

C-NERF: Representing Scene Changes as Directional Consistency Difference-based NeRF

Rui Huang¹ Binbin Jiang¹ Qingyi Zhao¹ William Wang² Yuxiang Zhang¹ Qing Guo^{3,4,*}

¹ College of Computer Science and Technology, Civil Aviation University of China, China

²University of South Carolina, The USA

³ IHPC, Agency for Science, Technology and Research, Singapore

⁴ CFAR, Agency for Science, Technology and Research, Singapore

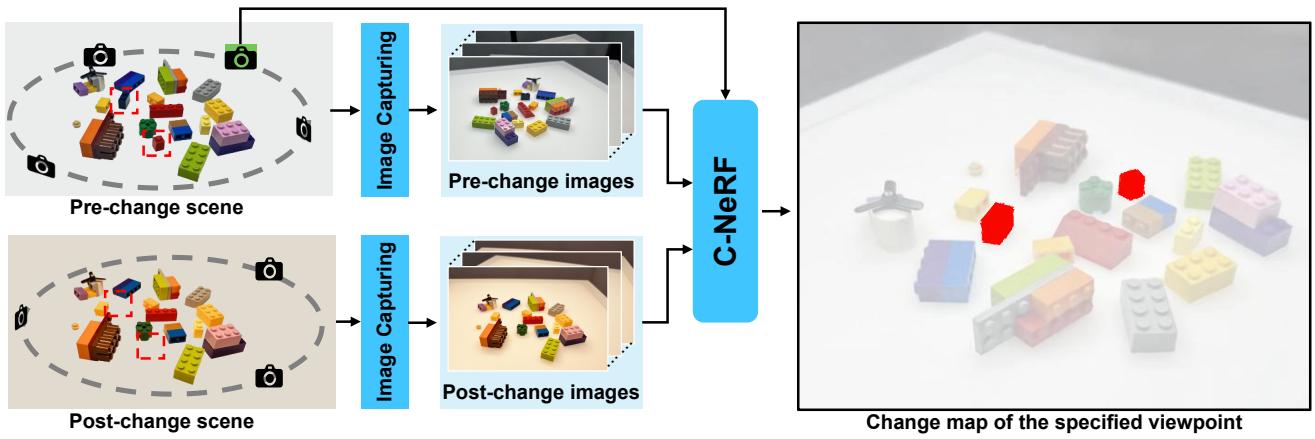


Figure 1. Overview of C-NERF. We aim to build a representation for 3D changes of a scene from multiview images captured before and after object changes. C-NERF can generate change maps based on arbitrarily specified views (see the left image). We highlight the changes in the scene via dashed red rectangles.

Abstract

In this work, we aim to detect the changes caused by object variations in a scene represented by the neural radiance fields (NeRFs). Given an arbitrary view and two sets of scene images captured at different timestamps, we can predict the scene changes in that view, which has significant potential applications in scene monitoring and measuring. We conducted preliminary studies and found that such an exciting task cannot be easily achieved by utilizing existing NeRFs and 2D change detection methods with many false or missing detections. The main reason is that the 2D change detection is based on the pixel appearance difference between spatial-aligned image pairs and neglects the stereo information in the NeRF. To address the limitations, we propose the C-NERF to represent scene changes as directional consistency difference-based NeRF, which mainly contains three modules. We first perform the spatial alignment of two NeRFs captured before and after changes. Then, we identify the change points based on the direction-consistent constraint; that is, real change points have similar change representa-

tions across view directions, but fake change points do not. Finally, we design the change map rendering process based on the built NeRFs and can generate the change map of an arbitrarily specified view direction. To validate the effectiveness, we build a new dataset containing ten scenes covering diverse scenarios with different changing objects. Our approach surpasses state-of-the-art 2D change detection and NeRF-based methods by a significant margin.

1. Introduction

Change detection aims to identify regions or objects that have changed in the scene, which has various applications in fields such as remote sensing [5, 19], surveillance [17, 50], and robot navigation [21, 43]. Previous works mainly detect scene changes based on spatially aligned image pairs that are captured by a camera at two timestamps. Although such 2D methods achieve significant progress [16], they have two inherent limitations: *First*, 2D change detections could only identify changes in one view direction. Some changes

hidden in other objects caused by occlusion could not be identified. *Second*, the 3D shape of the changing objects cannot be known, which limits the potential applications significantly. A 3D change detector should be developed.

In this work, we propose to extend the 2D change detection task to 3D based on the neural radiance fields (NeRFs) [8, 28]. Compared with other 3D representations (*e.g.*, point-clouds [2], meshes [37, 46], voxels [27]), NeRF is able to represent a 3D scene from multi-view images and synthesize photo-realistic novel views by learning a multi-layer perceptron (MLP) to map an input 5D coordinates (*i.e.*, spatial location and 2D view direction) to the volumen density and view-dependent emitted radiance color, which has been extended to diverse challenging scenarios [32, 44].

With the NeRF [28] for 3D representation, the change detection in 3D can be formulated as a change map synthesis task: given an arbitrary view and two sets of scene images captured at different timestamps, we aim to predict a binary change map indicating the changing regions in that view, as shown in Fig. 1. However, such a task is non-trivial and cannot be simply achieved through existing NeRF methods. In particular, a naive solution is to build two independent NeRFs based on the two image sets and synthesize two view images by feeding the two NeRFs with the same view direction. Then, we get two synthesized images and can calculate the difference between two images for change detection. As shown in Fig. 2, such a method easily leads to a lot of false or missing detections caused by: ❶ The two NeRFs of the 3D scene before and after change are not registered in the same coordinate. This means that even if the spatial location and viewing direction are identical for the NeRF models, the generated images are still not aligned, which results in a lot of false detections (see Fig. 2 (Left)). ❷ Even if we have two aligned images, the absolute difference of the image pair is also affected by rendering quality and environment differences (*e.g.*, light changes in Fig. 1 and Fig. 2 (Right)).

To address the limitations, we propose C-NeRF to represent scene changes as directional consistency difference-based NeRF, which mainly contains three modules. We first perform the spatial alignment of two NeRFs captured before and after changes. Then, we identify the change points based on the direction-consistent constraint; that is, real change points have similar change representations across view directions, but fake change points do not. Finally, we design the change map rendering process based on the built NeRFs and can generate the change map of an arbitrarily specified view direction. To validate the effectiveness, we build a new dataset containing ten scenes covering diverse scenarios with different changing objects. In summary, the main contributions of our paper are as follows:

- We extend the 2D change detection task to 3D, and propose C-NeRF by predicting the scene changes of an arbitrary view from two sets of scene images captured at different

timestamps. With such a novel method, we can synthesize change maps under arbitrary specified view directions, which indicates the change regions under those views.

- We observe that the real change points have similar change representations across different view directions but the fake change points do not. Inspired by this, we propose a real change points identification method based on the direction-consistent constraint, which is able to identify instance-level and fine-grained changes effectively.
- To validate the effectiveness, we build a new scene change detection dataset containing ten scenes covering diverse scenarios with different change objects and compare with a series of baselines.

2. Related Work

2D change detection methods. 2D change detection has been widely studied in image change captioning, remote sensing applications and street view scenes [1, 5, 16, 18, 41]. The basic requirement of these methods is that images of before and after scene change have limited viewpoint shifts.

Image change captioning aims to locate and describe changed objects by using natural language [31]. Huang et al. [15] propose a method to accurately locate changed objects using fine-grained features such as visual, semantic, and positional features at the instance-level. Kim et al. [20] present a viewpoint-agnostic change captioning network with cycle consistency to improve the robustness of the change detection captured with large camera perspectives. Qiu et al. [34] propose multi-change captioning transformers that identify change regions by densely correlating different regions in image pairs and dynamically determining the related change regions with words in sentences. Tu et al. [47] propose a semantic relation-aware difference representation learning network that explicitly learns the difference representation in the presence of illumination or viewpoint change. Unlike change captioning, change detection in remote sensing or street view scenes always outputs binary maps indicating the changed pixel. In [40], Sakurada et al. use a convolutional neural network in combination with superpixel segmentation to detect changes on a pair of street-view images. Lei et al. [24] design an effective feature fusion method to improve the accuracy of the corresponding change maps. Chen et al. [7] propose the temporal attention and explore the impact of the dependency-scope size of temporal attention on the performance of change detection. Sachdeva et al. [38] aims to detect the “object-level” change in an image pair with different illumination at different viewpoints. In their latest work [39], the authors’ further study change detection with a significant shift in camera pose. 2D change detection requires image alignment, training with a large volume of change image pairs, and can only generate change maps at the viewpoint of the input images. Obtaining change maps at different viewpoints can be challenging.

