

Unsupervised Flood Detection on SAR Time Series

Ritu Yadav, Andrea Nascetti, Hossein Azizpour, Yifang Ban, *Member, IEEE*,

Abstract—Human civilization has an increasingly powerful influence on the earth system. Affected by climate change and land-use change, natural disasters such as flooding have been increasing in recent years. Earth observations are an invaluable source for assessing and mitigating negative impacts. Detecting changes from Earth observation data is one way to monitor the possible impact. Effective and reliable Change Detection (CD) methods can help in identifying the risk of disaster events at an early stage. In this work, we propose a novel unsupervised CD method on time series Synthetic Aperture Radar (SAR) data. Our proposed method is a probabilistic model trained with unsupervised learning techniques, reconstruction, and contrastive learning. The change map is generated with the help of the distribution difference between pre-incident and post-incident data. Our proposed CD model is evaluated for flood detection task. We verified the efficacy of our model on 8 different flood sites, including three recent flood events from Copernicus Emergency Management Services and six from the Sen1Floods11 dataset. Our proposed model achieved an average of 64.53% Intersection Over Union (IoU) value and 75.43% F1 score. Our achieved IoU score is approximately 6-27% and F1 score is approximately 7-22% better than the compared unsupervised and supervised existing CD methods. Based on our CD method, we also proposed an automatic change point detection framework where time series data is processed through the model to identify percentage change and the date on which significant change started to reflect on SAR data. This can help in early detection of floods giving more time for response. Our proposed model and change point detection framework are lightweight and easy to deploy. We conducted a range of experiments and ablation on our model. The results and extensive discussion presented in the study show the effectiveness of the proposed unsupervised CD method.

Index Terms—SAR, Change Detection, Time Series, VAE, LSTM, Contrastive Learning, Flood Detection.

I. INTRODUCTION

ACCORDING to a report from the Centre for Research on the Epidemiology of Disasters (CRED) [1], In 2021, a total of 432 catastrophic events were recorded, which is considerably higher than the average of 357 annual catastrophic events for 2001-2020. Floods dominated these events, with 223 occurrences, up from an average of 163 annual flood occurrences recorded across the 2001-2020 period. Countries such as India, China, Afghanistan, and Germany faced the loss of thousands of lives and billions of dollars [1]. Current flood predictions and evacuation services are gradually improving and are not fully reliable to handle the situation before a flood. In most flooding events, rescue services are launched

This research is part of the project ‘EO-AI4GlobalChange’ funded by Digital Future.

Ritu Yadav, Andrea Nascetti and Yifang Ban is with Division of Geoinformatics, KTH Royal Institute of Technology, Sweden. (e-mail: rituy@kth.se, nascetti@kth.se, yifang@kth.se)

Hossein Azizpour is with Department of Robotics, Perception and Learning (RPL), KTH Royal Institute of Technology, Sweden. (e-mail: azizpour@kth.se)

afterward. In such a scenario, accurate and reliable flood maps reflecting the damaged areas can help in efficient emergency response. The maps can be used for rescue missions, re-routing traffic, delivering aid, and many more. In the case of large-scale floods, on-ground evaluation for the identification of affected areas can be risky due to unfavorable weather conditions and collapsed transportation systems. Whereas, satellites can help in quick access to ground information over a large geographical area. The data can be used in detecting and mapping flooded areas and their severity. Satellites are a leading technology in gathering quick information on a large scale. A rapid increase in remote sensing technology leads to an immense amount of earth observation sensors providing data at different spectral, spatial, and temporal resolutions. Compared to optical data, Synthetic Aperture Radar (SAR) imagery is preferred for flood mapping [2]. Unlike optical sensors, SAR has the capability of imaging day and night, irrespective of the weather conditions.

Water surface can be detected using SAR because the water surface is smooth and SAR backscatter from a smooth surface is very low [3]. As a result, the water surface appears in a darker tone whereas the land surface with rough soil texture, building, vegetation, and others appears in bright tones. During floods, land surfaces are partially covered with water causing a significant change in backscatter. Therefore, floods can be detected with a CD framework that is proficient in detecting these backscatter changes.

II. RELATED WORK

In recent years, deep learning in Earth observation has received significant attention. Before deep learning there were classical unsupervised CD methods such as ImageDiff which is simply the difference between bi-temporal images, ImageRatio [4] uses the ratio of two bands, ImageRegr [5], CVA [6] which is a conceptual extension of image differencing. Several super-pixels and spatial neighborhood-based variants of CVA have been proposed, such as parcel change vector analysis (PCVA) [7] and robust change vector analysis (RCVA) [8]. DPCA [9], and PCDA [10] are examples of principal component analysis used for land cover CD.

Although there exist several classical methods to detect changes in multi-temporal images, deep learning gained new achievements due to its powerful discriminative ability [11]. FC-EF [12] and FC-Siam-diff [13], [12] are fully connected siamese network variations developed for CD. DS-IFN [14] proposed a fusion network for bitemporal CD on high-resolution optical images. There are multiple works on attentive siamese networks, DASNet [15] proposed a change map using L2 distance between the attentive feature maps from a dual siamese encoder, ADS-Net [16] proposed a

multiscale siamese encoder followed by an attentive decoder and DAUSAR [17] proposed a dual stream attentive U-Net architecture to detect changes. DASNet and ADS-Net operate on optical imagery whereas DAUSAR detects changes in SAR imagery. There are few works of CD with a generative network such as BIT [18] proposed a tokenized transformer network embedded with a deep difference-based CD framework. BIT detects changes in bitemporal high-resolution optical imagery. Deep learning CD methods in remote sensing are predominantly supervised. The above-mentioned methods are some of well-known supervised deep learning methods for CD.

Deep neural networks harness their great feature learning power from a large amount of labeled data. Unfortunately, in earth observation labeling data is a time taking task and requires domain expertise. The challenge is further elevated by the low-resolution data causing difficulty in feature discrimination for data labeling. Due to these challenges, there are not many large-scale datasets in earth observation. Apart from urban/building monitoring and land cover classification other earth observation task such as flooding, landslide, and wildfire suffers severely due to lack of labeled data. Training supervised networks on small datasets raises questions about their generalizability to other sites. In fact, many studies in earth observation are being conducted on a single site [19], [20], [21]. On the other hand, we have a large amount of unlabeled earth observation data which is readily available for use. An unsupervised CD method can be trained on these unlabeled data and can give more generalized results in comparison to supervised methods trained on small labeled datasets [22], [23], [24]. More recently, unsupervised deep learning methods are proposed for CD on remote sensing data such as [25] where a denoising autoencoder (DAE) is proposed and [26] proposed a highly coupled convolutional network for detecting changes between SAR and optical images. These unsupervised methods are tested on scenes with limited spatial complexity and didn't explore time series data. One of the recent work [24] trained simple variational autoencoder on reconstruction task and used distance metric on latent parameters to get low resolution change maps. This network is trained on time series data and designed to detect changes between two Sentinel-2 multispectral images.

In past few years, unsupervised learning techniques like SimCLR [27], MoCo [28], BYOL [29] and DeepCluster [30] has shown tremendous success. SimCLR and MoCo proposed contrastive learning from 'positive pairs' (augmented version of the same image) and 'negative pairs' (augmented version of a different image). These methods need careful treatment for negative pairs by relying on large batches or memory banks. The need for negative pairs was eliminated by BYOL, relying on learning from positive pairs. DeepCluster is a clustering method that jointly learns the parameters of a neural network and cluster assignments of the resulting features. [31] used deepcluster method and implemented unsupervised clustering with CNN to learn clustering-friendly feature representations of SAR data. There are multiple CD works on heterogeneous remote sensing data such as [32] used both deep cluster and contrastive learning to train a multisensor siamese CD network. Another recent work is Code Aligned Autoencoder

(CAA) [33] where an encoder-decoder network learns features from cross modality with the help of contrastive learning. The network generate output which merge features from the two modalities and looks like somewhere in-between the two modalities. The change map is produced by calculating difference image between the pre change input and generated output followed by manual thresholding. One of the recent work RaVAEN [24] trained a simple variational autoencoder on reconstruction task and used distance metric on latent parameters to obtain low-resolution change maps. This network is trained on time series data and designed to detect changes between two Sentinel-2 multispectral images.

Inspired by unsupervised deep learning techniques, we targeted our challenging problem of CD in a fully self-supervised manner. In this study, we introduce a generative network for CD on Sentinel-1 SAR data. We named our CD method as **Contrastive ConvLSTM Variational Autoencoder (CLVAE)**. Our method utilizes unlabeled time series data to train our network without external supervision at any stage. The key contributions of our work are as follows.

- 1) We propose a novel self-supervised CD method that gains its ability predominantly from the strong latent representations learned by the probabilistic reconstruction architecture of the variational autoencoder. We acknowledge the ability of time series data in CD task and to accommodate the benefits we embrace our proposed network with convolutional long short-term memory.
- 2) We show how the learned latent parameters of a reconstruction network can be employed to generate change maps. See subsection IV-F and framework 4. We further empowered our reconstruction network with cross connections between encoder and decoder branches similar to U-net.
- 3) We present network trained in a fully self-supervised manner. Along with reconstruction loss, the network is trained using contrastive learning where layers can learn to reconstruct SAR input so well that they can differentiate between dissimilar patches. We followed the contrastive learning idea from MoCo and simplified it for our remote sensing CD task. See training pipeline 3 and subsection IV-E for explanation.
- 4) Our training network is lightweight with 577,239 total parameters. The inference network is even smaller as it only uses the encoder part of the trained network. Both training and inference networks are memory efficient, making it easier for testing and deployment.
- 5) We display adaptability of unsupervised learning on sparse spatiotemporal SAR satellite data that are substantially different from the natural images commonly used in computer vision.

