

Topology-Constrained Graph-Mamba with Logit Diffusion for Change Detection

Anonymous CVPR submission

Paper ID 7

Abstract

001 Remote sensing change detection requires distinguishing
 002 meaningful structural changes from appearance varia-
 003 tions such as illumination or seasonal shifts. We pro-
 004 pose **GMD-CD** (Graph-Mamba Diffusion for Change De-
 005 tection), a framework that combines topology-constrained
 006 cross-temporal modeling with efficient refinement. Bi-
 007 temporal features are organized as a graph-structured bi-
 008 partite representation, where regions from each time step
 009 are explicitly separated and jointly processed. Global in-
 010 teraction is performed using a bidirectional state-space
 011 model, enabling structured propagation with linear com-
 012 plexity. To improve robustness, we introduce orthogonal
 013 feature disentanglement that separates change-related and
 014 invariant components, reducing interference from appear-
 015 ance variations. Building on this representation, we design
 016 a conditional diffusion decoder in logit space that performs
 017 uncertainty-aware residual refinement, improving boundary
 018 quality while preserving global consistency. Experiments
 019 on multiple benchmarks show consistent gains in both ac-
 020 curacy and efficiency, yielding a strong accuracy–efficiency
 021 trade-off for high-resolution change detection.

022 1. Introduction

023 Change detection (CD) in remote sensing aims to identify
 024 meaningful structural differences between images acquired
 025 at different times, and is fundamental to applications such
 026 as urban monitoring [4], disaster assessment [21], and envi-
 027 ronmental analysis. Despite significant progress, three chal-
 028 lenges remain: (i) **pseudo-changes** caused by illumination,
 029 seasonal effects, or sensor variation, (ii) **boundary degra-**
 030 **dation** for thin and irregular regions, and (iii) **the need**
 031 **for global context** to distinguish true changes efficiently.
 032 Existing methods address these challenges with different
 033 trade-offs. CNN-based approaches [16, 19] are efficient but
 034 rely on local receptive fields, while transformer-based meth-
 035 ods [2, 3, 7] capture long-range dependencies at quadratic
 036 cost. State-space models (SSMs) [20] provide global in-
 037 teraction with linear complexity, and recent Mamba-based

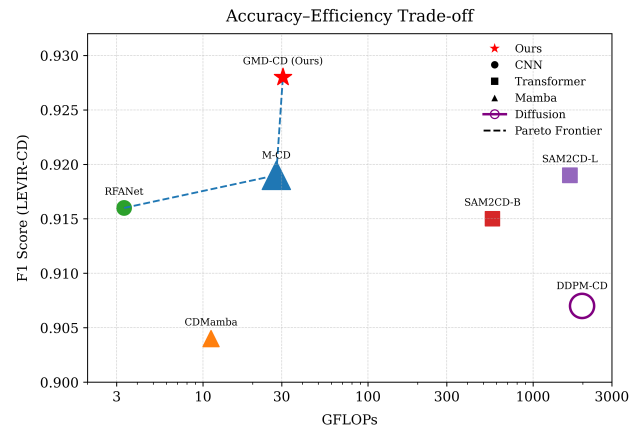


Figure 1. Accuracy–efficiency trade-off across representative CD methods.

038 CD methods [10, 34, 44] improve efficiency. However, they
 039 typically process features as a single flattened sequence, ig-
 040 noring cross-temporal structure and leading to unstructured
 041 information mixing. Diffusion-based approaches such as
 042 DDPM-CD [1] improve prediction quality through stochas-
 043 tic refinement but incur high computational cost. As shown
 044 in Fig. 1, existing methods either favor efficiency (e.g., CD-
 045 Mamba [44]) or accuracy at high cost (e.g., DDPM-CD,
 046 SAM2-CD [35]). This motivates a formulation that enables
 047 structured cross-temporal reasoning with efficiency. We
 048 propose **GMD-CD**, a unified framework that models cross-
 049 temporal interaction as a *structured state-space propagation*
 050 *problem*. We impose a *bipartite structural constraint* over
 051 bi-temporal features, representing them as an ordered bipar-
 052 tite sequence and performing bidirectional state-space prop-
 053 agation. This induces topology-aware cross-temporal cou-
 054 pling and enables global interaction with linear complex-
 055 ity, avoiding unstructured mixing in prior Mamba-based ap-
 056 proaches. Building on this, we introduce **Orthogonal Fea-**
 057 **ture Disentanglement (OFD)** to separate change-sensitive
 058 and invariant components, improving robustness to pseudo-
 059 changes. To address boundary ambiguity, we design a **con-**
 060 **ditional diffusion decoder** operating in logit space, which

061 performs stochastic refinement to correct uncertain predic- 113
 062 tions while preserving efficiency. 114

063 Contributions. 115

- 064 • We introduce a **structured bipartite state-space for-** 116
 065 **mulation** that models cross-temporal interaction as 117
 066 topology-constrained propagation. 118
- 067 • We propose a **BGMO module** that performs partition- 119
 068 aware bidirectional state-space propagation with linear 120
 069 complexity. 121
- 070 • We design **Orthogonal Feature Disentanglement** 122
 071 **(OFD)** to decouple change and invariant components, im- 123
 072 proving robustness to pseudo-changes. 124
- 073 • We develop a **conditional diffusion decoder in logit** 125
 074 **space** for uncertainty-aware refinement with moderate 126
 075 overhead. 127
- 076 • We demonstrate consistent improvements across bench- 128
 077 marks while maintaining a favorable accuracy–efficiency 129
 078 trade-off. 130

079 2. Related Work 131

080 **CNN and transformer-based change detection.** Early 132
 081 CD methods use Siamese CNN architectures with fea- 133
 082 ture differencing and encoder–decoder designs [16, 19]. 134
 083 Attention-based extensions improve spatial and temporal 135
 084 reasoning [4, 7]. Transformer models such as BIT [8] 136
 085 and ChangeFormer [3] capture long-range dependencies 137
 086 via self- and cross-attention, but incur quadratic complex- 138
 087 ity. Recent approaches scale performance with larger back- 139
 088 bones [41], while foundation-based methods like SAM2- 140
 089 CD [35] achieve strong generalization at high computa- 141
 090 tional cost. In contrast, we target efficient cross-temporal 142
 091 modeling with linear complexity. 143

092 **State-space models for vision and change detec-** 144
 093 **tion.** Selective state-space models (SSMs), particularly 145
 094 Mamba [20], enable global interaction with linear complex- 146
 095 ity. Vision adaptations serialize spatial features into se- 147
 096 quences [47], and CD variants such as ChangeMamba [10], 148
 097 CDMamba [44], and RSMamba [12] achieve strong effi- 149
 098 ciency. However, they treat bi-temporal features as a sin- 150
 099 gle unstructured sequence, limiting explicit cross-temporal 151
 100 modeling. We instead introduce a structured bipartite for- 152
 101 mulation with partition-aware state-space propagation. 153

