

SDC-GAE: Structural Difference Compensation Graph Autoencoder for Unsupervised Multimodal Change Detection

Te Han^{id}, Yuqi Tang^{id}, Yuzeng Chen^{id}, Xin Yang, Yuqiang Guo, and Shujing Jiang

Abstract—Multimodal change detection (MCD) is a crucial technology for applications in natural resource monitoring, disaster assessment, and urban planning. To address the reliance on labeled data and enhance the robustness of structural features in the existing methods, we propose a structure difference compensation graph autoencoder (SDC-GAE) for unsupervised MCD. It is recognized that the registered multimodal images exhibit consistency in structural features in unchanged areas, while the structural features in changed areas are distinct. SDC-GAE utilizes a graph convolutional network (GCN) to extract deep structural features from multimodal images. It uses the structural features of one time-phase image to reconstruct its spectral features in the spectral feature space of the target image. Through structural difference compensation, SDC-GAE learns the structural disparities between different images, with the compensation value directly reflecting the intensity of the changes. The SDC-GAE loss function consists of three components: image reconstruction loss, which evaluates the spectral feature discrepancy between the reconstructed and target images, guiding the model to reduce these differences via structural difference compensation; sparse constraint loss, which accounts for the fact that changes are typically confined to a few areas, ensuring the sparsity of the detected changes; and structural consistency loss, which aligns the structural features of the reconstructed image closely with those of the target image. The efficacy of our method is validated through experiments on eight multimodal datasets, where it is compared with the state-of-the-art methods.

Index Terms—Compensation, graph convolutional network (GCN), multimodal change detection (MCD), multisource data, structural difference, structural feature, structured graph.

Manuscript received 4 March 2024; revised 13 April 2024; accepted 25 April 2024. Date of publication 2 May 2024; date of current version 13 May 2024. This work was supported in part by the National Natural Science Foundation of China under Grant 42271411, in part by the Scientific Research Innovation Project for Graduate Students in Hunan Province under Grant CX20220169, in part by the Research Project on Monitoring and Early Warning Technologies for Implementation of Land Use Planning in Guangzhou City under Grant 2020B0101130009, and in part by the Collaborative Innovation Center for Natural Resources Planning and Marine Technology of Guangzhou under Grant 2023B04J0326. (Corresponding author: Yuqi Tang.)

Te Han, Yuqi Tang, Xin Yang, and Yuqiang Guo are with the School of Geosciences and Info-Physics, Central South University, Changsha 410083, China (e-mail: tehanrs@163.com; yqtang@csu.edu.cn; yang_x@csu.edu.cn; 235012186@csu.edu.cn).

Yuzeng Chen is with the School of Geodesy and Geomatics, Wuhan University, Wuhan 430072, China (e-mail: yuzeng_chen@whu.edu.cn).

Shujing Jiang is with Guangzhou Urban Planning & Design Survey Research Institute and the Collaborative Innovation Center for Natural Resources Planning and Marine Technology of Guangzhou, Guangzhou 510060, China (e-mail: jiangshujing128@126.com).

Digital Object Identifier 10.1109/TGRS.2024.3396141

I. INTRODUCTION

A. Background

REMOTE sensing image change detection (CD) is a technique that has garnered attention in fields, such as natural resource monitoring [1], disaster assessment [2], and urban planning [3], [4]. This technique enables the detection and analysis of changes on Earth's surface by comparing remote sensing images of the same geographical area captured at different times [5]. The advent of various remote sensing satellites, including multispectral, hyperspectral, and synthetic aperture radar (SAR) ones, has increased the variety and availability of remote sensing images, thereby advancing the development of CD. Based on the attributes of the remote sensing data, CD can be categorized into unimodal CD (UCD) and multimodal CD (MCD).

UCD primarily relies on data from a single type of satellite sensor. However, this approach often encounters challenges due to low data quality and potential data loss, influenced by factors, such as satellite performance and environmental conditions. For instance, optical satellite images are prone to cloud cover and solar illumination issues.

In contrast, MCD offers several advantages over UCD.

- 1) It leverages the observational strengths of various sensors, allowing for a combination of data types. For instance, optical remote sensing satellites provide high-resolution surface information, while SAR satellites can observe under any weather or lighting conditions.
- 2) It adapts to complex spatial-temporal characteristics, as surface changes are influenced by numerous natural and human factors, exhibiting intricate spatial-temporal patterns. MCD can better accommodate these complexities by integrating information from different remote sensing data sources, enhancing the accuracy and robustness of CD.
- 3) It improves the temporal frequency and coverage of CD, as different satellites have distinct revisit cycles and coverage capabilities. By utilizing a diverse range of satellites, more frequent and extensive remote sensing data can be collected for continuous monitoring of surface changes.

However, the varying spatial, spectral, radiometric, and temporal resolutions of different remote sensing data pose

challenges for applying traditional UCD methods. Consequently, there is an urgent need to develop specialized methods for MCD.

B. Related Work

Multimodal imagery presents challenges due to imaging differences that cause the same surface feature to exhibit varying characteristics across different images, complicating direct comparison. To address this, researchers have developed methods to construct comparable features, enabling a uniform analysis of images from various sensors. These efforts have led to the development of four main types of MCD methods: postclassification comparison, similarity measurement, feature space mapping, and image space transformation.

1) *Postclassification Comparison-Based Methods*: These methods involve selecting appropriate classification algorithms for individual processing and then comparing the results to detect changes. Representative methods include the kernel-based framework (KBF) [6], multitemporal segmentation and compound classification (MS-CC) [7], cooperative MS and hierarchical compound classification (CMS-HCC) [8], and hierarchical extreme learning machine classification (HELMC) [9]. These methods are straightforward and can be tailored to specific applications. They also facilitate the understanding and interpretation of detection results by categorizing image data. Despite the advantages of postclassification comparison methods, several challenges remain: first, classification errors inherent in these algorithms tend to accumulate during the comparison process, potentially diminishing the accuracy of CD. Second, these methods often necessitate extensive training datasets to effectively learn the characteristics of different feature types or changes. However, acquiring high-quality labeled data for MCD is particularly arduous, especially for intricate tasks. Third, multimodal images present diverse data features. It is crucial to select classification algorithms or feature extraction methods that are compatible with these varying data types. However, ensuring the consistency of classification criteria across different methods is a laborious and complex endeavor that can undermine the methods' efficiency.

2) *Similarity Measurement-Based Methods*: These methods posit a pattern correlation between multimodal images and leverage this correlation to construct invariant operators for measuring image similarity. For instance, the multidimensional statistical model (MSM) [10] employs statistical methods to model multimodal images, evaluating changes by comparing pixel-level statistical features. The MCD Markov model (M3CD) [11] identifies changed regions by establishing a Markov model to describe pixel relationships across different modalities. Other methods, such as energy-based model (EBM) [12], use energy distribution or difference metrics to assess similarity or changes. These approaches determine changes by comparing pixel value distributions or statistical information, bypassing the need for complex training processes. However, these methods may underutilize spatial information and are vulnerable to image noise. To utilize spatial information from imagery, sorted histogram (SH) [13]

assesses pixel similarity by sorting and comparing image histograms, and Mignotte [14] proposes a novel Bayesian statistical approach for MCD, involving a two-stage process that begins with preliminary estimation of spatially adaptive class conditional likelihoods specific to the imaging modality pair, followed by segmentation based on these likelihoods for each pixel and modality. The advantage of these methods is their ability to perform unsupervised CD (USCD) without specific data type dependencies, offering flexibility across various change scenarios. Nonetheless, the reliance on hand-designed mode operators, which depend on prior knowledge and expert insights, limits their ability to capture the complex dependencies between multimodal data. This challenge persists in creating an invariant operator that accurately reflects the correlation between multimodal images.

3) *Feature Space Mapping-Based Methods*: These methods project multimodal remote sensing images into a shared feature space, ensuring that similar objects are represented similarly within this space. Techniques, such as symmetric convolutional coupling network (SCCN) [15], approximately symmetrical deep neural network (ASDNN) [16], two-stage joint feature learning (TSJFL) [17], multicue contrastive self-supervised learning (MC-CSSL) [18], deep sparse residual model (DSRM) [19], commonality autoencoder (CAE) [20], and log-based transformation feature learning (LTFL) [21]. For multiscale feature learning, methods, such as deep pyramid feature learning networks (DPFL-Nets) [22], deep homogeneous feature fusion (DHFF) [23], iterative joint global-local translation (IJGLT) [24], topological structure coupling network (TSCNet) [25], and structural relationship graph convolutional autoencoder (SRGCAE) [26], are proposed. These approaches enhance consistency across different modalities by mapping multimodal images into a unified feature space, making them adaptable to various types of remote sensing data. However, the varying noise levels and imaging features of multimodal images can lead to differences in the relationships between similar object features within the shared feature space.

4) *Image Space Transformation-Based Methods*: These methods establish image transformation models between multimodal images, enabling the transformation of multimodal images from their original image space to another image space. This means that the transformed images are closer to the original images in terms of imaging features, reducing the impact of modality differences on CD. For example, homogeneous pixel transformation (HPT) [27] and deep translation-based CD network (DTCDN) [28] construct spatial transformation relationships between multimodal images using label data. To enhance the autonomy of the algorithm, methods, such as unsupervised image regression (UIR) [29] and coupled dictionary learning (CDL) [30], have been proposed. To utilize the structural information of image space, some graph-based methods have been introduced, such as patch similarity graph matrix (PSGM) [31], sparse-constrained adaptive structure consistency (SCASC) [32], and graph-based image regression and Markov random field (GIR-MRF) [33]. Additionally, some scholars have used deep learning methods to achieve spatial transformation of multimodal images, such as

generative adversarial networks under cutmix transformations (GANCT) [34], USCD [35], conditional adversarial network (CAN) [36], image translation network and postprocessing (ITNPP) [37], hierarchical extreme learning machine image transformation (HELMIT) [38], code-aligned autoencoders (CAA) [39], and adversarial cyclic encoder network (ACE-Net) [40]. Through image transformation, the feature representation of one temporal image in the image space of another temporal image can be obtained, enhancing the diversity and richness of the image data before and after the change event and providing more comprehensive change information. To further enhance the performance of such algorithms, it is necessary to consider how to establish accurate multimodal image space transformation models.

C. Motivations and Contributions

1) *Advancements of Unsupervised MCD Methods:* Recent advancements in MCD have seen the introduction of several supervised learning methods that achieve remarkable results by training models with annotated ground truth changes. Notably, multitask CD network (MTCDN) [41] and deep translation-based CD network (DTCDN) [28] have developed end-to-end image conversion frameworks utilizing UNet++ and generative adversarial networks (GAN), respectively. Hierarchical attention feature fusion (HAFF) [42] has enhanced CD capabilities through a hierarchical attention mechanism and feature fusion. Furthermore, domain adaptive cross reconstruction (DACR) [43] has facilitated domain adaptation between heterogeneous remote sensing images through feedback guidance. Leveraging spatial structural information, the dual neighborhood hypergraph neural network (DHGNN) [44] has introduced a novel network structure for high-resolution CD using a dual-neighbor hypergraph neural network. In this article, we aim to propose an unsupervised MCD method to reduce reliance on label data. Unsupervised methods offer several advantages over supervised methods: 1) they eliminate the need for training data, reducing the labor and subjectivity associated with manual data label, which is especially beneficial for large-scale datasets or when label data are scarce; 2) they are versatile, capable of handling image data from various sensors, bands, or time points, making them adaptable to a wide range of CD tasks; and 3) their lack of dependence on predefined change patterns or prior knowledge grants them robustness against unknown or complex changes.

2) *Structural Feature Consistency in Multimodal Images:* Despite the substantial differences in imaging features among multimodal images, they share consistent structural features in regions that have not changed [45], [46]. Detecting changes in multimodal images involves encoding these structural features into structural graphs and assessing the differences between them. Fractal projection and Markovian segmentation (FPMS) [47] leverages the spatial self-similarity of images. It projects patterns from one time phase to another using fractal encoding and employs pixel-level difference map binarization and Markov segmentation strategies within an unsupervised Bayesian framework to detect the changes between multimodal images. Convolution model-based mapping (CMM)

[48] captures local structural information through convolutional operations. Improved nonlocal patch-based graph (INLPG) [49] focuses on generating nonlocal structural features by considering the relationships between distant image patches, which are then quantified by mapping them into a common image domain for change measurement. On the other hand, graph-based fusion (GBF) [50] treats multimodal images as graph data, leveraging their intrinsic similarities to detect changes by fusing graph data and minimizing graph similarity. Graph learning based on signal smoothness representation (GLSSR) [51] integrates signal smoothness using graph structures to enhance detection accuracy. To further refine structural features, the iterative structure transformation and conditional random field (IST-CRF) [52] combines iterative optimization of structural transformations with CRF models for USCD. To reduce the influence of changed areas on image structural features, methods, such as enhanced graph structure representation (EGSR) [53], iterative robust graph and Markov co-segmentation (IRG-McS) [45], and adaptive optimization of structured graph (AOSG) [46], improve detection accuracy by iteratively optimizing graph structure. Structure graph-based methods offer several advantages over pixel, image patch, or superpixel-based methods in CD.

