

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

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

Te Han, Yuqi Tang, Yuzeng Chen, Xin Yang, Yuqiang Guo, 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 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 state-of-the-art methods.

Index Terms—Multimodal change detection, multi-source data, structural feature, structured graph, graph convolutional network, structural difference, compensation.

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

Manuscript received **** This work was Supported by the National Natural Science Foundation of China (Grant 42271411); the Scientific Research Innovation Project for Graduate Students in Hunan Province (No. CX20220169); Research Project on Monitoring and Early Warning Technologies for Implementation of Land Use Planning in Guangzhou City (2020B0101130009); Collaborative Innovation Center for Natural Resources Planning and Marine Technology of Guangzhou (No. 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. (yuzeng_chen@whu.edu.cn)

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

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

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: post-classification comparison, similarity measurement, feature space mapping, and image space transformation.

1) *Post-classification 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 multitemporal segmentation 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 post-classification comparison methods, several challenges remain: Firstly, classification errors inherent in these algorithms tend to accumulate during the comparison process, potentially diminishing the accuracy of CD. Secondly, 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. Thirdly, 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 multimodal change detection 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 [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 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 like deep pyramid feature learning networks (DPFL-Nets) [22], deep homogeneous feature fusion (DHFF) [23], iterative joint global-local translation (IJGLT) [24] and Topological structure coupling network (TSCNet) [25], structural relationship graph convolutional autoencoder (SR-GCAE) [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 change detection 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], unsupervised change detection (USCD) [35], conditional adversarial network (CAN) [36], image translation network and post-processing (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

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

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 change detection network (MTCNDN) [41] and deep translation based change detection 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 paper, we aim to propose an unsupervised MCD method to reduce reliance on label data. Unsupervised methods offer several advantages over supervised methods: a) 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; b) 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; c) 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 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 non-local 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 conditional random field models for unsupervised CD. To reduce the influence of changed areas on image structural features, methods like 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: a) They can overcome imaging differences in multimodal images by mining and comparing consistent structural features in unchanged regions, enhancing detection precision and robustness; b) 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; c) Encoding structural features into structure graphs can mitigate the effects of image noise on CD. However, these methods rely on traditional K-nearest neighbors (KNN) 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) *The potential of graph convolutional network (GCN) in extracting structural features of multimodal images.* GCN [54] are 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 address 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 paper 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, 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,

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

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} , 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 paper are:

- 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 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) Validation of the proposed method's effectiveness through experiments on 8 datasets and comparisons with 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.

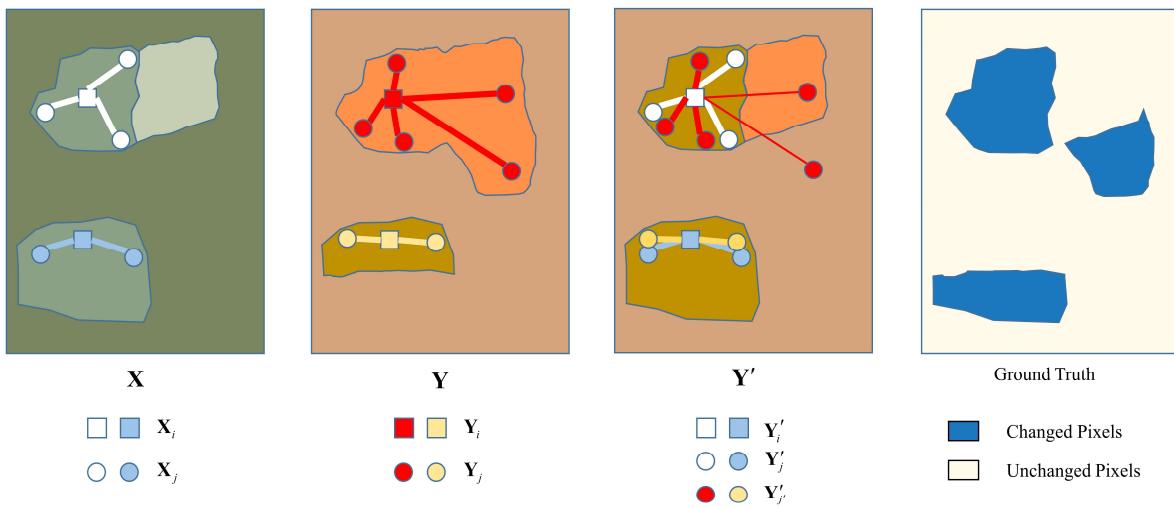
The method proposed in this paper consists of three primary components (Fig. 2): structural graph construction, SDC-GAE learning, and change map (CM) generation. The following sections will detail these steps.