102 **Graph and structured modeling.** Graph-based ap- 154
 103 proaches model spatial relationships via explicit node inter- 155
 104 actions and message passing, but incur additional computa- 156
 105 tional overhead. We adopt a lightweight alternative by or- 157
 106 ganizing pooled regions into an ordered bipartite structure, 158
 107 where interactions are implicitly realized through state- 159
 108 space transitions without explicit graph construction. 160

109 **Diffusion-based change detection.** Diffusion models 157
 110 have been applied to CD as generative or refinement mech- 158
 111 anisms. DDPM-CD [1] performs full diffusion-based gen- 159
 112 eration with high cost, while other works explore generative 160

modeling for representation learning [25]. In contrast, we 113
 use diffusion as a conditional decoder in logit space, en- 114
 abling targeted refinement with a small number of DDIM 115
 steps and moderate overhead. 116

3. Method 117

3.1. Overall Architecture 118

Given bi-temporal images (I^{t_1}, I^{t_2}) , we extract hierarchical 119
 features using a shared Swin-Tiny backbone [29], produc- 120
 ing multi-scale representations $\{f_l^a, f_l^b\}_{l=1}^4$, as illustrated 121
 in Fig. 2. The model consists of three stages: (i) *multi-* 122
scale state encoding, (ii) *structured state propagation*, and 123
 (iii) *logit refinement*. At shallow levels ($l = 1, 2$), we 124
 apply *LightFusion* to capture local differences using fea- 125
 ture interaction and residual aggregation. These stages fo- 126
 cus on fine-grained appearance variations. At deeper lev- 127
 els ($l = 3, 4$), we perform global cross-temporal modeling 128
 using the proposed *Graph-structured State-Space Operator* 129
(BGMO). Features are pooled into a bipartite node repres- 130
 entation and arranged in a partitioned sequence, which pre- 131
 serves temporal identity. Bidirectional state-space propa- 132
 gation [20] then enables structured cross-temporal interac- 133
 tion with linear complexity. The resulting representation is 134
 refined using *Orthogonal Feature Disentanglement (OFD)*, 135
 which separates change-related and invariant components. 136
 The multi-scale features are fused using a lightweight top- 137
 down decoder to produce coarse change logits L_0 . Finally, 138
 a *conditional diffusion decoder* operates in logit space and 139
 refines L_0 using a small number of denoising steps, improv- 140
 ing boundary quality while preserving global consistency. 141

3.2. Multi-scale State Encoding 142

We extract hierarchical features from bi-temporal inputs 143
 (I^{t_1}, I^{t_2}) using a shared encoder, producing $\{f_l^a, f_l^b\}_{l=1}^4$ at 144
 multiple resolutions. These features capture spatial context 145
 at increasing receptive fields and serve as inputs for sub- 146
 sequent state modeling. At shallow levels ($l \in \{1, 2\}$), 147
 changes are mostly local. Following prior CD meth- 148
 ods [16, 19], we encode both differences and agreements: 149

$$x_l = \text{Conv}_{1 \times 1}([\lvert f_l^a - f_l^b \rvert; f_l^a \odot f_l^b]), \quad (1) \quad 150$$

where $\lvert \cdot \rvert$ captures change magnitude and \odot captures shared 151
 responses. This provides a compact local representation 152
 sensitive to both variation and consistency. To improve spa- 153
 tial selectivity, we apply lightweight channel and spatial at- 154
 tention [38]: 155

$$S_l = x_l \cdot \sigma(\text{Conv}_{7 \times 7}([\text{AvgPool}(x_l); \text{MaxPool}(x_l)])), \quad (2) \quad 156$$

which suppresses background noise and enhances change- 157
 relevant regions with minimal overhead. At deeper levels 158
 $(l \in \{3, 4\})$, features encode higher-level semantics. In- 159
 stead of explicit differencing, we preserve their structure 160

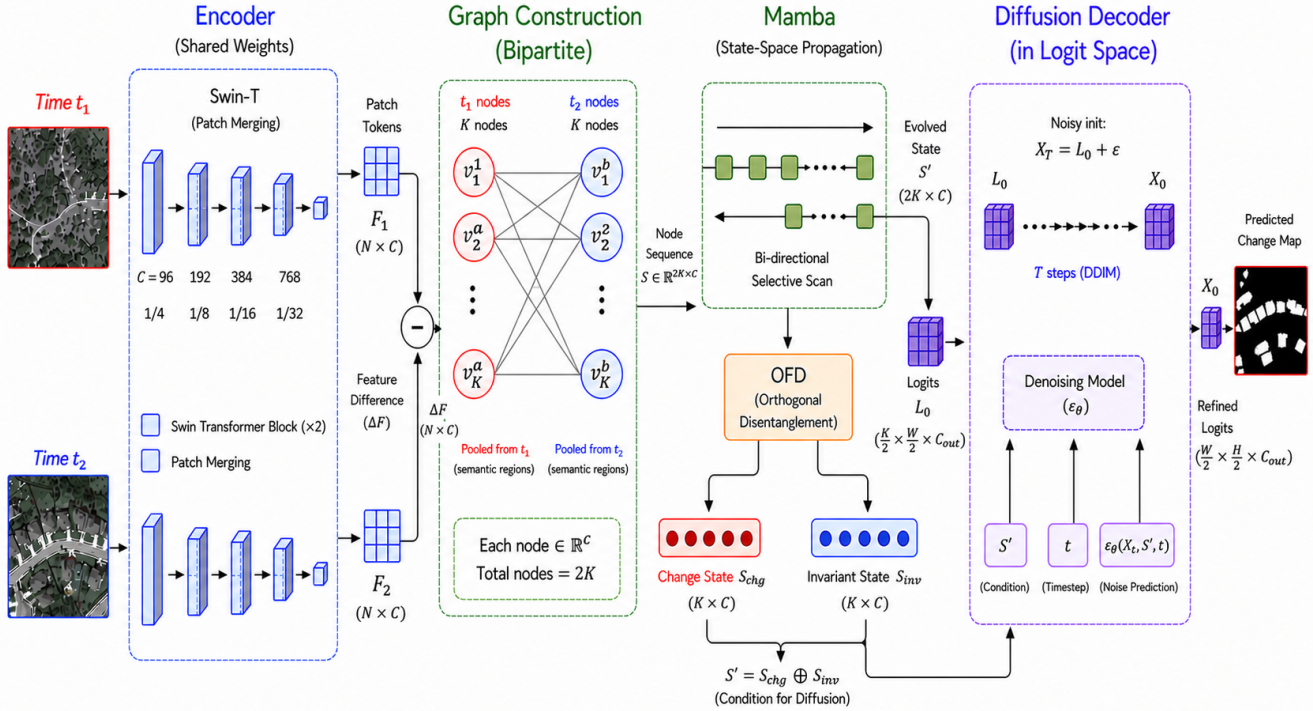


Figure 2. Overview of the proposed GMD-CD framework. Bi-temporal images are encoded into multi-scale features, followed by structured cross-temporal modeling using the Bipartite Graph-Mamba Operator (BGMO). The resulting representations are disentangled into change and invariant components via OFD and fused to produce coarse logits. A conditional diffusion decoder then refines the logits in logit space, improving boundary quality and reducing false detections.