1) They can overcome imaging differences in multimodal images by mining and comparing consistent structural features in unchanged regions, enhancing detection precision and robustness.

2) The vertices in the structure graph represent image objects, and the edges between them reflect the similarity and correlation between these objects, providing valuable contextual information. This comprehensive consideration of vertex and edge attributes allows for more accurate differentiation between changed and unchanged areas.

3) Encoding structural features into structure graphs can mitigate the effects of image noise on CD. However, these methods rely on traditional K-nearest neighbors (KNNs) graphs to establish structural relationships, which typically consider only the direct proximity between pixels or objects, failing to capture more complex spatial structures and contextual information.

3) *Potential of GCN in Extracting Structural Features of Multimodal Images:* Graph convolutional network (GCN) [54] is capable of uncovering deeper structural features in images, which is particularly advantageous in MCD. Traditional methods, relying on pixel or superpixel analysis, often struggle to capture global structural information due to imaging disparities and a focus on local features. GCN addresses this limitation by acting as a robust tool for graph structure learning. They perform convolutional operations on graph structures, integrating not only the local information of nodes but also learning the intricate relationships between them. This process unveils the images' global structural features, aiding models in comprehending complex changes. Furthermore, GCN's high-dimensional feature representation enriches the contextual information of images, enabling a clearer distinction between actual changes and those falsely induced

by imaging differences, thereby improving the reliability of detection.

Therefore, this article proposes a structure difference compensation graph autoencoder (SDC-GAE) as an unsupervised method for MCD. The rationale behind this method is the assumption that, in the absence of changes, the structural features of registered multimodal images \mathbf{X} and \mathbf{Y} should align perfectly. By leveraging the structural features of image \mathbf{X} , we can reconstruct an image \mathbf{Y}' in the domain of \mathbf{Y} , ensuring that the spectral features of \mathbf{Y} and \mathbf{Y}' are identical. However, the changes in the imagery introduce structural discrepancies, resulting in spectral feature differences between \mathbf{Y} and \mathbf{Y}' . To reconcile these differences, spectral difference compensation is employed, reflecting the intensity of the changes in the multimodal images. SDC-GAE constructs a graph model with superpixels as vertices and employs a GCN to extract deep structural features from the multimodal images. This process also involves learning the structural differences between images through structural difference compensation. The encoder component of SDC-GAE maps image \mathbf{X} into a latent space, utilizing its structural features and the spectral features of image \mathbf{Y} . The decoder then reconstructs image \mathbf{Y}' from this latent space, striving for structural consistency with image \mathbf{Y} . To address spectral feature discrepancies due to changes, SDC-GAE incorporates a structural difference compensation to align the spectral features of \mathbf{Y}' with those of \mathbf{Y} . The loss function of SDC-GAE is composed of three components: image reconstruction loss, which quantifies the spectral feature differences between \mathbf{Y}' and \mathbf{Y} and guides the model to minimize these through structural difference compensation; sparse constraint loss, which is designed based on the fact that changes in images are typically confined to a few areas; and structural consistency loss, which ensures that \mathbf{Y}' closely mirrors the structural features of \mathbf{Y} . The contributions of this article are as follows.

1) The proposed SDC-GAE for MCD has been developed to eliminate the need for additional supervision signals. SDC-GAE uses the structural features of imagery from one time phase to reconstruct the spectral spatial features in another image, establishing a spectral mapping relationship between the same objects in multimodal imagery, thus obtaining the imagery of different modalities at the same time.

2) Unlike traditional methods that rely on shallow features, SDC-GAE extracts deep structural features from multimodal imagery, taking into account the complex spatial contextual relationships within the imagery.

3) SDC-GAE introduces a structural difference compensation mechanism, which optimizes the compensation value to make the reconstructed imagery structurally closer to the target imagery. The loss function design of SDC-GAE considers the characteristics of MCD, employing three types of loss functions to achieve the reconstruction of multimodal imagery, enabling the accurate identification of changed areas through structural difference compensation.

4) The validation of the proposed method's effectiveness through experiments on eight datasets and comparisons with the state-of-the-art methods.

II. METHODOLOGY

Given a pair of registered multimodal images $\mathbf{X} \in \mathbb{R}^{M \times N \times B_X}$ and $\mathbf{Y} \in \mathbb{R}^{M \times N \times B_Y}$, where M , N , and $B_X(B_Y)$ denote the length, width, and number of bands of image $\mathbf{X}(\mathbf{Y})$, respectively, and the pixels are represented as $x(m, n, b_X)$ and $y(m, n, b_Y)$. Despite the significant imaging differences, the structural features in regions that have not changed remain consistent. This consistency enables the detection of changed areas by measuring the differences in the structural features between the multimodal images. As depicted in Fig. 1, squares and circles symbolize image objects, with the line thickness indicating the degree of similarity between them. Image \mathbf{Y}' represents the spectral expression of image \mathbf{X} within the image domain \mathcal{Y} . The structural features of multimodal images are manifested through the similarity relationships among these image objects. If images \mathbf{X} and \mathbf{Y} have not any changes, the structural features of the corresponding regions should be identical, meaning that the spectral features of image \mathbf{Y}' can be represented by those of the same objects in image \mathbf{Y} . However, if changes are present in the multimodal images, their structural features will differ (as illustrated in Fig. 1 by the changed similarity relationships between objects in images \mathbf{X} and \mathbf{Y}), resulting in the spectral features of image \mathbf{Y}' in the changed areas being unable to be represented by the spectral features of the corresponding objects in image \mathbf{Y} . To address this, structural difference compensation can be employed to rectify structural discrepancies in areas experiencing changes. This approach refines the spectral properties of the reconstructed image \mathbf{Y}' to correspond with those of image \mathbf{Y} . The compensation process adeptly detects changes in image intensity across different modalities, thereby enhancing the detection of changes.

The method proposed in this article consists of three primary components (Fig. 2): structural graph construction, SDC-GAE learning, and change map (CM) generation. Sections II-A–II-C will detail these steps.

A. Structural Graph Construction

Structural graphs intuitively represent the structural features of images by using vertices and edges to illustrate the connections between them. In this study, we utilize structural graphs to depict the inherent structural features of images, with superpixels acting as graph vertices. Superpixels are pixel clusters within an image that share similar color and texture, and each superpixel is represented by a vertex in the graph. These vertices not only capture the local image features but also indicate the similarity between superpixels through their connections. This approach is more efficient than the traditional methods that rely on image patches or individual pixels, as it better preserves the image's structural information, accurately captures regional boundaries, and minimizes fragmentation and noise issues associated with small processing units. Moreover, since superpixels consist of multiple pixels, the overall data processing volume is reduced.

To segment superpixels in images, we apply the simple linear iterative clustering (SLIC) [55], which delineates

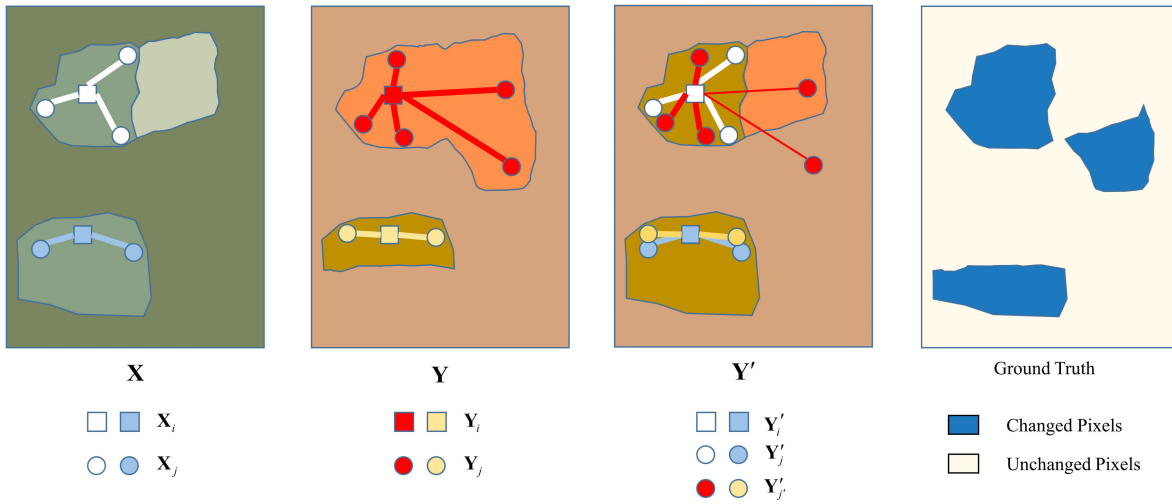


Fig. 1. Schematic of structural difference compensation in multimodal images. Squares and circles represent image objects, with the thickness of the connecting lines indicating the strength of their similarity. This similarity reflects the structural characteristics of the images. Image \mathbf{Y} represents the spectral expression of image \mathbf{X} in the image domain \mathcal{Y} . In the unchanged areas, the consistent structural features of image \mathbf{X} and \mathbf{Y} in that region allow the objects in image \mathbf{Y} to be characterized by the spectral features of the same objects in image \mathbf{X} . However, in the changed areas, the structural features of images \mathbf{X} and \mathbf{Y} will differ, preventing image \mathbf{Y} from being characterized by the spectral features of the corresponding objects in image \mathbf{X} . Therefore, structural difference compensation can be applied to the changed areas, making the reconstructed image \mathbf{Y}' have the same spectral features as the original image \mathbf{Y} . This compensation value reflects the intensity of the changes between multimodal images.

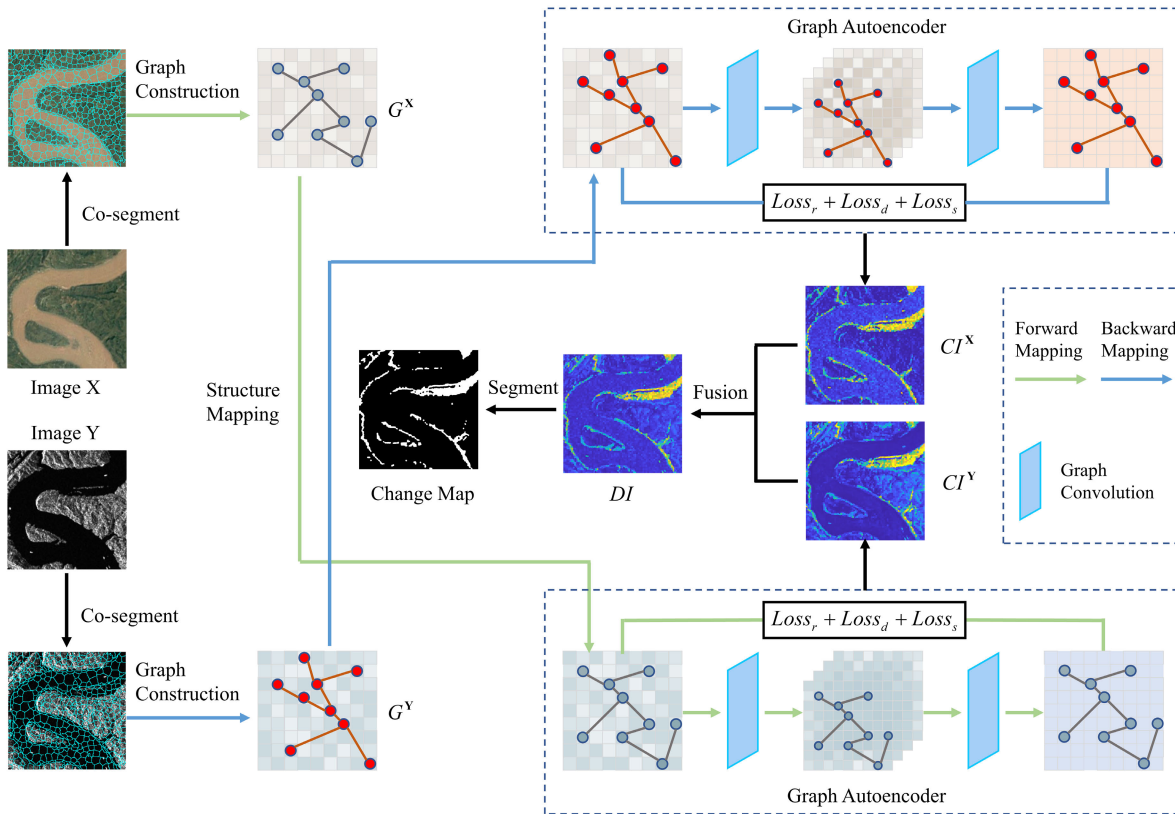


Fig. 2. Framework of SDC-GAE-based MCD method. Following the superpixel segmentation of images \mathbf{X} and \mathbf{Y} , an initial structural graph is constructed, where superpixels serve as vertices. Utilizing the structural features of image \mathbf{X} (\mathbf{Y}), the SDC-GAE models the spectral feature expressions of image \mathbf{X} (\mathbf{Y}) within image \mathbf{Y} (\mathbf{X}). The process yields a measure of the intensity of changes within the images through structural difference compensation.

