# Implicit neural representation for change detection

Peter Naylor<sup>1†</sup>, Diego Di Carlo<sup>1</sup>, Arianna Traviglia<sup>2</sup>, Makoto Yamada<sup>1,3</sup> and Marco Fiorucci<sup>2</sup>

<sup>1</sup> RIKEN AIP, Kyoto, Japan firstname.lastname@riken.jp

<sup>2</sup> Istituto Italiano di Tecnologia, Venice, Italy firstname.lastname@iit.it

<sup>3</sup> OIST, Okinawa, Japan makoto.yamada@oist.jp

# Abstract

Identifying changes in a pair of 3D aerial LiDAR point clouds, obtained during two distinct time periods over the same geographic region presents a significant challenge due to the disparities in spatial coverage and the presence of noise in the acquisition system. The most commonly used approaches to detecting changes in point clouds are based on supervised methods which necessitate extensive labelled data often unavailable in real-world applications. To address these issues, we propose an unsupervised approach that comprises two components: Implcit Neural Represenation (INR) for continuous shape reconstruction and a Gaussian Mixture Model for categorising changes. INR offers a grid-agnostic representation for encoding bi-temporal point clouds, with unmatched spatial support that can be regularised to enhance high-frequency details and reduce noise. The reconstructions at each timestamp are compared at arbitrary spatial scales, leading to a significant increase in detection capabilities. We apply our method to a benchmark dataset comprising simulated LiDAR point clouds for urban sprawling. This dataset encompasses diverse challenging scenarios, varying in resolutions, input modalities and noise levels. This enables a comprehensive multi-scenario evaluation, comparing our method with the current state-ofthe-art approach. We outperform the previous methods by a margin of 10% in the intersection over union metric. In addition, we put our techniques to practical use by applying them in a real-world scenario to identify instances of illicit excavation of archaeological sites and validate our results by comparing them with findings from field experts.

# 1. Introduction

In contemporary times, we observe the Earth through various sensors at unprecedented spatial and temporal resolutions. One of the most popular Earth Observation technologies is LiDAR: by using laser light to measure distances, it creates detailed three-dimensional maps or point cloud representations of the Earth's surface and objects, see Fig. 1. LiDAR has been applied to autonomous driving [26], robotics [31], digital terrain mapping [51], city planning [36], urban sprawling [9, 10, 23], and cultural heritage [18, 43, 50]. The surge in popularity can be attributed to three main factors. Firstly, LiDAR data exhibit a remarkable level of precision, with spatial resolutions typically falling below 1 m in most applications (though they can be larger based on the distance to the scanned scene or specific detector characteristics), effectively capturing the intricate details of a 3D environment. Secondly, LiDAR acquisition systems remain unaffected by varying lighting conditions. Lastly, LiDAR possesses the capability to map terrain and unveil structures concealed by vegetation canopies [4, 6].

The process of comparing two or more slightly coregistered EO data to identify and analyse discrepancies that have emerged between them is called Change Detection (CD). In this work, we focus on the application of CD to urban sprawl, shown in Fig. 1, and on the identification of illicit excavations of archaeological sites (looting) by detecting relevant changes in height. Urban sprawl monitoring, which entails the identification of recently erected and demolished structures within multi-temporal LiDAR point cloud datasets, has been identified as a method to assist landscape and city managers in promoting sustainable development [9]. The identification of illicit excavations of archaeological sites holds paramount significance due to the potential for looting to cause damage, displacement, or irrevocable loss of priceless archaeological artefacts.

The majority of existing techniques handle data that are defined on discrete regular grids (e.g., images) with supervised learning frameworks. Over the past decade, Deep Neural Networks (DNN) have risen to prominence as the state-of-the-art solution for CD on LiDAR data [11,23,54]. However, these approaches require a preprocessing step to cleanse and project a raw 3D LiDAR point cloud onto a consistent 2D regular grid with predefined spatial resolution. The prevailing projection techniques encompass the use of the Digital Elevation Model (DEM) [32] and Digi-

<sup>&</sup>lt;sup>†</sup>Now at ESA/ESRIN, Φ-lab, Italy. peter.naylor@esa.int



Figure 1. Simulated airborne LiDAR data for change detection of a clipped test data area.



Figure 2. Proposed unsupervised CD of PCs based on implicit neural representation and clustering.

tal Surface Model (DSM) [15, 28, 54]. Consequently, the accuracy of these models hinges on the accuracy of the projection, and memory complexity increases proportionately with the desired resolution, limiting their application to large point clouds (PCs). In contrast, our investigation delves into the application of CD directly onto the off-grid, raw 3D LiDAR PCs in an unsupervised manner, which suits better real-world applications.

In summary, our approach involves two main steps, illustrated in Fig. 2. In the first step, we use INR [53] to encode the bi-temporal PC as a continuous function of both time and space. This function is estimated blindly from data using the total variation (TV) norm [3] to enforce discontinuities along the time dimension to model sharp temporal changes and increase robustness to noisy measurements. In the second step, we categorise the change in altitude given by decoding the surface at both timestamps at an arbitrary resolution. To validate the proposed approach, we consider two EO datasets. The first is an open simulated airborne Li-DAR dataset comprising 15 distinct PCs. The goal here is to identify both newly built and demolished buildings. The second dataset consists of a pair of PCs over the Kulen, a region of Cambodia featuring numerous archaeological structures lying under the forest where instances of illegal excavations have been recorded. Our goal is to precisely pinpoint locations where looting has taken place and to record the shape and attributes of the looting pits.