161 and project them into a shared state-space:

162
$$S_l^a = W_l f_l^a, \quad S_l^b = W_l f_l^b, \quad (3)$$

163 where W_l aligns feature dimensions across scales. This design yields multi-scale states $\{S_l\}$ that capture both local variations and global semantics, providing stable inputs for structured cross-temporal state propagation.

167 3.3. Bipartite Graph-Mamba Operator (BGMO)

168 As illustrated in Fig. 3, BGMO organizes bi-temporal features into a structured bipartite representation and performs bidirectional state-space propagation over the resulting node sequence. At deeper scales ($l \in \{3, 4\}$), we model cross-temporal interaction as a *structured state-space evolution over a bipartite representation*. Given projected features $q^a = W_a f_l^a$ and $q^b = W_b f_l^b$ (Alg. 1, 1.1), we aggregate spatial regions into node embeddings $N^a, N^b \in \mathbb{R}^{B \times K \times C}$ via adaptive pooling (1.2–3). The two partitions correspond to t_1 and t_2 , and are concatenated into an ordered sequence $S = [N^a; N^b]$ (1.4), which explicitly preserves temporal identity while enabling joint processing. Interaction is performed using a bidirectional selective state-space model (Mamba) [20] (1.6–8). Unlike prior SSM-

Algorithm 1 Bipartite Graph-Mamba Operator (BGMO)

Input: $f_l^a, f_l^b \in \mathbb{R}^{C \times H_l \times W_l}$
Output: $S_l' \in \mathbb{R}^{C \times H_l \times W_l}$

- 1: $q^a \leftarrow W_a f_l^a, \quad q^b \leftarrow W_b f_l^b$
- 2: $N^a \leftarrow \text{AdaptivePool}(q^a, \sqrt{K} \times \sqrt{K})$
- 3: $N^b \leftarrow \text{AdaptivePool}(q^b, \sqrt{K} \times \sqrt{K})$
- 4: $S \leftarrow [N^a; N^b]$
- 5: $Z \leftarrow \text{Linear}(\text{SiLU}(\text{Linear}(\text{LayerNorm}(S))))$
- 6: $Z^{\rightarrow} \leftarrow \text{SSM}(Z)$
- 7: $Z^{\leftarrow} \leftarrow \text{reverse}(\text{SSM}(\text{reverse}(Z)))$
- 8: $S^{(1)} \leftarrow Z^{\rightarrow} + Z^{\leftarrow} + S$
- 9: $R \leftarrow \text{Linear}(\text{GELU}(\text{Linear}(\text{LayerNorm}(S^{(1)}))))$
- 10: $S^{(2)} \leftarrow R + S^{(1)}$
- 11: $(S_{\text{chg}}, S_{\text{inv}}) \leftarrow \text{ChannelSplit}(S^{(2)})$
- 12: $U \leftarrow \text{Upsample}(\text{LinearExpand}(S_{\text{chg}}))$
- 13: $U \leftarrow \text{SiLU}(\text{LayerNorm}(\text{Conv}_{3 \times 3}(U)))$
- 14: $S_l' \leftarrow f_l^a + f_l^b + \gamma \cdot U$

based CD methods that process flattened, unstructured sequences, the proposed bipartite ordering induces structured cross-temporal coupling: forward and backward scans repeatedly traverse both partitions, enabling each node to incorporate information from both temporal partitions under bidirectional traversal. This does not explicitly mask intra-

partition interactions, but biases propagation toward cross-temporal mixing through repeated traversal under bidirectional dynamics. As a result, global interaction is achieved with linear complexity, without pairwise attention or explicit graph construction. The propagated sequence is refined through residual feed-forward transformation (1.9–10) and decomposed into complementary subspaces (1.11). The change component is projected back to the spatial domain (1.12–13) and fused with the original features via a gated residual (1.14), enabling stable optimization and progressive integration of propagated context. Overall, BGMO reformulates cross-temporal reasoning as structured state evolution over a bipartite sequence, providing a structured alternative to unstructured sequence modeling in SSM-based change detection while maintaining efficiency.

3.4. Orthogonal Feature Disentanglement (OFD)

The BGMO representation encodes both change-sensitive variations and invariant contextual structure, which can interfere during prediction. We introduce *Orthogonal Feature Disentanglement (OFD)* to explicitly separate these components at the node level, enabling stable state evolution and reducing pseudo-change leakage. Given $\mathbf{S} \in \mathbb{R}^{B \times 2K \times C}$, we partition channels into complementary subspaces (\mathbf{S}_{chg} , \mathbf{S}_{inv}) with a fixed channel partition. To enforce separation, we apply channel-wise ℓ_2 normalization and impose an orthogonality constraint computed per node and averaged over the batch:

$$\mathcal{L}_{\text{orth}} = \|\mathbf{S}_{\text{chg}}^\top \mathbf{S}_{\text{inv}}\|_F^2, \quad (4)$$

which suppresses shared directions and reduces feature leakage between subspaces. To avoid degenerate solutions, we regularize the change component with $\mathcal{L}_{\text{var}} = \|\mathbf{S}_{\text{chg}}\|_F^2$, ensuring it remains informative. The subspaces are used asymmetrically: \mathbf{S}_{chg} is projected back to the spatial domain for prediction, while \mathbf{S}_{inv} is retained in the residual pathway to preserve context. This separation is necessary as BGMO propagates features across time; without disentanglement, invariant responses can contaminate change signals during propagation. By enforcing orthogonal subspaces on graph-evolved features, OFD stabilizes representation dynamics and improves robustness to appearance variations in subsequent refinement.