A. Structural Graph Construction

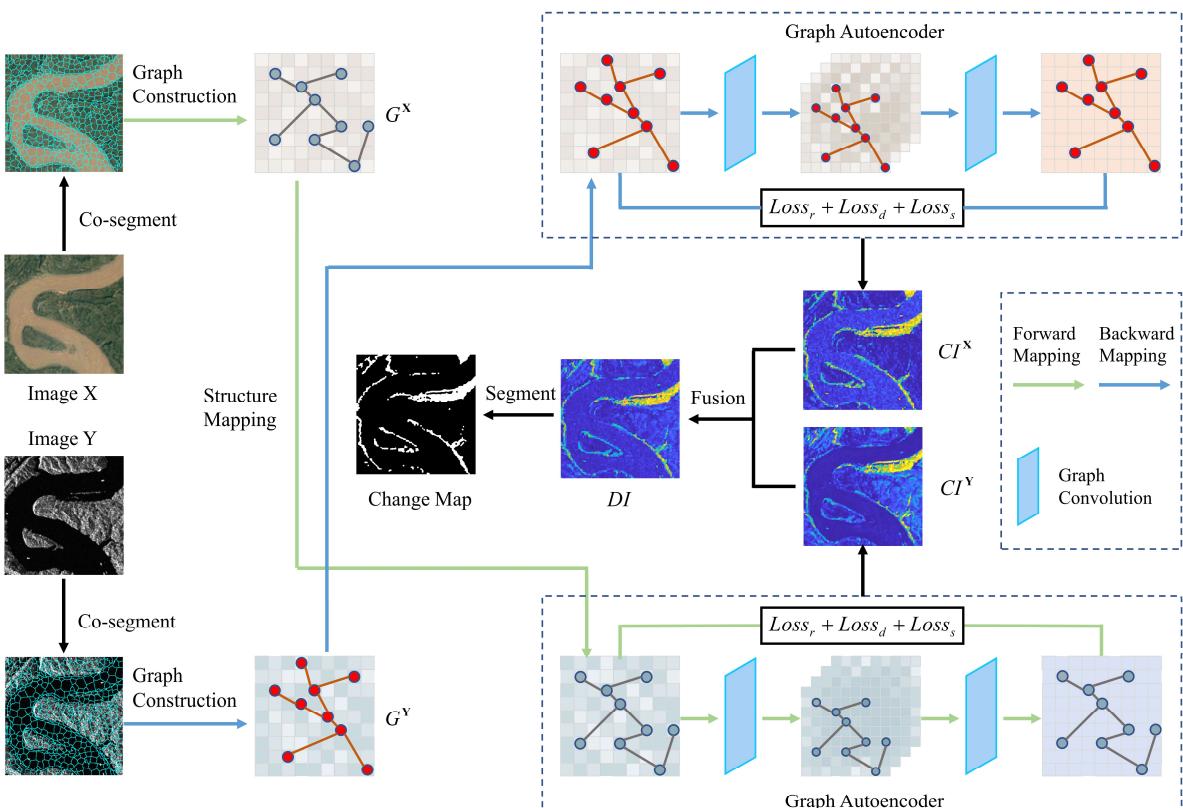
Structural graphs intuitively represent the structural features of images by using vertices and edges to illustrate the connections between them. In this study, we utilize structural graphs to depict the inherent structural features of images, with superpixels acting as graph vertices. Superpixels are 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 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 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 $P = \{P_i | i = 1, 2, \dots, N_p\}$. 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 $\tilde{\mathbf{X}} \in \mathbb{R}^{N_p \times 3B_X}$ and $\tilde{\mathbf{Y}} \in \mathbb{R}^{N_p \times 3B_Y}$ for images \mathbf{X} and \mathbf{Y} , respectively.