Contributions. This paper tackles the lack of extensively

labelled 3D PCs for CD in real-world applications by introducing an effective unsupervised pipeline. Specifically:

- Using a single and scalable Implicit Neural Representation to encode the height of two geo-referenced 3D point clouds as a continuous quantity.
- Denoising and regularising the encoded surface with total variation norm over the time dimension.
- Proposing an end-to-end unsupervised change detection pipeline that leverages on automatic hyperparameter tuning and a simple clustering algorithm.

#### 2. Related works

The off-grid and noisy nature of PCs collected by a Li-DAR system detrimentally impacts the performances of CD methods. Moreover, weather conditions and remote sensor trajectory may differ for two measurements at different timestamps, leading to spatial unmatching support. Despite this, supervised methods for CD can achieve great performances even in challenging noisy data. However, their usage is still limited to a few real-world applications due to the prohibitive cost of collecting and labelling datasets.

For these reasons, unsupervised methods represent an attractive alternative. Such methods can be broadly classified into three categories [9]: those based on distance computation [19, 45], those based on optimal transport [8, 17] and those adapting PC DNN for unsupervised learning [10, 12]. Distance-based methods, such as C2C [19], and M3C2 [25, 45] divide the PC via octrees, estimate the surface normal and orientation to calculate pair-wise euclidean distances. Alternatively, optimal transport-based methods estimate a distance based on the projection matrix of the first PC onto the second. In both cases, the actual changes are then classified via empirical thresholding, or the OTSU method [33]. These methods have been developed and applied to the only available airborne LiDAR dataset for building CD [9]. However, datasets usually contain millions of points because they are acquired with high-spatial resolution acquisition systems [42]. The previous methods do not scale well with the data size and have to subdivide the PC

for analysis. The adaptation of PC DNN for unsupervised learning are very recent and were absent at the initial time of writing (they have yet to be peer-reviewed) [10, 12]. PC DNN use Kernel Points [49] to compare co-referenced subsets of the two PCs. SSL-DCVA uses intermediate representation of the cylinders in the DNN to build an estimator for change [12] inspired by DCVA [41]. DC3DCD inspires itself from DeepCluster [5] with pseudo-labels for training. Based on recent reviews [47, 52], the aforementioned methods are the only unsupervised methods and fully automatic state-of-the-art approaches for PC CD.

Other methods for unsupervised CD do not use raw 3D LiDAR data directly. Instead, they used 2D images obtained by projecting the 3D PCs on a 2D regular grid, e.g., for DEM and DSM [56]. Due to the grid's regularity, ordering and consistency, these projected 2D images are ideally suited for convolutional operations. Hence, convolutional neural network-based architectures can be applied to the 2D digital models for building CD [44, 54, 55]. However, these projections lead to precision loss of the LiDAR data as the height measurements are interpolated to output the DSM or the DEM [9]. In contrast, we use INR to model the surface and obtain an intermediate 2D images at any desired resolution of the surfaces height.

Novel DNN models, called INRs [53], have tackled off-grid PC analysis [23, 38] and 3D surface reconstruction [21, 46, 48] in both supervided and unsupervised way. The INRs are *coordinate*-based continuous deep learning models that map coordinates to a target value, called field (e.g. from pixel position to colour in an image). This function can be estimated by fitting observation and inducing desired properties by application-based regularisation terms. The major property of these parametric models is their capability to natively interpolate the target field.

A common limitation of INR is their poor capability of capturing high-frequency details of the surface, referred to as spectral bias. To address this, positional encoding and Random Fourier Features (RFF) [39, 48] of the input coordinates have become standard practice. This allowed to apply INR to a new plethora of applications and methodologies, such as 3D shape reconstruction from sparse images, animation of human bodies and faces, as well as video cod-ing [53]. Moreover, RFF have improved Physics-Informerd Neural Networks (PINNs) [40], which are coordinate-based DNNs trained to fit observations while solving partial differential equation evaluated with network's backpropagation algorithm [13, 22].

An alternative approach to surpass the spectral bias is SIREN [46], where the standard nonlinearities are replaced with periodic sine functions. Such architectures show a very low reconstruction error, both in fitting the target field and its gradient with respect to the input. SIRENs have become standard backbone networks for both INR and PINNs [53]. Despite their success, SIRENs suffer from difficulties in training as they are prone to overfit or being stuck in local minima, for which careful parameter tuning is required. These issues have been addressed in subsequent works [16, 27], where novel architectures have been proposed claiming easier training and faster convergence. Nevertheless, SIREN's performance depends on the applications. As reported in this work, SIREN doesn't consistently outperform RFF-based architecture.

Most of the INR-based models for 3D reconstruction recover the object shape by minimising the signed distance between any given point and the closest surface [35]. Then, the associated INR is a function of the three spatial coordinates. A positive sign implies the point lies outside the object and vice versa. The shape contours are then found by checking the 0 iso-line of the learned function. Neural Unsigned Distance Fields [7] remove the sign from the distance function and encode continuous locations for a stronger regularisation, leading to better reconstruction. Nonetheless, these methods require knowledge of the normals vectors of each point within the PC [21]. While the rudimentary normal estimation of LiDAR data might be achieved using conventional geometrical methods, our preliminary investigations have demonstrated that errors stemming from such preprocessing steps significantly undermine the performances. As elucidated in Section 3, the phenomenon of urban sprawl can be reasonably approximated through a continuous function across a 2D spatial domain, while the exploration of 3D aspects will be deferred to forthcoming research endeavors.