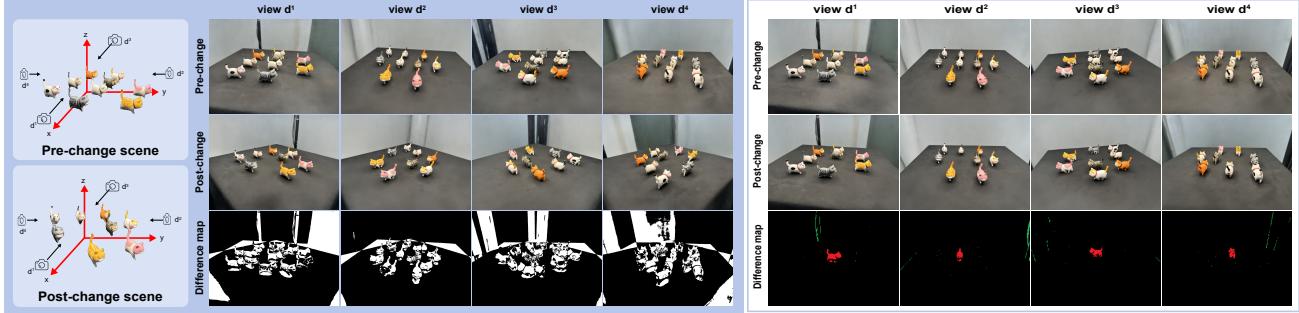


Figure 2. Left: images generated from two independent NeRF models with identical spatial location and viewing direction. Right: images generated from two spatially aligned NeRFs. The red points represent real changes. The green points indicate fake changes.

3D change detection. 3D change detection can reflect change in 3D space with multi-view images of before and after scene change and 3D scans [3, 33, 45, 48]. Previous research [13] has utilized dense color and depth information to detect changes between 3D maps. Fehr et al. [11] develop a 3D reconstruction algorithm based on an extended Truncated Signed Distance Function to more accurately solve the scene differencing problem. Ku et al. [22] have contributed a street-scene dataset for 3D point cloud change detection that can detect changes in a complex street environment from multi-temporal point clouds. Qiu et al. [35] propose a method to explicitly localize changes in 3D bounding boxes from two point clouds, describing detailed scene changes. In contrast with existing methods, our method only requires multi-view images before and after scene change, without camera pose alignment. In addition, our method is capable of capturing the change in 3D space and rendering change maps at any given viewpoint.

NeRF for scene representation. Neural Radiance Field (NeRF) [28] is a 3D scene representation technology that uses a deep neural network to generate novel views of a 3D scene based on the corresponding view directions. Several methods have been proposed to improve the quality and speed of NeRF models. Barron et al., [4] propose a method to represent the scene at a continuously-valued scale and render anti-aliased conical frustums instead of rays, which improves the fine details of the generated image. Niemeyer et al., [30] regularize the geometry and appearance of patches rendered from unobserved viewpoints to enhance the NeRF. Zhang et al., [49] propose a random ray casting policy that enables training unseen views using seen views to produce high-quality renderings under novel viewpoints. Slow rendering times always hinder the use of NeRF. To achieve high rendering speed, FastNeRF [12] and KiloNeRF [36] propose using multiple MLPs. Instant-NGP[29] propose to use multi-resolution hash encoding as network input to reduce the number of floating point and memory access operations. EfficientNeRF[14] proposes using multi-resolution hash encoding as network input to reduce the number of floating-point and memory access operations. MobileNeRF[9] intro-

duces a new NeRF representation based on textured polygons that can synthesize novel images efficiently with standard rendering pipelines. Different from these works improving the rendering quality and speed of NeRF models, our method is built on NeRF and solves change detection without camera pose alignment.

3. Preliminary: Neural Radiance Fields

Neural radiance fields (NeRF) [28] can synthesize new views via a MLP (multilayer perception) Ψ trained with a set of sparse images $\mathcal{I} = \{\mathbf{I}_1, \mathbf{I}_2, \dots, \mathbf{I}_N\}$ captured at different views, which can be denoted as $\Psi = \Gamma_{\text{train}}(\mathcal{I}, \mathcal{D})$. The set $\mathcal{D} = \{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_N\}$ contains view directions of images in \mathcal{I} , which are calculated by the method [42]. During the inference process, NeRF takes a 5D coordinate including the spatial location $\mathbf{x} = (x, y, z)$ and the viewing direction $\mathbf{d} = (\theta, \phi)$ as inputs and can estimate the volume density σ and view-dependent emitted radiance color $\mathbf{c} = (r, g, b)$ of the specified coordinate, that is $(\mathbf{c}, \sigma) = \Psi(\mathbf{x}, \mathbf{d})$.

When we render an image under the view direction \mathbf{d} , for a pixel \mathbf{p} , we have a ray that is shot from the camera center through the pixel, and we can estimate its color by $\mathbf{I}[\mathbf{p}] = \sum_{i=1}^K T_i \alpha_i \mathbf{c}_i$, where $\{\mathbf{x}_i\}_{i=1}^K$ contains K points from the ray. The variables α_i and T_i are the alpha values for blending and the accumulated transmittance, which are determined by the estimated volume density. We denote the whole render process for the given view direction \mathbf{d} as $\mathbf{I} = \Gamma_{\text{render}}(\Psi, \mathbf{d})$. Please refer to [28] for more details.

4. Scene Change Representation

4.1. Problem Formulation

For a scene we want to monitor, we collect a set of images $\mathcal{I}^a = \{\mathbf{I}_1^a, \mathbf{I}_2^a, \dots, \mathbf{I}_{N_a}^a\}$ at the first time observation (*i.e.*, pre-change images). After some time, we observe the scene again and get another set of images $\mathcal{I}^b = \{\mathbf{I}_1^b, \mathbf{I}_2^b, \dots, \mathbf{I}_{N_b}^b\}$ (*i.e.*, post-change images). Our goal is to learn a mapping that can represent the changes in the scene at arbitrary viewing direction, in particular, the fine-grained changes, between two observations. Such a requirement is expected in the real

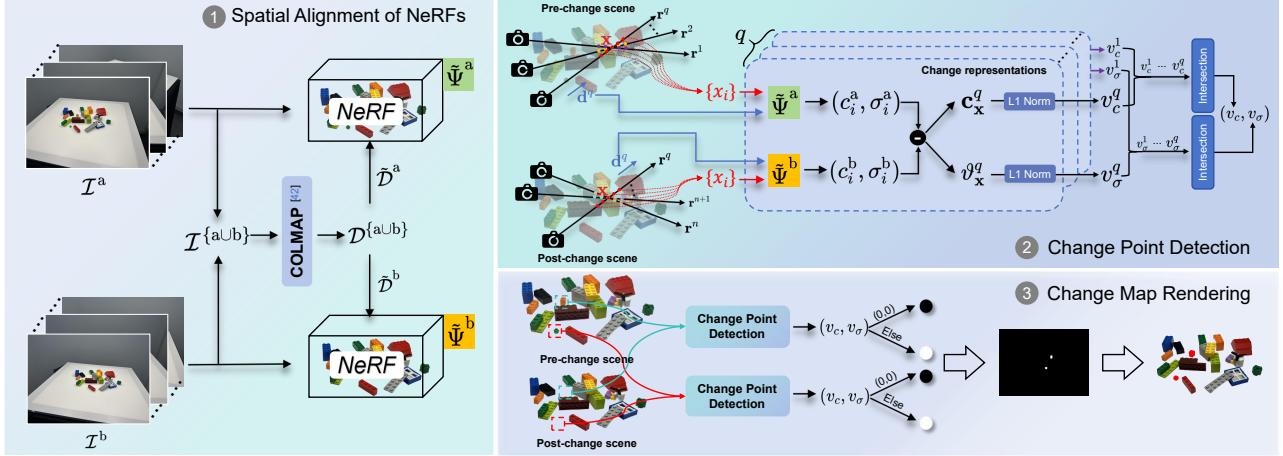


Figure 3. The main framework of the proposed C-NERF.

world. For example, we would like to monitor and measure the scene changes around high-value objects, such as cultural relics and jewelry, to protect them. However, existing 2D change detections cannot achieve the goal since they usually require the images captured two times to be spatially aligned.

Specifically, given a 3D point $\mathbf{x} = (x, y, z)$ and a view direction $\mathbf{d} = (\theta, \phi)$, a learned representation can estimate whether the point is a changed point or not

$$(v_c, v_\sigma) = \Psi^v(\mathbf{x}, \mathbf{d}), \quad (1)$$

where $v_c \in \{0, 1\}$ and $v_\sigma \in \{0, 1\}$ are two change indicators indicating whether the color and volume density of \mathbf{x} change or not dependent on \mathbf{d} . Different from learning Ψ within Sec. 3, Ψ^v should depend on the two observations with two sparse sets, *i.e.*, \mathcal{I}^a and \mathcal{I}^b .