3.5. Hierarchical Fusion and Logit Projection

The graph-evolved features $\{S'_l\}_{l=1}^4$ form a multi-scale hierarchy, where deeper levels encode global context and shallower levels retain spatial detail. We adopt a top-down fusion scheme to propagate high-level information to higher-resolution features, similar to feature pyramid designs [27]. Fusion is performed recursively:

$$\tilde{S}_l = S'_l + \phi(\text{Up}(\tilde{S}_{l+1})), \quad l \in \{3, 2, 1\}, \quad (5)$$

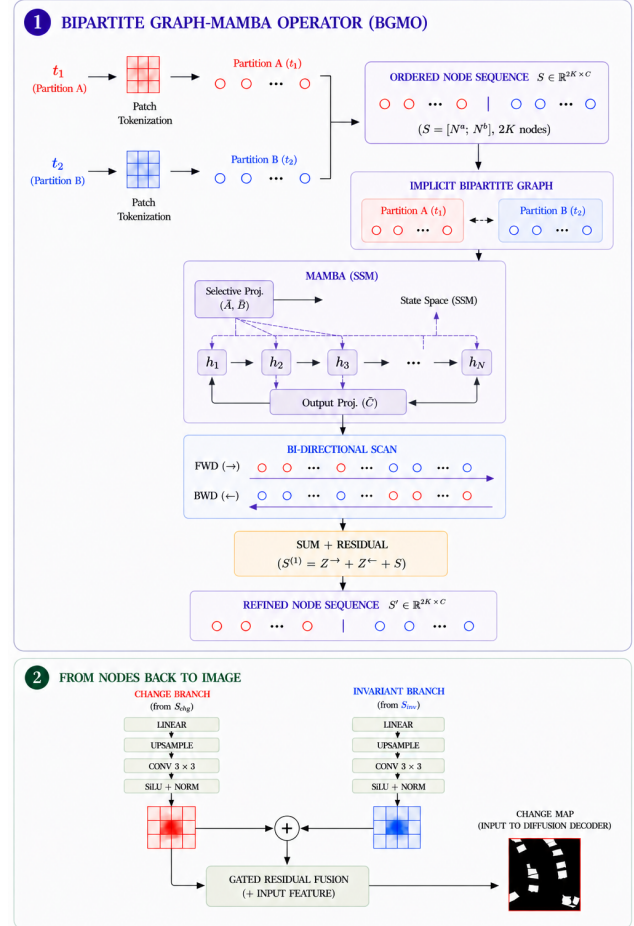


Figure 3. Architecture of Bipartite Graph-Mamba Operator (BGMO).

with $\tilde{S}_4 = S'_4$. Here, $\text{Up}(\cdot)$ denotes bilinear upsampling, and $\phi(\cdot)$ is a 3×3 convolution for channel alignment and local smoothing. This yields a consistent multi-scale representation where coarse features provide global priors for finer scales. The final fused feature $\tilde{S}_1 \in \mathbb{R}^{C \times H \times W}$ is projected to logit space:

$$L_0 = W_o \tilde{S}_1, \quad W_o \in \mathbb{R}^{1 \times C \times 3 \times 3}, \quad (6)$$

producing a coarse change logit map $L_0 \in \mathbb{R}^{1 \times H \times W}$. This hierarchical fusion integrates global context with local detail, providing a stable and informative initialization for subsequent diffusion-based refinement.

3.6. Conditional Diffusion Decoder in Logit Space

As shown in Fig. 4, we formulate refinement as a conditional diffusion process operating directly in logit space. Given the fused representation \tilde{S}_1 , we obtain a coarse logit map $L_0 \in \mathbb{R}^{1 \times H \times W}$ that captures globally consistent change structure from BGMO. However, L_0 is deter-

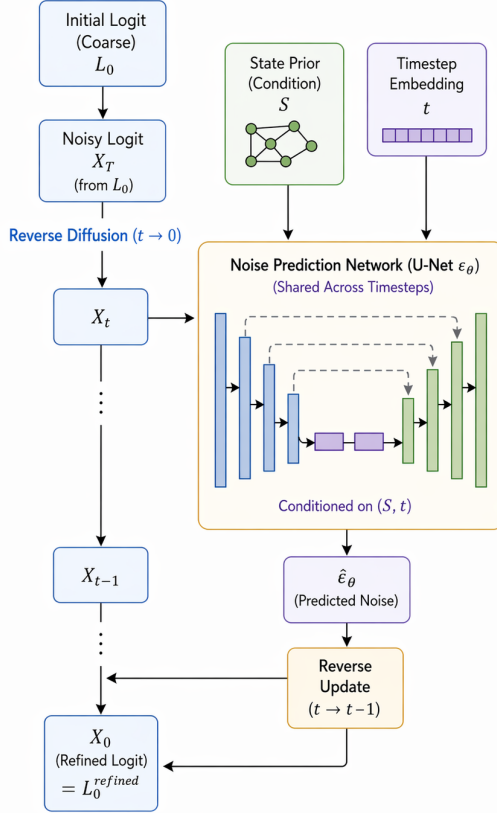


Figure 4. Architecture of Conditional Diffusion Decoding.

ministic and tends to be over-confident in ambiguous regions such as boundaries and small objects. To address this, we formulate refinement as a *conditional diffusion process* directly in logit space. We adopt a DDPM formulation [23, 33] conditioned on the graph-evolved state prior $S \in \mathbb{R}^{C \times H \times W}$. Unlike prior diffusion-based CD methods that perform full generative modeling [1], we restrict diffusion to a *residual correction regime* around L_0 . During training, the model learns to predict noise from perturbed logits; at inference, we initialize from $x_T = L_0 + \epsilon$, ensuring the trajectory remains close to a meaningful solution. The denoiser $\epsilon_\theta(x_t, S, t)$ is implemented as a lightweight U-Net and conditioned via feature concatenation with a projected state prior, enabling each step to be guided by globally structured context. Training uses noisy logits derived from the same forward pipeline, reducing train–test mismatch when initialized from L_0 . We employ a deterministic DDIM sampler [37] with a small number of steps ($T = 5$), resulting in efficient inference. In contrast to deterministic decoders that produce point estimates, diffusion captures uncertainty through stochastic perturbation and iterative denoising, enabling targeted correction in regions where predictions are unreliable. Operating in logit space constrains the process to refinement ($x_0 = L_0 + \Delta L$), where x_0 de-

notes the ground-truth logit obtained via inverse sigmoid of the binary mask. This helps preserve global semantics while correcting high-frequency errors. This design complements BGMO: structured state propagation ensures global consistency, while diffusion performs uncertainty-aware local refinement, improving boundary fidelity and reducing false positives without the cost of full generative diffusion.

4. Experiments

4.1. Datasets

We evaluate our method on four widely used remote sensing change detection benchmarks: *LEVIR-CD* [45], *DSIFN-CD* [39], *WHU-CD* [24], and *CDD* [6], which collectively cover building-level changes, multi-class land-cover variations, and seasonal transformations. *LEVIR-CD* consists of 637 high-resolution (1024×1024 , 0.5 m/pixel) Google Earth image pairs with annotated building changes, which are partitioned into non-overlapping 256×256 patches following standard practice ($7120/1024/2048$ for training/validation/testing). *DSIFN-CD* provides large-scale multi-class urban change data with 512×512 images, similarly divided into 256×256 patches using the official split ($14,400/1,360/192$). *WHU-CD* contains bi-temporal aerial imagery for building change detection, where we adopt the commonly used patch-based preprocessing and standard splits. *CDD* includes diverse change scenarios such as buildings, vehicles, and seasonal variations from Google Earth imagery, and is processed using the same 256×256 patch protocol with standard evaluation splits. This preprocessing ensures consistent training and fair comparison.