### 3. Method

We want to detect changes in two geo-referenced LiDAR PCs with unsupervised methods. A PC at time t will be denoted by  $\mathcal{X}_t \subset \mathbb{R}^3$  where each element is a 3D coordinate, i.e.  $(x, y, z) \in \mathcal{X}_t$ . We will denote by  $t_0 = 0$  and  $t_1 = 1$  the two timestamps for which we wish to detect change. If the support of  $\mathcal{X}_0$  and  $\mathcal{X}_1$  match, we naturally define the addition of an element by a positive difference, i.e. the 2D point (x, y) is of the label "Addition" if the associated altitudes:  $z_1 - z_0 > \alpha$ , where  $z_t$  corresponds to the altitude at time t and with  $\alpha$  a fixed scalar. Similarly, we can define the "Deletion" class by a negative difference.

In general, the above operations are not directly applicable to LiDAR PCs as the supports do not match due to different acquisition conditions. To fix the support, some have used projection methods [9] and optimal transport [17]. In the current paper, we fix the support matching by estimating a surface from the PC, allowing us to interpolate and query any spatial point. In other words, we reconstruct the surface at a given time point. We can then compute the difference to find additions or deletions. It is important to note that we are estimating a function that maps a position to an altitude, which is different from the actual PC. Some points will share the same (x, y) but have different altitudes due to the inclination of the LiDAR emitter and the verticality of elements in the maps, especially with buildings.

#### 3.1. Regression model

We denote by  $f_{\theta}$  a DNN model with learnable parameters  $\theta$ . Given an input vector  $\mathbf{v}$ , we estimate the density by  $f_{\theta}$  with inputs in V and values in  $\mathbb{R}$ . As a baseline, we can independently reconstruct the first and second PC with  $V = \mathbb{R}^2$ . We learn two functions and name this model (D). The formulae for CD with the baseline method is:

$$\Delta z(x,y) = (f_{\theta_1} - f_{\theta_0})(x,y). \tag{1}$$

Here  $\theta_t$  corresponds to the set of DNN parameters to reconstruct the surface at time *t* that are optimised by minimising the mean squared error (MSE) between the estimated and observed altitude.

A more compact and efficient representation is possible where we modify the input  $\mathbf{v} \in \mathbb{R}^3$  of the network to incorporate time. In particular, we learn a single model (S) with parameters  $\theta$ , and the detection formula is modified to:

$$\Delta z(x,y) = f_{\theta}(x,y,t_1) - f_{\theta}(x,y,t_0) \tag{2}$$

A priori, it is not evident to know in advance which method is better suited for a given dataset. In some configurations of the simulated datasets,  $\mathcal{X}_0$  and  $\mathcal{X}_1$  are not drawn from the same distribution, which could potentially be harmful to the single model. Therefore we investigate the benefits of a single model as opposed to two. In the following sections, we will describe in more detail the different regularisations and model specificities we apply to improve reconstruction.

#### **3.2. Random Fourier Features**

As common in INR models, we map the input v to a higher dimensional space with RFF. It has been shown that this projection is crucial for estimating high frequencies for a better reconstruction [48]. For a fixed model, RFF is defined as follows:

$$B \in \mathbb{R}^{M \times 3}, \ \forall (i, j), B_{ij} \sim \mathcal{N}(0, \sigma)$$
$$\gamma_B(\mathbf{v}) = [\cos(2\pi B \mathbf{v}), \ \sin(2\pi B \mathbf{v})]$$

The size of B depends on the size of  $\mathbf{v}$ , and here we supposed  $\mathbf{v} \in \mathbb{R}^3$ . The mapping size M and the scale  $\sigma$  are two hyper-parameters that need tuning.

#### 3.3. Network architecture

We propose several network architectures comprising MLP layers, different activations, and skip layers in Table A.2 of the Appendix. We show in Fig. A.1 a particular model that we name 'skip-ten-only'. In particular, the network size will depend on the data complexity. We use skip

layers allowing a better gradient flow [20]. The activation functions will be ReLu or hyperbolic tangent (tanh), as tanh allows the resulting  $f_{\theta}$  to be  $C^{\infty}$ , we do not use any activation layers for the last layer. We use a similar architecture to the one presented in Fig. A.1 for the SIREN methods. We allow the final model to fine-tune the architecture for both methods by making it a hyperparameter.

### 3.4. Total variation norm (TVN)

The total variation norm, defined as  $\sum |u_{i+1} - u_i|$  is a standard regularisation scheme for sequential like data [3]. It is also helpful to smooth spatial patterns, as a sudden altitude change should be penalised. Such a scheme has been generalised to a continuous version by enforcing that the resulting gradient of the function  $f_{\theta}$  be sparse over the spatial coordinates [40]. It is possible to add to the loss function the regularisation term  $\mathcal{R}_{TV} = |\frac{\partial f_{\theta}}{\partial x}| + |\frac{\partial f_{\theta}}{\partial y}|$ .

### **3.5.** Time difference (TD)

Similarly to the discrete total variation norm, we can enforce that the change over time be sparse, which indicates to the DNN that most points do not change over time, but we will allow some to change. To enforce such a constraint, we add the following regularisation to the loss function:  $\mathcal{R}_{TD} = |f_{\theta}(x, y, t_1) - f_{\theta}(x, y, t_0)|$ . This regularisation is only possible when the input v contains time, and therefore it cannot be applied where we reconstruct the surface for two models (1). This regularisation over time is very similar to the total variation norm over the temporal domain. We stress that shape reconstruction does not need this regularisation term, and this is only a CD regulariser. In particular, adding this term to the loss allows us to fuse the information from both PCs more efficiently and enables the PCs to benefit from each other mutually.