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <



1 Fig. 1. Schematic diagram of structural difference compensation in multimodal images. Squares and circles represent image objects, with the thickness of the
 2 connecting lines indicating the strength of their similarity. This similarity reflects the structural characteristics of the images. Image \mathbf{Y}' represents the spectral
 3 expression of image \mathbf{X} in the image domain \mathcal{V} . In the unchanged areas, the consistent structural features of image \mathbf{X} and \mathbf{Y} in that region allow the objects
 4 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}
 5 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 dif-
 6 ference 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
 7 compensation value reflects the intensity of the changes between multimodal images.
 8



10 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,
 11 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 $\mathbf{X}(\mathbf{Y})$
 12 within image $\mathbf{Y}(\mathbf{X})$. The process yields a measure of the intensity of changes within the images through structural difference compensation.
 13

14 Based on the obtained superpixels, we can construct a struc-
 tural graph G_X for image \mathbf{X} to represent its structural

features. We define $G_X = \{V_X, E_X\}$ as follows:

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

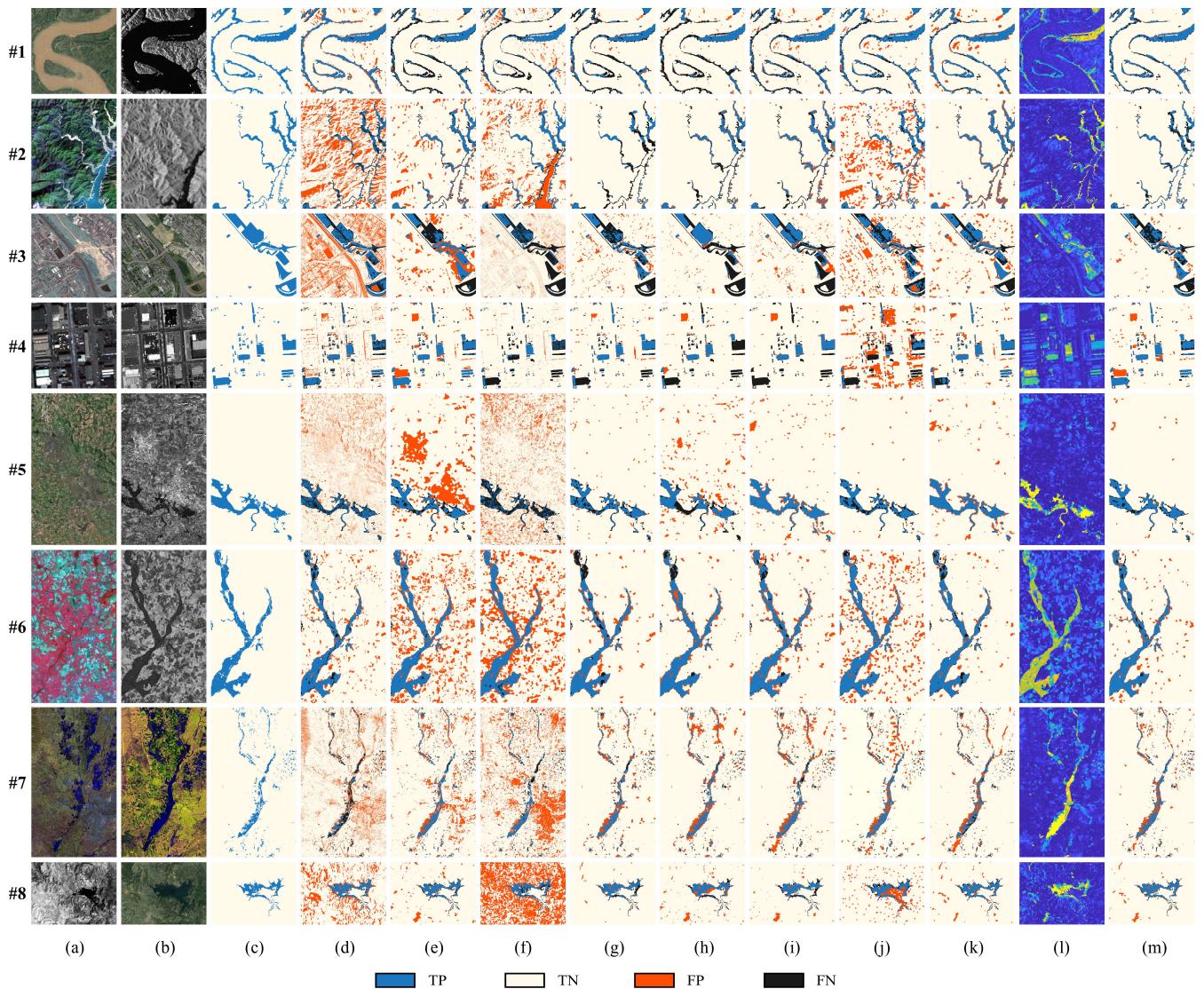


Fig. 3. CMs of different methods and the CIMs for SDC-GAE. From top to bottom, they correspond to datasets #1 to #8. From left to right, they are as follows: (a) image X ; (b) image Y ; (c) reference image; (d) CMs of LTFL; (e) CMs of INLPG; (f) CMs of GBF; (g) CMs of IRG-McS; (h) CMs of SCASC; (i) CMs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

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, 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 changes in buildings and vehicles, saw LTFL 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 in the lower left corner. Datasets #5 to #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