superpixel boundaries based on local color similarity and spatial continuity, yielding well-defined superpixels. To ensure that multimodal images have consistent superpixel boundaries for comparison, we overlay images \mathbf{X} and \mathbf{Y} and segment the combined image using SLIC to create a superpixel map

$\mathbf{P} = \{\mathbf{P}_i | i = 1, 2, \dots, N_p\}$. This map is then mapped back to the original images \mathbf{X} and \mathbf{Y} , yielding superpixel sets $\mathbf{X} = \{\mathbf{X}_i | i = 1, 2, \dots, N_p\}$ and $\mathbf{Y} = \{\mathbf{Y}_i | i = 1, 2, \dots, N_p\}$, with N_p indicating the total number of superpixels. To characterize each superpixel, we calculate the mean and median

of its constituent pixels, providing insights into the color distribution. The mean indicates the average color trend, while the median is less affected by outliers and better represents the central tendency of the color spectrum. This process results in superpixel feature matrices $\tilde{X} \in \mathbb{R}^{N_p \times 3B_x}$ and $\tilde{Y} \in \mathbb{R}^{N_p \times 3B_y}$ for images \mathbf{X} and \mathbf{Y} , respectively.

Based on the obtained superpixels, we can construct a structural graph G_X for image \mathbf{X} to represent its structural features. We define $G_X = \{V_X, E_X\}$ as follows:

$$\begin{aligned} V_X &= \{\mathbf{X}_i | i = 1, 2, \dots, N_p\} \\ E_X &= \{(\mathbf{X}_i, \mathbf{X}_j) | i = 1, 2, \dots, N_p, j \in \Omega_{\mathbf{X}_i}, j \neq i\} \end{aligned} \quad (1)$$

where V_X and E_X represent the vertices and edges of the graph G_X , respectively, while $\Omega_{\mathbf{X}_i}$ denotes the set of indices of vertices \mathbf{X}_j connected to vertex \mathbf{X}_i . Similarly, we can construct the structural graph G_Y for image \mathbf{Y} .

To more conveniently describe the connection relationships between vertices in graph G_X , we introduce the adjacency matrix $\mathbf{A}^X \in \mathbb{R}^{N_p \times N_p}$ of graph G_X

$$\mathbf{A}_{i,j}^X = \begin{cases} 1, & (\mathbf{X}_i, \mathbf{X}_j) \in E_X \\ 0, & (\mathbf{X}_i, \mathbf{X}_j) \notin E_X. \end{cases} \quad (2)$$

This article introduces the degree matrix \mathbf{D}^X for graph G_X to analyze the connection strength between vertices and to uncover the graph's clustering features and the significance of individual vertices. The degree matrix, a diagonal matrix where $\mathbf{D}_{i,i}^X = \sum_j \mathbf{A}_{i,j}^X \in \mathbb{R}^{N_s \times N_s}$ represents the frequency of vertex connections, serves as a direct mapping of each vertex's connectivity. In parallel, a structural graph G_Y is constructed for image \mathbf{Y} to apply similar analysis.

The selection of an appropriate number of associated vertices K_i for vertex \mathbf{X}_i is pivotal for accurately capturing the graph's structural features. A suboptimal K_i , either too small or too large, can lead to an inaccurate representation of the graph's structure. To determine K_i , we adopt a method similar to that described in [32], which involves setting a range for the number of neighbors with $k_{\max} = \lfloor K_{\text{ratio}} \times N_p \rfloor$ as the upper limit and $k_{\min} = \lfloor k_{\max}/10 \rfloor$ as the lower limit, where $\lfloor \cdot \rfloor$ denotes the floor function and K_{ratio} is the neighbor ratio. We calculate the feature distance $\text{dist}_{i,j}^X = \|\mathbf{X}_i - \mathbf{X}_j\|_2^2$ between vertex \mathbf{X}_i and other vertices to identify the k_{\max} closest neighbors. The degree value \mathbf{D}^X of vertex \mathbf{X}_i is then computed. Subsequently, K_i is set to the minimum of $\max\{\mathbf{D}_{i,i}^X, k_{\min}\}$ and k_{\max} to ensure a balanced representation of the graph's structure.

B. SDC-GAE Learning

1) *Structure Difference Compensation Graph Autoencoder:* GCNs are specialized network models designed for processing graph data, adept at capturing and learning the structural characteristics inherent in graphs. Their strong interpretability and expressive capabilities have made them a staple in fields ranging from computer vision to recommendation systems. To precisely determine spectral difference compensation, which measure the changes intensity in multimodal images, the proposed SDC-GAE is composed of an encoder and a decoder. Both components are equipped with graph convolutional layers

to facilitate efficient feature extraction and reconstruction, akin to traditional autoencoders.

The SDC-GAE reconstructs the original image to obtain its spectral feature representation in a different temporal image domain, while also calculating the structural difference compensation value *dif*. As established earlier, if multimodal images remain unchanged, their structural features should be consistent. Leveraging this consistency, we preserve the structural features of image \mathbf{X} and reconstruct it within the image domain \mathcal{Y} , guided by the spectral features of image \mathbf{Y} . In the reconstructed image \mathbf{Y}' , the feature representation $h_{\mathbf{Y}'_i}^{(l+1)}$ of superpixel \mathbf{Y}'_i at layer $l+1$ is refined through a weighted aggregation of the features of its l -layer neighboring vertices $h_{\mathbf{Y}'_j}^{(l)}$

$$h_{\mathbf{Y}'_i}^{(l+1)} = \sigma \left(\sum_{j \in \Omega_{\mathbf{X}_i}} \alpha_{i,j}^{\mathbf{X}^{(l)}} W^{(l)} h_{\mathbf{Y}'_j}^{(l)} \right) \quad (3)$$

where $\sigma(\cdot)$ is the activation function, $W^{(l)}$ represents the learnable weight matrix, $\alpha_{i,j}^{\mathbf{X}^{(l)}}$ denotes the dynamically allocated contribution degree for each vertex, which reflects the importance of neighboring vertex in updating the current vertex, and $j \in \Omega_{\mathbf{X}_i}$ denotes that in the superpixel \mathbf{Y}'_i update process, the spatial index utilized for neighboring vertices corresponds to that of the adjacent vertices in graph G_X . Consequently, this implies that the feature learning for reconstructing image \mathbf{Y}' is performed using the structural features of image \mathbf{X} . It is worth noting that the input layer features of vertex $\tilde{\mathbf{Y}}'_j$ are $h_{\mathbf{Y}'_j}^{(0)} = \tilde{\mathbf{Y}}'_j$.

To enhance the accuracy of establishing connection relationships between graph vertices and to improve the graph neural network's ability to model image structural features, this article introduces the graph attention mechanism (GAM) [56], [57]. This mechanism is used to learn the contribution degree $\alpha_{i,j}^{\mathbf{X}^{(l)}}$ for each vertex. Specifically, for each vertex \mathbf{X}_i in graph $G_X = \{V_X, E_X\}$, the attention weight $\alpha_{i,j}^{\mathbf{X}^{(l)}}$ with its neighboring vertices can be calculated as follows:

$$\alpha_{i,j}^{\mathbf{X}^{(l)}} = \frac{\exp\left(\text{LeakyReLU}\left(a^{(l)T} \left[W^{(l)} h_{\mathbf{X}_i}^{(l)} \parallel W^{(l)} h_{\mathbf{X}_j}^{(l)} \right] \right)\right)}{\sum_{k \in \Omega_{\mathbf{X}_i}} \exp\left(\text{LeakyReLU}\left(a^{(l)T} \left[W^{(l)} h_{\mathbf{X}_i}^{(l)} \parallel W^{(l)} h_{\mathbf{X}_k}^{(l)} \right] \right)\right)} \quad (4)$$

where $a^{(l)}$ represents a parameterized vector of a learnable attention mechanism, \parallel denotes the concatenation of feature vectors, and $\text{LeakyReLU}(\cdot)$ is the activation function.

Additionally, to achieve a more stable update process, we introduce multihead attention. Assuming there are R attention heads for vertex feature updates, (3) can be rewritten as

$$h_{\mathbf{Y}'_i}^{(l+1)} = \sigma \left(\frac{1}{R} \sum_{r=1}^R \sum_{j \in \Omega_{\mathbf{X}_i}} \alpha_{i,j}^{\mathbf{X}^{(l)r}} W^{(l)r} h_{\mathbf{Y}'_j}^{(l)} \right) \quad (5)$$

where $\alpha_{i,j}^{\mathbf{X}^{(l)r}}$ represents the weight of the r th ($r = 1, 2, \dots, R$) attention mechanism a^r , while $W^{(l)r}$ is the corresponding learnable weight matrix.

SDC-GAE employs graph convolutional layers to progressively learn and reconstruct the image \mathbf{Y}' of \mathbf{X} within the image domain \mathcal{Y} . By leveraging the structural features of image \mathbf{Y} , we can similarly construct the reconstructed image \mathbf{X}' within the image domain \mathcal{X} .

2) *Loss Function*: Influenced by changed areas, images $\mathbf{Y}'(\mathbf{X}')$ and $\mathbf{Y}(\mathbf{X})$ may exhibit differences, necessitating compensation with values $dif^{\mathbf{X}}$ and $dif^{\mathbf{Y}}$. These values, along with other SDC-GAE parameters, are optimized through the backpropagation algorithm [58] and gradient descent. To ensure that the reconstructed images \mathbf{Y}' and \mathbf{X}' closely resemble images \mathbf{Y} and \mathbf{X} in spectral features, respectively, we introduce a reconstruction loss function to achieve this objective

$$\text{loss}_r = \sum_{i=1}^{N_p} \|\tilde{X}_i - \tilde{X}'_i + dif_i^{\mathbf{X}}\|_2^2 + \sum_{i=1}^{N_p} \|\tilde{Y}_i - \tilde{Y}'_i + dif_i^{\mathbf{Y}}\|_2^2 \quad (6)$$

where $\tilde{Y}' \in \mathbb{R}^{N_p \times 3B_Y}$ represents the feature matrix of the reconstructed image \mathbf{Y}' .

The sparse constraint loss is engineered to minimize the structural difference compensation values $dif^{\mathbf{X}}$ and $dif^{\mathbf{Y}}$, acknowledging that changed regions are generally smaller compared to unchanged regions. Consequently, the majority of pixels are categorized as unchanged (assigned values of zero or close to zero), whereas a small subset of pixels is designated as changed (assigned a larger value values)

$$\text{loss}_d = \|dif^{\mathbf{X}}\|_2^2 + \|dif^{\mathbf{Y}}\|_2^2. \quad (7)$$

As shown in Fig. 1, the image \mathbf{X} and its reconstructed image \mathbf{Y}' share consistent structural features. Specifically, the interconnections and interactions between \mathbf{Y}'_i and \mathbf{Y}'_j should reflect the same pattern and intensity as those between \mathbf{X}_i and \mathbf{X}_j , which means

$$\min \sum_{i,j=1}^{N_p} \|\tilde{Y}'_i - \tilde{Y}'_j\|_2^2 \mathbf{A}_{i,j}^{\mathbf{X}} = 2Tr(\tilde{Y}'^T \mathbf{L}^{\mathbf{X}} \tilde{Y}') \quad (8)$$

where $\mathbf{L}^{\mathbf{X}} = \mathbf{D}^{\mathbf{X}} - \mathbf{A}^{\mathbf{X}}$ is the Laplacian matrix of the graph $G_{\mathbf{X}}$.