4.2. Implementation Details

We use a **Swin Transformer-Tiny (Swin-T)** backbone [29], pretrained on ImageNet-1K, as the shared multi-scale encoder to extract hierarchical features for graph-based state evolution. Optimization is performed using **AdamW** [30] with ($\beta_1 = 0.9$, $\beta_2 = 0.999$) and weight decay 1×10^{-4} . A **OneCycle learning rate schedule** is adopted with a maximum learning rate of 3×10^{-4} , 10% warmup, and cosine annealing. Training is conducted with a batch size of 8 using **automatic mixed precision (FP16)** for 125 epochs. The segmentation objective combines binary cross-entropy and Dice loss, $\mathcal{L}_{\text{seg}} = \mathcal{L}_{\text{BCE}} + \lambda_{\text{Dice}} \mathcal{L}_{\text{Dice}}$, and the overall objective is $\mathcal{L} = \mathcal{L}_{\text{seg}} + \lambda_{\text{orth}} \mathcal{L}_{\text{orth}} + \lambda_{\text{var}} \mathcal{L}_{\text{var}} + \lambda_{\text{diff}} \mathcal{L}_{\text{diff}}$. For the diffusion decoder, we follow the DDPM formulation [22, 33] with $T = 100$ forward steps and a cosine noise schedule, and perform efficient **5-step DDIM sampling** [36] during inference. All experiments are conducted on a single **NVIDIA RTX 4090 GPU (24GB)**. We report four standard evaluation metrics: F1-score, Intersection-over-Union (IoU), Overall Accuracy

Table 1. Quantitative comparison on WHU-CD, DSIFN-CD, LEVIR-CD and CDD datasets. Best results in **bold**, second-best underlined.

Method	Extra data	WHU-CD			DSIFN-CD			LEVIR-CD			CDD		
		F1	IoU	OA	F1	IoU	OA	F1	IoU	OA	F1	IoU	OA
CNN-based Methods													
FC-Siam-conc [15]	None	0.798	0.665	98.5	0.597	0.426	87.6	0.837	0.720	98.5	0.751	0.601	94.9
SNUNet [18]	None	0.835	0.717	98.7	0.662	0.495	87.3	0.882	0.788	98.8	0.839	0.721	96.2
IFNet [40]	IN1k	0.834	0.715	98.8	0.601	0.430	87.8	0.881	0.788	98.9	0.840	0.719	96.03
CNN + Attention Methods													
DT-SCN [28]	IN1k	0.914	0.842	99.3	0.706	0.545	82.9	0.877	0.781	98.8	0.921	0.853	98.2
STANet [5]	IN1k	0.823	0.700	98.5	0.645	0.478	88.5	0.873	0.774	98.7	0.841	0.722	96.1
Transformer-based Methods													
BIT [8]	IN1k	0.905	0.834	99.3	0.876	0.780	92.3	0.893	0.807	98.92	0.889	0.800	97.5
ChangeFormer [3]	None	0.886	0.795	99.12	0.947	0.887	93.2	0.904	0.825	99.0	0.946	0.898	98.7
SAM2-CD [35]	SAM-pretrained	0.926	0.889	<u>99.77</u>	-	-	-	0.919	0.855	99.19	-	-	-
PeftCD [17]	IN1k	<u>0.959</u>	<u>0.9205</u>	99.6	-	-	-	0.923	0.856	<u>99.2</u>	<u>0.985</u>	<u>0.970</u>	<u>99.63</u>
Self-supervised Pretraining													
SiamSiam [42]	IN1k, IBSD, GoogleEarth	0.847	0.734	-	-	-	-	0.880	0.786	-	-	-	-
MoCo-v2 [13]	IN1k, IBSD, GoogleEarth	0.882	0.789	-	-	-	-	0.879	0.784	-	-	-	-
SeCo [32]	IN1k, IBSD, GoogleEarth	0.883	0.790	-	-	-	-	0.881	0.787	-	-	-	-
SaDL-CD [26]	IN1k, IBSD, GoogleEarth	0.909	0.833	-	-	-	-	0.899	0.818	-	-	-	-
Diffusion-based Methods													
DDPM-CD [1]	GoogleEarth	0.927	0.863	99.4	0.967	0.913	97.1	0.909	0.833	99.1	0.956	0.916	99.0
DiffRegCD [31]	GoogleEarth	0.934	0.883	99.0	0.940	0.890	96.7	<u>0.929</u>	<u>0.881</u>	98.7	-	-	-
Mamba-based Methods													
RSMamba [11]	IN1k	0.927	0.865	99.4	0.965	0.913	97.0	0.897	0.814	98.9	0.943	0.902	98.8
ChangeMamba [9]	IN1k	0.925	0.861	99.4	0.875	0.778	95.8	0.902	0.821	99.0	0.944	0.920	99.0
CDMamba [43]	IN1k	0.937	0.882	99.5	0.966	0.914	97.0	0.907	0.831	99.0	0.960	0.919	99.1
M-CD [34]	IN1k	0.953	0.911	99.6	<u>0.970</u>	<u>0.935</u>	<u>98.9</u>	0.921	0.850	<u>99.2</u>	0.982	0.963	99.5
Boundary-focused Methods													
LRNet [46]	IN1k	0.925	0.861	99.47	-	-	-	0.911	0.836	99.10	-	-	-
<i>Ours</i>													
GMD-CD	IN1k	0.962	0.921	99.78	0.978	0.948	99.12	0.9375	0.882	99.35	0.990	0.975	99.64

(OA), and Boundary F1 (B-F1). B-F1 is computed using a 3-pixel tolerance following [14], measuring the alignment between predicted and ground-truth boundaries.

4.3. Comparison with State-of-the-Art

Tab. 1 reports quantitative comparisons across four benchmarks. CNN-based methods (e.g., FC-Siam-conc, SNUNet [15, 18]) are limited by local receptive fields and underperform on complex scenes. Attention-based models (STANet [5], BIT [8], ChangeFormer [3]) improve global context modeling but incur higher complexity. Recent approaches explore complementary directions: foundation-based models such as SAM2-CD [35] and adaptation-based methods such as PeftCD [17] leverage large pretrained priors, while state-space models (RSMamba, ChangeMamba, CDMamba, M-CD [10, 12, 34, 44]) provide efficient long-range modeling. Diffusion-based methods (DDPM-CD, DiffRegCD [1, 31]) improve refinement quality, and LRNet [46] focuses on boundary modeling. However, these approaches typically improve either global interaction or refinement in isolation. In contrast, GMD-CD jointly models structured cross-temporal interaction and uncertainty-aware refinement. The BGMO module induces structured prop-

agation across temporal partitions with linear complexity, while the diffusion decoder performs logit-space correction guided by global context. This combination improves both region accuracy and boundary quality. As shown in Tab. 1, GMD-CD achieves the best performance across all datasets, improving both F1 and IoU over CNN and transformer baselines and outperforming recent state-space and diffusion-based methods. Compared to strong baselines such as M-CD and PeftCD, the proposed model yields consistent gains, indicating that structured propagation and logit-space refinement are complementary. Qualitative results in Fig. 5 further support these findings: compared to M-CD and DDPM-CD, GMD-CD produces more coherent regions with sharper boundaries and fewer false positives, particularly in densely structured scenes.