4.2. Naive Solution

We can employ existing NeRF methods and calculate appearance and density difference to realize the task defined in Eq. (1). Specifically, we first use an existing NeRF method to build NeRF models based on the two sets \mathcal{I}^a and \mathcal{I}^b , respectively, *i.e.*,

$$\Psi^a = \Gamma_{\text{train}}(\mathcal{I}^a, \mathcal{D}^a), \text{ and } \Psi^b = \Gamma_{\text{train}}(\mathcal{I}^b, \mathcal{D}^b), \quad (2)$$

where \mathcal{D}^a and \mathcal{D}^b contain directions of images in \mathcal{I}^a and \mathcal{I}^b , respectively. We calculate \mathcal{D}^a and \mathcal{D}^b by feeding \mathcal{I}^a and \mathcal{I}^b to [42], respectively. Then, we estimate the color of the same point under a view direction via the two NeRF models, respectively, and get

$$(c^a, \sigma^a) = \Psi^a(\mathbf{x}, \mathbf{d}), \quad (c^b, \sigma^b) = \Psi^b(\mathbf{x}, \mathbf{d}). \quad (3)$$

After that, we can calculate the difference between the two colors and get

$$c^{a-b} = |c^a - c^b|, \quad \sigma^{a-b} = |\sigma^a - \sigma^b|. \quad (4)$$

Intuitively, we can calculate v via a threshold for c^{a-b} ,

$$v_c = \begin{cases} 1, & c^{a-b} > \epsilon_c \\ 0, & \text{otherwise} \end{cases}, \quad v_\sigma = \begin{cases} 1, & \sigma^{a-b} > \epsilon_\sigma \\ 0, & \text{otherwise} \end{cases}. \quad (5)$$

Although easy, this strategy requires the two built NeRF models to be well-aligned in both spatial space and appearance. That is, given the same direction, the Ψ^a and Ψ^b should render the same view and have the same color at the same point. Nevertheless, Ψ^a and Ψ^b are built independently and usually cannot generate the same view with a given direction. We show an example in Fig. 2 (Left) where we build two NeRFs for a scene respectively and render a view by specifying the same direction. The rendered views are different and lead to numerous fake changes. In addition, the difference caused by environmental changes (*e.g.*, light changes) and rendering quality instead of object changes will also lead to a high c^{a-b} and affect the change detection results significantly. As the case is shown in Fig. 2 (Right), even though the two NeRF models are well aligned, there are a lot of fake changes (*i.e.*, green points), and real changes are missed.

5. Methodology: C-NERF

Instead of detecting changes based on differences between two rendered images, we argue that a real change point should have distinct color differences in all view directions. To this end, we propose aligning two NeRF models in the same coordinate system through structure-from-motion (SfM), as detailed in Sec. 5.1. Then, we propose the *direction-consistent change point detection* in Sec. 5.2 to filter fake changes and identify fine-grained changes automatically. Finally, we introduce how to render the change map with the specified direction in Sec. 5.3. The main framework of C-NERF is shown in Fig. 3.

5.1. Spatial Alignment of NeRFs

Given two image sets captured at different timestamps (*i.e.*, \mathcal{I}^a and \mathcal{I}^b), we merge them and get a set $\mathcal{I}^{\{a\cup b\}} = \{\mathcal{I}^a, \mathcal{I}^b\}$.

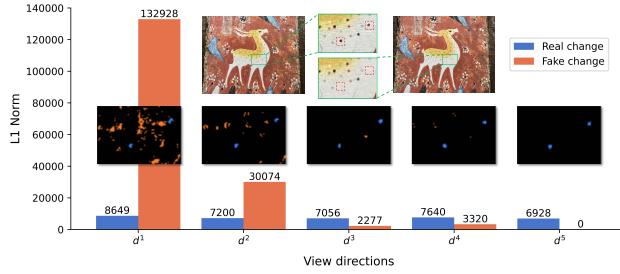


Figure 4. Total L1 norms of real and fake change points of the green box on five view directions.

Different from the Eq. (2) where the two NeRF models are trained with independently calculated \mathcal{D}^a and \mathcal{D}^b , we use the method [42] to assign directions for all images in $\mathcal{I}^{\{a \cup b\}}$ as a whole and get $\mathcal{D}^{\{a \cup b\}} = \{\tilde{\mathcal{D}}^a, \tilde{\mathcal{D}}^b\}$. As a result, the two image sets are aligned in the same coordinate system. With the new direction sets $\{\tilde{\mathcal{D}}^a, \tilde{\mathcal{D}}^b\}$, we can use Eq. (2) to rebuild the two NeRF models denoted as $\tilde{\Psi}^a$ and $\tilde{\Psi}^b$.

With such a simple strategy, the two NeRF models captured at two different timestamps are well-aligned in the spatial space, which eliminates the fake changes significantly. However, the fake changes caused by environmental changes and rendering quality cannot be appropriately addressed.

5.2. Direction-Consistent Change Point Detection

A real change point \mathbf{x} should have the following properties: **①** Its appearances or densities estimated from the two NeRF models are usually distinct. **②** The first property should be consistent across different view directions. To achieve the two properties, we design a change point detection method based on the built $\tilde{\Psi}^a$ and $\tilde{\Psi}^b$.

Change representations. Given a point \mathbf{x} and q th view direction \mathbf{d}^q that can see the point, we have a ray \mathbf{r}^q that is shot from the camera center to the point \mathbf{x} . We sample points near \mathbf{x} along the ray \mathbf{r}^q and get $\mathcal{N} = \{\mathbf{x}_i\}$. Then, we calculate the color difference of $\mathbf{x}_i \in \mathcal{N}$ based on $\tilde{\Psi}^a$ and $\tilde{\Psi}^b$, and use their weighted difference as the representation of the appearance change of the point \mathbf{x} under the view \mathbf{d}^q ,

$$\mathbf{c}_{\mathbf{x}}^q = [c_1^{a-b}, c_2^{a-b}, \dots, c_{|\mathcal{N}|}^{a-b}], \quad (6)$$

$$\vartheta_{\mathbf{x}}^q = [\sigma_1^{a-b}, \sigma_2^{a-b}, \dots, \sigma_{|\mathcal{N}|}^{a-b}], \quad (7)$$

$$\begin{aligned} \text{s.t. } c_i^{a-b} &= |T_i^a \alpha_i^a c_i^a - T_i^b \alpha_i^b c_i^b|, \\ \sigma_i^{a-b} &= |T_i^a \alpha_i^a \sigma_i^a - T_i^b \alpha_i^b \sigma_i^b|, \end{aligned}$$

where $(c_i^a, \sigma_i^a) = \tilde{\Psi}^a(\mathbf{x}_i, \mathbf{d}^q)$, and $(c_i^b, \sigma_i^b) = \tilde{\Psi}^b(\mathbf{x}_i, \mathbf{d}^q)$ with $i \in [1, \dots, |\mathcal{N}|]$. In addition, we follow the sampling strategy used in the original NeRF and have $T_i^* = \exp(-\sum_{j=1}^{i-1} \sigma_j^* \delta_j)$, $\alpha_i^* = 1 - \exp(\sigma_i^* \delta_i)$, $* \in \{a, b\}$, and δ_i denotes the distance between two neighboring sampled points. Intuitively, a real change point \mathbf{x} easily leads to distinct pre-post changes of neighboring points around the \mathbf{x} itself, and we consider the influence of point density by

involving $\{T_i^* \alpha_i^* | * \in \{a, b\}\}$ in Eq. (6). Then, we use the weighted pre-post differences of neighboring points to represent appearance changes of \mathbf{x} and calculate the change indicators of \mathbf{x} in Eq. (1) based on the q th view by

$$v_c^q = \begin{cases} 1, & \|\mathbf{c}_{\mathbf{x}}^q\|_1 > \epsilon_c \\ 0, & \text{otherwise} \end{cases}, \quad v_{\sigma}^q = \begin{cases} 1, & \|\vartheta_{\mathbf{x}}^q\|_1 > \epsilon_{\sigma} \\ 0, & \text{otherwise} \end{cases}. \quad (8)$$

However, based on a single view, the representation is easily disturbed by rendering noise or environmental changes. As an example shown in Fig. 2 (Right), there are numerous fake changes based on the Eq. (8).

Direction-consistent constrain. We observe that the real change points' representations always have high L_1 norms and get consistent v_c or v_{σ} across different views. In contrast, the fake change points' representations only have high L_1 norms in a few views and low L_1 norms across other views. An example is shown in Fig. 4 where we calculate the total L_1 norms of real and fake change points of the green box on five view directions. Clearly, the real change points have similar L_1 norms across different directions while the fake ones have significantly large L_1 norm on d_1 and zero L_1 norm on d_5 . Specifically, given a set of view directions $\tilde{\mathcal{D}} = \{\mathbf{d}^q\}$, we can compute the representation of \mathbf{x} under each view direction and get $\{(v_c^q, v_{\sigma}^q) | \mathbf{d}^q \in \tilde{\mathcal{D}}\}$ via Eq. (8). According to the direction consistent fact, if \mathbf{x} is a real change point caused by color variation (or density variation), all $\{v_c^q\}$ (or $\{v_{\sigma}^q\}$) should be one and we can calculate the two indicators of \mathbf{x} as

$$(v_c, v_{\sigma}) = \Psi^v(\mathbf{x}, \mathbf{d}) = \left(\bigcap_{q=1}^{|\tilde{\mathcal{D}}|} v_c^q, \bigcap_{q=1}^{|\tilde{\mathcal{D}}|} v_{\sigma}^q \right). \quad (9)$$

Then, the key problem becomes how to collect view directions to build $\tilde{\mathcal{D}}$ to calculate v .

View direction sampling strategies. We follow two principles to select view directions: **①** the desired point \mathbf{x} could be seen along the selected views. **②** As many views as possible should be included. Moreover, we need to consider two scenarios: *Forward-facing scene*, the camera always faces the scene, and each point can be seen in arbitrary view directions. Hence, we can uniformly sample V view directions to form $\tilde{\mathcal{D}}$. In Fig. 5a, we show the sampled view directions to calculate v . *Surround scene*, for a point \mathbf{x} , the view directions that cannot see the point due to the self-occlusion should be excluded from the set $\tilde{\mathcal{D}}$. We take the view direction \mathbf{d} as the center line and rotate the direction from -90 degree to 90 degree centered on the point \mathbf{x} as the sampling range. This range is able to cover most of the view directions that can see the point \mathbf{x} . Then, within this range, we uniformly sample V view directions to form the set $\tilde{\mathcal{D}}$. We show an example in Fig. 5b.