Additionally, we have performed normalization on the Laplacian matrix, which offers three significant benefits. First, normalization stabilizes the degree values of vertices, essential for maintaining nodal force equilibrium during analysis. This step also minimizes computational errors, boosting the reliability and accuracy of numerical calculations. Moreover, the normalized matrix adapts to different network structures, making the algorithm versatile for a variety of scenarios. Thus, we can obtain the structural consistency loss

$$\text{loss}_s^{\mathbf{X}} = 2Tr(\tilde{Y}'^T \tilde{L}^{\mathbf{X}} \tilde{Y}') \quad (9)$$

where $\tilde{L}^{\mathbf{X}} = \mathbf{I} - (\mathbf{D}^{\mathbf{X}})^{-1/2} \mathbf{A}^{\mathbf{X}} (\mathbf{D}^{\mathbf{X}})^{-1/2}$ represents the normalized Laplacian matrix, while $\mathbf{I} \in \mathbb{R}^{N_p \times N_p}$ represents the identity matrix. Similarly, for the reconstructed image \mathbf{X}' , we can obtain

$$\text{loss}_s^{\mathbf{Y}} = 2Tr(\tilde{X}'^T \tilde{L}^{\mathbf{Y}} \tilde{X}') \quad (10)$$

TABLE I
FRAMEWORK OF SDC-GAE

SDC-GAE
Input: images \mathbf{X} and \mathbf{Y} , parameters of SDC-GAE
1. Structure Graph Construction:
Image superpixel segmentation and superpixel feature extraction;
Calculation of similarity between graph vertices;
Determination of the number of neighboring vertices for a graph vertex;
2. SDC-GAE Learning:
Construct the structure of SDC-GAE;
Construct the loss function;
Update vertex features according to (3);
3. CM Generation:
Compute the binary CM with (13), (14) and Otsu's threshold segmentation.

where $\tilde{X}' \in \mathbb{R}^{N_p \times 3B_X}$ represents the feature matrix of the reconstructed image \mathbf{X}' , and $\tilde{L}^{\mathbf{Y}} = \mathbf{I} - (\mathbf{D}^{\mathbf{Y}})^{-1/2} \mathbf{A}^{\mathbf{Y}} (\mathbf{D}^{\mathbf{Y}})^{-1/2}$ represents the normalized Laplacian matrix of graph $G_{\mathbf{Y}}$. The overall structural consistency loss function is then given by

$$\text{loss}_s = \text{loss}_s^{\mathbf{X}} + \text{loss}_s^{\mathbf{Y}}. \quad (11)$$

Then, the overall loss function for SDC-GAE is

$$\text{loss} = \text{loss}_r + \text{loss}_d + \text{loss}_s. \quad (12)$$

C. CM Generation

Structural difference compensation $dif^{\mathbf{X}}$ and $dif^{\mathbf{Y}}$ reflect the change intensity of image; thus, the final change intensity is a combination of both

$$dif^{\text{final}} = dif^{\mathbf{X}} / \text{mean}(dif^{\mathbf{X}}) + dif^{\mathbf{Y}} / \text{mean}(dif^{\mathbf{Y}}). \quad (13)$$

The change intensity map (CIM) is obtained as follows:

$$CI_{(m,n)} = dif_i^{\text{final}}; (m, n) \in P_i. \quad (14)$$

The generation of the final CM can be visualized as a binary segmentation problem, and we obtain the CM using Otsu's threshold segmentation [59]. The summary of the proposed SDC-GAE is presented in Table I.

III. EXPERIMENTS

In this section, we will provide a detailed introduction to the experimental setup, including comparative methods, evaluation metrics, experimental parameters, and so on, following with the experimental datasets and the results of the experiments.

A. Comparative Methods and Evaluation Metrics

To validate the effectiveness of the methods presented in this article, we employ several state-of-the-art methods as comparative methods.

1) LTFL [21] utilizes a stacked denoising autoencoder to extract deep features from multimodal images and applies difference metrics to these features. High-confidence samples are selected for training a classifier, which performs binary classification on feature difference maps to generate CM.

2) INLPG [49] constructs nonlocal structural features for multimodal images by considering nonlocal correlations

between image patches. These features are mapped into a consistent image domain for comparison, enabling the quantification of changes within multimodal images.

3) GBF [50] treats multimodal images as graph data, capturing their intrinsic similarities. It fuses multitemporal graph data and minimizes graph similarity to detect changes in multimodal images.

- 1) IRG-McS [45] develops an adaptive nonlocal structural graph based on superpixels to represent image structure features. It refines the graph structure of unchanged regions using structural difference metrics and McS, enhancing the accuracy of CD.
- 2) SCASC [32] retains the structural features of the source image and applies sparse constraints to transform these features into the target image domain. Change information is extracted through a comparison between the original and transformed images.
- 3) GIR-MRF [33] employs a UIR method that integrates global and local graph structure learning to capture image features and relationships. It uses a Markov segmentation model to segment difference maps, resulting in CM.
- 4) SRGCAE [26] leverages a GCN to learn graph structure relationships in multimodal images, expressing their structural features. It extracts change information by comparing these structural features.
- 5) AOSG [46] constructs an adaptive structured graph by combining the spatial scale of image objects with the feature distance between image patches. It optimizes the graph by considering the change attributes of neighboring regions, thereby improving the precision of structural feature representation.

In order to evaluate the effectiveness of various methods quantitatively, we have employed a suite of evaluation metrics, including overall accuracy (OA), kappa coefficient (KC), and $F1$ -measure ($F1$). These metrics provide a comprehensive assessment of model performance. The detailed descriptions of these evaluation metrics are given in the following.

The OA serves as a key indicator when assessing the efficacy of a classification model. It measures the model's ability to correctly classify samples, expressed as a ratio of the correctly classified instances to the total sample count N . The OA is calculated as follows:

$$OA = (TP + TN)/(TP + TN + FP + FN) \quad (15)$$

where TP, FP, TN, and FN denote true positives, false positives, true negatives, and false negatives, respectively.

The KC is used to measure the degree of agreement among data, and it is calculated as follows:

$$KC = (OA - PRE)/(1 - PRE) \quad (16)$$

where

$$PRE = ((TP + FN)(TP + FP) + (TN + FP)(TN + FN))/N^2. \quad (17)$$

The $F1$ score is the harmonic mean of precision and recall, used to measure the comprehensive performance of a

classification model on a specific category, especially in cases of class imbalance. It is calculated as follows:

$$F1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (18)$$

where

$$\text{Precision} = TP/(TP + FP) \quad (19)$$

$$\text{Recall} = TP/(TP + FN). \quad (20)$$

In this article, the proposed SDC-GAE employs the Adam optimizer [60] with a learning rate set to 0.01, weight decay to 0.0001, and the number of attention heads R is set to 4. Both layers of the graph convolutional layers use convolutional kernels of 16. The epochs are set to 300. In addition, the main parameters of SDC-GAE, the neighbor ratio k_{ratio} and the number of superpixels N_p , are set to 0.1 and 5000, respectively, and these parameters will be analyzed in Section IV. All experiments of SDC-GAE are conducted in the following environment: Pytorch, CPU is AMD Ryzen 7 3800X 8-core processor, 3.89 GHz, Windows 11, 32-GB RAM, and an NVIDIA GeForce GTX 2070 SUPER.

B. Experimental Datasets

This article validates the effectiveness of the proposed multimodal image CD methods using eight diverse remote sensing datasets, as depicted in Fig. 3(a) and (b). The datasets include a range of sensor images from optical satellites, such as Sentinel-1, QuickBird, and Landsat-8, as well as SAR satellites, capturing both short-term and long-term changes. This variety ensures a comprehensive assessment of the methods' performance across different sensors and time frames. Spanning various geographical regions, the datasets feature changes at multiple scales, such as urban development and river expansion, which are essential for testing the algorithms' applicability and generalization capabilities. The reference images, derived from expert knowledge and high-resolution, temporally and spatially proximate images, provide a reliable benchmark for validation. The selection of these datasets aims to rigorously evaluate the proposed methods' effectiveness in diverse environmental contexts. All datasets underwent preprocessing, such as radiometric correction, atmospheric correction, and geometric correction, with each dataset's images resampled to the same spatial resolution. For further details, refer to Table II.

C. Experimental Results

Fig. 3 illustrates the CMs for different methods applied to datasets #1 through #8, alongside the CIMs of SDC-GAE. Datasets #1 and #2, which depict river changes, present a challenge due to "pseudo-changes" resulting from varying shadow distributions over land. Visual inspection reveals that all methods successfully identified the primary change areas in dataset #1. However, GBF and SRGCAE exhibited notable FPs, while INLPG, IRG-McS, SCASC, and GIR-MRF, though having fewer FPs, missed several detections. AOSG managed to detect more comprehensive change areas but encountered isolated FPs in the image's center. In dataset #2, LTFL,

TABLE II
DESCRIPTION OF THE DATASETS

Dataset	Sensor	Size(pixels)	Date	Location	Event (& Spatial resolution)
#1	Google Earth/Sentinel-1	600×600×3(1)	Dec. 1999-Nov. 2017	Chongqing, China	River expansion (10 m)
#2	Sentinel-2/ Sentinel-1	444×571×3(1)	Apr. 2017 - Oct. 2020	Lake Poyang, China	Lake expansion (10 m)
#3	Pleiades/WorldView2	2000 × 2000 × 3(3)	May 2012 – July 2013	Toulouse, France	Urban construction (0.52m)
#4	QuickBird/LiDAR	700×700×4(1)	Nov. 2007 – June 2011	San Francisco, USA	Urban construction (0.5m)
#5	QuickBird-2/TerraSAR-X	4135×2325×3(1)	July 2006 – July 2007	Gloucester, England	Flooding (\approx 0.65m)
#6	Spot/NDVI	990 × 554 × 3(1)	1999 - 2000	Gloucester, England	Flooding (\approx 25m)
#7	Landsat-8/Sentinel-1	3500×2000×11(3)	Jan. 2017 – Feb. 2017	Sutter County, USA	Flooding (\approx 15m)
#8	Landsat-5/Google Earth	300 × 412 × 1(3)	Sept. 1995 - July 1996	Sardinia, Italy	Lake expansion (30m.)

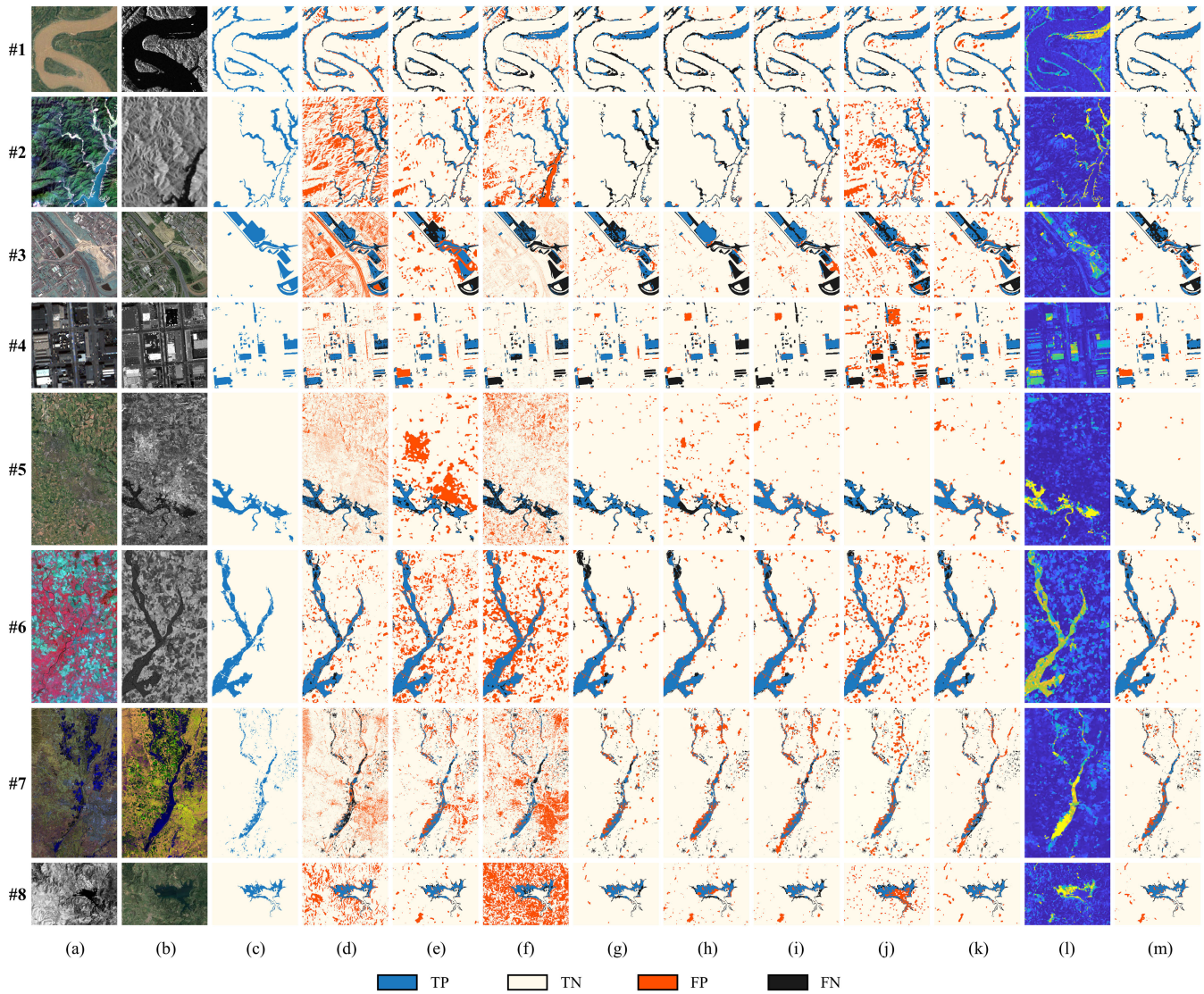


Fig. 3. CMs of different methods and the CIMs for SDC-GAE. From top to bottom, they correspond to datasets #1–#8. From left to right, they are as follows. (a) Image X, (b) image Y, (c) reference image, (d) CMs of LTFM, (e) CMs of INLPG, (f) CMs of GBF, (g) CMs of IRG-McS, (h) CMs of SCASC, (i) CMs of GIR-MRF, (j) CMs of SRGCAE, (k) CMs of AOSG, (l) CIMs of SDC-GAE, and (m) CMs of SDC-GAE.