#### **3.6. SIREN**

From an architectural point of view, the SIREN network [46] is a simple modification of an MLP where standard activation functions, e.g. ReLU and tanh, are continuous sine functions. This substitution enables the modelling of a continuous complicated signal without the need for explicit upsampling in various domains. The SIREN network can then be described by

$$f(\mathbf{x}) = \mathbf{W}_n(g_{n-1} \circ g_{n-2} \circ \dots \circ g_0 \mathbf{x}) + \mathbf{h}_n,$$
  
where  $g_i(\mathbf{h}_i) = \sin(\mathbf{W}_i \mathbf{h}_i + \mathbf{h}_i).$  (3)

SIREN operates as a composition of sinusoidal transformations recalling the principle behind RFF. In fact, it has been shown that positional encoding with RFF is equivalent to periodic nonlinearities with one hidden layer as the first DNN layer [2]. Good initialisation of the weights is critical for their successful training. To avoid saturation of the sine activations, a scalar hyperparameter  $\sigma$  is introduced to scale the layers' weights, like for RFF.

#### **3.7.** Unsupervised labelling of $\Delta z$

Once the surface reconstruction is performed, we have access to  $\Delta z$  given by equation (1) or (2). The OTSU threshold method has been used to separate binary sources [30,33] for CD. However, in our case, we have three sources to distinguish and use a Gaussian Mixture Model (GMM) with three components [29]. We show in Fig. A.2 the application of GMM to  $\Delta z$  for a small clipped sample. The success of the GMM depends on the distribution of  $\Delta z$ . GMM will divide the distribution into three, regardless of the shape of the distribution and lead to a random score.

### 4. Experimental setting

#### 4.1. Simulated dataset

We use the publicly simulated airborne LiDAR dataset for CD: Urb3DCD [9]. Even if this dataset is simulated, it mimics true data with different noise levels and sensors used in practice. Five simulation configurations are given: 1 - low resolution – low noise, 2 - high resolution – low noise, 3 - low resolution – high noise, 4 - photogrammetry and 5 - multi-sensor. The photogrammetry setting is low resolution, high noise and tight scan angle for each timestamp. They mimic satellite acquisition. The multi-sensor simulation is characterised by  $pc_0$  and  $pc_1$  having different resolutions and noise levels.  $pc_0$  is low resolution and high noise, whereas  $pc_1$  is high resolution and low noise. In this situation  $\#pc_1 >> \#pc_0$  which is different to other subsets where  $\#pc_1 \approx \#pc_0$ . # denotes the cardinality of the set. Each configuration is divided into training and testing datasets. We will only apply the method to the testing sets as the methods used are unsupervised. For each testing configuration, three simulated datasets exist where the ground truth is different for each. The testing set comprises three different geographical areas of the city of Lyon in France [9]. Only the second PC  $pc_1$  is annotated with additions and deletion changes and will be used to evaluate the methods. In Fig. 1c, we overlay  $pc_1$  with the annotation.

This dataset was updated (Urb3DCD-v2) to include vegetation changes and mobile objects in two simulation settings, low density and multi-sensor LiDAR acquisition. This version is used to compare the DC3DCD model [10].

#### 4.2. Metrics

The minimised MSE used to cross-validate training will not be used for evaluation. Due to the noise level, a good performance on this metric will not imply a good reconstruction. Indeed an MSE of 0 implies that the model perfectly reconstructs the data and the noise. **Intersection over union** uses predicted and true labels:  $IoU(P,G) = \frac{P \cap G}{P \cup G}$ . This metric is very sensitive to small changes, especially when the ground truth is small, like in our situation. The IoU will be measured after applying the GMM to  $\Delta z$ . In particular, a low score could mean that the GMM is unfit for converting the differences into labels or that the surface reconstruction failed.

Average AUC will be computed to highlight goodperforming methods irrespective of the GMM results. We compute the standard AUC over three settings: addition vs no addition, deletion vs no deletion and change vs no change. In the last setting, we use  $|\Delta z|$ .

#### 4.3. Training procedure for surface reconstruction

In detail, we will describe how we train our network  $f_{\theta}$ given a PC. When a single network is used,  $pc_0$  and  $pc_1$ are concatenated, and the input dimension is three. With no loss of generality, we will consider that we have a single PC pc of dimension two or three. We normalise the PC to be in [-1; 1] on each axis which is a requirement for the methods RFF and SIREN [46]. We randomly split the dataset in two where 80% is retained for training and the other 20% is used for validating the surface reconstruction. We minimise and backpropagate through the training loss and evaluate the MSE performance on the validation. We use Optuna [1] to find the best set of hyperparameters via bayesian optimisation that minimises the validation MSE. The tuned hyper-parameters are the model architecture, the learning rate, the batch size, the scale of the gaussian mapping, the scalars associated with the regularisation terms  $\lambda_{TD}$ , and  $\lambda_{TV}$ . When we use the SIREN model, we also optimise the scale size by multiplying the signal in the sinusoidal activation function, the number of layers, and the number of hidden units for each layer. We use the optimisation method Adam [24] wrapped with the Layer-wise Adaptive Rate Scaling (LARS) [34] that enables the use of an enormous batch size that for us is essential to carry out our experiments in a reasonable time. To speed up computation, we set the number of epochs to 50, use learning rate decay and early stopping. To compute the TV norm, we sample random elements from pc that we corrupt with noise and backpropagate their prediction to the input [40].