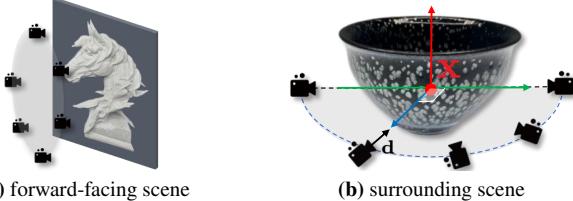


Figure 5. Examples of $\tilde{\mathcal{D}}$ for forward-facing and surrounding scenes.

5.3. Change Map Rendering under New view

Given a view direction \mathbf{d} , we aim to render a binary change map C that indicates the changed pixels under the view direction. Specifically, for the pixel \mathbf{p} of the image captured under view \mathbf{d} , we have a ray \mathbf{r} shooting from the camera center through the pixel center. Then, we sample points along the ray as done in NeRF and find the max values $T_i^a \alpha_i^a$ and $T_j^b \alpha_j^b$ from $\{T_i^a \alpha_i^a\}_{i=1}^K$ and $\{T_j^b \alpha_j^b\}_{j=1}^K$, respectively.

We select the point closest to the camera from $T_i^a \alpha_i^a$ and $T_j^b \alpha_j^b$ as the center point denoted as \mathbf{x} and can determine the binary value of pixel \mathbf{p} by

$$C[\mathbf{p}] = \max(v_c, v_\sigma), \text{ s.t., } (v_c, v_\sigma) = \Psi^v(\mathbf{x}, \mathbf{d}). \quad (10)$$

5.4. Implementation Details

Given two image sets observed at two timestamps (*i.e.*, \mathcal{T}^a and \mathcal{T}^b) and a view direction \mathbf{d} , we adopt the following steps to generate the change map under \mathbf{d} : ① We use the method in Sec. 5.1 to generate two spatially aligned NeRFs, *i.e.*, $\hat{\Psi}^a$ and $\hat{\Psi}^b$. ② We build the view direction set $\tilde{\mathcal{D}}$ according to different scenarios as detailed in ‘view direction sampling strategies.’ ③ We calculate the change map for each pixel via Eq. (10) in Sec. 5.3 and Eq. (9) in Sec. 5.2.

Configuration. In this paper, we adopt the NeRF model [28] with default parameters. The image size is set to 1008×756 . All the experiments are conducted on Nvidia RTX 3090Ti. We set $\epsilon_c \in [60, 180]$ and $\epsilon_\sigma \in [100, 600]$.

6. Experimental Results

6.1. Setups

Dataset. Most existing change detection datasets are built from coupled images, which cannot support NeRF model training. To validate our method, we build a new dataset **NeRFCD**. It consists of 2 *surrounding* scenes (*e.g.*, Cats and Block) and 8 *forward-facing* scenes (*e.g.*, Go pieces, Mural, Card, Text, Potting, Desk, sculpture and Desk).

As shown in Fig. 6, we adopt a “Zig-zag” path to shoot about 20 images for the forward-facing scene. To capture the surrounding scene, we use three horizontal circles around the objects: upper, middle and lower. We sample 40-60 shooting positions on each circle. Note that we keep the height of the camera when shooting at the middle circle. The camera always faces the center of the scene for shooting.

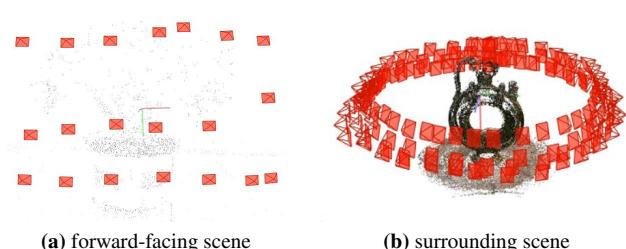


Figure 6. Shooting positions of forward-facing and surrounding scenes.

Table 1. Quantitative comparison of CYWS-3D [39] with C-NERF. The best MAP values are highlighted in **bolded**.

Scene	Cats	Blocks	Potting	Go pieces
CYWS-3D [39]	15.9	0	0	1.5
C-NERF	85.1	67.9	34.5	87.3

The motion trajectories of the two shots before and after the change should be as similar as possible.

To train the fine-grained change detectors, we render 120 image pairs for each scene and use our C-NERF to generate the corresponding change maps. Afterward, we select 10 image pairs from the captured image pairs from each scene to build the testing dataset. The training set consists of 1, 200 image pairs, and the testing set consists of 100 image pairs.

Evaluation metrics. We evaluate instance change detection using mean average precision (MAP) and pixel-level change detection using Precision (P), Recall (R), F1-measure (F1), and Intersection over Union (IoU).

6.2. Results

Instance change detection. As CYWS-3D [39] cannot detect fine-grained changes, we conducted the instance change detection experiments on instance change scenes, including Cats, Blocks, Potting and Go Pieces. To generate our instance change detection results, we use the tightest bounding box of the change region detected by C-NERF.

As shown in Fig. 7, CYWS-3D can detect large change objects such as the white cat and the bag hanging on the tree but fails to detect small objects such as go pieces and blocks. In contrast, our method, C-NERF, can accurately detect all changes in these four scenes.

In Table 1, we can see that CYWS-3D achieves a MAP value of 15.9 on the Cats scene and 1.5 on the Go pieces scene, which is not as high as the MAP values achieved by C-NERF. For instance, C-NERF achieves much higher MAP values of 85.1 and 87.3 on these respective scenes. Moreover, CYWS-3D attains 0 MAP values on the Blocks and Potting scenes, whereas C-NERF achieves 67.9 and 34.5 MAP values on these two scenes, respectively. These quantitative results clearly demonstrate the superiority of the C-NERF model over CYWS-3D.

Fine-grained change detection. For fine-grained change detection, we choose 2D change detectors DTCDSN [26], BIT [6], A2Net [25] and USSFC-Net [23], and TinyCD [10].

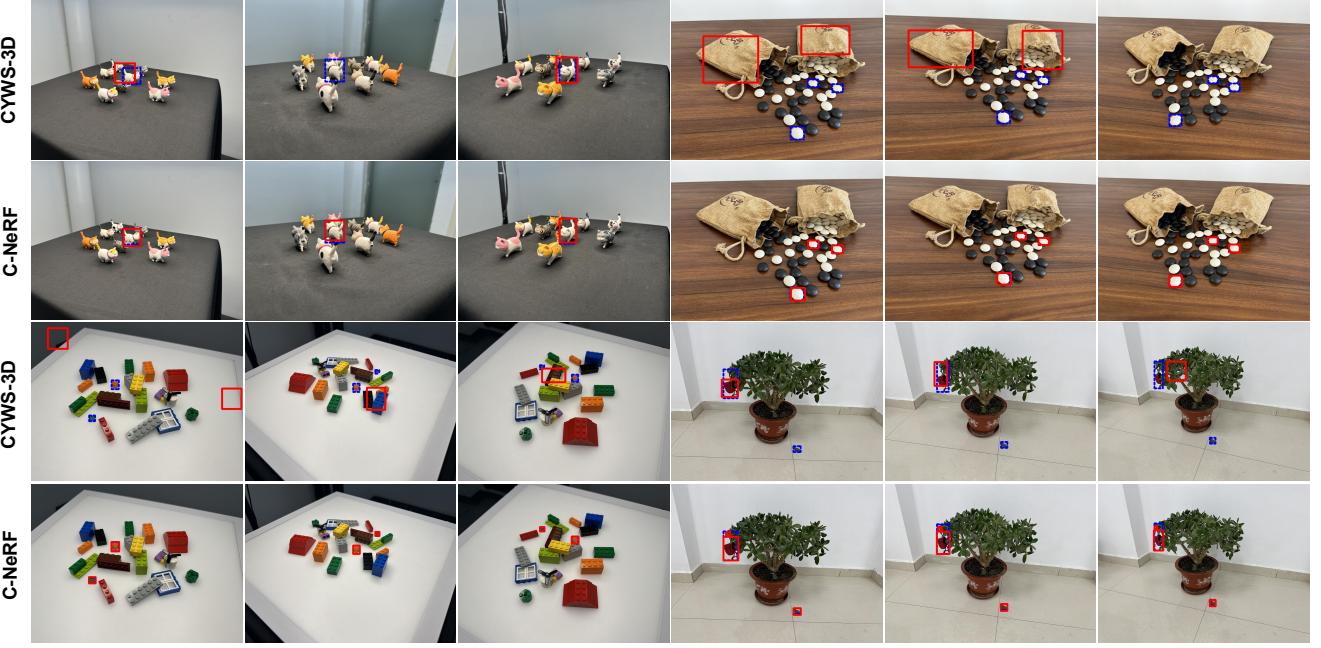


Figure 7. Comparison of instance change detection method CYWS-3D [39] with C-NeRF on different scenes. Red boxes are the predictions and blue boxes are the ground truth.