INLPG, GBF, and SRGCAE showed significant FPs, whereas IRG-McS, SCASC, GIR-MRF, and AOSG had fewer FPs but more missed detections. Dataset #3, with its complex change scenario involving bare land, grassland, buildings, and roads, saw most methods, except SCASC, producing many false positives. SCASC, however, failed to identify the change area

in the lower right corner of the image. Dataset #4, which reflects the changes in buildings and vehicles, saw LTFM and SRGCAE generating more FPs than other methods. Meanwhile, INLPG, IRG-McS, SCASC, and GIR-MRF missed three distinct building change areas. AOSG achieved a balance between FPs and FNs but still overlooked a building change

TABLE III
ACCURACY EVALUATIONS ON DATASETS #1–#4. THE MAXIMUM AND SECOND MAXIMUM VALUES ARE REPRESENTED BY BOLD AND UNDERLINED TEXT, RESPECTIVELY

Methods	#1			#2			#3			#4		
	OA	KC	F1	OA	KC	F1	OA	KC	F1	OA	KC	F1
LTFL	0.9204	0.7119	0.7579	0.7016	0.1490	0.2600	0.6800	0.2181	0.3862	0.9109	0.5367	0.5862
INLPG	0.9083	0.5635	0.6092	0.9148	0.5426	0.5879	0.8171	0.3395	0.4482	0.9279	0.6445	0.6843
GBF	0.8989	0.5530	0.6109	0.8300	0.3458	0.4231	0.8261	0.2155	0.3105	0.9202	0.4167	0.4544
IRG-McS	0.9128	0.6022	0.6483	0.9450	0.4948	0.5182	0.8685	0.4239	0.4973	0.9387	0.6220	0.6552
SCASC	0.8955	0.5069	0.5599	<u>0.9537</u>	0.6443	0.6684	<u>0.8918</u>	0.4711	0.5247	0.9148	0.3929	0.4345
SRGCAE	0.9223	0.6812	0.7259	0.8354	0.3585	0.4339	0.8231	0.3817	0.4867	0.7708	0.0821	0.2005
GIR-MRF	0.9037	0.5913	0.6457	0.9481	0.6398	0.6678	0.8960	0.4840	0.5350	0.9300	0.4957	0.5294
AOSG	<u>0.9270</u>	<u>0.7291</u>	<u>0.7725</u>	0.9465	<u>0.7061</u>	0.7347	0.8745	0.5131	0.5871	<u>0.9445</u>	<u>0.6870</u>	<u>0.7177</u>
SDC-GAE	0.9429	0.7590	0.7914	0.9627	0.7067	<u>0.7259</u>	0.8860	<u>0.4932</u>	<u>0.5561</u>	0.9615	0.6978	0.7178

TABLE IV
ACCURACY EVALUATIONS ON DATASETS #5–#8. THE MAXIMUM AND SECOND MAXIMUM VALUES ARE REPRESENTED BY BOLD AND UNDERLINED TEXT, RESPECTIVELY

Methods	#5			#6			#7			#8		
	OA	KC	F1	OA	KC	F1	OA	KC	F1	OA	KC	F1
LTFL	0.8658	0.2879	0.3509	0.9195	0.6689	0.7145	0.8504	0.0631	0.1240	0.8078	0.2849	0.3549
INLPG	0.8316	0.2615	0.3321	0.7829	0.4134	0.5162	0.9059	0.3730	0.4128	0.9503	0.6128	0.6392
GBF	0.8284	0.1032	0.1832	0.8272	0.4827	0.5695	0.7915	0.1093	0.1737	0.4667	0.0344	0.1418
IRG-McS	0.9672	0.7029	0.7201	0.9358	0.7136	0.7502	<u>0.9469</u>	<u>0.4703</u>	<u>0.4975</u>	0.9700	0.7242	0.7401
SCASC	0.9295	0.5010	0.5381	0.9503	0.7762	<u>0.8046</u>	0.9381	0.4585	0.4888	0.9472	0.5958	0.6238
SRGCAE	<u>0.9707</u>	0.7184	0.7335	0.8728	0.5724	0.6400	0.9376	0.4246	0.4557	0.9129	0.4737	0.5164
GIR-MRF	0.9639	0.7524	0.7713	0.9332	0.7249	0.7627	0.9446	0.4674	0.4954	0.9587	0.6332	0.6551
AOSG	0.9644	<u>0.7587</u>	<u>0.7773</u>	0.9553	<u>0.7788</u>	0.8040	0.9450	0.4259	0.4544	0.9605	0.6733	0.6944
SDC-GAE	0.9774	0.7945	0.8063	<u>0.9530</u>	0.7902	0.8169	0.9525	0.5134	0.5378	<u>0.9645</u>	<u>0.7185</u>	<u>0.7374</u>

in the lower left corner. Datasets #5–#7, all showing river changes, faced challenges due to extensive image coverage and complex object textures. In dataset #5, GBF's performance was compromised by numerous false and FNs. LTFL and INLPG had many FPs, yet other methods managed to detect change areas more thoroughly, albeit with some FPs. Dataset #6 saw INLPG, GBF, and SRGCAE with significant FPs, while other methods detected change areas more completely with minimal FPs. In dataset #7, LTFL, INLPG, and GBF had more FPs, but other methods detected the main change areas with relatively fewer FPs. Finally, dataset #8, which captured subtle lake changes, found LTFL and GBF with many FPs in land areas, and SRGCAE incorrectly identified the lake's interior as a change area. In contrast, other methods detected a more accurate change area. In summary, the proposed SDC-GAE method not only detected more complete change areas but also had the fewest FPs. This is attributed to its capability to extract high-level image features via a graph encoder and accurately pinpoint change areas through structural difference compensation.

Fig. 3(i) showcases the CIMs produced by the proposed SDC-GAE across various datasets. To evaluate the performance of these CIMs, we utilize the area under the receiver operating characteristics (ROC) curve (AUC) as a metric. As depicted in Fig. 4, the ROC curves for SDC-GAE on datasets #1 through #8 yield the AUC values of 0.9207, 0.9018, 0.8383, 0.9070, 0.9460, 0.9835, 0.8899, and 0.9331,

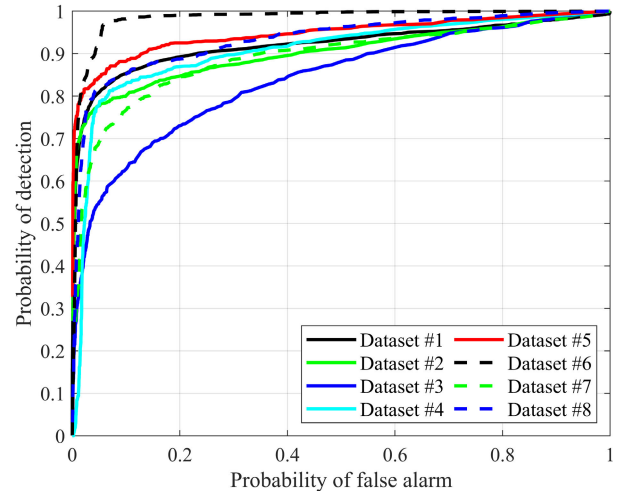


Fig. 4. ROC curves of CIMs generated by SDC-GAE across all datasets.

respectively. These results indicate that SDC-GAE is adept at generating high-quality CIMs. Furthermore, these maps can be efficiently converted into CMs using straightforward threshold segmentation techniques.

The accuracy assessment results for various methods on datasets #1–#4 and #5–#8 are detailed in Tables III and IV, respectively. In these tables, the boldface denotes the maximum values, while underlines highlight the second maximum

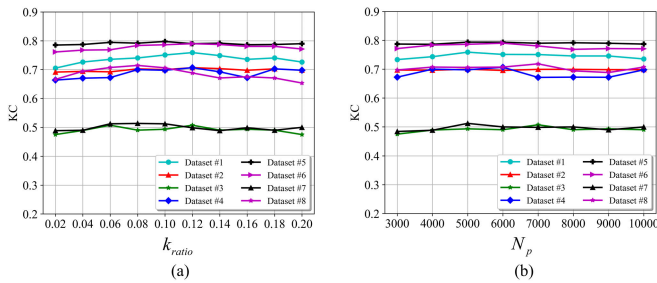


Fig. 5. Sensitivity analysis of parameters in SDC-GAE. (a) k_{ratio} -KC curves. (b) N_p -KC curves.

values. With the exception of dataset #3, the proposed SDC-GAE consistently ranks within the top two for OA across the remaining datasets. Moreover, SDC-GAE secures a position in the top two for both KC and $F1$ scores in all datasets, achieving the highest KC in six out of the eight datasets. These results underscore the effectiveness and robustness of the SDC-GAE.

IV. DISCUSSION

A. Parameter Analysis

1) *Neighbor Ratio K_{ratio}* : This section examines the sensitivity of the proposed SDC-GAE to the parameter K_{ratio} . To do this, we maintained N_p at 5000 and varied K from 0.02 to 0.20 in increments of 0.02. The resulting accuracy changes for SDC-GAE are depicted in Fig. 5(a). Analyzing Fig. 5(a) reveals that SDC-GAE's accuracy generally ascends with K_{ratio} values from 0.02 to 0.08, peaks and stabilizes when K_{ratio} is between 0.08 and 0.12, and subsequently declines for K_{ratio} values from 0.12 to 0.20. This trend underscores the impact of the number of neighbors on the graph's structural representation quality. The rationale is twofold: first, a vertex with an excessive number of neighbors may be inundated with noisy information, compromising the graph's accuracy. On the other hand, a vertex with insufficient neighbors might lack the context necessary for precise feature representation. Second, graph neural networks are prone to oversmoothing, where an abundance of neighbors can lead to each vertex being overly influenced by similar ones during information aggregation, thus diminishing the distinctiveness of vertex features. Given these considerations, we suggest an optimal K_{ratio} value of 0.1 in this article.

2) *Number of Superpixels N_p* : To investigate the effect of N_p on the proposed SDC-GAE, we varied N_p from 3000 to 10000 in increments of 1000, with the neighbor ratio K_{ratio} held constant at 0.1. Fig. 5(b) illustrates the accuracy of SDC-GAE across this range of N_p values. The KC for SDC-GAE initially increases as N_p rises from 3000 to 4000, stabilizes between 4000 and 6000, and then exhibits a slight decline for N_p values from 6000 to 10 000. This is because the number of superpixels determines the spatial resolution of the segmented image. Too few superpixels may fail to capture subtle changes, while too many can lead to oversegmentation, making it difficult to distinguish real changes from noise in CD. However, overall, the impact on the proposed SDC-GAE is not significant. Moreover, an excessive number

of superpixels will increase the algorithm's running time. Therefore, considering both the accuracy and efficiency of the algorithm, we recommend setting N_p to 5000.

B. Similarity Between the Reconstructed Image and the Target Image

Fig. 6 displays the reconstructed images generated by SDC-GAE across datasets #1-#8. It can be observed from Fig. 6 that the reconstructed image $\mathbf{X}'(\mathbf{Y}')$ shares similar spectral characteristics with image $\mathbf{X}(\mathbf{Y})$ and structural features with image $\mathbf{Y}(\mathbf{X})$. To further illustrate the differences in spectral characteristics before and after image reconstruction, we selected unaltered regions within each dataset (indicated by the red boxes) and calculated the histogram distribution of each band of the images [red, green, and blue bands are labeled with numbers (1), (2), and (3), respectively, in Fig. 7]. It is evident that there are significant differences in the shape, peak position, and distribution range of the probability density function (pdf) curves between the original images \mathbf{X} and \mathbf{Y} . In contrast, the pdf curves of the reconstructed image $\mathbf{X}'(\mathbf{Y}')$ and image $\mathbf{X}(\mathbf{Y})$ are much closer. This indicates that SDC-GAE can encode image $\mathbf{X}(\mathbf{Y})$ into the spectral space of image $\mathbf{Y}(\mathbf{X})$ while preserving the structural features of image $\mathbf{X}(\mathbf{Y})$.