### 5. Results and discussion

#### 5.1. Feature mapping

In Fig. 3, we show the IoU (in %) results between the different mapping methods: no feature mapping, RFF and SIREN. In Fig. A.3 of the Appendix, we show the results for the AUC metric. In Fig. 4, we show some resulting crops trained with no feature mapping, RFF and SIREN. We show additional crops as well as the whole map for data (3) in Fig.A.4, A.5, A.6 and A.7 of the Appendix. For both



Figure 3. *IoU* results (in %) for different feature mapping methods for every LiDAR airborne simulated dataset.

metrics, IoU and AUC, RFF outperforms SIREN and the default configuration by a fair margin. In terms of average performance (given in %) with standard deviation, RFF reaches an IoU of  $53.0 \pm 12$ . and an AUC of  $97.6 \pm 0.7$  and, SIREN  $39.0 \pm 10$ . and  $95.9 \pm 2.8$ , and not using any mapping reaches  $30. \pm 15$ . and  $95.9 \pm 3.8$ . A deconvolution of the reported values is given in Table A.1 in the Appendix. From the visualisation in Fig. 4, no feature mapping leads to a reconstruction where the building delimitation is unclear and fuzzy. SIREN gives the sharpest reconstruction with close to no noise between the buildings. Conversely, SIREN's projection onto the support of the second timestamp, given in the final row, is subject to many false positives along building boundaries. The RFF method produces distinctive buildings, like SIREN, but with a noisier output. However, the number of false positives in the final row is smaller. This ablation study shows the necessity of feature mapping to achieve good IoU and, therefore, a good reconstruction. Capturing high-frequencies is essential to our current problem due to the verticality of buildings.

#### 5.2. Hyper-parameter influence

In Fig. 5, we show a study on the regularisation parameter  $\lambda_{TD}$  and  $\lambda_{TVN}$  for both RFF, in darker and SIREN, in lighter colours. We measure the *IoU* and RMSE for both and compare them to the setting without penalty. Naturally, a too-strong penalty damages the performance, and a toolow value will render the penalisation negligible. Only for the method using RFF and TD regularisation do we see a 6% improvement in *IoU* compared to the baseline. This optimal  $\lambda_{TD}$  does not necessarily correspond to a minimal RMSE, metric used for the validation scheme. Similarly to Section 5.1, RFF features obtain better performance on both metrics and a more stable reconstruction noticeable by a lower RMSE and smaller confidence intervals.

The literature reports better reconstruction for SIREN

over the RFF methods. However, in our situation, SIREN gives lower performances in terms of IoU and MSE. SIREN suffers more from the mathematical formulation given in Section 3, which is ill-posed because of many points on the sides of the buildings. As SIREN induces qualitatively a better reconstruction, i.e. sharper edges and less noise. Having samples from the PC sharing similar geographical coordinates but radically different altitudes (along the verticality of the building) leads SIREN to slightly misplace the boundaries of the buildings, leading to many false positives along their edges, and hence a lower IoU. SIREN's MSE is higher, as the MSE penalises false positives more strongly as they correspond to larger differences between the ground truth and the prediction. The sharper edges of SIREN, compared to RFF, are harmful in terms of IoU and MSE. In other words, SIREN's reconstruction is sharper and, due to noise, wrongly estimates the building size. In contrast, the RFF method has softer edges, i.e. less overfitting induced by capturing fewer high-frequencies. This leads to a better mean error along the building edges and, to a lower minimisation of the MSE and fewer false positives.

#### **5.3.** Comparison to state of the art

In Table 1, we benchmark our methods on Urb3DCD and compare them to previous state-of-the-art in unsupervised detection, M3C2 [25] and OT [17]. RFF outperforms the other methods by a large margin, about 13% and 8% in *IoU* over the previous state of the art. In addition, the previous state-of-the-art maximised the *IoU* with respect to a predefined threshold, whereas our method is completely unsupervised. For example, SIREN still improves over the previous state-of-the-art even if the results do not show this because the metric report for SIREN is unbiased.

The experiments favour the use of one single function for both timestamps. We have a 2% difference in *IoU* between using a single model (S) and two models (D) when no regularisation is applied. The best model uses regularisation.

The final performed comparison is with DC3DCD on Urb3DCD-v2, and the results are shown in Table 2. To compare fairly, we use the same weakly supervised setting as the authors to map unlabelled classes, which for us corresponds to our binned distance, to labelled classes [10]. In Table 2, we compare, DNN to DNN, and show that INR outperforms the only other unsupervised DNN for CD on the building classes. We explain the poor scores on the vegetation classes because the hyper-parameters and postprocessing were optimised on Urb3DCD, i.e., for building change. In particular, the TVN penalty (which penalises the detection of small objects) and the TD penalty (detection of small growth) should be relaxed for the specific task of vegetation changes. DC3DCD shows better performance than our method when coupled with specific manual input features [10]. However, it should be noted that



Figure 4. Visualisation of a crop where in each column we show a different method comprising a single DNN trained on dataset (3). In the two first rows, we reconstruct the surface along a regular grid for timestamp  $t_0$  and  $t_1$ . In the third row, we show the difference  $\Delta z$  on the support of  $\mathcal{X}_1$  with it's predicted labels from the GMM, we filter out points where  $|\Delta z| < 2m$ . Each column shows a different method. In the final row and in the first column we show the true cloud point overlaid with the ground truth. To compare fairely, the color map ranges from dark purple, 160m altitude, to yellow, 205m, for the first two rows and from -30m to 30m for the visualisation of  $\Delta z$ .



