ultraser: Spatial encoding is a missing key for implicit image function-based arbitrary-scale super-resolution," *CoRR*, vol. abs/2103.12716, 2021. [Online]. Available: <https://arxiv.org/abs/2103.12716>
- [84] T. Shen, Y. Zhang, L. Qi, J. Kuen, X. Xie, J. Wu, Z. Lin, and J. Jia, "High quality segmentation for ultra high-resolution images," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 2022, pp. 1300–1309.
- [85] H. Hu, Y. Chen, J. Xu, S. Borse, H. Cai, F. Porikli, and X. Wang, "Learning implicit feature alignment function for semantic segmentation," in *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XXIX*, ser. Lecture Notes in Computer Science, S. Avidan, G. J. Brostow, M. Cissé, G. M. Farinella, and T. Hassner, Eds., vol. 13689. Springer, 2022, pp. 487–505. [Online]. Available: https://doi.org/10.1007/978-3-031-19818-2_28
- [86] B. Cheng, O. Parkhi, and A. Kirillov, "Pointly-supervised instance segmentation," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 2022, pp. 2607–2616.
- [87] Y. Wu, Z. Zou, and Z. Shi, "Remote sensing novel view synthesis with implicit multiplane representations," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, 2022.
- [88] Z. Qi, H. Chen, C. Liu, Z. Shi, and Z. Zou, "Implicit ray-transformers for multi-view remote sensing image segmentation," *CoRR*, vol. abs/2303.08401, 2023.
- [89] Z. Qi, Z. Zou, H. Chen, and Z. Shi, "Remote-sensing image segmentation based on implicit 3-d scene representation," *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2022.
- [90] L. Liu, Z. Zou, and Z. Shi, "Hyperspectral remote sensing image synthesis based on implicit neural spectral mixing models," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–14, 2023.
- [91] K. Chen, W. Li, S. Lei, J. Chen, X. Jiang, Z. Zou, and Z. Shi, "Continuous remote sensing image super-resolution based on context interaction in implicit function space," *CoRR*, vol. abs/2302.08046, 2023.
- [92] J. Luo, L. Han, X. Gao, X. Liu, and W. Wang, "Sr-feinr: Continuous remote sensing image super-resolution using feature-enhanced implicit neural representation," *Sensors*, vol. 23, no. 7, 2023. [Online]. Available: <https://www.mdpi.com/1424-8220/23/7/3573>
- [93] K. Chen, W. Li, J. Chen, Z. Zou, and Z. Shi, "Resolution-agnostic remote sensing scene classification with implicit neural representations," *IEEE Geoscience and Remote Sensing Letters*, vol. 20, pp. 1–5, 2023.
- [94] H. Zheng, M. Gong, T. Liu, F. Jiang, T. Zhan, D. Lu, and M. Zhang, "Hfa-net: High frequency attention siamese network for building change detection in vhr remote sensing images," *Pattern Recognition*, vol. 129, p. 108717, 2022.
- [95] Y. Shangguan, J. Li, and Z. Hua, "Contour-enhanced densely connected siamese network for change detection," *Journal of Applied Remote Sensing*, vol. 17, no. 1, pp. 016515–016515, 2023.
- [96] H. Liu, Z. Hu, Q. Ding, and X. Chen, "Idan: Image difference attention network for change detection," *arXiv preprint arXiv:2208.08292*, 2022.
- [97] J. Canny, "A computational approach to edge detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PAMI-8, no. 6, pp. 679–698, 1986.
- [98] M. A. Lebedev, Y. V. Vizilter, O. V. Vygolov, V. A. Knyaz, and A. Y. Rubis, "Change detection in remote sensing images using conditional adversarial networks," vol. XLII-2. Copernicus GmbH, may 2018, pp. 565–571.
- [99] A. Toker, L. Kondmann, M. Weber, M. Eisenberger, A. Camero, J. Hu, A. P. Hoderlein, Ç. Senaras, T. Davis, D. Cremers, G. Marchisio, X. X. Zhu, and L. Leal-Taixé, "Dynamicearthnet: Daily multi-spectral satellite dataset for semantic change segmentation," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 2022, pp. 21 126–21 135.
- [100] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang, "Unet++: A nested u-net architecture for medical image segmentation," in *Deep Learning in Medical Image Analysis - and - Multimodal Learning for Clinical Decision Support - 4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings*, ser. Lecture Notes in Computer Science, D. Stoyanov, Z. Taylor, G. Carneiro, T. F. Syeda-Mahmood, A. L. Martel, L. Maier-Hein, J. M. R. S. Tavares, A. P. Bradley, J. P. Papa, V. Belagiannis, J. C. Nascimento, Z. Lu, S. Conjeti, M. Moradi, H. Greenspan, and A. Madabhushi, Eds., vol. 11045. Springer, 2018, pp. 3–11. [Online]. Available: https://doi.org/10.1007/978-3-030-00889-5_1
- [101] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature Pyramid Networks for Object Detection," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Jul. 2017, pp. 936–944.
- [102] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention - MICCAI 2015 - 18th International Conference Munich, Germany, October 5 - 9, 2015, Proceedings, Part III*, ser. Lecture Notes in Computer Science, N. Navab, J. Hornegger, W. M. W. III, and A. F. Frangi, Eds., vol. 9351. Springer, 2015, pp. 234–241. [Online]. Available: https://doi.org/10.1007/978-3-319-24574-4_28
- [103] H. Wang, Z. Wang, M. Du, F. Yang, Z. Zhang, S. Ding, P. Mardziel, and X. Hu, "Score-cam: Score-weighted visual explanations for convolutional neural networks," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. Los Alamitos, CA, USA: IEEE Computer Society, jun 2020, pp. 111–119. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1109/CVPRW50498.2020.00020>