C. Effectiveness of GAM

The GAM allows SDC-GAE to dynamically adjust the weights of graph vertices during processing, prioritizing those vertices deemed more important. This adaptive approach facilitates the learning of vertex representations. To assess the impact of GAM, we compared the accuracy of SDC-GAE both with and without GAM (Table V). The results, as presented in Table V, show that SDC-GAE achieved an average enhancement of 2.16% in OA, 7.30% in KC, and 4.71% in $F1$ when GAM was employed. These improvements underscore the GAM's ability to direct SDC-GAE's focus toward salient vertex features, which in turn optimizes the structural difference compensation values and enhances the accuracy of CD.

D. Ablation Study of the Loss Function

In the structural difference compensation learning process of SDC-GAE, three loss functions collaboratively guide the model: image reconstruction loss $loss_r$, sparsity constraint loss $loss_d$, and structural consistency loss $loss_s$. $loss_r$ includes the loss term of the reconstructed image and the structural difference compensation value and thus serves as the fundamental loss function of the proposed SDC-GAE. Therefore, this article focuses on the ablation of $loss_d$ and $loss_s$ to evaluate their individual impacts on algorithm performance, as detailed in Table VI. Fig. 6 demonstrates that the accuracy of SDC-GAE, when relying solely on $loss_r$, is markedly inferior to that of the comprehensive model incorporating all loss functions. This disparity highlights the insufficiency of $loss_r$ for precise CD and underscores the model's reliance on $loss_d$ and $loss_s$ for enhanced performance. By incorporating $loss_d$ and $loss_s$ into $loss_r$, SDC-GAE achieves an average improvement of

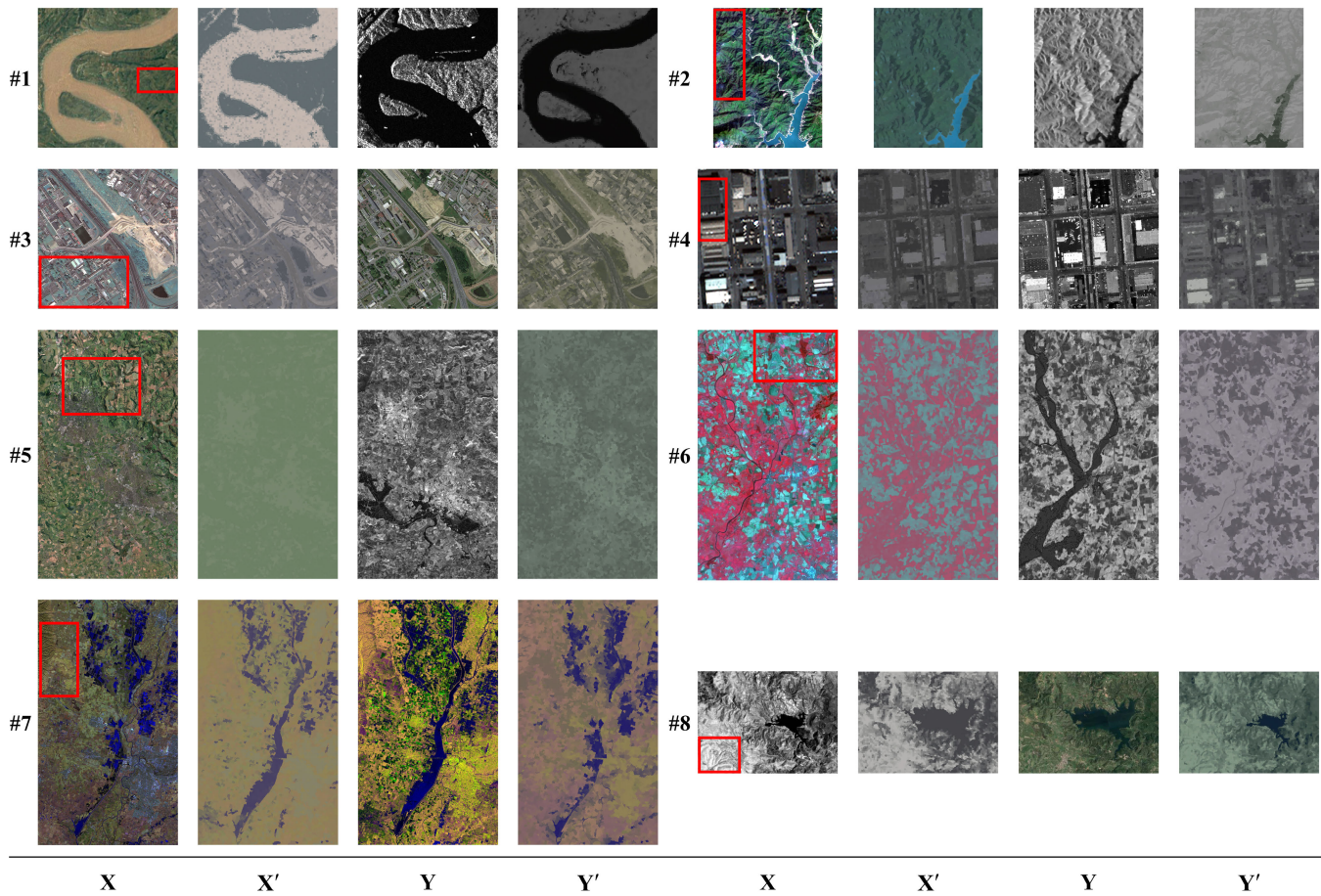


Fig. 6. Reconstructed images of SDC-GAE on datasets #1–#8 are presented. Within each dataset, from left to right, they are as follows: image X , the reconstructed image X' of image Y in the domain of \mathcal{X} , image Y , and the reconstructed image Y' of image X in the domain of \mathcal{Y} . The red boxes indicate the selected changed regions.

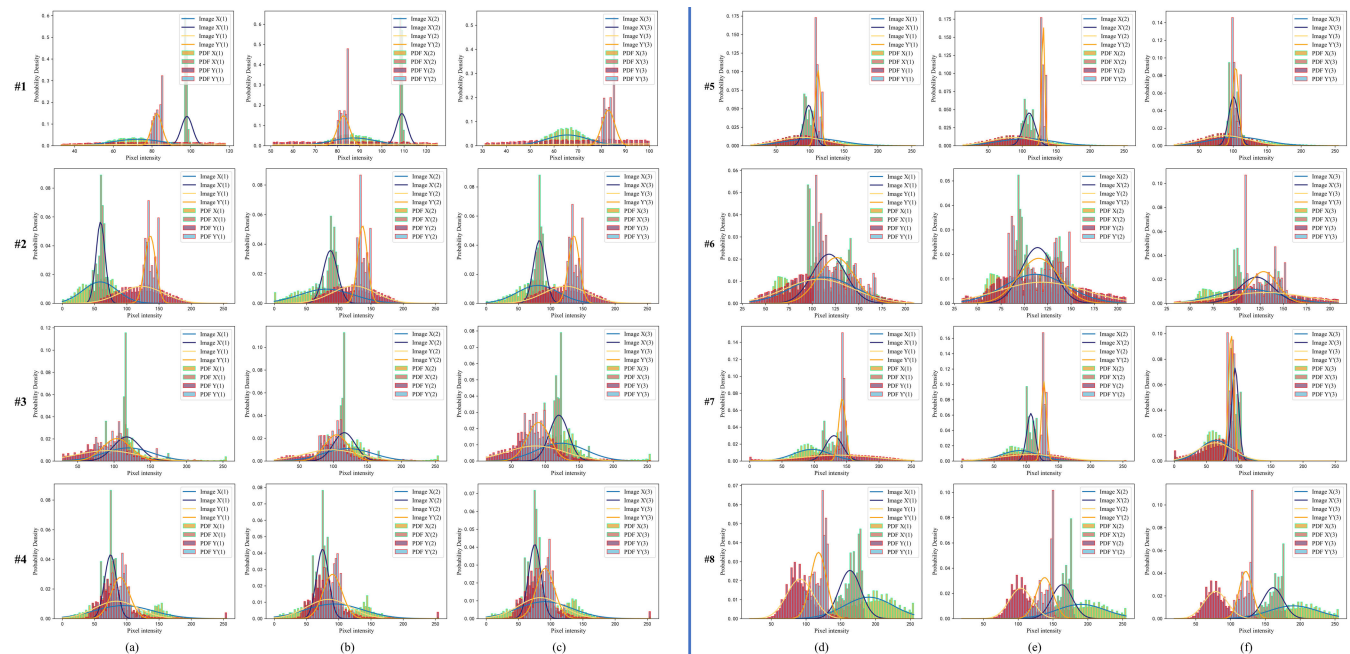


Fig. 7. Comparative analysis of histograms for the reconstructed image X' (Y') and the original image X (Y) across eight datasets. The datasets are divided into two groups. (a)–(c) Datasets #1–#4 and (d)–(f) datasets #5–#8. For each dataset, the histograms of the respective bands for both images are juxtaposed from left to right, allowing for a side-by-side comparison.

4.7% and 8.70% in OA, 55.28% and 67.10% in KC, and 39.26% and 46.28% in $F1$, respectively. Conversely, when

$loss_d$ and $loss_s$ are omitted from the total loss function, the model with the complete loss function still shows an average

TABLE V
EFFECTIVENESS OF GAM

	AT	#1	#2	#3	#4	#5	#6	#7	#8
OA	×	0.9379	0.9378	0.8786	0.9392	0.9397	0.9627	0.8722	0.9591
	√	0.9429	0.9627	0.8860	0.9615	0.9774	0.9530	0.9525	0.9645
KC	×	0.7263	0.6703	0.4748	0.6719	0.7404	0.7067	0.4513	0.6661
	√	0.7590	0.7067	0.4932	0.6978	0.7945	0.7902	0.5134	0.7185
F1	×	0.7604	0.7049	0.5432	0.7057	0.7741	0.7259	0.5236	0.6879
	√	0.7914	0.7259	0.5561	0.7178	0.8063	0.8169	0.5378	0.7374

TABLE VI
EFFECTIVENESS OF SPARSITY CONSTRAINT LOSS $loss_d$ AND STRUCTURAL CONSISTENCY LOSS $loss_s$

	$loss_r$	$loss_d$	$loss_s$	#1	#2	#3	#4	#5	#6	#7	#8
OA	√	×	×	0.9202	0.8309	0.8354	0.8447	0.9531	0.9197	0.8400	0.7979
	√	×	√	0.9133	0.9393	0.8361	0.9371	0.9629	0.9251	0.8722	0.9625
	√	√	×	0.9279	0.9591	0.8818	0.9346	0.9749	0.9468	0.9495	0.9484
	√	√	√	0.9429	0.9627	0.8860	0.9615	0.9774	0.9530	0.9525	0.9645
KC	√	×	×	0.7122	0.3158	0.2342	0.3516	0.6675	0.6842	0.2579	0.2910
	√	×	√	0.6952	0.6715	0.3701	0.6486	0.7193	0.7002	0.4513	0.5949
	√	√	×	0.7232	0.6661	0.4738	0.6659	0.7862	0.7757	0.4983	0.6061
	√	√	√	0.7590	0.7067	0.4932	0.6978	0.7945	0.7902	0.5134	0.7185
F1	√	×	×	0.7598	0.4070	0.3212	0.4318	0.6922	0.7292	0.3105	0.3615
	√	×	√	0.7470	0.7053	0.4670	0.6486	0.7391	0.7424	0.5236	0.6136
	√	√	×	0.7658	0.6879	0.5390	0.7022	0.7996	0.8060	0.5241	0.6333
	√	√	√	0.7914	0.7259	0.5561	0.7178	0.8063	0.8169	0.5378	0.7374

enhancement of 3.53% and 1.03% in OA, 14.16% and 5.55% in KC, and 10.08% and 4.41% in $F1$ score, respectively. These findings underscore the critical roles of $loss_d$ and $loss_s$ in refining the model's detection of change areas and preserving structural consistency. $loss_d$, by enforcing sparsity in change areas, aids in the precise localization of changes and serves as a regularization term to prevent overfitting, particularly in scenarios with subtle or minimal changes. $loss_s$, on the other hand, ensures the spatial structural consistency between the reconstructed and original images, which is crucial for identifying genuine changes and mitigating false positives due to noise, shadows, or other nonstructural elements.