Figure 5. Influence of parameter  $\lambda$ , refering to either  $\lambda_{TD}$  or  $\lambda_{TVN}$ , with respect to the *IoU* or reconstruction metric. We show 95% confidence intervals for each estimator. Each points followed other hyper-parameter selection procedures (20 attemps to minimise the reconstruction) and was repeated 20 times. In red dashes, we show the model with  $\lambda$  set to 0.

our model comprises approx. 600K trainable parameters, whereas DC3DCD has more than 100M parameters [49].<sup>1</sup>

#### 5.4. Application to Cultural heritage

The two-fold purpose of using LiDAR PCs to identify looting activities is to validate model 'S+RFF+TVN+TD'

on real (non-simulated) bi-temporal pairs of LiDAR PCs and assess its capability to detect looting, which is a pressing global-scale problem. We processed one bi-temporal pair of LiDAR PCs acquired over the Phnom Kulen region (Cambodia), where temples and ancient dams of the Angkor era are largely obscured by thick and closed canopies. Still, LiDAR deals well with this environment thanks to its capability to penetrate landscapes covered by continuous vege-

<sup>&</sup>lt;sup>1</sup>This is only an approximation as the precise number requires the number of kernel points, which is unknown.

Table 1. Comparison to state-of-the-art on the IoU metric (in %), we only report the best configuration when no feature mapping and SIREN. We show the best performing model in each row. D denotes the model with two DNN given by equation (1) and S the model with a single DNN given by equation (2). The complete table can be found in Table A.1 of the Appendix.

Data	M3C2 [25]	OT [17]	None	SIREN	RFF (Proposed Method)					
			S+TVN	S	D	D+TVN	S	S+TVN	S+TD	S+TVN+TD
(1)	29.87	40.65	37.22	40.14	50.13	44.68	52.49	55.87	54.73	49.68
(2)	53.73	55.20	45.22	53.98	57.39	59.33	56.45	61.17	60.34	59.52
(3)	38.72	39.26	33.11	38.57	46.54	43.17	51.87	46.94	54.00	53.33
(4)	35.01	39.89	33.54	39.01	48.62	49.70	51.38	51.10	53.40	53.99
(5)	37.78	48.17	37.97	40.48	42.55	43.16	42.95	43.26	40.55	47.17
Avg	39.02	44.63	37.41	42.43	49.04	48.00	51.02	51.67	52.60	52.74

Table 2. Comparaison to state-of-the-art on the IoU metric per class on Urb3DCD-v2 (in %). In each column, we highlight the best achieving model.

Method	DC3DCD	SIREN+S+TVN+TD	RFF+S+TVN+TD
Unchanged	$90.90 \pm 0.70$	$84.83 \pm 6.32$	$87.47 \pm 3.17$
New building	$64.06 \pm 5.13$	$62.62 \pm 11.4$	$71.81 \pm 2.76$
Demolition	$54.35 \pm 3.84$	$47.92 \pm 10.51$	$57.63 \pm 5.09$
New veg.	$58.14 \pm 20.03$	$4.26 \pm 4.87$	$5.44 \pm 5.88$
Veg. Growth	$1.45 \pm 2.05$	$0.62 \pm 0.29$	$1.54 \pm 1.42$
Missing veg.	$0.94 \pm 0.78$	$3.89 \pm 4.83$	$8.82 \pm 7.53$
Mobile Object	$47.57 \pm 2.58$	$0.26 \pm 0.19$	$0.58 \pm 0.26$

tation. The first,  $pc_0$ , was acquired in 2012, and the second,  $pc_1$ , in 2015. Fig. 6 shows detected looting. Archaeologists drew the red bounding box to identify an area where looting occurred. The archaeologists verified the predicted changes through visual inspection and confirmed that all the looting pits inside the bounding box were correctly identified. The false positive in the top right part can be easily filtered due to the 20 meter diameter, which is too big to be considered as a looting pit.

## 6. Conclusion

The amount of Earth observation acquired with 3D Li-DAR data is rising exponentially, which opens up the possibility of monitoring human activity through CD algorithms. In particular, we focused on urban planning and looting activities identification. Thanks to advances in DNN, we can now estimate and reconstruct large areas with high precision that allows passing the reconstructed surfaces to downstream tasks. However, the amount of training data and their discrete modelling limits their application to real scenarios. To address these issues, we propose a novel unsupervised grid-agnostic scheme for CD based on surface reconstruction and clustering, which achieves 52.74% IoU (in %), surpassing previous state-of-the-art M3C2 and OT by 10% on average on the Urb3DCD dataset. Moreover, we demonstrated in this paper that RRF mapping outperforms SIREN for CD on PC data acquired from airborne LiDAR sensors



Figure 6. Detection of looting pits on a bi-temporal pair of LiDAR point clouds. The ground truth bounding box identified the geographical area where looting occurred. The red points represent detected looting pits.

and allows us to identify looting activity correctly.

## **Code availability**

The code with a colab notebook is made fully available at the following URL: NN-4-change-detection. The code runs efficiently with hyperparameter selection thanks to the Optuna [1] package. We use GPU computation with Py-Torch [37] and combine each experiment into a pipeline with Nextflow [14] for easy reproducibility.

### Aknowledgement

DDC was supported by JSPS KAKENHI 23K16913. MY was supported by MEXT KAKENHI 20H04243 and partly by MEXT KAKENHI no. 21H04874. MF: This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No101027956. MF is thankful to Damian Evans for providing us with the LiDAR PC pairs for looting identification. We thank the RAIDEN computing system and its support team at the RIKEN AIP, which we heavily used.