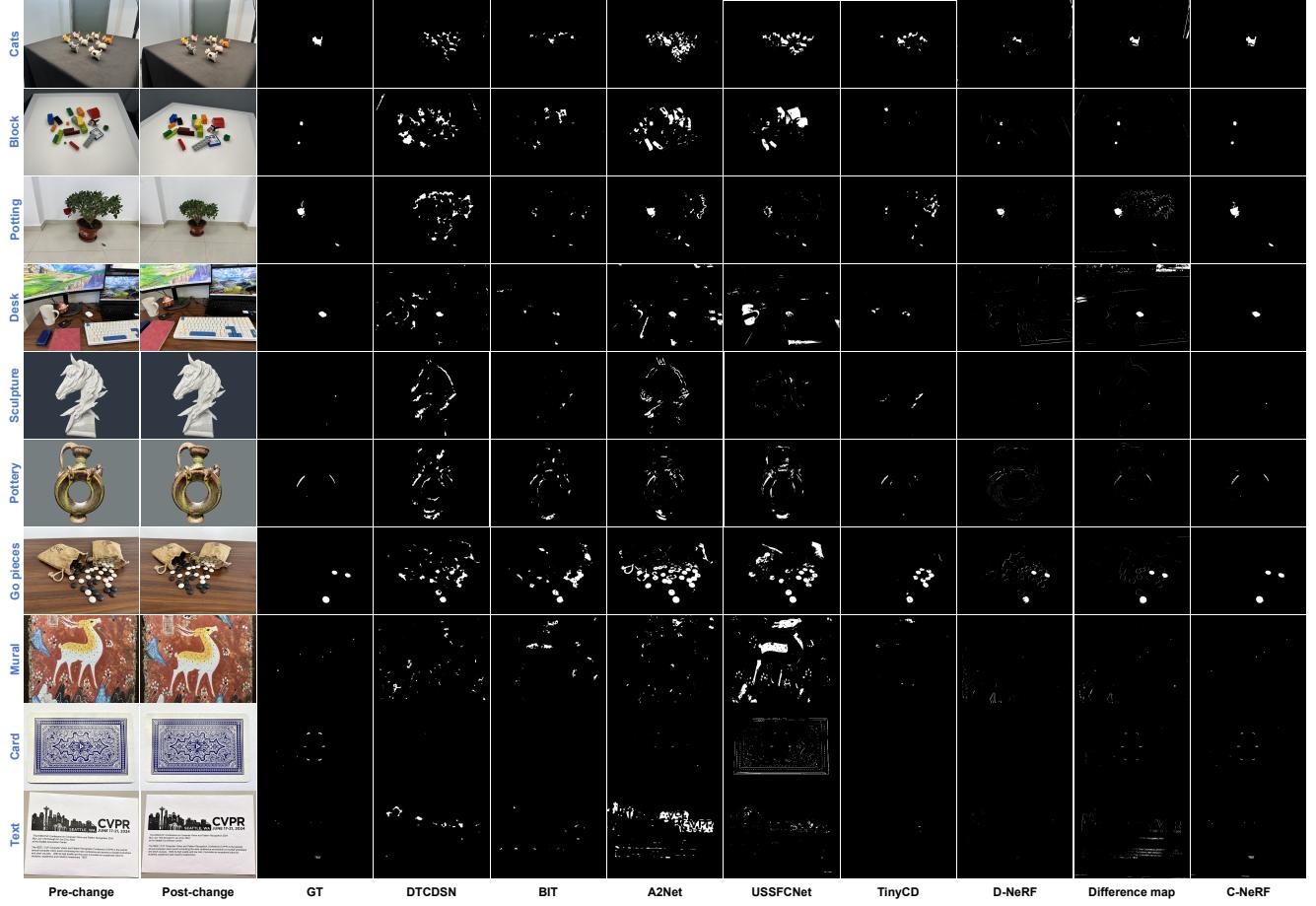


Figure 8. Some typical fine-grained change detection results of different methods on different scenes.

Table 2. Quantitative comparison of 2D change detection with C-NERF. The best values are highlighted in **bolded**.

Method	P↑	R↑	F1 ↑	IoU ↑
DTCDSN [26]	8.47	53.34	14.61	7.88
BIT [6]	20.97	48.81	29.34	17.19
A2Net [25]	11.41	78.29	19.92	11.06
USSFC-Net [23]	10.14	72.21	17.85	9.8
TinyCD [10]	34.49	47.62	40.02	25.01
D-NeRF [32]	35.68	41.92	35.22	24.32
Difference	42.61	67.46	46.52	33.81
C-NERF	81.2	75.1	74.48	62.18

We train all the change detectors with 1, 200 training image pairs. Besides, we modify D-NeRF [32] to generate aligned image pairs to calculate the change maps. We train D-NeRF with pre-change and post-change images as two time periods. Then we generate aligned image pairs at the given view directions of two time points. The final change maps of D-NeRF are generated by the absolute difference of the image pair with the best thresholds. We also generate the difference map with the aligned image pairs by C-NERF.

As shown in Fig. 8, we can find that almost all the compared methods produce fake changes on the object boundaries due to misaligned images. DTCDSN, BIT, A2Net, and TinyCD failed to detect tiny changes of the Mural (see the 8th row of Fig. 8) and Card (see the 9th row of Fig. 8). A2Net can identify large changes but with low precision (see the 2nd, 3rd and 4th rows of Fig. 8). TinyCD can identify changes, but they are not complete (see the 1st and 6th rows of Fig. 8). D-NeRF fails to detect the changes of the planner scenes in the 8th, 9th and 10th rows of Fig. 8. In contrast, C-NERF can identify real changes on different scenes while filtering out the fake ones.

As shown in Table 2, TinyCD achieves the best performances among the compared methods on NeRFCD. While more complex models, BIT and DTCDSN, show worse performance than TinyCD as they need more training data. Although D-NeRF can align pre- and post-change images, it obtains higher precision and lower recall than TinyCD. With the aligned image pairs, the simple absolute difference-based method achieves 46.52 and 33.81 F1 and IoU values, which are higher than other compared methods. The F1 and IoU values of C-NERF are 74.48 and 62.18, which rank at the first place. The quantitative results demonstrate that C-NERF is capable of capturing fine-grained changes in various scenes.

6.3. Discussion

Performance of C-NERF on different scenes. Fig. 9 shows the F1 and IoU values of C-NERF on different scenes. For general scenes with large changes, such as Cats, Desk, and Go pieces, C-NERF can achieve higher performances. While for the scenes with tiny changes, such as Mural, Card, and Text, C-NERF obtains lower F1 and IoU values. This might be caused by the rendering quality of NeRF. Low ren-

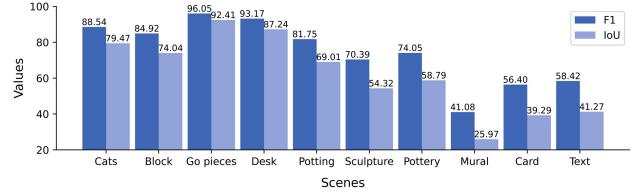


Figure 9. Performance of C-NERF on different scenes.

dering quality blurs the boundaries of the change regions, which reduces the areas of the changes. Thus, the change regions of C-NERF are small of the regions of ground truth.

Running time. We use 186 image pairs to train the NeRF models and render 120 change maps. As shown in able 3, aligning NeRFs is the most time-consuming process.

Table 3. Time consumption of main stages of C-NERF.

Stage	NeRF alignment	Change point detection	Rendering
Time (Hour)	16.25	0.75	2.6

Pros and cons of C-NERF. There are three main advantages of C-NERF: ① C-NERF eliminates the need for high-precision camera pose and image alignment, making scene shooting more accessible for ordinary people. ② C-NERF breaks the traditional pipeline of 2D change detection, introduces a novel change detection framework. ③ With C-NERF, we are able to synthesize change maps under arbitrary specified view directions, unlike 2D change detections which only identify changes in one view direction. However, there are some limitations: ① Large scenes like remote sensing and street view are not considered in this paper. Further study is needed on how to conduct change detection on these scenes. ② The performance of C-NERF is heavily influenced by NeRF [28]. For example, C-NERF cannot deal with the changes in the specular regions and has a long training time. Using different NeRF models might extend the capabilities of C-NERF, which will be studied in our future work.

7. Conclusion

In this work, we have identified a challenging but meaningful task, *i.e.*, 3D change detection based on NeRF representation, which aims to synthesize change maps of a scene under arbitrary specified view directions. To this end, we studied the naive solution based on existing NeRF methods and found that this method is not able to identify the change regions accurately. To address the limitations, we proposed a novel method denoted as C-NERF that contains three modules to address the challenges of the naive solutions. Specifically, we first performed spatial alignment to build two NeRFs for two image sets captured under two different timestamps. Then, we proposed to add a direction-consistent constraint to filter noise change points. Finally, we detailed the rendering method to generate the change map under a view direction. We build a new dataset, NERFCD, with 10 different

scenes. Extensive experiment results have demonstrated that our method is effective in detecting instance-level and fine-grained changes in various scenes, outperforming a series of NeRF-based 2D change detection methods significantly.