E. Computational Time

This article presents an analysis of the computational time for the proposed SDC-GAE model, focusing on the smallest (dataset #8) and largest (dataset #5) datasets by size, as detailed in Table VII. The data reveal a direct correlation between the number of superpixels N_p and the neighbor ratio K_{ratio} with the computational time required for SDC-GAE. Notably, when N_p is set to a higher value, such as 10 000, an increase in K_{ratio} significantly extends the computational time. This increase is attributed to the fact that a larger N_p necessitates the processing of more neighbor vertices for each vertex in SDC-GAE, complicating the aggregation operation. Furthermore, managing a higher volume of vertices and their relationships demands additional memory for storing vertex features, weight matrices, and intermediate computation results. Consequently, the augmented memory requirements

TABLE VII
COMPUTATIONAL TIME (s) OF SDC-GAE

		$N_p=3000$	$N_p=5000$	$N_p=10000$
Dataset #8 $300 \times 412 \times 1(3)$	$K_{ratio}=0.02$	9.12	12.18	22.76
	$K_{ratio}=0.10$	17.91	32.33	623.45
Dataset #5 $4135 \times 2325 \times 3(1)$	$K_{ratio}=0.02$	120.02	163.68	289.29
	$K_{ratio}=0.10$	125.55	200.06	810.66

can potentially restrict the model's efficiency when implemented on hardware.

V. CONCLUSION

This research tackles the challenge of imaging feature discrepancies in MCD by introducing an innovative SDC-GAE model. Utilizing an unsupervised learning framework, SDC-GAE delves into the extraction of deep structural features from images. This model, augmented with a structural difference compensation mechanism, enables precise detection of change areas. The loss function of SDC-GAE is meticulously designed with three key components: image reconstruction loss, sparsity constraint loss, and structural consistency loss. These components work in tandem to guide the model in minimizing spectral differences between reconstructed and target images, focusing on change areas, and ensuring the structural features of both images remain consistent. Comparative experiments on eight multimodal datasets with the state-of-the-art methods have validated the effectiveness and superiority of SDC-GAE in MCD tasks.

SDC-GAE demonstrates its potential in processing multimodal remote sensing data, providing valuable technical support for CD in domains, such as natural resource monitoring, disaster assessment, and urban planning. Future research will focus on the continued refinement of the SDC-GAE's performance. The planned enhancements aim to address the following key areas.

1) Advanced graph representation learning: we suggest employing cutting-edge GCN frameworks to more effectively extract and encode image structural features. Integrating multiscale graph learning will facilitate the simultaneous capture of both local and global structural elements, thereby increasing the model's responsiveness to changes across different scales.

2) Refined attention mechanisms: we advocate for the development of more nuanced attention mechanisms that prioritize CD areas and minimize overfitting to stable regions. An adaptive attention distribution should be implemented, which adjusts focus based on image content, enabling more precise identification of change features.

3) Multitask learning: we recommend integrating CD with complementary tasks, such as classification and segmentation, to benefit from mutual information and enhance the precision of structural difference compensation.

4) Loss function optimization: the design of novel loss functions or the enhancement of the existing ones should be pursued to better address structural discrepancies. This could involve the integration of gradient-based loss components to diminish reconstruction inaccuracies.

5) Data augmentation and preprocessing: the use of data augmentation strategies, including image rotation, scaling, and brightness adjustments, during training can improve the model's adaptability to diverse change scenarios. Additionally, preprocessing steps, such as noise reduction and contrast enhancement, should be applied to refine input data quality and, consequently, the accuracy of structural difference compensation.

In the subsequent research, efforts will also be directed toward improving the model's real-time detection capabilities. This can be achieved through the following approaches: 1) employing techniques, such as model pruning and knowledge distillation to reduce the complexity of the model while ensuring its performance is maintained; 2) improving the computation of the loss function to minimize redundant calculations during the model's gradient descent; and 3) further increasing computational efficiency by utilizing parallel computing frameworks and specialized hardware accelerators.

Moreover, we intend to investigate the application of SDC-GAE to a broader range of remote sensing image analysis tasks, with the goal of offering more extensive technical support to related fields.

REFERENCES

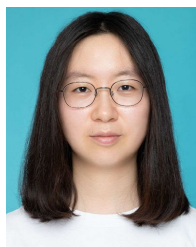
- [1] M. C. Hansen and T. R. Loveland, "A review of large area monitoring of land cover change using Landsat data," *Remote Sens. Environ.*, vol. 122, pp. 66–74, Jul. 2012, doi: [10.1016/j.rse.2011.08.024](https://doi.org/10.1016/j.rse.2011.08.024).
- [2] D. Brunner, G. Lemoine, and L. Bruzzone, "Earthquake damage assessment of buildings using VHR optical and SAR imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 48, no. 5, pp. 2403–2420, May 2010, doi: [10.1109/TGRS.2009.2038274](https://doi.org/10.1109/TGRS.2009.2038274).
- [3] Y. Tang, L. Zhang, and X. Huang, "Object-oriented change detection based on the Kolmogorov–Smirnov test using high-resolution multispectral imagery," *Int. J. Remote Sens.*, vol. 32, no. 20, pp. 5719–5740, Oct. 2011, doi: [10.1080/01431161.2010.507263](https://doi.org/10.1080/01431161.2010.507263).
- [4] Y. Tang and L. Zhang, "Urban change analysis with multi-sensor multispectral imagery," *Remote Sens.*, vol. 9, no. 3, p. 252, Mar. 2017, doi: [10.3390/rs9030252](https://doi.org/10.3390/rs9030252).
- [5] T. Bai et al., "Deep learning for change detection in remote sensing: A review," *Geo-Spatial Inf. Sci.*, vol. 26, no. 3, pp. 262–288, Jul. 2023, doi: [10.1009/10095020.2022.2085633](https://doi.org/10.1009/10095020.2022.2085633).
- [6] G. Camps-Valls, L. Gomez-Chova, J. Munoz-Mari, J. L. Rojo-Alvarez, and M. Martinez-Ramon, "Kernel-based framework for multitemporal and multisource remote sensing data classification and change detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 46, no. 6, pp. 1822–1835, Jun. 2008, doi: [10.1109/TGRS.2008.916201](https://doi.org/10.1109/TGRS.2008.916201).
- [7] L. Wan, Y. Xiang, and H. You, "A post-classification comparison method for SAR and optical images change detection," *IEEE Geosci. Remote Sens. Lett.*, vol. 16, no. 7, pp. 1026–1030, Jul. 2019, doi: [10.1109/LGRS.2019.2892432](https://doi.org/10.1109/LGRS.2019.2892432).
- [8] L. Wan, Y. Xiang, and H. You, "An object-based hierarchical compound classification method for change detection in heterogeneous optical and SAR images," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 12, pp. 9941–9959, Dec. 2019, doi: [10.1109/TGRS.2019.2930322](https://doi.org/10.1109/TGRS.2019.2930322).
- [9] T. Han, Y. Tang, X. Yang, Z. Lin, B. Zou, and H. Feng, "Change detection for heterogeneous remote sensing images with improved training of hierarchical extreme learning machine (HELM)," *Remote Sens.*, vol. 13, no. 23, p. 4918, Dec. 2021, doi: [10.3390/rs13234918](https://doi.org/10.3390/rs13234918).
- [10] J. Prendes, M. Chabert, F. Pascal, A. Giros, and J.-Y. Tourneret, "A new multivariate statistical model for change detection in images acquired by homogeneous and heterogeneous sensors," *IEEE Trans. Image Process.*, vol. 24, no. 3, pp. 799–812, Mar. 2015, doi: [10.1109/TIP.2014.2387013](https://doi.org/10.1109/TIP.2014.2387013).
- [11] R. Touati, M. Mignotte, and M. Dahmane, "Multimodal change detection in remote sensing images using an unsupervised pixel pairwise-based Markov random field model," *IEEE Trans. Image Process.*, vol. 29, pp. 757–767, 2020, doi: [10.1109/TIP.2019.2933747](https://doi.org/10.1109/TIP.2019.2933747).
- [12] R. Touati and M. Mignotte, "An energy-based model encoding nonlocal pairwise pixel interactions for multisensor change detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 2, pp. 1046–1058, Feb. 2018, doi: [10.1109/TGRS.2017.2758359](https://doi.org/10.1109/TGRS.2017.2758359).
- [13] L. Wan, T. Zhang, and H. J. You, "Multi-sensor remote sensing image change detection based on sorted histograms," *Int. J. Remote Sens.*, vol. 39, no. 11, pp. 3753–3775, Jun. 2018, doi: [10.1080/01431161.2018.1448481](https://doi.org/10.1080/01431161.2018.1448481).
- [14] M. Mignotte, "MRF models based on a neighborhood adaptive class conditional likelihood for multimodal change detection," *AI, Comput. Sci. Robot. Technol.*, vol. 2022, pp. 1–20, Mar. 2022, doi: [10.5772/acrt.02](https://doi.org/10.5772/acrt.02).
- [15] J. Liu, M. Gong, K. Qin, and P. Zhang, "A deep convolutional coupling network for change detection based on heterogeneous optical and radar images," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 29, no. 3, pp. 545–559, Mar. 2018, doi: [10.1109/TNNLS.2016.2636227](https://doi.org/10.1109/TNNLS.2016.2636227).
- [16] W. Zhao, Z. Wang, M. Gong, and J. Liu, "Discriminative feature learning for unsupervised change detection in heterogeneous images based on a coupled neural network," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 12, pp. 7066–7080, Dec. 2017, doi: [10.1109/TGRS.2017.2739800](https://doi.org/10.1109/TGRS.2017.2739800).
- [17] T. Han, Y. Tang, and Y. Chen, "Heterogeneous image change detection based on two-stage joint feature learning," in *Proc. IEEE Int. Geosci. Remote Sens. Symp.*, Jul. 2022, pp. 3215–3218, doi: [10.1109/IGARSS46834.2022.9883323](https://doi.org/10.1109/IGARSS46834.2022.9883323).
- [18] M. Yang et al., "Multicue contrastive self-supervised learning for change detection in remote sensing," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, 2023, Art. no. 5221014, doi: [10.1109/TGRS.2023.3330494](https://doi.org/10.1109/TGRS.2023.3330494).
- [19] R. Touati, M. Mignotte, and M. Dahmane, "Anomaly feature learning for unsupervised change detection in heterogeneous images: A deep sparse residual model," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 13, pp. 588–600, 2020, doi: [10.1109/JSTARS.2020.2964409](https://doi.org/10.1109/JSTARS.2020.2964409).
- [20] Y. Wu, J. Li, Y. Yuan, A. K. Qin, Q.-G. Miao, and M.-G. Gong, "Commonality autoencoder: Learning common features for change detection from heterogeneous images," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 33, no. 9, pp. 4257–4270, Sep. 2022, doi: [10.1109/TNNLS.2021.3056238](https://doi.org/10.1109/TNNLS.2021.3056238).
- [21] T. Zhan, M. Gong, X. Jiang, and S. Li, "Log-based transformation feature learning for change detection in heterogeneous images," *IEEE Geosci. Remote Sens. Lett.*, vol. 15, no. 9, pp. 1352–1356, Sep. 2018, doi: [10.1109/LGRS.2018.2843385](https://doi.org/10.1109/LGRS.2018.2843385).
- [22] M. Yang, L. Jiao, F. Liu, B. Hou, S. Yang, and M. Jian, "DPFL-nets: Deep pyramid feature learning networks for multiscale change detection," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 33, no. 11, pp. 6402–6416, Nov. 2022.