Additionally, we propose a change point detection framework (see Figure 5). Change point detection aims to locate abrupt property changes in time series data [33]. We can use the framework for continuous change monitoring, event detection and temporal anomaly detection such as detecting the point when the change started [34]. A significant change is an indicator of a major activity and might require human

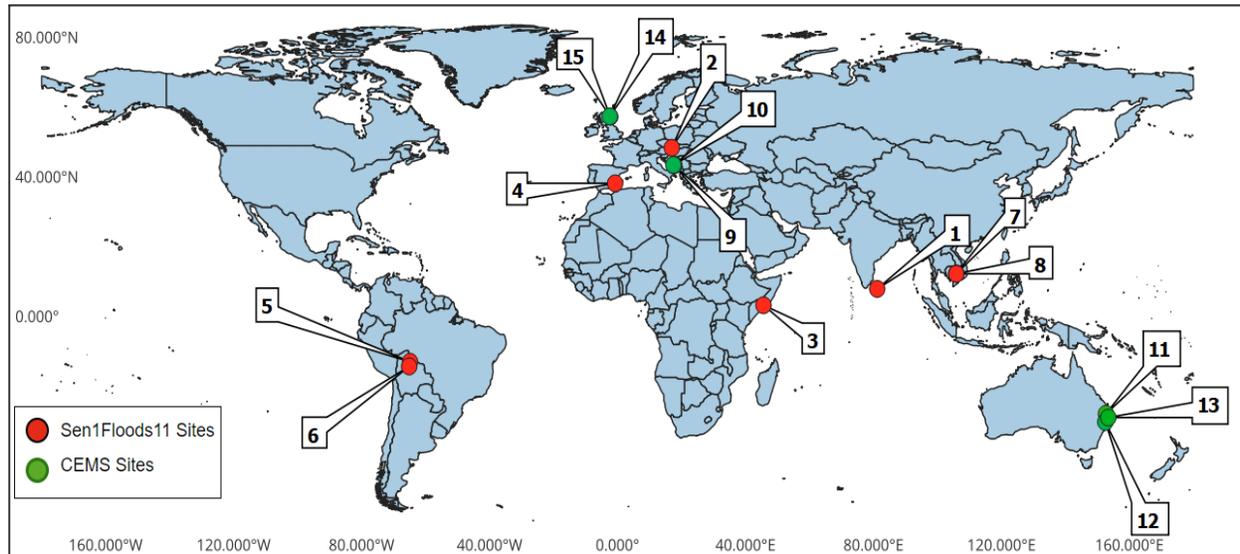


Fig. 1. Overview of the Study Sites. Colored dots represents tile locations and numbers are the references given to each tile. Red dots represents Sen1Floods11 sites and green dots represents CEMS sites.

attention.

III. DATA AND STUDY AREA

Our proposed unsupervised CD method was validated on Sentinel-1 SAR data. Two sources Sen1Floods11 [35] and Copernicus Emergency Management Service (CEMS)[36] were utilised to prepare the data. Collectively, our method was validated on data from 9 flood sites. The location and details of these flood events are presented in Figure 1 and Table I. The data collection process from both sources is explained in Subsection III-A followed by the data pre-processing steps in Subsection III-B

TABLE I

FLOOD EVENT METADATA. 'TILE REF.' IS THE REFERENCE NUMBER GIVEN TO EACH TILE, 'SITE' IS THE NAME OF THE FLOOD SITE, 'S1 POST DATE' IS THE DATE OF THE SENTINEL-1 POST-IMAGE, 'GT DATE' IS THE DATE OF THE SATELLITE IMAGE WHICH IS USED TO CREATE THE GROUND TRUTH, 'REL. ORBIT' IS THE RELATIVE ORBIT NUMBER OF THE SENTINEL-1 POST IMAGE AND 'ORBIT' IS ASCENDING(ASC) OR DESCENDING(DES) ORBIT INFORMATION OF THE SENTINEL-1 IMAGE.

Tile Ref.	Site	S1 Post Date	GT Date	Rel. Orbit	Orbit
Sen1Floods11					
1	Sri-Lanka	2017-05-30	2017-05-30	19	Des
2	Slovakia	2020-10-20	2020-10-20	73	Asc
3	Somalia	2020-05-07	2020-05-07	116	Asc
4	Spain	2019-09-17	2019-09-17	110	Des
5,6	Bolivia	2018-02-15	2018-02-15	156	Des
7,8	Mekong	2018-08-05	2018-08-05	26	Asc
CEMS					
9,10	Bosnia	2022-04-06	2022-04-03	51	Des
11,12,13	Australia	2022-03-31	2022-03-31	147	Des
14,15	Scotland	2022-11-18	2022-11-18	30	Asc

A. Data Collection

Sen1Floods11 Dataset consist of Sentinel-1 data from 11 different flood events covering a wide variety of geographical

area. In total, there are 446 non-overlapped Sentinel-1 tiles in the dataset and each tile is of 512x512 pixel size with a 20-meter ground resolution. Each data sample is composed of two bands VV (vertical transmit, vertical receive) and VH (vertical transmit, horizontal receive). The dataset also provides pixel-wise classification ground truth (flood segmentation maps). The dataset is hand labeled by experts by using information from Sentinel-1 and Sentinel-2 data followed by manual validation. Each pixel in the ground truth is classified into three categories, 0, 1, and -1. Class 0 represents the absence of water, class 1 represents water, and -1 indicates missing data. Since the ground truth was prepared using both Sentinel-1 and Sentinel-2, heavy clouds in the Sentinel-2 data affected the ground truth preparation. Whenever there is a cloudy pixel in the Sentinel-2, the corresponding pixel in ground truth is marked as missing data i.e., -1. Even though the dataset is big in terms of number of sites covered and number of tiles, a large part of the covered area has no(missing) corresponding ground truth. Therefore, this dataset is not sufficient and we still need a bigger dataset to efficiently train a deep network in a supervised setting. As of now, to the best of our knowledge Sen1Floods11 dataset is the biggest global dataset available on Sentinel-1, hence we decided to use some of good sites for evaluating our unsupervised CD method. To detect, evaluate and visualize flood on each pixel of the region, we choose reliable sites from the test data where there is no missing data in the ground truth. With this criteria, we ended up with tiles from six sites Bolivia, Spain, Cambodia, Slovakia, Somalia, and Sri Lanka.

In Sentinel-1 data samples in the Sen1Floods11 dataset were acquired after floods and are sufficient for a segmentation task. While the CD task requires samples from both pre and post-flood events. Therefore, we also collected pre-flood images and the data collection process is as follows. First,

we extracted geometry, relative orbit, and the passing orbit of post-flood images; Second, we downloaded Sentinel-1 images in two months window before the flood event date using geometry and orbit criteria. For each flood event, four pre-flood images (pre-images) were selected. Further data specifications for each site are provided in Table I. All Sentinel-1 data was collected from Google Earth Engine (GEE) [37]. Before downloading, data were pre-processed as described in subsection III-B.

Copernicus Emergency Management Service (CEMS) is one of the six worldwide services provided by the Copernicus program. It provides early warning, monitoring platforms, and mapping services for different natural and man-made disasters. CEMS helps countries with prevention, preparation, response, and recovery activities. We evaluated our proposed method on three recent flood events listed on the CEMS website under "List of EMS Rapid Mapping Activations"[38]. These floods occurred in current year (2022) in Mostar, Bosnia [39], Coraki, Australia [40] and Aberdeen, Scotland [41]. The CEMS provides the official flood maps and is publicly available on the CEMS website. According to the information given on the CEMS website, the flood maps for the mentioned flood events were derived from pre and post-event satellite images using a semi-automatic approach. The Bosnia flood map was generated using pre-image from Sentinel-2B and a postimage from the RADARSAT2 satellite. The Australia flood map was generated using pre-image from ESRI imagery and post-image from COSMO-SkyMed satellite. The Scotland flood map was generated using pre-image from ESRI imagery and post-image from Sentinel-1 satellite. The acquisition date of post image which is used in preparing the reference label is mentioned in Table I under column 'GT Date'.

In this study, the reference labels for the three flood event were collected from the CEMS website. These are the official flood maps and publicly accessible. We downloaded and pre-processed Sentinel-1 data for both pre and post-flood events. Multiple tiles were selected for each flood event covering urban, and surrounding agricultural areas. The tile size was kept as 512x512 pixels to maintain consistency with the Sen1Floods11 data. We selected four pre-images for each post-flood tile. Similar to Sen1Floods11 Dataset, all Sentinel-1 data are preprocessed, downloaded from GEE. For further details on collected Sentinel-1 data see Table I.

B. Data Preprocessing

All collected data were preprocessed and subsequently exported from the cloud-based platform GEE. It is becoming one of the most popular platforms for geospatial big data analysis. One of the biggest advantages of GEE is that Sentinel-1 SAR data is directly available as analysis-ready data cubes. Several studies have highlighted the potential of GEE platform to analyse large amount of geospatial data in a timely manner (e.g.[42], [43], [44], [45], [46]).

The Sentinel-1 mission collects C-band SAR images at 20 m resolution with dual polarization (HH+HV and VV+VH). Sentinel-1 images in GEE were preprocessed to Ground Range Detected (GRD) images using the Sentinel-1 Toolbox.

Preprocessing includes removal of thermal noise, radiometric calibration, and terrain correction. In addition, backscatter coefficients were converted to decibels using log scaling ($10 \log_{10} x$). We fetched dual-band VV+VH scene acquired in Interferometric Wide swath (IW) mode in a given period, orbit and location. While collecting scenes we also filtered them by Ascending and Descending passes due to the strong influence of incidence angle in the backscatter coefficient. We made sure that orbit pass and relative orbit of all pre and post-flood images are in agreement. Scenes were carefully selected, ensuring better data quality. Then we mask backscatter noise by clipping VV and VH channels in the range (-23, 0) dB and (-28, -5) dB respectively. Finally, both channels are normalized in the range [0, 1].