References

- [1] Pablo F Alcantarilla, Simon Stent, German Ros, Roberto Arroyo, and Riccardo Gherardi. Street-view change detection with deconvolutional networks. *Autonomous Robots*, 42:1301–1322, 2018. [2](#)
- [2] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII 16*, pages 696–712. Springer, 2020. [2](#)
- [3] Rareş Ambruş, Nils Bore, John Folkesson, and Patric Jensfelt. Meta-rooms: Building and maintaining long term spatial models in a dynamic world. In *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 1854–1861, 2014. [3](#)
- [4] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021. [3](#)
- [5] Hao Chen and Zhenwei Shi. A spatial-temporal attention-based method and a new dataset for remote sensing image change detection. *Remote Sensing*, 12(10):1662, 2020. [1](#), [2](#)
- [6] Hao Chen, Zipeng Qi, and Zhenwei Shi. Remote sensing image change detection with transformers. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–14, 2021. [6](#), [8](#)
- [7] Shuo Chen, Kailun Yang, and Rainer Stiefelhagen. Dr-tanet: Dynamic receptive temporal attention network for street scene change detection. In *2021 IEEE Intelligent Vehicles Symposium (IV)*, pages 502–509. IEEE, 2021. [2](#)
- [8] Tianlong Chen, Peihao Wang, Zhiwen Fan, and Zhangyang Wang. Aug-nerf: Training stronger neural radiance fields with triple-level physically-grounded augmentations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15191–15202, 2022. [2](#)
- [9] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16569–16578, 2023. [3](#)
- [10] Andrea Codegoni, Gabriele Lombardi, and Alessandro Ferrari. Tinyed: a (not so) deep learning model for change detection. *Neural Computing and Applications*, 35(11):8471–8486, 2023. [6](#), [8](#)
- [11] Marius Fehr, Fadri Furrer, Ivan Dryanovski, Jürgen Sturm, Igor Gilitschenski, Roland Siegwart, and Cesar Cadena. Tsdf-based change detection for consistent long-term dense reconstruction and dynamic object discovery. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5237–5244, 2017. [3](#)
- [12] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14346–14355, 2021. [3](#)
- [13] Evan Herbst, Peter Henry, Xiaofeng Ren, and Dieter Fox. Toward object discovery and modeling via 3-d scene comparison. In *2011 IEEE International Conference on Robotics and Automation*, pages 2623–2629, 2011. [3](#)
- [14] Tao Hu, Shu Liu, Yilun Chen, Tiancheng Shen, and Jiaya Jia. Efficientnerf efficient neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12902–12911, 2022. [3](#)
- [15] Qingbao Huang, Yu Liang, Jielong Wei, Yi Cai, Hanyu Liang, Ho-fung Leung, and Qing Li. Image difference captioning with instance-level fine-grained feature representation. *IEEE transactions on multimedia*, 24:2004–2017, 2021. [2](#)
- [16] Rui Huang, Ruofei Wang, Qing Guo, Jieda Wei, Yuxiang Zhang, Wei Fan, and Yang Liu. Background-mixed augmentation for weakly supervised change detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 7919–7927, 2023. [1](#), [2](#)
- [17] Stefan Huwer and Heinrich Niemann. Adaptive change detection for real-time surveillance applications. In *Proceedings Third IEEE International Workshop on Visual Surveillance*, pages 37–46. IEEE, 2000. [1](#)
- [18] Harsh Jhamtani and Taylor Berg-Kirkpatrick. Learning to describe differences between pairs of similar images. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4024–4034, 2018. [2](#)
- [19] Salman H Khan, Xuming He, Fatih Porikli, and Mohammed Bennamoun. Forest change detection in incomplete satellite images with deep neural networks. *IEEE Transactions on Geoscience and Remote Sensing*, 55(9):5407–5423, 2017. [1](#)
- [20] Hoeseong Kim, Jongseok Kim, Hyungseok Lee, Hyunsung Park, and Gunhee Kim. Agnostic change captioning with cycle consistency. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2095–2104, 2021. [2](#)
- [21] Alexander Krawciw, Jordy Sehn, and Timothy D Barfoot. Change of scenery: Unsupervised lidar change detection for mobile robots. *arXiv preprint arXiv:2309.10924*, 2023. [1](#)
- [22] Tao Ku, Sam Galanakis, Bas Boom, Remco C Veltkamp, Darshan Bangera, Shankar Gangisetty, Nikolaos Stagakis, Gerasimos Arvanitis, and Konstantinos Moustakas. Shrec 2021: 3d point cloud change detection for street scenes. *Computers & Graphics*, 99:192–200, 2021. [3](#)
- [23] Tao Lei, Xinze Geng, Hailong Ning, Zhiyong Lv, Maoguo Gong, Yaochu Jin, and Asoke K Nandi. Ultralightweight spatial–spectral feature cooperation network for change detection in remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–14, 2023. [6](#), [8](#)
- [24] Yinjie Lei, Duo Peng, Pingping Zhang, Qihong Ke, and Haifeng Li. Hierarchical paired channel fusion network for street scene change detection. *IEEE Transactions on Image Processing*, 30:55–67, 2020. [2](#)
- [25] Zhenglai Li, Chang Tang, Xinwang Liu, Wei Zhang, Jie Dou, Lizhe Wang, and Albert Y Zomaya. Lightweight remote

- sensing change detection with progressive feature aggregation and supervised attention. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–12, 2023. 6, 8
- [26] Yi Liu, Chao Pang, Zongqian Zhan, Xiaomeng Zhang, and Xue Yang. Building change detection for remote sensing images using a dual-task constrained deep siamese convolutional network model. *IEEE Geoscience and Remote Sensing Letters*, 18(5):811–815, 2020. 6, 8
- [27] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: learning dynamic renderable volumes from images. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019. 2
- [28] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2, 3, 6, 8
- [29] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 3
- [30] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022. 3
- [31] Ariyo Oluwasanmi, Enoch Frimpong, Muhammad Umar Aftab, Edward Y Baagyere, Zhiguang Qin, and Kifayat Ullah. Fully convolutional captionnet: Siamese difference captioning attention model. *IEEE access*, 7:175929–175939, 2019. 2
- [32] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021. 2, 8
- [33] Rongjun Qin, Jiaojiao Tian, and Peter Reinartz. 3d change detection—approaches and applications. *ISPRS Journal of Photogrammetry and Remote Sensing*, 122:41–56, 2016. 3
- [34] Yue Qiu, Shintaro Yamamoto, Kodai Nakashima, Ryota Suzuki, Kenji Iwata, Hirokatsu Kataoka, and Yutaka Satoh. Describing and localizing multiple changes with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1971–1980, 2021. 2
- [35] Yue Qiu, Shintaro Yamamoto, Ryosuke Yamada, Ryota Suzuki, Hirokatsu Kataoka, Kenji Iwata, and Yutaka Satoh. 3d change localization and captioning from dynamic scans of indoor scenes. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1176–1185, 2023. 3
- [36] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14335–14345, 2021. 3
- [37] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIX 16*, pages 623–640. Springer, 2020. 2
- [38] Ragav Sachdeva and Andrew Zisserman. The change you want to see. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3993–4002, 2023. 2
- [39] Ragav Sachdeva and Andrew Zisserman. The change you want to see (now in 3d). In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2060–2069, 2023. 2, 6, 7
- [40] Ken Sakurada and Takayuki Okatani. Change detection from a street image pair using cnn features and superpixel segmentation. In *British Machine Vision Conference*, 2015. 2
- [41] Ken Sakurada, Takayuki Okatani, and Koichiro Deguchi. Detecting changes in 3d structure of a scene from multi-view images captured by a vehicle-mounted camera. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 137–144, 2013. 2
- [42] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4104–4113, 2016. 3, 4, 5
- [43] Koji Takeda, Kanji Tanaka, and Yoshimasa Nakamura. Life-long change detection: Continuous domain adaptation for small object change detection in everyday robot navigation. In *2023 18th International Conference on Machine Vision and Applications (MVA)*, pages 1–5. IEEE, 2023. 1
- [44] Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022. 2
- [45] Aparna Taneja, Luca Ballan, and Marc Pollefeys. Image based detection of geometric changes in urban environments. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2336–2343. IEEE, 2011. 3
- [46] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *AcM Transactions on Graphics (TOG)*, 38(4):1–12, 2019. 2
- [47] Yunbin Tu, Tingting Yao, Liang Li, Jiedong Lou, Shengxiang Gao, Zhengtao Yu, and Chenggang Yan. Semantic relation-aware difference representation learning for change captioning. In *Findings of the association for computational linguistics: ACL-IJCNLP 2021*, pages 63–73, 2021. 2
- [48] Ali Osman Ulusoy and Joseph L Mundy. Image-based 4-d reconstruction using 3-d change detection. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part III 13*, pages 31–45. Springer, 2014. 3
- [49] Jian Zhang, Yuanqing Zhang, Huan Fu, Xiaowei Zhou, Bowen Cai, Jinchi Huang, Rongfei Jia, Binqiang Zhao, and Xing Tang. Ray priors through reprojection: Improving neural radiance fields for novel view extrapolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18376–18386, 2022. 3
- [50] Xiang Zhang, Honggang Wu, Min Wu, and Celimuge Wu. Extended motion diffusion-based change detection for airport

In this supplementary material, we provide the code link of **C-NERF**, dataset details, quantitative results of different change detectors on different scenes, more visualization results of C-NERF, and a video demo.

8. Dataset Details

We capture all scenes with iPhone 14 Pro and set the resolution of the images to 1008×756 to train D-NeRF and C-NERF. To train 2D change detectors, we resize the images into 512×512 . The detailed number of images for each scene is shown in Table. 4.

Table 4. Dataset details.