4.4. Ablation and Analysis

Table 2 quantifies the role of each component on LEVIR-CD and WHU-CD. The Mamba-only baseline captures global context (88.30/93.85 F1) but exhibits poor boundary quality (66.50/80.10 B-F1), indicating that unstructured sequence propagation is insufficient for precise localization. Introducing structured interaction (BGMO) im-

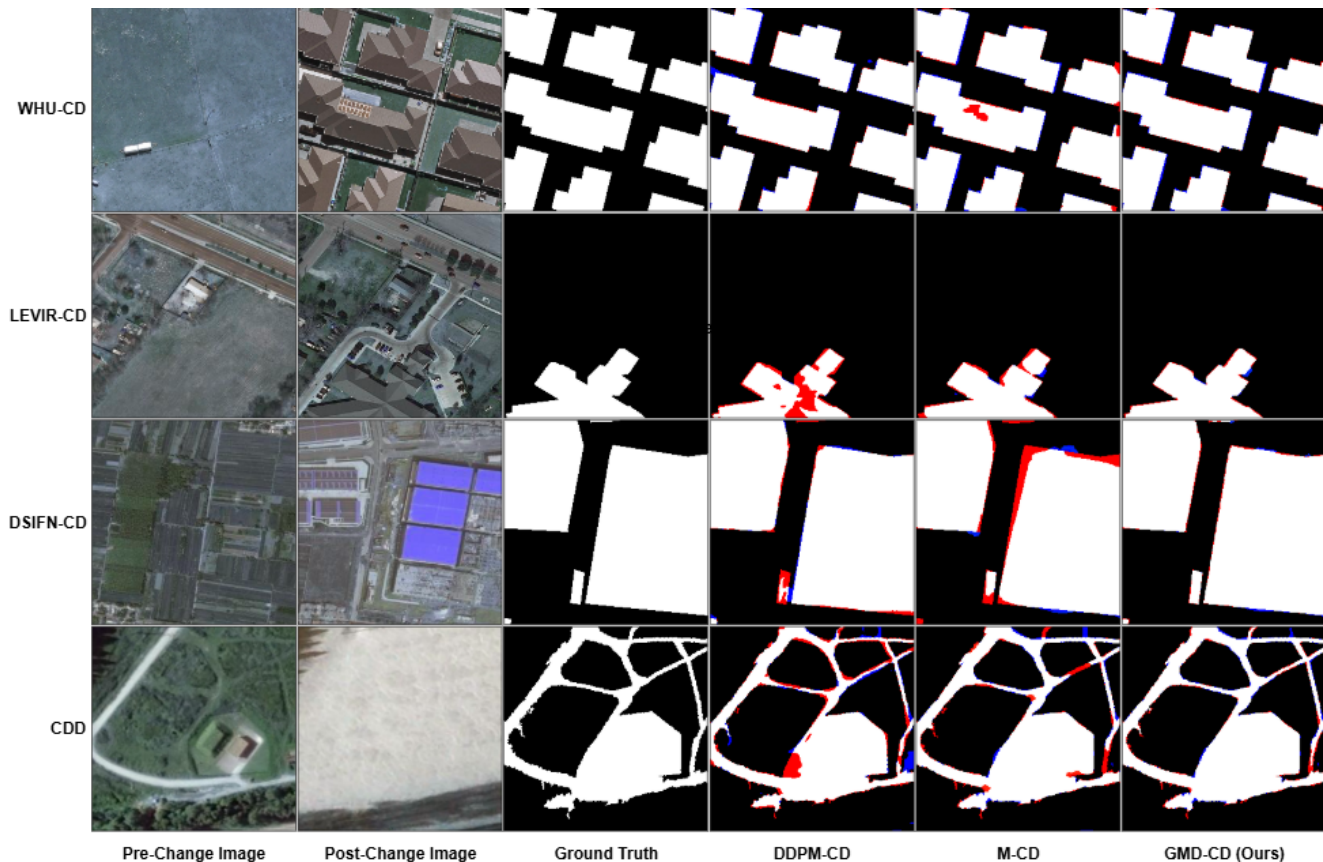


Figure 5. Change map comparison. The six columns show: Time 1 image, Time 2 image, Ground Truth, DDPM-CD, M-CD, and GMD-CD (Ours). In the error overlay, white denotes true positives (TP), black denotes true negatives (TN), red denotes false positives (FP), and blue denotes false negatives (FN).

372 proves recall (88.10 \rightarrow 90.95 on LEVIR; 93.50 \rightarrow 95.80
 373 on WHU) and increases F1 to 90.04/95.45, showing that
 374 partition-aware ordering improves coverage of spatially
 375 distributed changes. Diffusion, in contrast, yields large gains
 376 in boundary quality (66.50 \rightarrow 84.10 on LEVIR; 80.10 \rightarrow
 377 88.60 on WHU) with smaller F1 improvements, indicating
 378 uncertainty-aware refinement rather than primary predic-
 379 tion. Combining both leads to consistent gains (92.34/96.40
 380 F1 and 85.90/89.40 B-F1), demonstrating complementary
 381 behavior. Finally, adding OFD further improves all metrics
 382 (93.750/96.20 F1 and 87.40/91.20 B-F1), confirming that
 383 disentangling change and invariant features reduces inter-
 384 ference during propagation. Overall, structured propagation
 385 improves semantic completeness, diffusion refines uncer-
 386 tain regions, and OFD stabilizes feature dynamics, leading
 387 to joint gains in region accuracy and boundary quality.

388 4.5. Graph & Diffusion Analysis

389 Table 3 analyzes the sensitivity of the structured state-space
 390 module and diffusion refinement. Performance peaks at
 391 $K = 64$, indicating a balance between spatial abstraction

Table 2. Component-wise ablation on LEVIR-CD and WHU-CD.

Mamba	Graph	Diff	OFD	LEVIR-CD				WHU-CD			
				Prec.	Rec.	F1	B-F1	Prec.	Rec.	F1	B-F1
✓	×	×	×	88.50	88.10	88.30	66.50	94.20	93.50	93.85	80.10
✓	✓	×	×	89.15	90.95	90.04	69.10	95.10	95.80	95.45	82.30
✓	×	✓	×	90.50	90.10	90.30	84.10	95.80	94.70	95.25	88.60
✓	✓	✓	×	90.98	93.750	92.34	85.90	94.10	94.70	95.40	89.40
✓	✓	✓	✓	93.41	94.04	93.75	87.40	95.80	96.17	96.20	91.20

and representation capacity; smaller K reduces spatial cov- 392
 erage, while larger K approaches dense tokenization with 393
 diminishing returns. Average pooling consistently outper- 394
 forms max pooling, suggesting that stable region aggrega- 395
 tion is preferable to sparse activation selection. The $[a; b]$ 396
 ordering improves performance over interleaved and ran- 397
 dom layouts, confirming that preserving temporal partition- 398
 ing influences cross-temporal interaction within the state- 399
 space formulation. Bidirectional scanning further improves 400
 results over forward-only scanning, indicating that symmet- 401
 ric traversal enhances information exchange across parti- 402
 tions. For diffusion, increasing DDIM steps from 1 to 5 403

Table 3. BGMO sensitivity analysis on LEVIR-CD. Default: $K = 64$, avg pooling, $[a; b]$, bidirectional scan, 5-step DDIM.