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

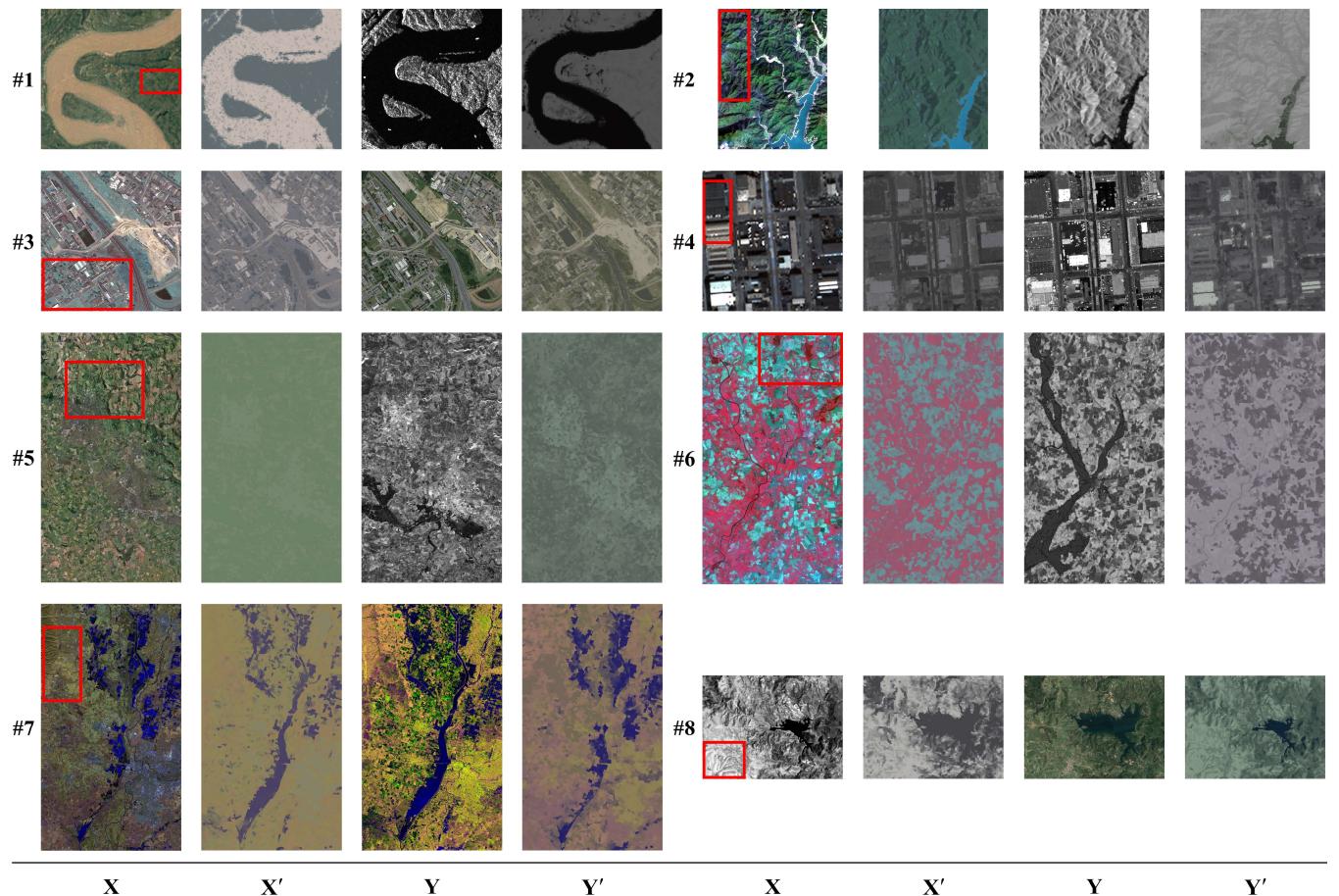


Fig. 6. The reconstructed images of SDC-GAE on datasets #1 to #8 are presented. Within each dataset, from left to right, they are as follows: image \mathbf{X} , the reconstructed image \mathbf{X}' of image \mathbf{Y} in the domain of \mathcal{X} , image \mathbf{Y} , and the reconstructed image \mathbf{Y}'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

GAM's ability to direct SDC-GAE's focus towards 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 paper 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 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 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 non-structural elements.

E. Computational Time

This paper 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 reveals a direct correlation between the

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

- [54] T. N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks." arXiv, Feb. 22, 2017. Accessed: Feb. 25, 2024. [Online]. Available: <http://arxiv.org/abs/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.
- [56] A. Vaswani *et al.*, "Attention Is All You Need." arXiv, Aug. 01, 2023. Accessed: Feb. 24, 2024. [Online]. Available: <http://arxiv.org/abs/1706.03762>
- [57] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph Attention Networks." arXiv, Feb. 04, 2018. Accessed: Feb. 24, 2024. [Online]. Available: <http://arxiv.org/abs/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.
- [59] N. Otsu, "A Threshold Selection Method from Gray-Level Histograms," *IEEE Trans. Syst., Man, Cybern.*, vol. 9, no. 1, pp. 62–66, Jan. 1979, doi: 10.1109/TSMC.1979.4310076.
- [60] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization." arXiv, Jan. 29, 2017. Accessed: Feb. 24, 2024. [Online]. Available: <http://arxiv.org/abs/1412.6980>



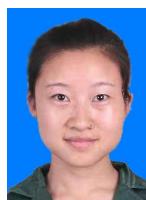
Xin Yang received the B.S. degree in Geographic Information Systems (GIS) from Nanjing University of Posts and Telecommunications, Nanjing, China, in 2021. She is currently pursuing her M.Eng. in Surveying Engineering within the Department of Surveying and Remote Sensing at the School of Geosciences and Info-Physics, Central South University, China. Her research interest lies in change detection using multimodal remote sensing imagery.