IV. METHODOLOGY

A. Autoencoder

In supervised settings, a neural network uses labels to learn features of input data. Labels guide the network to learn specific features depending on the target task, such as classification, segmentation, CD, and others. Input features can also be learned in an unsupervised manner and to do so, autoencoders are one of the widely used network categories. Autoencoders are pixel-wise reconstruction networks, which try to reconstruct their input x from a learned representation z . Unlike supervised CNNs, autoencoders generally learn input features for the reconstruction task and can be used for anomaly detection as a downstream task. The architecture of an autoencoder composed of an encoder $e(\cdot)$ [$z = e(x)$], which tries to capture input features x and encodes them into a smaller feature representation z ; a decoder $d(\cdot)$ [$\hat{x} = d(z)$], which decodes the representation z to reconstruct the input x . The network is trained with a loss derived by comparing the input x and the reconstructed output \hat{x} .

For our proposed CD network we used a probabilistic autoencoder known as variational autoencoder (VAE)[47]. It is a type of generative model that, unlike standard autoencoders, uses probabilistic encoding and decoding and learns to output a distribution over the latent representation $\hat{z} \sim P(z|x)$ and the reconstruction $\hat{x} \sim P(x|z)$. Both encoder and decoder networks in a VAE learn to output the parameters of the learnt distributions, $P(z|x)$ and $P(x|z)$ (e.g., mean and variance in case of a normality assumption). Thus, a trained VAE can be used for both latent representation of an observed input (using the probabilistic encoder) and for generating new unseen data (using the probabilistic decoder). In this work, we propose to detect the changes using the latent representation of a VAE, since it concisely summarizes the content of an input. In particular, we used a divergence measure between the latent distributions (or parameters thereof) of two corresponding image patches to determine whether there is a substantial change from one patch to the other.

B. Convolutional LSTM

Long Short-Term Memory (LSTM) is a commonly used method to learn temporal features of a time series data. LSTM however operates on $1 \times N$ dimensional vector and does not

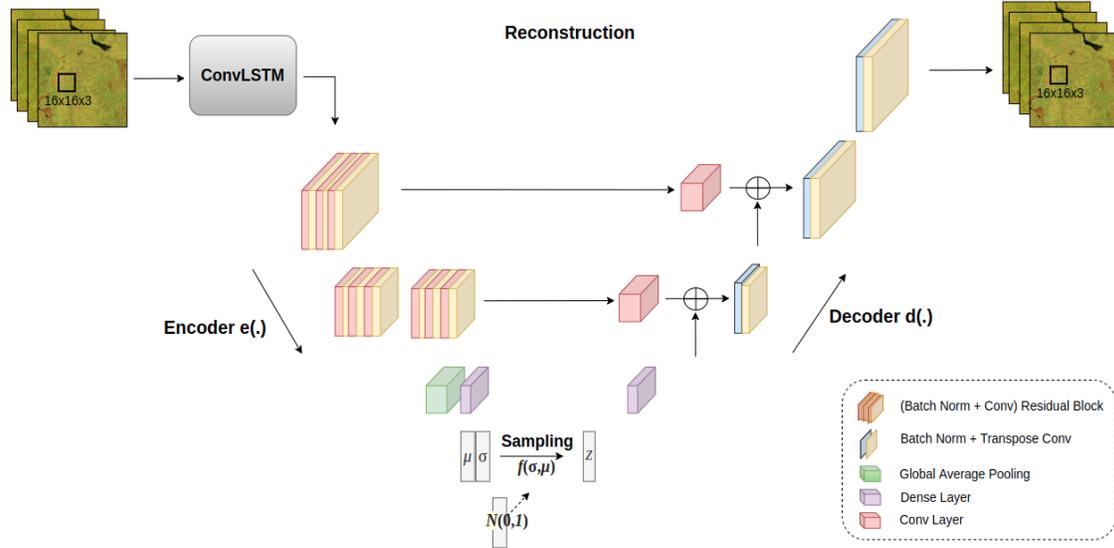


Fig. 2. Overview of the proposed Network Architecture for unsupervised CD. The network is trained on small 16x16x3 patches to learn the distribution of the small region at a time.

explore spatial feature learning. For CD tasks on time series data, both spatial and temporal characteristics are significant. Therefore we embedded our proposed network with Convolutional LSTM (ConvLSTM) [48], an LSTM which captures spatiotemporal correlation. ConvLSTM replaced all 1-d matrix multiplications by convolutional operation thus taking care of spatial neighboring features [49].

C. Proposed Network Architecture

Our training network is shown in Figure 2. The architecture is composed of an encoder, intermediate layers, and a decoder. The encoder consists of a convolutional LSTM layer, two residual blocks, a GlobalAveragePooling3D layer, and a dense layer. The convolutional LSTM layer takes a time series of input patches and extracts both temporal and spatial features. Each residual block has three sets of 3D convolutional layer and batch normalization layer. All convolutional layers used kernel of size 3 and stride 2. Non-linearity is added using the relu activation function. The residual block is followed by a GlobalAveragePooling3D layer, which calculates the spatial average value for each channel and reduces the dimensionality effectively. At last, a bottleneck dense layer of 8 channels is added. The output of the dense layer goes through two intermediate layers which are dense layers of a size equivalent to latent space. We fixed the size of the latent space to 128. These two dense layers output mean (μ_x) and log-variance (σ_x) values of the latent distribution corresponding to an input x . The output μ and σ of the dense layers are 1D vectors, which are then used to sample a latent vector $z \sim \mathcal{N}(\mu_x, \sigma_x)$ with the help of reparameterization trick [47], [50] for the forward pass to remain differentiable w.r.t. to μ and σ .

The decoder takes sampled z as input and passes through a dense layer. These features are then fed into three sets of transpose convolution and batch normalization layers. For transpose convolution Conv3DTranspose layer is implemented. First Conv3DTranspose layer is implemented with kernel size 3 and stride 2. Second Conv3DTranspose layer with 3 filters of kernel size 3x3, stride 3, and 'same' padding is employed to reconstruct the input stack of patches. In our implementation, we follow the common setup of outputting and using only the mean from the decoder (assuming a fixed unit variance). Furthermore, the decoder network's capacity is improved by employing cross-connections from the encoder network. Note that such skip layers deviate from the standard VAE by having the reconstruction conditioned not only on the latent representation but also on intermediate representations of the encoder. We found this change to be helpful.

The architecture of our proposed method is lightweight since we are using a limited number of time-series images to train the network. Another reason is that Sentinel-1 SAR data is low resolution in comparison to computer vision images and sparse spatial features can be easily learned by a shallow network.

D. Training Pipeline

An overview of the network training pipeline is shown in Figure 3. In the training pipeline, two proposed networks are placed in parallel forming two streams. The inputs to the two streams are time series patches $P1$ and $P2$. The input patches were selected randomly to ensure that they refer to different locations. Both streams were trained to reconstruct their corresponding inputs. Since the two inputs were from different locations, the two networks were also trained to increase the distance between the reconstructed outputs.

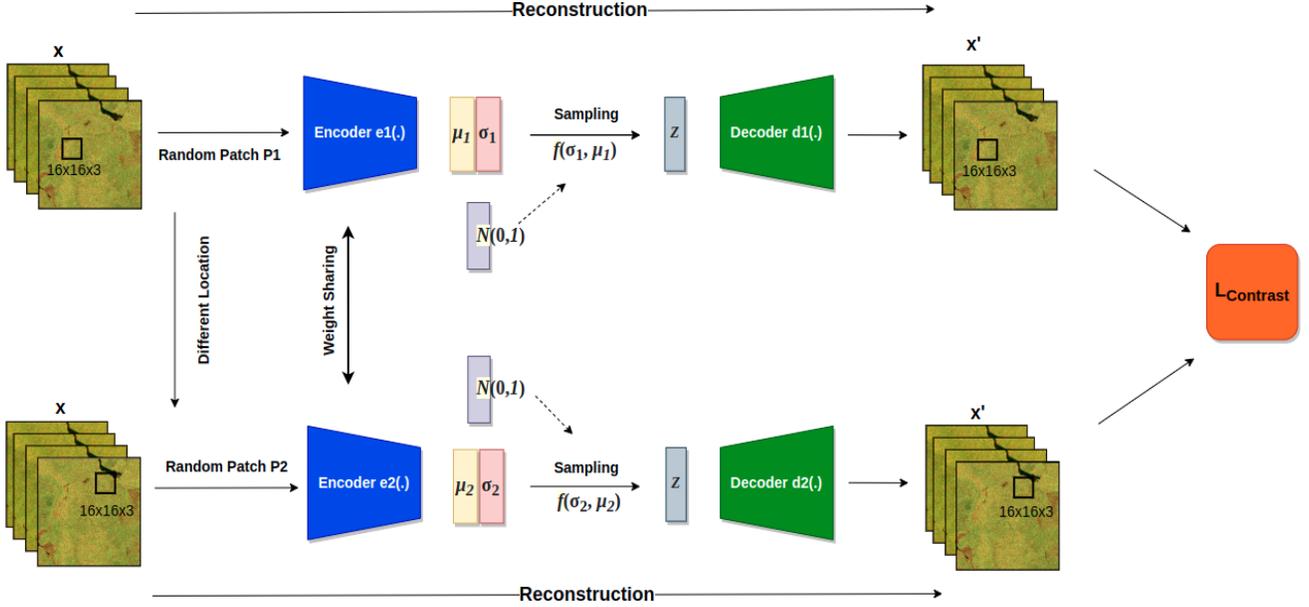


Fig. 3. The overview of our unsupervised training pipeline. The training pipeline is to train the model on reconstruction task. Figure Best viewed in color.