Configuration	F1	Δ F1
<i>Number of graph nodes K</i>		
$K = 16$	92.35	-0.75
$K = 64$ [default]	93.750	—
$K = 256$	92.88	-0.22
<i>Node construction</i>		
Max pooling	92.60	-0.50
Average pooling [default]	93.750	—
<i>Sequence ordering</i>		
Interleaved	92.45	-0.65
$[a; b]$ [default]	93.750	—
Random	91.80	-1.30
<i>Scan direction</i>		
Forward	92.30	-0.80
Bidirectional [default]	93.750	—
<i>Diffusion steps (DDIM)</i>		
1 step	92.62	-0.48
3 steps	92.95	-0.15
5 steps [default]	93.750	—
10 steps	93.754	+0.04

404 improves F1, while further increases yield marginal gains.
 405 This behavior indicates that a small number of denoising
 406 iterations is sufficient for correcting high-frequency errors.
 407 We therefore use 5 steps as an effective trade-off between
 408 accuracy and efficiency. Overall, the results show that struc-
 409 tured state propagation governs global consistency, while
 410 diffusion provides targeted refinement, and OFD stabilizes
 411 feature interactions under varying configurations. Unlike
 412 deterministic refiners, diffusion corrects predictions condi-
 413 tioned on uncertainty, which is reflected in larger gains in
 414 B-F1 compared to F1.

415 4.6. Efficiency Analysis

416 We evaluate computational efficiency using three metrics:
 417 number of trainable parameters, GFLOPs, and average infer-
 418 ence time per image pair. Results on LEVIR-CD and
 419 WHU-CD are reported in Table 4. All measurements are
 420 conducted on a single NVIDIA RTX 4090 GPU with input
 421 resolution 256×256 . For methods without reported la-
 422 tency, inference time is estimated from GFLOPs under the
 423 same hardware setting, and may vary across implementa-
 424 tions. The proposed model without diffusion requires 13.2
 425 GFLOPs and 22 ms inference time, achieving 0.913 F1 on
 426 LEVIR-CD. Incorporating diffusion increases computation
 427 to 32.09 GFLOPs and 47 ms, while improving performance
 428 to 0.9375 (LEVIR-CD) and 0.962 (WHU-CD). GFLOPs in-
 429 clude all components, including the 5-step DDIM refine-
 430 ment. Compared to ChangeFormer, the proposed model
 431 achieves higher accuracy with substantially lower computa-
 432 tional cost and reduced latency. Lightweight models such
 433 as CDMamba and PefiCD achieve lower FLOPs but ex-
 434 hibit reduced accuracy, indicating limited modeling capac-

Table 4. Computational complexity comparison on LEVIR-CD and WHU-CD. Models are ordered by increasing GFLOPs.

Method	Params(M)	GFLOPs	Avg. Inf. Time	F1 _L	F1 _w
CDMamba [43]	11.90	10.4	15 ms	0.907	0.937
PefiCD [17]	10.15	22.0	25 ms	0.923	0.959
M-CD [34]	69.80	29.58	160 ms	0.921	0.953
LRNet [46]	48.71	92.23	60 ms	0.911	0.925
ChangeFormer [3]	41.02	254.8	85 ms	0.904	0.886
SAM2-CD [35]	2.59	1193.83	95 ms	0.919	0.926
DDPM-CD [1]	46.41	2175.46	88.10 ms	0.909	0.927
Ours (w/o Diff)	30.7	13.2	22 ms	0.913	0.9589
Ours (full)	33.55	32.09	47 ms	0.9375	0.962

ity. In contrast, M-CD attains competitive performance at
 significantly higher latency (160 ms), reflecting an unfavor-
 able efficiency trade-off. High-capacity approaches such as
 SAM2-CD and DDPM-CD incur substantially higher compu-
 tational cost due to large-scale pretrained backbones or
 full diffusion processes. Overall, GMD-CD operates at
 moderate computational cost while achieving the best accu-
 racy across datasets. The gains from diffusion are obtained
 at the cost of moderate additional latency, consistent with
 its role as a lightweight refinement stage rather than a full
 generative process. This results in a favorable accuracy-
 efficiency trade-off, combining structured state-space
 interaction with efficient logit-space refinement.

5. Conclusion

We introduced **GMD-CD**, a unified framework for remote
 sensing change detection that combines structured state-
 space interaction with logit-space diffusion refinement. The
 proposed BGMO module models cross-temporal relation-
 ships through a bipartite structured sequence and bidirec-
 tional state-space propagation, enabling global interaction
 with linear complexity. Complementing this, the condi-
 tional diffusion decoder refines coarse predictions directly
 in logit space, improving boundary quality without the over-
 head of full generative diffusion. Extensive experiments
 across four benchmarks demonstrate consistent gains over
 CNN, transformer, and recent state-space and diffusion-
 based methods. Ablation results further show that struc-
 tured state propagation improves semantic completeness,
 while diffusion provides targeted boundary refinement, in-
 dicating complementary behavior. Overall, the results sug-
 gest that combining structured state-space modeling with
 lightweight refinement is an effective direction for high-
 resolution change detection. Future work will explore ex-
 tensions to multi-class settings, adaptive node construction,
 and stronger pretraining strategies for state representations.

References

- [1] Chaminda Bandara and Vishal Patel. Ddpm-cd: Diffusion-
 based change detection. In *CVPR*, 2022. 1, 2, 5, 6, 8