- [23] X. Jiang, G. Li, Y. Liu, X.-P. Zhang, and Y. He, "Change detection in heterogeneous optical and SAR remote sensing images via deep homogeneous feature fusion," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 13, pp. 1551–1566, 2020, doi: [10.1109/JSTARS.2020.2983993](https://doi.org/10.1109/JSTARS.2020.2983993).
- [24] H. Chen, F. He, and J. Liu, "Heterogeneous images change detection based on iterative joint global–local translation," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 15, pp. 9680–9698, 2022, doi: [10.1109/JSTARS.2022.3192251](https://doi.org/10.1109/JSTARS.2022.3192251).
- [25] X. Wang, W. Cheng, Y. Feng, and R. Song, "TSCNet: Topological structure coupling network for change detection of heterogeneous remote sensing images," *Remote Sens.*, vol. 15, no. 3, p. 621, Jan. 2023, doi: [10.3390/rs15030621](https://doi.org/10.3390/rs15030621).
- [26] H. Chen, N. Yokoya, C. Wu, and B. Du, "Unsupervised multimodal change detection based on structural relationship graph representation learning," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5635318, doi: [10.1109/TGRS.2022.3229027](https://doi.org/10.1109/TGRS.2022.3229027).
- [27] Z. Liu, G. Li, G. Mercier, Y. He, and Q. Pan, "Change detection in heterogeneous remote sensing images via homogeneous pixel transformation," *IEEE Trans. Image Process.*, vol. 27, no. 4, pp. 1822–1834, Apr. 2018, doi: [10.1109/TIP.2017.2784560](https://doi.org/10.1109/TIP.2017.2784560).
- [28] X. Li, Z. Du, Y. Huang, and Z. Tan, "A deep translation (GAN) based change detection network for optical and SAR remote sensing images," *ISPRS J. Photogramm. Remote Sens.*, vol. 179, pp. 14–34, Sep. 2021, doi: [10.1016/j.isprsjprs.2021.07.007](https://doi.org/10.1016/j.isprsjprs.2021.07.007).
- [29] L. T. Luppino, F. M. Bianchi, G. Moser, and S. N. Anfinsen, "Unsupervised image regression for heterogeneous change detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 12, pp. 9960–9975, Dec. 2019, doi: [10.1109/TGRS.2019.2930348](https://doi.org/10.1109/TGRS.2019.2930348).
- [30] M. Gong, P. Zhang, L. Su, and J. Liu, "Coupled dictionary learning for change detection from multisource data," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 12, pp. 7077–7091, Dec. 2016, doi: [10.1109/TGRS.2016.2594952](https://doi.org/10.1109/TGRS.2016.2594952).
- [31] Y. Sun, L. Lei, X. Li, X. Tan, and G. Kuang, "Patch similarity graph matrix-based unsupervised remote sensing change detection with homogeneous and heterogeneous sensors," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 6, pp. 4841–4861, Jun. 2021, doi: [10.1109/TGRS.2020.3013673](https://doi.org/10.1109/TGRS.2020.3013673).
- [32] Y. Sun, L. Lei, D. Guan, M. Li, and G. Kuang, "Sparse-constrained adaptive structure consistency-based unsupervised image regression for heterogeneous remote-sensing change detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 4405814, doi: [10.1109/TGRS.2021.3110998](https://doi.org/10.1109/TGRS.2021.3110998).
- [33] Y. Sun, L. Lei, X. Tan, D. Guan, J. Wu, and G. Kuang, "Structured graph based image regression for unsupervised multimodal change detection," *ISPRS J. Photogramm. Remote Sens.*, vol. 185, pp. 16–31, Mar. 2022, doi: [10.1016/j.isprsjprs.2022.01.004](https://doi.org/10.1016/j.isprsjprs.2022.01.004).
- [34] A. Radoi, "Generative adversarial networks under CutMix transformations for multimodal change detection," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022, doi: [10.1109/LGRS.2022.3201003](https://doi.org/10.1109/LGRS.2022.3201003).
- [35] Z.-G. Liu, Z.-W. Zhang, Q. Pan, and L.-B. Ning, "Unsupervised change detection from heterogeneous data based on image translation," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 4403413, doi: [10.1109/TGRS.2021.3097717](https://doi.org/10.1109/TGRS.2021.3097717).
- [36] X. Niu, M. Gong, T. Zhan, and Y. Yang, "A conditional adversarial network for change detection in heterogeneous images," *IEEE Geosci. Remote Sens. Lett.*, vol. 16, no. 1, pp. 45–49, Jan. 2019, doi: [10.1109/LGRS.2018.2868704](https://doi.org/10.1109/LGRS.2018.2868704).
- [37] D. Wang, F. Zhao, H. Yi, Y. Li, and X. Chen, "An unsupervised heterogeneous change detection method based on image translation network and post-processing algorithm," *Int. J. Digit. Earth*, vol. 15, no. 1, pp. 1056–1080, Dec. 2022, doi: [10.1080/17538947.2022.2092658](https://doi.org/10.1080/17538947.2022.2092658).
- [38] T. Han, Y. Tang, B. Zou, H. Feng, and F. Zhang, "Heterogeneous images change detection method based on hierarchical extreme learning machine image transformation," *J. Geo-Inf. Sci.*, vol. 24, no. 11, pp. 2212–2224, 2022, doi: [10.12082/dqxxkx.2022.220089](https://doi.org/10.12082/dqxxkx.2022.220089).
- [39] L. T. Luppino et al., "Code-aligned autoencoders for unsupervised change detection in multimodal remote sensing images," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 35, no. 1, pp. 60–72, Jan. 2024, doi: [10.1109/TNNLS.2022.3172183](https://doi.org/10.1109/TNNLS.2022.3172183).
- [40] L. T. Luppino et al., "Deep image translation with an affinity-based change prior for unsupervised multimodal change detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 4700422, doi: [10.1109/TGRS.2021.3056196](https://doi.org/10.1109/TGRS.2021.3056196).
- [41] Z. Du, X. Li, J. Miao, Y. Huang, H. Shen, and L. Zhang, "Concatenated deep-learning framework for multitask change detection of optical and SAR images," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 17, pp. 719–731, 2024, doi: [10.1109/JSTARS.2023.3333959](https://doi.org/10.1109/JSTARS.2023.3333959).
- [42] Z. Lv, J. Liu, W. Sun, T. Lei, J. A. Benediktsson, and X. Jia, "Hierarchical attention feature fusion-based network for land cover change detection with homogeneous and heterogeneous remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, 2023, Art. no. 4411115, doi: [10.1109/TGRS.2023.3334521](https://doi.org/10.1109/TGRS.2023.3334521).
- [43] Q. Liu, K. Ren, X. Meng, and F. Shao, "Domain adaptive cross reconstruction for change detection of heterogeneous remote sensing images via a feedback guidance mechanism," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, 2023, Art. no. 4507216, doi: [10.1109/TGRS.2023.3320805](https://doi.org/10.1109/TGRS.2023.3320805).
- [44] J. Wu et al., "A dual neighborhood hypergraph neural network for change detection in VHR remote sensing images," *Remote Sens.*, vol. 15, no. 3, p. 694, Jan. 2023, doi: [10.3390/rs15030694](https://doi.org/10.3390/rs15030694).
- [45] Y. Sun, L. Lei, D. Guan, and G. Kuang, "Iterative robust graph for unsupervised change detection of heterogeneous remote sensing images," *IEEE Trans. Image Process.*, vol. 30, pp. 6277–6291, 2021, doi: [10.1109/TIP.2021.3093766](https://doi.org/10.1109/TIP.2021.3093766).
- [46] T. Han, Y. Tang, B. Zou, and H. Feng, "Unsupervised multimodal change detection based on adaptive optimization of structured graph," *Int. J. Appl. Earth Observ. Geoinf.*, vol. 126, Feb. 2024, Art. no. 103630, doi: [10.1016/j.jag.2023.103630](https://doi.org/10.1016/j.jag.2023.103630).
- [47] M. Mignotte, "A fractal projection and Markovian segmentation-based approach for multimodal change detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 11, pp. 8046–8058, Nov. 2020, doi: [10.1109/TGRS.2020.2986239](https://doi.org/10.1109/TGRS.2020.2986239).
- [48] R. Touati, M. Mignotte, and M. Dahmane, "Multimodal change detection using a convolution model-based mapping," in *Proc. 9th Int. Conf. Image Process. Theory, Tools Appl. (IPTA)*, Nov. 2019, pp. 1–6, doi: [10.1109/IPTA.2019.8936127](https://doi.org/10.1109/IPTA.2019.8936127).
- [49] Y. Sun, L. Lei, X. Li, X. Tan, and G. Kuang, "Structure consistency-based graph for unsupervised change detection with homogeneous and heterogeneous remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 4700221, doi: [10.1109/TGRS.2021.3053571](https://doi.org/10.1109/TGRS.2021.3053571).
- [50] D. A. Jimenez-Sierra, H. D. Benítez-Restrepo, H. D. Vargas-Cardona, and J. Chanussot, "Graph-based data fusion applied to: Change detection and biomass estimation in Rice crops," *Remote Sens.*, vol. 12, no. 17, p. 2683, Aug. 2020, doi: [10.3390/rs12172683](https://doi.org/10.3390/rs12172683).
- [51] D. A. Jimenez-Sierra, D. A. Quintero-Olaya, J. C. Alvear-Muñoz, H. D. Benítez-Restrepo, J. F. Florez-Ospina, and J. Chanussot, "Graph learning based on signal smoothness representation for homogeneous and heterogeneous change detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 4410416, doi: [10.1109/TGRS.2022.3168126](https://doi.org/10.1109/TGRS.2022.3168126).
- [52] Y. Sun, L. Lei, D. Guan, J. Wu, and G. Kuang, "Iterative structure transformation and conditional random field based method for unsupervised multimodal change detection," *Pattern Recognit.*, vol. 131, Nov. 2022, Art. no. 108845, doi: [10.1016/j.patcog.2022.108845](https://doi.org/10.1016/j.patcog.2022.108845).
- [53] Y. Tang, X. Yang, T. Han, F. Zhang, B. Zou, and H. Feng, "Enhanced graph structure representation for unsupervised heterogeneous change detection," *Remote Sens.*, vol. 16, no. 4, p. 721, Feb. 2024, doi: [10.3390/rs16040721](https://doi.org/10.3390/rs16040721).
- [54] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," 2016, *arXiv:1609.02907*.
- [55] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Süsstrunk, "SLIC superpixels compared to state-of-the-art superpixel methods," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 11, pp. 2274–2282, Nov. 2012, doi: [10.1109/TPAMI.2012.120](https://doi.org/10.1109/TPAMI.2012.120).
- [56] A. Vaswani et al., "Attention is all you need," 2017, *arXiv:1706.03762*.
- [57] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," 2017, *arXiv:1710.10903*.
- [58] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, no. 6088, pp. 533–536, Oct. 1986, doi: [10.1038/323533a0](https://doi.org/10.1038/323533a0).
- [59] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Trans. Syst. Man, Cybern.*, vol. SMC-9, no. 1, pp. 62–66, Jan. 1979, doi: [10.1109/TSMC.1979.4310076](https://doi.org/10.1109/TSMC.1979.4310076).
- [60] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014, *arXiv:1412.6980*.



Te Han received the bachelor's and master's degrees in surveying and mapping science and technology from the School of Geosciences and Info-Physics, Central South University, Changsha, China, in 2017 and 2020, respectively, where he is currently pursuing the Ph.D. degree.

His research interests include machine learning and deep learning, as well as their application in remote sensing image change detection.



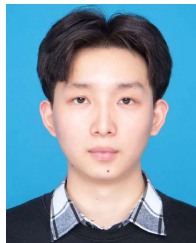
Xin Yang received the B.S. degree in geographic information systems (GIS) from Nanjing University of Posts and Telecommunications, Nanjing, China, in 2021. She is currently pursuing the M.Eng. degree in surveying engineering with the Department of Surveying and Remote Sensing, School of Geosciences and Info-Physics, Central South University, Changsha, China.

Her research interests include change detection using multimodal remote sensing imagery.



Yuqi Tang received the master's degree from Wuhan University of Technology, Wuhan, China, in 2008, and the Ph.D. degree in engineering from Wuhan University, Wuhan, in 2013.

Since 2013, she has been working at the School of Geosciences and Info-Physics, Central South University, Changsha, China, where she is an Associate Professor and a Doctoral Supervisor. She has published over 30 articles in SCI-indexed journals. Her research has long been focused on the intelligent interpretation of multispectral/hyperspectral remote sensing data and the analysis and monitoring of natural resources. Her research interests include multimodal remote sensing image change detection, satellite remote sensing video motion target detection, and hyperspectral remote sensing for water resource monitoring, among other applications in natural resource monitoring.



Yuqiang Guo received the B.S. degree in geographic information science from Zhengzhou University, Zhengzhou, China, in 2023. He is currently pursuing the M.Eng. degree in surveying engineering with the Department of Surveying and Remote Sensing, School of Geosciences and Info-Physics, Central South University, Changsha, China.

His research interests include land-cover/use change detection with homogeneous/heterogeneous remote sensing images.



Yuzeng Chen received the B.S. degree in geographic information science from the Southwest University of Science and Technology, Mianyang, China, in 2020, and the M.S. degree from Central South University, Changsha, China, in 2023. He is currently pursuing the Ph.D. degree with the School of Geodesy and Geomatics, Wuhan University, Wuhan, China.

His research interests include remote sensing/hyperspectral video object detection and tracking and change detection.



Shujing Jiang received the master's degree from Wuhan University School, Wuhan, China, in 2012.

She is currently working as a Senior Engineer at Guangzhou Urban Planning & Design Survey Research Institute, Guangzhou, China. Her research interests include surveying natural resources, monitoring and evaluation, early warning systems for land spatial planning, and the processing and application of real estate data.