Scene	Pre-change	Post-change	Render	Training	Testing
Cats	118	125	120	120	10
Block	186	187	120	120	10
Potting	21	21	120	120	10
Desk	20	20	120	120	10
Sculpture	26	30	120	120	10
Pottery	26	22	120	120	10
Go pieces	22	26	120	120	10
Mural	22	18	120	120	10
Card	27	27	120	120	10
Text	17	20	120	120	10

9. Quantitative results of different change detectors on different scenes

We report a quantitative comparison of D-NeRF, results of difference maps with aligned image pairs generated by C-NERF, and C-NERF in different scenes in Table 5. We can find that C-NERF achieves the highest F1 and IoU values in all ten scenes. This demonstrates that C-NERF is more suitable for detecting changes in various scenes.

10. More visualization results of C-NERF

To demonstrate the effectiveness of C-NERF in detecting changes in different viewpoints, we provide more visualization results in Figs. 10, 11, 12, 13, and 14. For each scene, we provide detection results of five view directions with intervals of ten. From the figures, we find that C-NERF can detect the changes of each scene under different viewpoints.

11. Video demo

In order to better showcase our results, we have provided images including pre-change image, post-change image, difference maps, and results of C-NERF in each video frame synchronously. Please see the attached video.

Table 5. Quantitative comparison of D-NeRF, results of difference maps with aligned image pairs generated by C-NERF, and C-NERF in different scenes. The best values are **bolded**.

Scene	Method	P \uparrow	R \uparrow	F1 \uparrow	IoU \uparrow
Cats	D-NeRF	59.24	60.57	59.68	42.71
	Difference	58.80	83.5	68.86	52.7
	C-NERF	86.70	90.46	88.54	79.47
Blocks	D-NeRF	25.83	79.48	38.06	23.88
	Difference	50.1	91.17	64.3	48.22
	C-NERF	77.70	93.62	84.92	74.04
Go pieces	D-NeRF	46.34	80.27	58.41	41.42
	Difference	99.97	17.69	30.04	17.68
	C-NERF	92.66	99.71	96.05	92.41
Desk	D-NeRF	20.81	41.72	27.74	16.42
	Difference	19.62	86.09	31.67	18.96
	C-NERF	96.02	90.48	93.17	87.24
Potting	D-NeRF	83.87	70.77	76.71	62.23
	Difference	44.97	89.42	59.81	42.68
	C-NERF	71.59	94.66	81.75	69.01
Sculpture	D-NeRF	69.5	38.78	49.68	33.08
	Difference	19.47	53.42	28.53	16.6
	C-NERF	92.48	56.81	70.39	54.32
Pottery	D-NeRF	7.51	9.3	8.25	4.3
	Difference	57.99	37.29	45.36	29.33
	C-NERF	90.08	62.86	74.05	58.79
Mural	D-NeRF	4.56	7.7	5.71	2.95
	Difference	5.92	29.32	9.74	5.15
	C-NERF	84.08	27.2	41.08	25.97
Card	D-NeRF	36.84	22.05	26.02	15.2
	Difference	79.19	39.61	52.8	35.87
	C-NERF	77.69	44.28	56.40	39.29
Text	D-NeRF	2.26	8.6	1.92	0.98
	Difference	8.54	66.56	15.05	8.18
	C-NERF	43.04	90.92	58.42	41.27

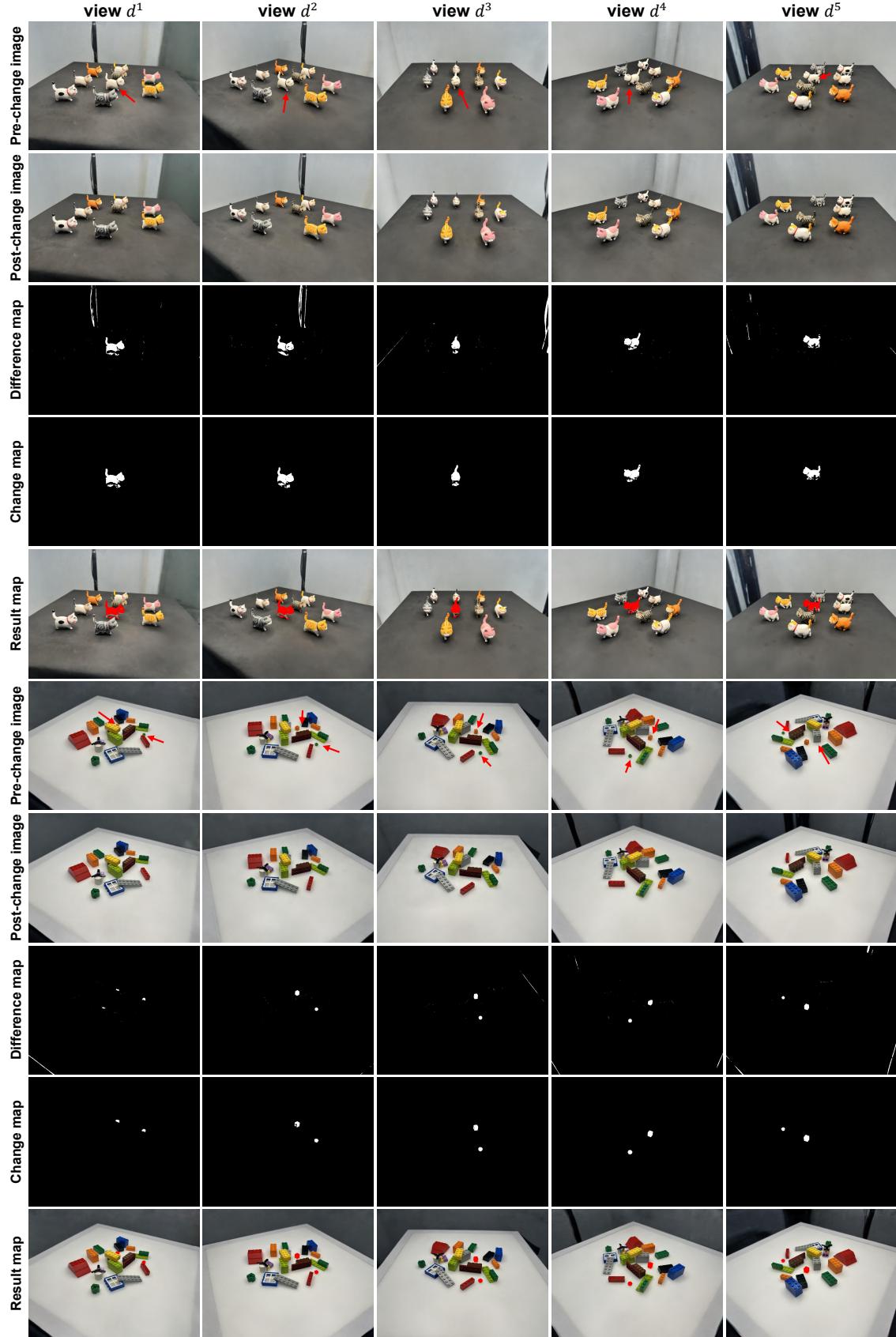


Figure 10. Change detection results of C-Nerf of different viewpoints of Cats and Blocks. The change objects are pointed by the red arrows.

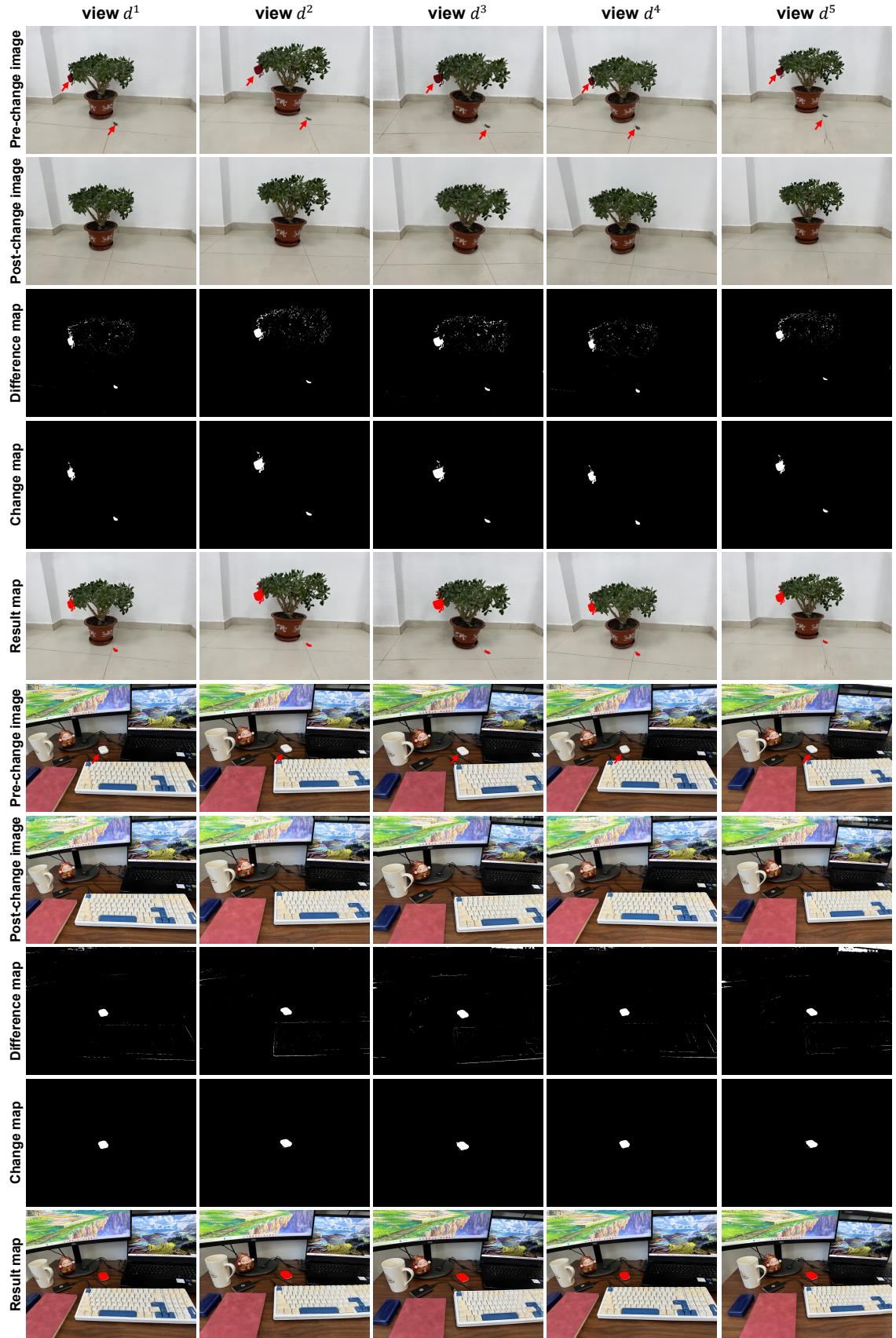


Figure 11. Change detection results of C-NERF of different viewpoints of Potting and Desk. The change objects are pointed by the red arrows