E. Training Objective Function

The network is optimized using three unsupervised loss functions: a reconstruction loss [47] (L_{Recon}), a Kullback-Leibler (KL) divergence loss [47] (L_{KL}), and a contrastive loss [51] ($L_{Contrast}$). Note that the first two loss functions constitute the standard VAE objective.

The reconstruction loss encourages the latent representation to contain adequate input information for an accurate reconstruction. The KL divergence loss pushes the latent distribution to be decorrelated and closer to a standard Gaussian.

The contrastive loss ensures diversity in the reconstructions of the two independent patches $P1$ and $P2$. The patches are from different locations capturing different areas and therefore should most frequently contain dissimilar features. Such a contrastive loss enables the network to learn latent representations that can differentiate features of separate patches. Moreover, it can learn uniformly distributed noise, resulting in a denoising architecture [52]. In this work we are proposing an architecture for change detection on SAR data. It is well known that SAR data contains peculiar speckle noise and removing such noise is one of the major challenge in remote sensing. Therefore, adding denoising ability to the architecture is of high importance. The combined objective of the training network is given in the below equation.

$$\begin{aligned}
 L_{Total} &= \alpha * [L_{KL}(\mu_{P1}, \sigma_{P1}, N(0, 1)) \\
 &+ L_{KL}(\mu_{P2}, \sigma_{P2}, N(0, 1))] \\
 &+ \beta * [L_{Recon}(P1, \hat{P}1) + L_{Recon}(P2, \hat{P}2)] \\
 &+ (1 - \alpha - \beta) * L_{Contrast}(\hat{P}1, \hat{P}2)
 \end{aligned} \quad (1)$$

where L_{KL} is the KL divergence loss, L_{Recon} is the reconstruction loss. Parameters α , β are the weight parameter for prioritizing losses.

F. Inference

Importantly, the proposed network is trained only on pre-event time series data. This means it has adapted its parameters to the distribution of features of pre-event data. Based on this, we assume that the distribution of latent variables for a patch $P1$ should undergo a significant change when affected by an extreme event. This assumption leads us to adopt a simple mechanism to detect change.

The proposed mechanism for CD is shown in Figure 4 and algorithm 1. The inference pipeline used two encoders $e1(.)$ and $e2(.)$ from the trained network. Since the encoders are trained on time series data of length four, we need to provide data of same length while taking the inference. The pre-event input data is prepared by stacking four sequential pre-images (time series), whereas post-event data is formed by stacking one post image four times. This is because Sentinel-1 provides one image every 6 days and flood extension change (increase or decrease) every day. So we use the latest possible image (stacked four times) to get best estimate of the flood extension rather than using four post flood image where the area might not be flooded anymore. Both pre and post-event data is now divided into small patches of size $16 \times 16 \times 3$ with stride 1.

From the trained network, the learned distribution can be retrieved as 1D vectors of mean μ and log-variance for latent representations. The variance σ is obtained by taking the exponent of log-variance. We pass pre-event patches through the encoder $e1(.)$ and obtain μ_1 , σ_1 . Similarly, we pass the post-event patches through encoder $e2(.)$ and retrieve μ_2 and σ_2 . We can use a number of different measures of divergence or difference between the two distributions of $\mathcal{N}(\mu_1, \sigma_1)$, and $\mathcal{N}(\mu_2, \sigma_2)$ to calculate the change. A few distribution difference measures are tested and compared (see experiment section VI-A). We opt for using a Cosine difference (CosD)

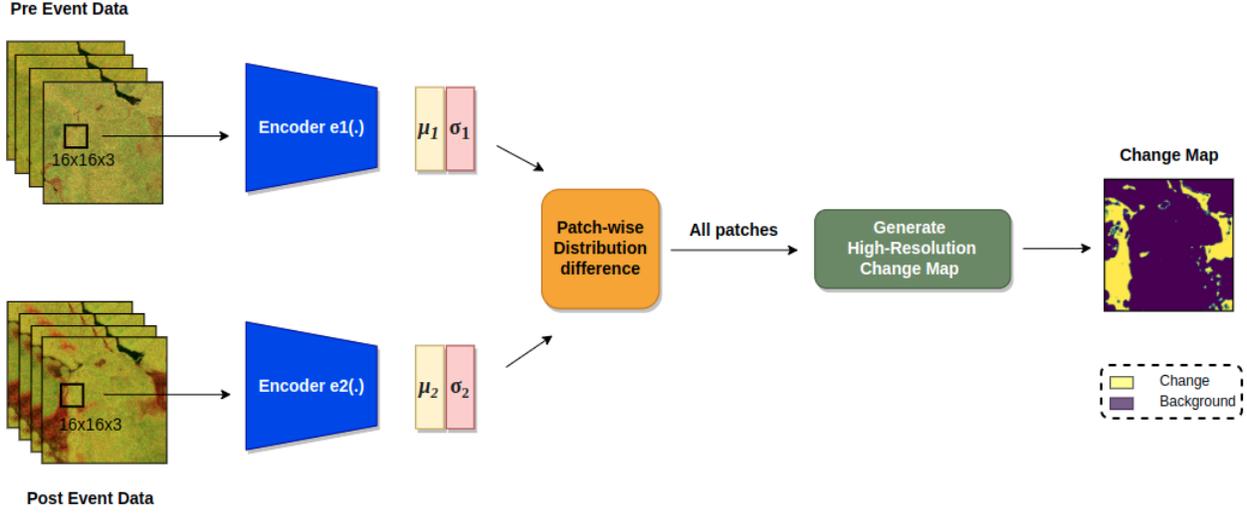


Fig. 4. The overview of our inference pipeline. The inference pipeline is for generating change maps between two time stamps. Figure Best viewed in color.

between the μ_1 and μ_2 , defined as follows

$$\text{CosD}(\mu_1, \mu_2) = -\frac{\mu_1}{\|\mu_1\|} \cdot \frac{\mu_2}{\|\mu_2\|} \quad (2)$$

Unlike training, input patches in the two encoders are from the same location. Input patches are generated with stride one and processed with a batch of size 512. The patch-wise distribution difference is calculated using CosD, resulting in a change map (see COSINE_DIFF_MAP in algorithm 1). Since, the output distribution difference is one value for each 16x16 size patch, the size of change map is smaller in comparison to width and height (W, H) of the input image. We tackled this problem by padding the input pre and post-event images. We used reflect padding or mirroring with length 8 which means reflect padding by 8 pixels on all four sides. This problem could also be resolved by zero padding, but zero padding introduced boundary errors to the detection results [53]. At last, a threshold of -0.9 is applied to get the binary change map (see BIN_CMAP in algorithm 1).

G. Change Point Detection

Change over an area can occur slowly over time such as slow flooding or it can be a sudden event. In case of slow changes, we should identify when the change starts becoming significant. The change can be verified on ground. If valid, the area can then be monitored with priority for further event prediction and warnings. With this motive, we developed a change point detection framework where we identify the point(i.e date) when the change is first started appearing on the available SAR data. Change point detection is performed over a long time series and the change is detected compared to a reference image x_{ref} . Any pre-image from the training data can be selected as a reference image and the corresponding acquisition date is the reference date $t_{1,ref}$. The reference image x_{ref} is the input to the first encoder $e1(.)$. For the time series data, we select a time window by providing the start date and length of the time window. The length of the

Algorithm 1: Binary Change Map Inference.

Input: Time series pre flood images of length 4(PRE_IMAGES), post flood image(POST_IMAGE),
 PATCH_SIZE=nx16x16x3, PAD_SIZE= 8,
 MODEL

Output: Binary change map BIN_CMAP.

- 1 Pad PRE_IMAGES and POST_IMAGE using reflect mode and PAD_SIZE.
 - 2 PRE_EVENT_DATA = stack all PRE_IMAGES .
 - 3 POST_EVENT_DATA = stack POST_IMAGE four times.
 - 4 PRE_PATCH = patches from PRE_EVENT_DATA of patchsize and stride 1.
 - 5 POST_PATCH = patches from POST_EVENT_DATA of patchsize and stride 1.
 - 6 **for** PRE_PATCH and POST_PATCH **do**
 - 7 PRE_MEAN, PRE_STD =
 MODEL.encoder1.predict(PRE_PATCH)
 POST_MEAN, POST_STD =
 MODEL.encoder2.predict(POST_PATCH)
 Calculate Cosine difference between PRE_MEAN and POST_MEAN.
 - 8 From patch-wise cosine difference get the change map COSINE_DIFF_MAP.
 - 9 Apply threshold of -0.9 to get binary change map BIN_CMAP= COSINE_DIFF_MAP>-0.9
-

time series can be adjusted as per the requirement, for the demonstration we selected the length as 4 and the time series is referred as $t_x \dots t_{x+4}$. All images from the calculated time window are fetched and processed one by one through encoder $e2(.)$. Change maps are generated for each image from the time window. This is done by following algorithm 1, where inputs are the x_{ref} and image from the time window (one at a time following the sequence). As our proposed network is