- 473 [2] Wele Gedara Chaminda Bandara and Vishal M. Patel. A
474 transformer-based siamese network for change detection. In
475 *ICIP*, 2022. 1
- 476 [3] Wele Gedara Chaminda Bandara and Vishal M. Patel.
477 Changeformer. In *CVPR*, 2022. 1, 2, 6, 8
- 478 [4] Hao Chen and Zhenwei Shi. A spatial-temporal attention-
479 based method and a new dataset for remote sensing image
480 change detection. In *CVPR*, 2020. 1, 2
- 481 [5] Hao Chen and Zhenwei Shi. Stanet: Spatial-temporal atten-
482 tion network. In *CVPR*, 2020. 6
- 483 [6] Hao Chen, Zhi Zhang, et al. Change detection dataset for
484 remote sensing images. *Remote Sensing*, 2020. 5
- 485 [7] Hao Chen, Zipeng Qi, and Zhenwei Shi. Remote sensing
486 image change detection with transformers. In *ICCV*, 2021.
487 1, 2
- 488 [8] Hao Chen, Zipeng Qi, and Zhenwei Shi. Bit: Bitemporal
489 image transformer. In *CVPR*, 2022. 2, 6
- 490 [9] Hongruixuan Chen, Jian Song, Chengxi Han, Junshi Xia, and
491 Naoto Yokoya. Changemamba. In *CVPR*, 2024. 6
- 492 [10] Hongruixuan Chen, Jian Song, Chengxi Han, Junshi Xia,
493 and Naoto Yokoya. Changemamba: State space models for
494 change detection. In *CVPR*, 2024. 1, 2, 6
- 495 [11] Keyan Chen, Bowen Chen, Chenyang Liu, Wenyuan Li,
496 Zhengxia Zou, and Zhenwei Shi. Rsmamba. In *CVPR*, 2024.
497 6
- 498 [12] Keyan Chen, Bowen Chen, Chenyang Liu, Wenyuan Li,
499 Zhengxia Zou, and Zhenwei Shi. Rsmamba: Multi-
500 directional state space models for remote sensing. In *CVPR*,
501 2024. 2, 6
- 502 [13] Xinlei Chen. Moco v2. In *CVPR*, 2020. 6
- 503 [14] Gabriela Csurka, Diane Larlus, Florent Perronnin, and Fa-
504 bien Meylan. What is a good evaluation measure for seman-
505 tic segmentation? In *BMVC*, 2013. 6
- 506 [15] Rodrigo Daudt, Bertrand Le Saux, and Alexandre Boulch.
507 Fully convolutional siamese networks for change detection.
508 In *ICIP*, 2018. 6
- 509 [16] Rodrigo Caye Daudt, Bertrand Le Saux, Alexandre Boulch,
510 and Yann Gousseau. Fully convolutional siamese networks
511 for change detection. In *ICIP*, 2018. 1, 2
- 512 [17] Sijun Dong, Yuxuan Hu, Libo Wang Wang, Geng Chen, and
513 Xiaoliang Meng. Peftcd: Leveraging vision foundation mod-
514 els with parameter-efficient fine-tuning for remote sensing
515 change detection. *TGRS*, 2026. 6, 8
- 516 [18] Shuai Fang and Kun Fu. Snunet-cd: A densely connected
517 siamese network for change detection. In *IEEE TIP*, 2021. 6
- 518 [19] Shuai Fang, Kaiyu Li, Jinyuan Shao, and Zhe Li. Snunet-cd:
519 A densely connected siamese network for change detection.
520 *IEEE Geoscience and Remote Sensing Letters*, 2022. 1, 2
- 521 [20] Albert Gu and Tri Dao. Mamba: Linear-time sequence mod-
522 eling with selective state spaces. In *ICLR*, 2024. 1, 2, 3
- 523 [21] Rishabh Gupta, Manish Shah, et al. Creating xbd: A dataset
524 for assessing building damage. *CVPR*, 2019. 1
- 525 [22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffu-
526 sion probabilistic models. In *NeurIPS*, 2020. 5
- 527 [23] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffu-
528 sion probabilistic models. In *NeurIPS*, 2020. 5
- [24] Shunping Ji, Shuang Jiang, and Yubing Lu. Fully convo-
lutional networks for multisource building change detection
using very high resolution images. *Pattern Recognition*,
2019. 5
- [25] Yilmaz Korkmaz and Vishal M. Patel. Remotevar: Autore-
gressive visual modeling for remote sensing. In *CVPR*, 2024.
2
- [26] Yunsong Li, Chao Zhang, Kai Chen, and Aojie Li. Sadl-cd.
In *CVPR*, 2023. 6
- [27] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S.
Belongie. Feature pyramid networks for object detection. In
CVPR, 2017. 4
- [28] Jing Liu, Lei Liu, and Yu Zheng. Dt-scnn: Deep temporal
siamese network. In *CVPR*, 2020. 6
- [29] Ze Liu, Yutong Lin, Yue Cao, et al. Swin transformer: Hier-
archical vision transformer using shifted windows. In *ICCV*,
2021. 2, 5
- [30] Ilya Loshchilov and Frank Hutter. Decoupled weight decay
regularization. In *ICLR*, 2017. 5
- [31] Seyedehanita Madani, Rama Chellappa, and Vishal M. Patel.
Diffregcd: Integrated registration and change detection with
diffusion features. In *CVPR*, 2024. 6
- [32] Oğuzhan Fatih Manas, Jean-Francois Todorovic, Shreyas
Kulkarni, and Julien Cornebise. Seco: Seco self-supervised
learning. In *CVPR*, 2021. 6
- [33] Alexander Quinn Nichol and Prafulla Dhariwal. Improved
denoising diffusion probabilistic models. In *ICML*, 2021. 5
- [34] Jay N. Paranjape, Celso M. de Melo, and Vishal M. Patel.
M-cd. In *WACV*, 2026. 1, 6, 8
- [35] Yuan Qin, Chaoting Wang, Yuanyuan Fan, and Chanling
Pan. Sam2-cd: Remote sensing image change detection with
sam2. *TGRS*, 2025. 1, 2, 6, 8
- [36] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denois-
ing diffusion implicit models. In *ICLR*, 2021. 5
- [37] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denois-
ing diffusion implicit models. In *ICLR*, 2021. 5
- [38] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So
Kweon. Cbam: Convolutional block attention module. In
ECCV, 2018. 2
- [39] Chen Wu, Bin Du, and Liangpei Zhang. A benchmark
dataset for change detection in remote sensing images.
TGRS, 2021. 5
- [40] Jian Zhang, Qingyun Li, and Guohua Cao. Ifnet for change
detection. In *CVPR*, 2023. 6
- [41] Jian Zhang, Qingyun Li, and Guohua Cao. Next2former-
cd: Convnext meets mask2former for change detection. In
CVPR, 2024. 2
- [42] Ming Zhang, Xiang He, and Zhaohui Li. Siamself-
supervised change detection. In *CVPR*, 2022. 6
- [43] Y. Zhang and et al. Cdmamba. In *CVPR*, 2024. 6, 8
- [44] Y. Zhang and et al. Cdmamba: Bidirectional state space
models for change detection. In *CVPR*, 2024. 1, 2, 6
- [45] Zhengxin Zhang, Qingjie Liu, and Yunhong Wang. Levir-
cd: A large-scale remote sensing dataset for building change
detection. In *CVPR*, 2020. 5

- 584 [46] Huan Zhong, Chen Wu, and Ziqi Xiao. Lrnet: Change
585 detection in high-resolution remote sensing imagery via a
586 localization-then-refinement strategy. In *Remote Sensing*,
587 2023. 6, 8
- 588 [47] Lianghai Zhu, Biao Wang, Tang Tang, Xiaofan Zhu, and et
589 al. Vision mamba: Efficient vision backbone with state space
590 model. *arXiv*, 2024. 2