Yuqiang Guo received the B.S. degree in Geographic Information Science from Zhengzhou University, Zhengzhou, China, in 2023. He is studying for the M.Eng. degree in Surveying Engineering with the Department of Surveying and Remote Sensing, School of Geosciences and Info-Physics, Central South University, China. His research interest is land-cover/use change detection with homogeneous/heterogeneous remote sensing images.



Te Han received the Bachelor's and Master's degrees in Surveying and Mapping Science and Technology from School of Geosciences and Info-Physics at Central South University in 2017 and 2020, respectively. He is currently pursuing his Ph.D. at the same institution. His research interests include machine learning and deep learning, as well as their application in remote sensing image change detection.



Shujing Jiang graduated from Wuhan University School in 2012 with a Master's degree. She is currently working as a Senior Engineer at the Guangzhou Urban Planning & Design Survey Research Institute, 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.



Yuqi Tang received the Master's degree from Wuhan University of Technology in 2008, and a Ph.D. in Engineering from Wuhan University in 2013. Since 2013, she has been working at the School of Geosciences and Info-Physics at Central South University, where she holds the position of Associate Professor and is a doctoral supervisor. She has published over 30 papers 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 includes 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.



Yuzeng Chen received the B.S. degree in geographic information science from Southwest University of Science and Technology, Mianyang, China, in 2020 and the M.S. degree from Central South University, in 2023, Changsha, China. He is pursuing the Ph.D. degree with the School of Geodesy and Geomatics, Wuhan University, Wuhan. His research interests include remote sensing/hyperspectral video object detection and tracking, change detection. More details can be found at <https://yzcu.github.io>. e-mail: yuzeng_chen@whu.edu.cn