## References

- [1] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of* the 25th ACM SIGKDD international conference on knowledge discovery & data mining, 2019. 5, 8
- [2] Nuri Benbarka, Timon Höfer, Andreas Zell, et al. Seeing implicit neural representations as fourier series. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pages 2041–2050, 2022. 4
- [3] Kevin Bleakley and Jean-Philippe Vert. The group fused lasso for multiple change-point detection. *arXiv preprint arXiv:1106.4199*, 2011. 2, 4
- [4] Marcello A Canuto, Francisco Estrada-Belli, Thomas G Garrison, Stephen D Houston, Mary Jane Acuña, Milan Kováč, Damien Marken, Philippe Nondédéo, Luke Auld-Thomas, Cyril Castanet, et al. Ancient lowland maya complexity as revealed by airborne laser scanning of northern guatemala. *Science*, 361(6409):eaau0137, 2018. 1
- [5] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European conference on computer vision (ECCV)*, pages 132–149, 2018. 3
- [6] Evian Pui Yan Chan, Tung Fung, and Frankie Kwan Kit Wong. Estimating above-ground biomass of subtropical forest using airborne lidar in hong kong. *Scientific Reports*, 11(1):1751, 2021. 1
- [7] Julian Chibane, Gerard Pons-Moll, et al. Neural unsigned distance fields for implicit function learning. Advances in Neural Information Processing Systems, 33:21638–21652, 2020. 3
- [8] Nicolas Courty, Rémi Flamary, Devis Tuia, and Thomas Corpetti. Optimal transport for data fusion in remote sensing. In 2016 IEEE international geoscience and remote sensing symposium (IGARSS), pages 3571–3574. IEEE, 2016. 2
- [9] Iris de Gélis, Sébastien Lefèvre, and Thomas Corpetti. Change detection in urban point clouds: An experimental comparison with simulated 3d datasets. *Remote Sensing*, 13(13):2629, 2021. 1, 2, 3, 5
- [10] Iris de Gélis, Sébastien Lefèvre, and Thomas Corpetti. Dc3dcd: unsupervised learning for multiclass 3d point cloud change detection. *arXiv preprint arXiv:2305.05421*, 2023. 1, 2, 3, 5, 6
- [11] Iris De Gélis, Sébastien Lefèvre, Thomas Corpetti, Thomas Ristorcelli, Chloé Thénoz, and Pierre Lassalle. Benchmarking change detection in urban 3d point clouds. In 2021 IEEE International Geoscience and Remote Sensing Symposium IGARSS, pages 3352–3355. IEEE, 2021. 1
- [12] Iris de Gélis, Sudipan Saha, Muhammad Shahzad, Thomas Corpetti, Sébastien Lefèvre, and Xiao Xiang Zhu. Deep unsupervised learning for 3d als point clouds change detection. *arXiv preprint arXiv:2305.03529*, 2023. 2, 3
- [13] Diego Di Carlo, Dominique Heitz, and Thomas Corpetti. Post processing sparse and instantaneous 2d velocity fields using physics-informed neural networks. In 20th International Symposium On Application Of Laser And Imaging Techniques To Fluid Mechanics, 2022. 3

- [14] Paolo Di Tommaso, Maria Chatzou, Evan W Floden, Pablo Prieto Barja, Emilio Palumbo, and Cedric Notredame. Nextflow enables reproducible computational workflows. *Nature biotechnology*, 35(4):316–319, 2017. 8
- [15] Mustafa Erdogan and Altan Yilmaz. Detection of building damage caused by van earthquake using image and digital surface model (dsm) difference. *International Journal of Remote Sensing*, 40(10):3772–3786, 2019. 2
- [16] Rizal Fathony, Anit Kumar Sahu, Devin Willmott, and J Zico Kolter. Multiplicative filter networks. In *International Conference on Learning Representations*, 2021. 3
- [17] Marco Fiorucci, Peter Naylor, and Makoto Yamada. Optimal Transport for Change Detection on LiDAR Point Clouds. *arXiv e-prints*, Feb. 2023. 2, 3, 6, 8
- [18] Marco Fiorucci, Wouter Verschoof-van der Vaart, Paolo Soleni, Bertrand Saux, and Arianna Traviglia. Deep learning for archaeological object detection on lidar: New evaluation measures and insights. *Remote Sensing*, 14:1694, 03 2022. 1
- [19] Daniel Girardeau-Montaut, Michel Roux, Raphaël Marc, and Guillaume Thibault. Change detection on points cloud data acquired with a ground laser scanner. *International Archives of Photogrammetry, Remote Sensing and Spatial Information Sciences*, 36(part 3):W19, 2005. 2
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016. 4
- [21] Zhangjin Huang, Yuxin Wen, Zihao Wang, Jinjuan Ren, and Kui Jia. Surface reconstruction from point clouds: A survey and a benchmark. arXiv preprint arXiv:2205.02413, 2022. 3
- [22] George Em Karniadakis, Ioannis G Kevrekidis, Lu Lu, Paris Perdikaris, Sifan Wang, and Liu Yang. Physics-informed machine learning. *Nature Reviews Physics*, 3(6):422–440, 2021. 3
- [23] Abderrazzaq Kharroubi, Florent Poux, Zouhair Ballouch, Rafika Hajji, and Roland Billen. Three dimensional change detection using point clouds: A review. *Geomatics*, 2(4):457–485, 2022. 1, 3
- [24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014. 5
- [25] Dimitri Lague, Nicolas Brodu, and Jérôme Leroux. Accurate 3d comparison of complex topography with terrestrial laser scanner: Application to the rangitikei canyon (nz). *ISPRS journal of photogrammetry and remote sensing*, 82:10–26, 2013. 2, 6, 8
- [26] You Li and Javier Ibanez-Guzman. Lidar for autonomous driving: The principles, challenges, and trends for automotive lidar and perception systems. *IEEE Signal Processing Magazine*, 37(4):50–61, 2020. 1
- [27] David B Lindell, Dave Van Veen, Jeong Joon Park, and Gordon Wetzstein. Bacon: Band-limited coordinate networks for multiscale scene representation. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 16252–16262, 2022. 3
- [28] Xuzhe Lyu, Ming Hao, and Wenzhong Shi. Building change detection using a shape context similarity model for li-