Figure 12. Change detection results of C-NERF of different viewpoints of Pottery and Sculpture. The change objects are pointed by the red arrows

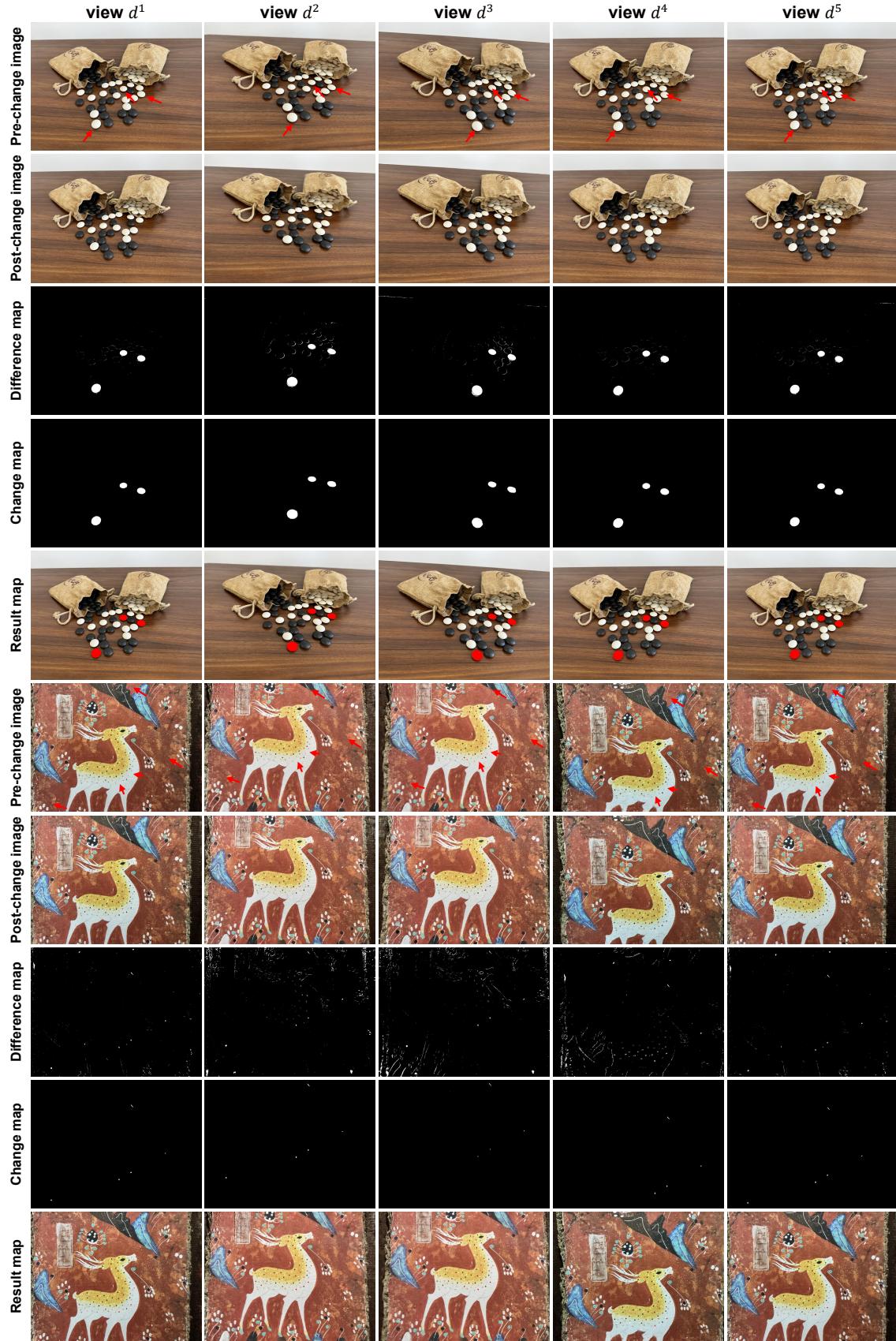


Figure 13. Change detection results of C-NERF of different viewpoints of Go pieces and Mural. The change objects are pointed by the red arrows

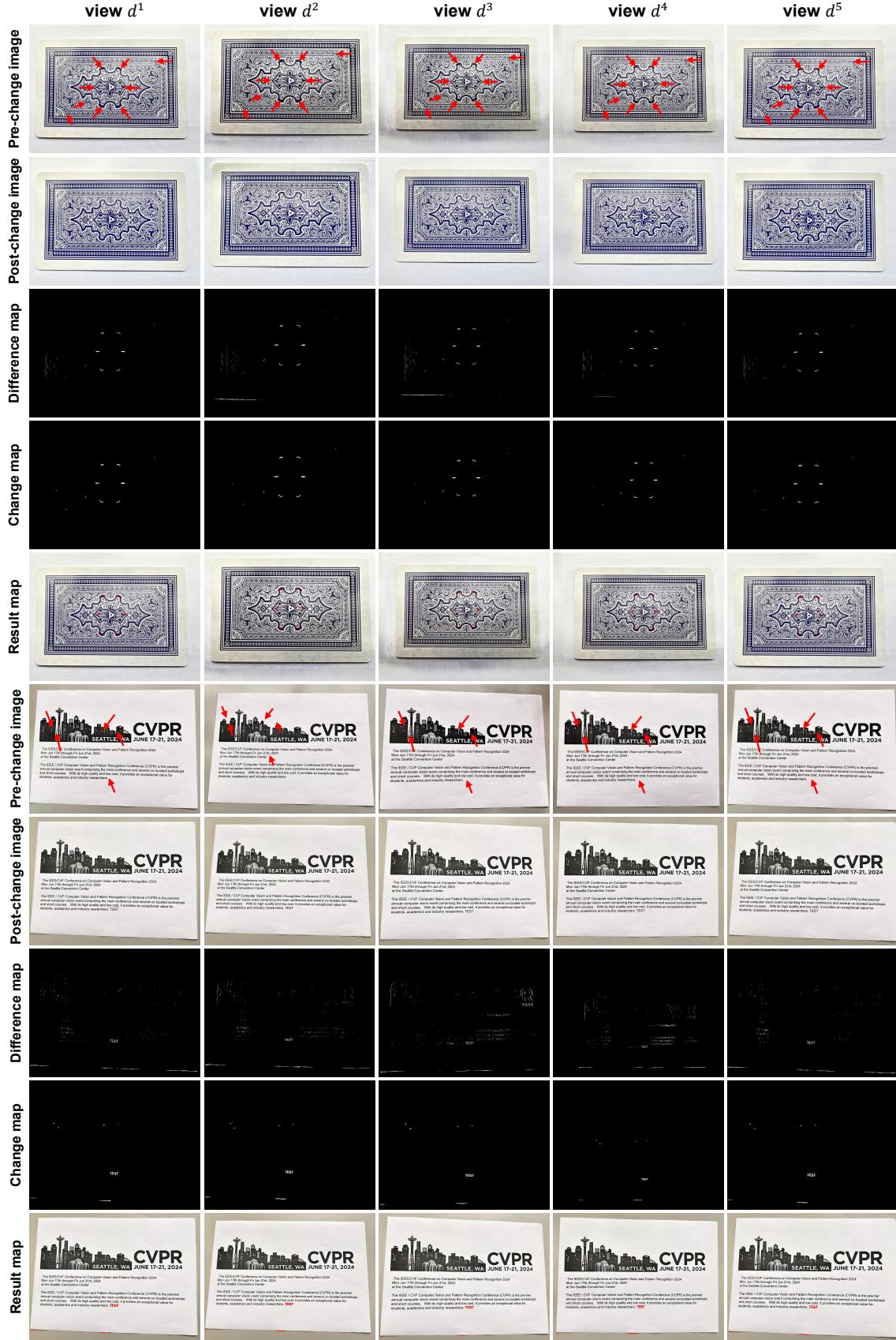


Figure 14. Change detection results of C-NERF of different viewpoints of Card and Text. The change objects are pointed by the red arrows