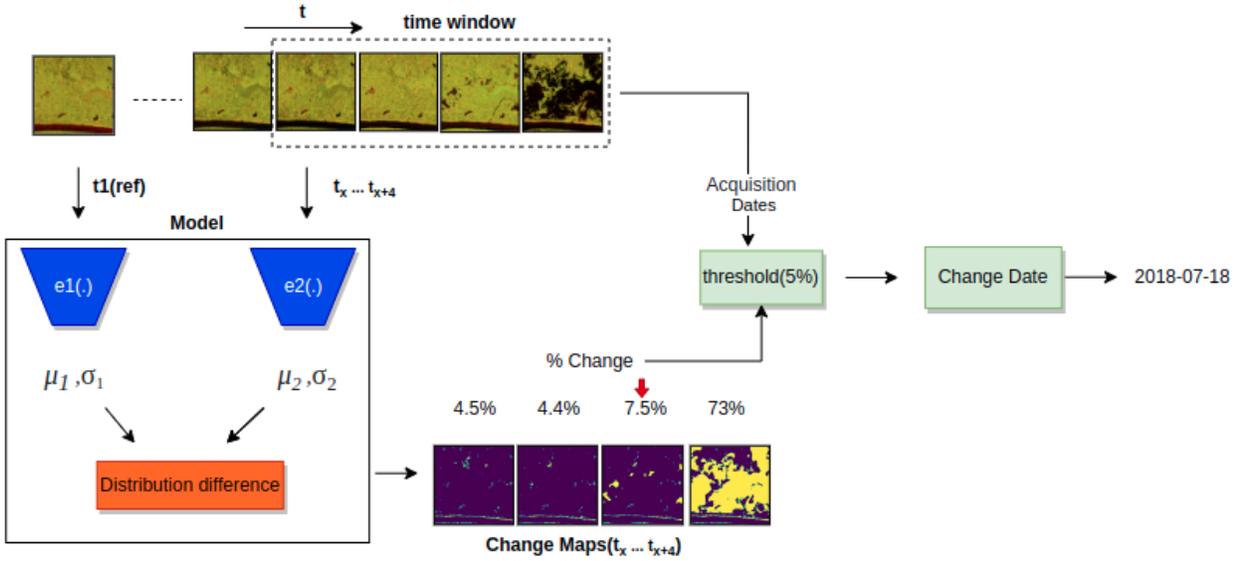


Fig. 5. Framework for the time-series Change Point Detection. It takes a selected pre-image as reference x_{ref} and generates change maps for all the images in a selected time window and then finds the change point in the time window.

currently set to 4 timestamps, both the inputs are stacked four times before feeding into encoder $e1(\cdot)$ and $e2(\cdot)$.

After getting the change maps, a pixel-based percentage change is calculated for each map. The change is considered significant if the percentage change is above a set threshold value. When the threshold is crossed for the first time, the framework fetches the acquisition date of the corresponding input image. The retrieved date is named as change point, i.e. the date when the change started appearing on the SAR data and probably the actual starting point of the change. The threshold can be different according to the sensitivity of the targeted change and can be adjusted by changing the parameter. In this study we set the threshold value as median of percentage changes. The overview of the proposed change point detection is depicted in Figure 5.

H. Implementation Details

One encoder of our proposed network takes four SAR time series images and each SAR here contains three channels. The first 2 channels are VV and VH, whereas the third channel is empty. We are using the three channel network because of two reasons. First, with three channels we can use imagenet pre-trained weights and second, the proposed network architecture can be reused for detecting changes on optical or RGB images.

Our model was trained on a pre-images to learn the distribution of the area and before feeding the images to the training model, input images were split into small patches of size $16 \times 16 \times 3$. We employed data augmentation technique to introduce more variations in the data which in turn increase the robustness of the model. This technique is widely used in classification, segmentation, change detection, and other tasks.

We utilized four types of augmentation for our dataset: gaussian blur, gammaContrast, flips and rotation. Enhancing the training dataset with these basic operations improves the performance of CNNs in remote sensing scene classification

compared to training on the original dataset [54]. Gaussian blur is a non-geometric augmentation that is applied to our input data with a kernel size of 3×3 . We use gammaContrast with range (0.25, 2.0) to adjust the image contrast. Both flips and rotation are geometric augmentation methods. Flips were applied left-right with a probability of 0.5 and up-down flips were applied with a probability of 0.2. Rotation was implemented randomly between -90 and 90 degrees. On each input sample, a random combination of Gaussian blur, flips, and rotation was applied before feeding them to the model.

The weight parameter of the objective function α is set to 0.1 which is the weight for KL divergence loss term to bring latent distribution closer to standard normal distribution. β is set to 0.7 and given higher weight to reconstruction loss term for accurate reconstruction learning. The remaining 0.2 weight is assigned to contrastive loss term. For better convergence, the model was trained with a decaying learning rate. The initial learning rate is 0.001 and decayed until it was at 0.00001. The decay steps were controlled with the "reduce on plateau" method, which decays the rate when the learning curve is stuck at a plateau. The learning rate decay after 2 steps of no learning (no change in loss) and the training terminates after four steps of no learning. The training network is shallow and lightweight. It contains total 576,395 trainable parameters, therefore the network is faster to train. The network was trained for 10 epochs. All the experiments were implemented in Keras and the training was conducted on one Google Colab GPU. The Code will be released as free and open source and publicly available on our GitHub account soon.

I. Evaluation metric

The output of the inference network is a pixel-level binary change map. So, the results are evaluated using pixel-level metrics. We used four accuracy metrics namely precision (P) and recall (R), F1 score and IoU. The formulas of the

metrics are given in Equation 3 – 6, where TP represents true positives, FP represents false positives, and FN represents false negatives.

$$\text{Precision (P)} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (3)$$

$$\text{Recall (R)} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (4)$$

$$\text{F1 score} = \frac{\text{TP}}{\text{TP} + \frac{1}{2}(\text{FP} + \text{FN})} \quad (5)$$

$$\text{IoU} = \frac{\text{TP}}{\text{TP} + \text{FP} + \text{FN}} \quad (6)$$

F1 score metric is the harmonic mean of precision(P) and recall(R) where P and R are given by 3 and 4 respectively. The IoU metric measures the intersection over union where TP is the intersection term and union term is (TP + FP + FN). Both F1 score and IoU range from 0 to 1.

V. RESULTS

A. Compared Methods

To demonstrate the benefits of our proposed unsupervised CD method (CLVAE), we compared it with log-ratio and Change Vector Analysis (CVA) [6] which are two well-known methods for CD on SAR images. We also show comparison with a recent VAE based unsupervised CD method RaVAEn [24]. It is important to note that our method is unsupervised and should be compared with only unsupervised methods. But a good part (6 out of 9 sites) of our tested sites comes from the Sen1Floods11 dataset test set. Therefore, we choose to compare our results with the results produced by the benchmark method on the Sen1Floods11 dataset i.e., the published work with highest IoU score on Sen1Floods11 dataset. To the best of our knowledge DAUSAR [17] network has provided the highest score on the Sen1Floods11 dataset, therefore in our work we will refer to this work as the benchmark method. We also compared with a recent deeply supervised network ADS-Net. We want to emphasize that the motive of the comparison is not only to find the best performing method, but also to study generalizability of our proposed novel unsupervised CD network. Below is an overview of the compared methods.

- 1) Log-ratio is commonly used to highlight changes in pairs of bi-temporal SAR images (e.g., [55]), and is formally defined as follows:

$$\text{LR} = 10 \log_{10} \left(\frac{S_{tn}}{S_{t1}} \right) = 10 \log_{10}(S_{tn}) - 10 \log_{10}(S_{t1}) \quad (7)$$

where S_{t1} and S_{tn} are the SAR images acquired over the same geographical area at the beginning and end of the time series, respectively. However, prior to computing the log-ratio, we used the Lee filter to suppress speckle noise from both images [56]. Binary change maps from log-ratio were generated using both otsu[57] and Yen's thresholding methods [58].

- 2) CVA is a popular CD method for multispectral optical images and SAR images. CVA generates change magnitude and change direction separately, which can be useful in determining change areas and change types. In this work, we focus on the magnitude of the change. Therefore, we calculated the magnitude change using CVA and then the changes are binarized using otsu thresholding.
- 3) RaVAEn is an unsupervised method recently proposed to detect changes in Sentinel-2 multispectral images instead of SAR. For comparison, we adapted RaVAEn on SAR data. This method uses a simple VAE network with residual encoder trained using default VAE losses. The method uses 32x32 size patches and change on each patch is calculated using cosine difference. Generated changes maps contain pixelated changes (see Figures 8 and 9) that result in low-resolution change maps.
- 4) DAUSAR is a supervised benchmark network [17] on Sen1Floods11 SAR dataset. It is a deep convolutional network for pixel-wise detailed CD between pre and post event images. Since the original Sen1Floods11 contains only post flood images, the authors extended the dataset by adding pre flood Sentinel-1 images. The pre-flood images were collected from GEE. The network is dual-stream Siamese U-Net enhanced with spatial and channel-wise attention.
- 5) ADS-Net is a supervised change detection network. It is a deep convolutional network which uses multi-scale features to extract changes between bitemporal remote sensing images. ADS-Net proved better change detection compared to existing deeply supervised networks including the famous FC-Siamese networks among others. This Network is proposed for 3 channel RGB optical images. We implemented the network for SAR images, where we used VV, VH, and VV/VH as three channels. We trained the network on pre-flood and post-flood images.

All mentioned methods are implemented from scratch. The two supervised network DAUSAR and ADS-Net are trained on the training set of Sen1Floods11 dataset and corresponding Sentinel-1 pre flood images from GEE. We use these trained network to take inference on our test sites listed in Table I. The inference results will further be used to see how well a supervised method trained on Sen1Floods11 dataset can generalize to new flood sites from CEMS website.

B. Quantitative Results

The quantitative comparison of the above-mentioned methods with our CLVAE is presented in Table II and III. The first table shows the comparison with all unsupervised methods and the second table with supervised methods. The comparison is displayed on each site in terms of four accuracy metrics; Recall, Precision, F1 score, and IoU. The log-ratio method with otsu thresholding is performing significantly low in comparison to the log-ratio with Yen's thresholding. Hence in Table II we presented the results from the best performing threshold method. However, the comparison of log-ratio with

TABLE II
 QUANTITATIVE COMPARISON WITH UNSUPERVISED METHODS: COMPARISON OF OUR CLVAE METHOD WITH UNSUPERVISED METHODS LOG-RATIO(YEN), CVA AND RAVAEN. THE COMPARISON IS PRESENTED IN TERMS OF PERCENTAGE RECALL(R), PRECISION(P), F1 SCORE(F1), IOU METRIC. THE PRESENTED METRIC VALUES WERE AVERAGED OVER 3 RUNS.