dar data. ISPRS International Journal of Geo-Information, 9(11), 2020. 2

- [29] Geoffrey J McLachlan and Kaye E Basford. *Mixture models: Inference and applications to clustering*, volume 38. M. Dekker New York, 1988. 5
- [30] Peter Naylor, Marick Laé, Fabien Reyal, and Thomas Walter. Segmentation of nuclei in histopathology images by deep regression of the distance map. *IEEE transactions on medical imaging*, 38(2):448–459, 2018. 5
- [31] Julian Nubert, Etienne Walther, Shehryar Khattak, and Marco Hutter. Learning-based localizability estimation for robust lidar localization. In 2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 17–24. IEEE, 2022. 1
- [32] Unal Okyay, Jennifer Telling, Craig L. Glennie, and William E. Dietrich. Airborne lidar change detection: An overview of earth sciences applications. *Earth-Science Reviews*, 198:102929, 2019. 1
- [33] Nobuyuki Otsu. A threshold selection method from graylevel histograms. *IEEE transactions on systems, man, and cybernetics*, 9(1):62–66, 1979. 2, 5
- [34] Chunmyong Park, Heungsub Lee, Myungryong Jeong, Woonhyuk Baek, and Chiheon Kim. torchlars, A LARS implementation in PyTorch. https://github.com/ kakaobrain/torchlars, 2019. 5
- [35] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 165–174, 2019. 3
- [36] Yujin Park and Jean-Michel Guldmann. Creating 3d city models with building footprints and lidar point cloud classification: A machine learning approach. *Computers, envi*ronment and urban systems, 75:76–89, 2019. 1
- [37] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In Advances in Neural Information Processing Systems 32, pages 8024–8035. Curran Associates, Inc., 2019. 8
- [38] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference* on computer vision and pattern recognition, pages 652–660, 2017. 3
- [39] Ali Rahimi and Benjamin Recht. Random features for largescale kernel machines. Advances in neural information processing systems, 20, 2007. 3
- [40] Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019. 3, 4, 5

- [41] Sudipan Saha, Francesca Bovolo, and Lorenzo Bruzzone. Unsupervised deep change vector analysis for multiplechange detection in vhr images. *IEEE Transactions on Geoscience and Remote Sensing*, 57(6):3677–3693, 2019. 3
- [42] Brent Schwarz. Lidar: Mapping the world in 3D. Nature Photonics, 4(7):429–430, July 2010. 2
- [43] Gregory Sech, Paolo Soleni, Wouter B. Verschoof-van der Vaart, Žiga Kokalj, Arianna Traviglia, and Marco Fiorucci. Tranfer Learning of Semantic Segmentation Methods for Identifying Buried Archaeological Structures on LiDAR Data. arXiv e-prints, July 2023. 1
- [44] Wenzhong Shi, Min Zhang, Rui Zhang, Shanxiong Chen, and Zhao Zhan. Change detection based on artificial intelligence: State-of-the-art and challenges. *Remote Sensing*, 12(10), 2020. 3
- [45] Sara Shirowzhan, Samad ME Sepasgozar, Heng Li, John Trinder, and Pingbo Tang. Comparative analysis of machine learning and point-based algorithms for detecting 3d changes in buildings over time using bi-temporal lidar data. *Automation in Construction*, 105:102841, 2019. 2
- [46] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *Proc. NeurIPS*, 2020. 3, 4, 5
- [47] Uwe Stilla and Yusheng Xu. Change detection of urban objects using 3d point clouds: A review. *ISPRS Journal of Photogrammetry and Remote Sensing*, 197:228–255, 2023.
  3
- [48] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. Advances in Neural Information Processing Systems, 33:7537–7547, 2020. 3, 4
- [49] Hugues Thomas, Charles R Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6411–6420, 2019. 3, 7
- [50] Wouter B. Verschoof van der Vaart. Learning to look at Li-DAR: combining CNN-based object detection and GIS for archaeological prospection in remotely-sensed data. PhD thesis, Leiden University, 2022. 1
- [51] George Vosselman and Hans-Gerd Maas. *Airborne and terrestrial laser scanning*. CRC press, 2010. 1
- [52] Wen Xiao, Hui Cao, Miao Tang, Zhenchao Zhang, and Nengcheng Chen. 3d urban object change detection from aerial and terrestrial point clouds: A review. *International Journal of Applied Earth Observation and Geoinformation*, 118:103258, 2023. 3
- [53] Yiheng Xie, Towaki Takikawa, Shunsuke Saito, Or Litany, Shiqin Yan, Numair Khan, Federico Tombari, James Tompkin, Vincent Sitzmann, and Srinath Sridhar. Neural fields in visual computing and beyond. In *Computer Graphics Forum*, volume 41, pages 641–676. Wiley Online Library, 2022. 2, 3

- [54] Zhenchao Zhang, George Vosselman, Markus Gerke, Claudio Persello, Devis Tuia, and Michael Ying Yang. Detecting building changes between airborne laser scanning and photogrammetric data. *Remote Sensing*, 11(20), 2019. 1, 2, 3
- [55] Zhenchao Zhang, George Vosselman, Markus Gerke, Devis Tuia, and Michael Ying Yang. Change detection between multimodal remote sensing data using siamese cnn. *ArXiv*, abs/1807.09562, 2018. 3
- [56] Benjamin Štular, Žiga Kokalj, Krištof Oštir, and Laure Nuninger. Visualization of lidar-derived relief models for detection of archaeological features. *Journal of Archaeological Science*, 39(11):3354–3360, 2012. 3