Site	Log-Ratio(Yen's)				CVA				RaVAEn				CLVAE(Ours)			
	R	P	F1	IoU	R	P	F1	IoU	R	P	F1	IoU	R	P	F1	IoU
Sri-Lanka (1)	31.4	11.0	16.3	8.9	57.8	5.3	9.7	5.1	55.65	3.57	6.71	3.47	25.0	8.5	12.7	6.8
Slovakia (2)	60.1	72.9	65.9	49.1	77.0	84.4	80.5	67.4	76.53	50.31	60.71	43.59	77.9	93.8	85.1	74.1
Somalia (3)	65.3	71.9	68.4	52.0	64.2	70.4	67.2	50.6	71.36	55.67	62.55	45.51	82.9	72.1	77.1	62.7
Spain (4)	57.4	60.4	58.9	41.7	56.4	70.4	62.6	45.6	34.97	41.91	38.13	23.55	73.8	74.7	74.2	59.0
Bolivia (5)	57.7	96.0	72.1	56.4	75.5	95.7	84.4	73.0	94.10	78.46	85.57	74.78	92.8	91.9	92.3	85.8
Bolivia (6)	34.2	92.6	50.0	33.3	63.7	60.1	61.8	44.8	65.44	78.60	71.42	55.54	81.9	81.8	81.8	69.3
Mekong (7)	59.8	93.4	72.9	57.4	76.8	99.7	86.8	76.6	67.07	91.00	77.22	62.90	89.5	96.2	92.7	86.4
Mekong (8)	41.7	94.9	57.9	40.8	78.3	99.6	87.7	78.0	72.06	86.80	78.75	64.95	94.9	95.8	95.3	91.1
Bosnia (9)	67.7	23.6	35.0	21.3	60.6	11.0	18.6	10.2	33.00	26.53	29.41	17.24	50.0	51.3	50.6	33.9
Bosnia (10)	58.5	54.1	56.2	39.1	57.7	65.7	61.4	44.3	76.61	52.06	61.99	44.92	70.5	81.1	75.4	60.5
Australia (11)	53.2	69.4	60.2	43.1	75.6	90.8	82.5	70.2	86.85	76.98	81.62	68.94	93.8	91.2	92.5	86.0
Australia (12)	55.7	58.5	57.1	39.9	75.1	83.2	78.9	65.2	82.75	64.44	72.46	56.81	88.9	86.2	87.5	77.8
Australia (13)	47.1	71.8	56.9	39.7	79.7	93.7	86.1	75.6	75.06	83.47	79.04	65.35	96.4	92.2	94.3	89.2
Scotland (14)	67.47	18.37	28.88	16.87	63.9	13.25	21.95	12.33	28.38	50.57	36.36	22.22	71.68	50.95	59.56	42.41
Scotland (15)	71.65	24.19	36.17	22.08	60.43	9.68	16.69	9.1	18.32	56.68	27.69	16.07	65.34	55.8	60.19	43.05
Average	55.26	60.87	52.86	37.44	68.18	63.53	60.46	48.54	62.54	59.80	57.98	44.39	77.01	74.90	75.43	64.53

TABLE III
 QUANTITATIVE COMPARISON WITH SUPERVISED METHODS: COMPARISON OF OUR CLVAE METHOD WITH SUPERVISED CHANGE DETECTION METHODS ADS-NET AND DAUSAR. THE COMPARISON IS PRESENTED IN TERMS OF PERCENTAGE RECALL(R), PRECISION(P), F1 SCORE(F1), IOU METRIC. THE PRESENTED METRIC VALUES WERE AVERAGED OVER 3 RUNS.

Site	ADS-Net				DAUSAR				CLVAE(Ours)			
	R	P	F1	IoU	R	P	F1	IoU	R	P	F1	IoU
Sri-Lanka (1)	13.7	30.9	18.98	10.5	43.8	12.7	19.69	10.9	25.0	8.5	12.7	6.8
Slovakia (2)	94.3	80.2	86.68	76.5	97.7	69.7	81.4	68.6	77.9	93.8	85.1	74.1
Somalia (3)	64.4	84.6	73.13	57.6	99	58.6	73.6	58.3	82.9	72.1	77.1	62.7
Spain (4)	50.5	90.2	64.75	47.8	90.4	55.7	68.9	52.6	73.8	74.7	74.2	59.0
Bolivia (5)	73.7	84.8	78.86	65.1	92.8	78.5	85.1	73.9	92.8	91.9	92.3	85.8
Bolivia (6)	47.1	82.7	60.02	42.9	63.5	77.2	69.7	53.5	81.9	81.8	81.8	69.3
Mekong (7)	93.3	95.9	94.58	89.7	96.2	93.4	94.8	90	89.5	96.2	92.7	86.4
Mekong (8)	92.9	97.3	95.05	90.6	97.1	95.1	96.1	92.5	94.9	95.8	95.3	91.1
Bosnia (9)	23.2	75.8	35.53	21.6	98.2	15.8	27.2	15.7	50.0	51.3	50.6	33.9
Bosnia (10)	59.2	81.4	68.55	52.1	89.6	46.9	61.6	44.4	70.5	81.1	75.4	60.5
Australia (11)	99.7	76.3	86.4	76.1	96.8	84.5	90.23	82.2	93.8	91.2	92.5	86.0
Australia (12)	99	61.5	75.9	61.2	93.2	75.6	83.48	71.6	88.9	86.2	87.5	77.8
Australia (13)	99.8	79.6	88.6	79.5	97.9	87.4	92.35	85.8	96.4	92.2	94.3	89.2
Scotland (14)	38.7	45.4	41.78	26.4	81.9	35.1	49.14	32.6	71.68	50.95	59.56	42.41
Scotland (15)	15.9	68.2	25.79	14.8	61.8	31.7	41.9	26.5	65.34	55.8	60.19	43.05
Average	64.36	75.65	66.31	54.16	86.66	61.19	69.01	57.27	77.01	74.90	75.43	64.53

otsu thresholding is depicted later in Figures 8 and 9 for qualitative analysis.

On average, CLVAE achieved an F1 score of 75.43% and IoU score of 64.53% with 77.01% recall and 74.90% precision. Among all the compared methods ours achieved the best average precision, F1 score and IoU, whereas recall is best achieved by the DAUSAR. In comparison to the unsupervised CD methods log-ratio, CVA and RaVAEn, CLVAE outperformed in all four average metrics. The lead in recall ranges from 9-22%, in precision from 11-15%, in F1 score from 15-22% and in IoU from 16-27%. RaVAEn was originally proposed for CD on multi spectral data and seems to be not promising for SAR data. This method performed better than log-ratio but couldn't outperform other compared methods. The supervised method DAUSAR has a high recall percentage(86.66%) which is approximately 10% higher compared to CLVAE's recall. Also, ADS-Net 75.65% precision which

is 0.7% better than CLVAE's precision. But overall CLVAE maintains high precision and recall and outperformed the supervised methods by 6% in F1 and 7% in IoU score.

On individual sites, CLVAE yielded the best F1 score and IoU except for 'Slovakia', 'Mekong' and 'Sri-Lanka' sites. On 'Slovakia' site DAS-Net gave the best results and it's DAUSAR on 'Mekong' site. Compared to CLVAE scores, the difference is not much and ranges from 1-2% in F1 score and 1-4% in IoU metric. On 'Sri-Lanka' site, DAUSAR gave the best results but the scores are extremely low. The site is a rice field that was flooded right after harvesting months. Therefore the half-cut stems are good enough to hold flood water and change is reflected between pre, post images. However, these flooded fields were not considered flooded in the ground truth leading to disagreement between the detected change and the ground truth. This explains the low scores by CLVAE and all compared methods.

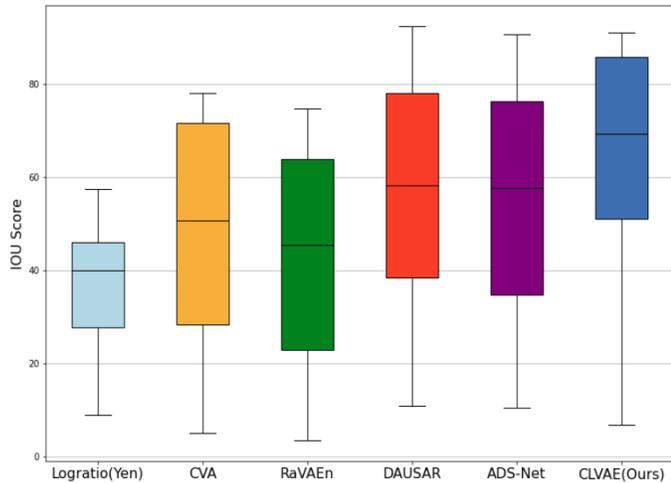


Fig. 6. Boxplot Graph: Graphical IoU comparison of CLVAE with log-ratio with Yen’s thresholding, CVA, RaVAEn, DAUSAR and ADS-Net change detection methods.

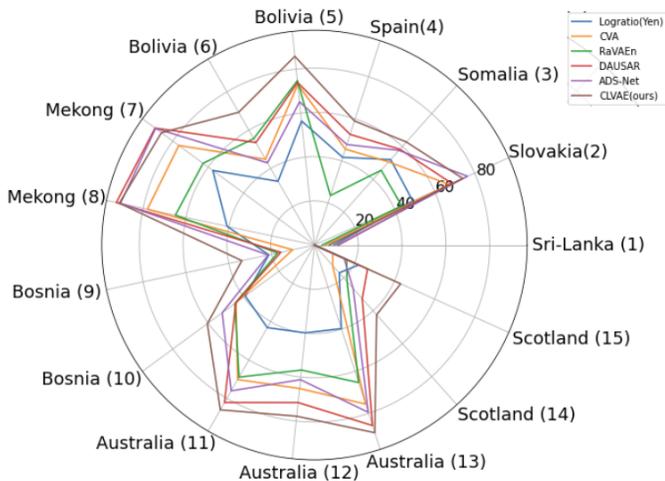


Fig. 7. Spider Graph: Graphical IoU comparison of our method with log-ratio with Yen’s thresholding, CVA, RaVAEn, DAUSAR and ADS-Net change detection methods.

On most of the sites recall is best achieved by the supervised method DAUSAR. Since high recall of DAUSAR is associated with a low precision, the IoU and F1-scores are also relatively low. Another minor deviation from the average metric results is in precision on ‘Bolivia’ and ‘Mekong’ sites. Unlike average results, the precision on ‘Bolivia’ sites is best by log-ratio and on ‘Mekong’ sites by CVA. But they also have low recall which leads them to significantly low IoU and F1 scores.

Further insights into the results are provided by Figure 6 and 7 where the IoU results are visualized in boxplot and spider graphs. In the boxplot, the x-axis represents the compared method and the y-axis represents the percentage IoU score. Notably, our CLVAE gave the highest IoU score median. In the spider graph, the axis represents the evaluation sites and the numbers on all the concentric circles represent the possible percentage IoU score from 0 at the center to 100 on the

outermost circle. The farther toward the end of the spike, the larger the value. Closest to the center means closer to zero. The outermost line represents the best performing model and in the current scenario it is our proposed method CLVAE.

C. Qualitative Results

For qualitative analysis, we selected three geographically different and challenging sites; Mekong, Bolivia, Slovakia and Bosnia. The visualization of change maps is presented in Figure 8 and 9, where (a) and (b) show the latest pre-flood image and post-flood image, (c) show the ground truth, (d) and (e) show the change map from supervised methods ADS-Net and DAUSAR, (f), (g), (h), (i) and (j) show the change map generated from log-ratio with otsu threshold, log-ratio with Yen’s threshold, CVA, RaVAEn and our proposed CLVAE respectively. The number below each change map is the corresponding percentage IoU score.

The First two rows of Figure 8 show the detection results on ‘Bosnia’ site. Log-ratio detected a good part of the flooded area correctly. But the detection is not smooth rather grainy causing false negatives. The speckle noise in surrounding areas is also adding false positives. CVA detected a good portion of the flooded area but the change map contains high speckle noise. RaVAEn detected changes in 32x32 patches showing ill-defined pixelated change. The change map doesn’t contain major false detection but also failed to detect most part of the changes shown in ground truth. Our CLVAE also couldn’t detect the flooded area with huge success but provided significantly good detection in comparison to others. The generated change map contains very low false detection and did not suffer from speckle noise. It is noteworthy that the supervised method ADS-Net is missing a good part of the flooded area and DAUSAR suffers from a large amount of false detection whereas our unsupervised method CLVAE produced relatively good detection results. In terms of IoU, CLVAE gave 11 to 23% better score than the compared methods.

The third and fourth rows of Figure 8 shows the detection results on ‘Bolivia’ site. The detection on this site is comparatively better than the ‘Bosnia’ site. CVA shows good detection results and doesn’t suffer from speckle noise. However, both log-ratio and CVA missed a significant part of the flooded area(change). The detection from RaVAEn is also better but lacks details due to the coarse resolution of the output change maps. Both supervised methods show good detection results but ADS-Net suffers from false negatives and DAUSAR suffers from false positives. Compared to the six methods our CLVAE method resulted in a clear and speckle-free change map. On ‘Bolivia’ site, the IoU score of CLVAE is 10 to 30% better than the compared methods.

The last two rows of Figure 8 show the detection results on ‘Mekong’ site. All the CD methods performed well on this site. The supervised methods ADS-Net and DAUSAR outperformed unsupervised CD methods. Among unsupervised results, log-ratio has problem of grainy detections, both CVA and RaVAEn missed some of the changed areas, RaVAEn suffers from pixelated detection and CLVAE missed small flooded streams. On an average all methods detected majority of the flooded areas.

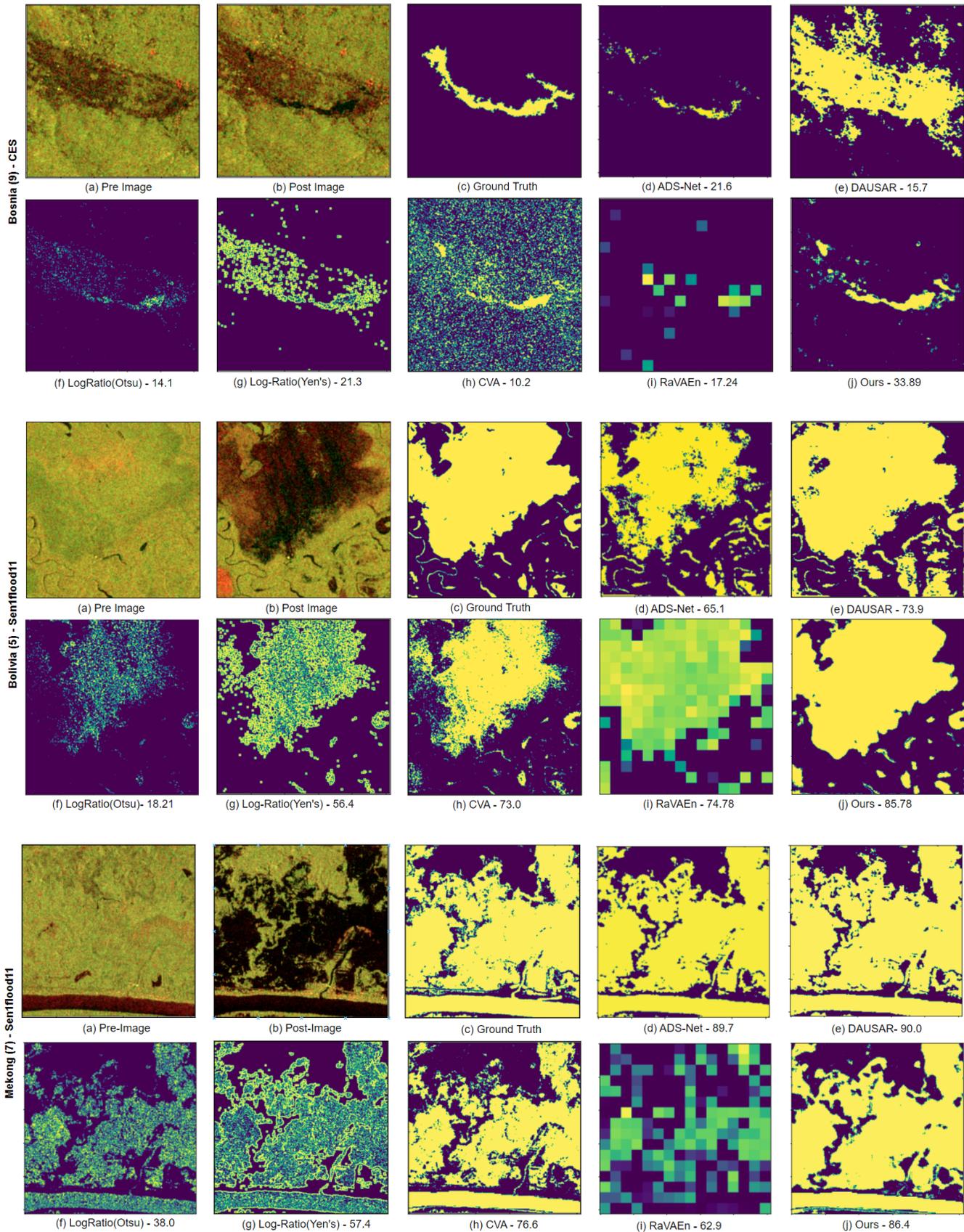


Fig. 8. Qualitative Comparison: (a), (b) and (c) represents latest pre flood image, post flood image and binary ground truth. The remaining images shows comparison of (d) ADS-Net, (e) DAUSAR, (f) Log-ratio with otsu threshold, (g) Log-ratio with yen's threshold, (h) CVA, and (i) RaVAEn with (j) our proposed CLVAE CD method. The number below each change map is corresponding percentage IoU score.

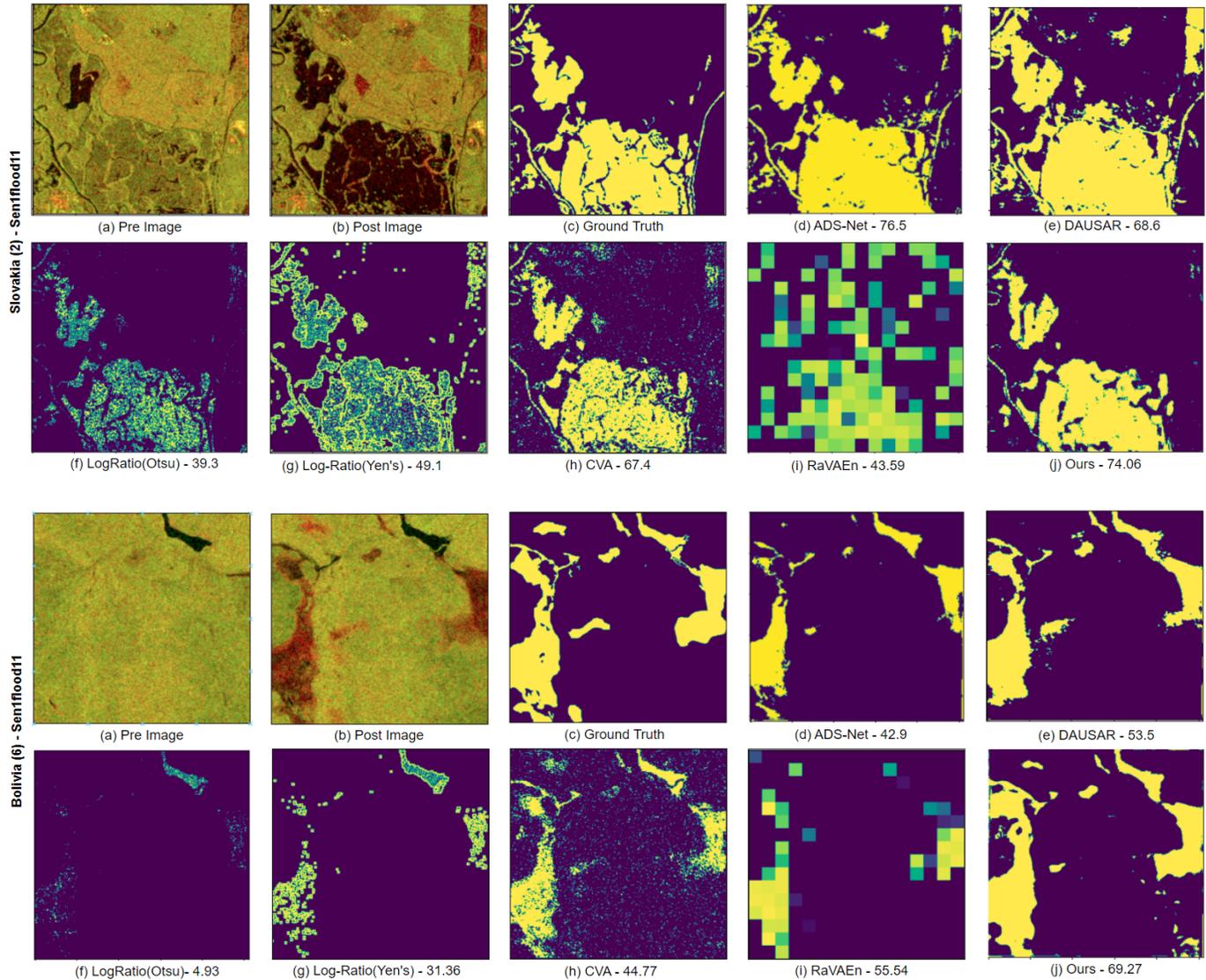


Fig. 9. Qualitative Comparison: (a), (b) and (c) represents latest pre flood image, post flood image and binary ground truth. The remaining images shows comparison of (d) ADS-Net, (e) DAUSAR, (f) Log-ratio with otsu threshold, (g) Log-ratio with yen's threshold, (h) CVA, and (i) RaVAEn with (j) our proposed CLVAE CD method. The number below each change map is corresponding percentage IoU score.

Two more samples from 'Slovakia' and 'Bolivia' sites are shown for evaluation in Figure 9. The first two rows shows the detection results on 'Slovakia' site, where ADS-Net gave better(2% better IoU score) detection results compared to CLVAE. Whereas, the second sample from 'Bolivia' site shows that best results are provided by our CLVAE.

VI. EXPERIMENTS

A. Performance With Respect to different Distribution Difference Methods

As discussed in Section IV-F, with our proposed CD architecture, different choices of difference between $\mathcal{N}_1 = \mathcal{N}(\mu_1, \sigma_1)$ and $\mathcal{N}_2 = \mathcal{N}(\mu_2, \sigma_2)$ are possible for change detection. Here, we tested four different methods for calculating distribution difference, namely Kullback-Leibler Divergence (KLD), Jensen-Shannon Divergence (JSD), Euclidean

Distance (ED) and Cosine Distance (CosD). The KL and JSD methods operate on full distribution (using both mean and variance), whereas we use ED and CosD only on the mean parameter of the distribution. The formulas for all four distance functions is given in Eq. 8 – 11, where P, Q are the distributions, μ represents mean, σ is variance and $\| \cdot \|$ represents L2 normalization function.

$$KLD(\mathcal{N}_1||\mathcal{N}_2) = \sum \log\left(\frac{\sigma_2}{\sigma_1}\right) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} \quad (8)$$

$$JSD(\mathcal{N}_1||\mathcal{N}_2) = \frac{KLD(\mathcal{N}_1||\mathcal{N}_m)}{2} + \frac{KLD(\mathcal{N}_2||\mathcal{N}_m)}{2} \quad (9)$$

$$\text{where, } \mathcal{N}_m = \mathcal{N}\left(\frac{\mu_1 + \mu_2}{2}, \frac{(\sigma_1 + \sigma_2)}{2}\right)$$

$$ED(\mathcal{N}_1, \mathcal{N}_2) = \sqrt{\sum (\mu_1 - \mu_2)^2} \quad (10)$$

$$CosD(\mathcal{N}_1, \mathcal{N}_2) = -\frac{\mu_1}{\|\mu_1\|} \cdot \frac{\mu_2}{\|\mu_2\|} \quad (11)$$

The average metric calculated by the mentioned four methods is shown in Table IV where CosD shows the best mean precision, F1 score and IoU value. The recall with CosD function is lower in comparison to other compared functions. It is important to point out that KLD, JSD, and ED resulted in similar values. This can be explained as follows, as explained in subsections IV-C and IV-E, in the training process, our network is guided to learn input data as standard normal distribution encouraging the variance to be 1. As a result, the majority of variance values from the trained encoder, at inference time, turn out to be 1 as well. This is true for both pre and post-flood images. If we equate σ_1 and σ_2 to 1 in KLD eq. 8, it comes down to $\frac{1}{2} \sum (\mu_1 - \mu_2)^2$. At this point, both ED and KL values are some positive fraction of $\sum (\mu_1 - \mu_2)^2$. In the calculation of distribution difference 1, the threshold for KLD, JSD, and ED is set to 0.0 indicating that only the existence of change is considered and not the magnitude of change. This condition equalizes the change maps from KLD and ED methods. The same explanation is valid for results with the JSD method as well. Therefore, we end up with ED and CosD as two different difference measures. On average, both distribution difference functions show similar results but these values are significantly different on individual sites.

TABLE IV
PERFORMANCE VARIATION WITH RESPECT TO DISTRIBUTION DIFFERENCE FUNCTIONS.

	Mean R	Mean P	Mean F1	Mean IoU
KLD	80.37	69.84	73.93	62.64
JSD	80.37	69.84	73.93	62.64
ED	80.37	69.84	73.93	62.64
CosD	77.01	74.90	75.43	64.53

B. Performance Variation With Respect to Number of Residual Blocks

In the encoder part of our proposed network, the first two residual blocks downsample the input data. The third block is for feature learning without downsampling. Before reaching our network architecture settings, we experimented with the number of non-downsampling residual blocks. The average results of the conducted experiments are shown in Table V. Best results were recorded with one non-downsampling residual block, therefore this setting is used in our proposed network. Increasing the number of residual blocks further

shows a slight decrease in the performance. Also, note that our approach is unsupervised and use limited training data. Therefore a shallow network is an appropriate choice hence the behavior.

TABLE V
PERFORMANCE VARIATION WITH RESPECT TO NUMBER OF RESIDUAL BLOCKS.

Blocks	Mean R	Mean P	Mean F1	Mean IoU
0	75.30	72.46	72.67	61.13
1	77.01	74.90	75.43	64.53
2	74.97	71.75	72.17	60.26

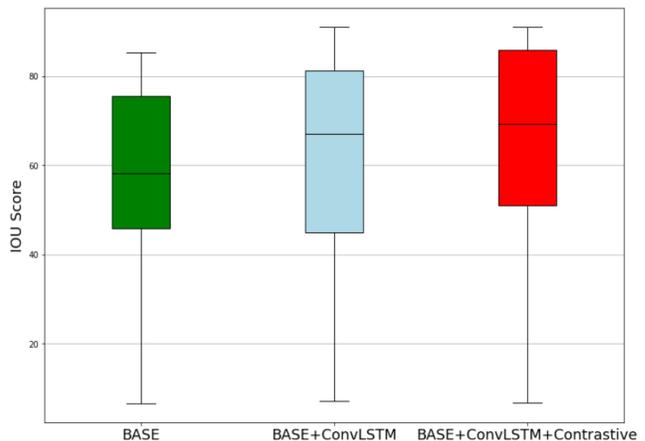


Fig. 10. Ablation results visualized with boxplot depicting mean IoU score across all sites. X-axis represents mean IoU score and y-axis represents boxplots for BASE network, BASE with Convolution LSTM network and our proposed CLAVE method which is BASE network with convolutional LSTM and trained with contrastive learning.

C. Ablation Experiments

The term ‘‘ablation study’’ is borrowed from the medical field that consist of removing parts of the nervous system of vertebrates to understand their purpose. This technique was originally introduced by the French physiologist M.J.P. Flourens [59]. In DL, ablation is removal of parts of the network and analysing the performance of the resulting networks. It helps in investigating the contribution of different parts or techniques used in the DL network. In this study we conducted an ablation of convLSTM and contrastive learning method. Quantitative results are shown in Table VI.

TABLE VI
ABLATION STUDY. COMPARISON OF BASE ARCHITECTURE WITH CONVOLUTIONAL LSTM, BASE ARCHITECTURE WITH CONTRASTIVE LEARNING AND OUR PROPOSED CLVAE NETWORK.

	Mean R	Mean P	Mean F1	Mean IoU
BASE	74.75	69.99	70.87	58.44
+ ConvLSTM	76.90	73.95	73.68	61.89
+ Contrastive Learning	77.01	74.90	75.43	64.53

The ‘BASE’ network refers to the encoder-decoder based VAE reconstruction network. The four metric values given in the table are averaged over all the sites. The results

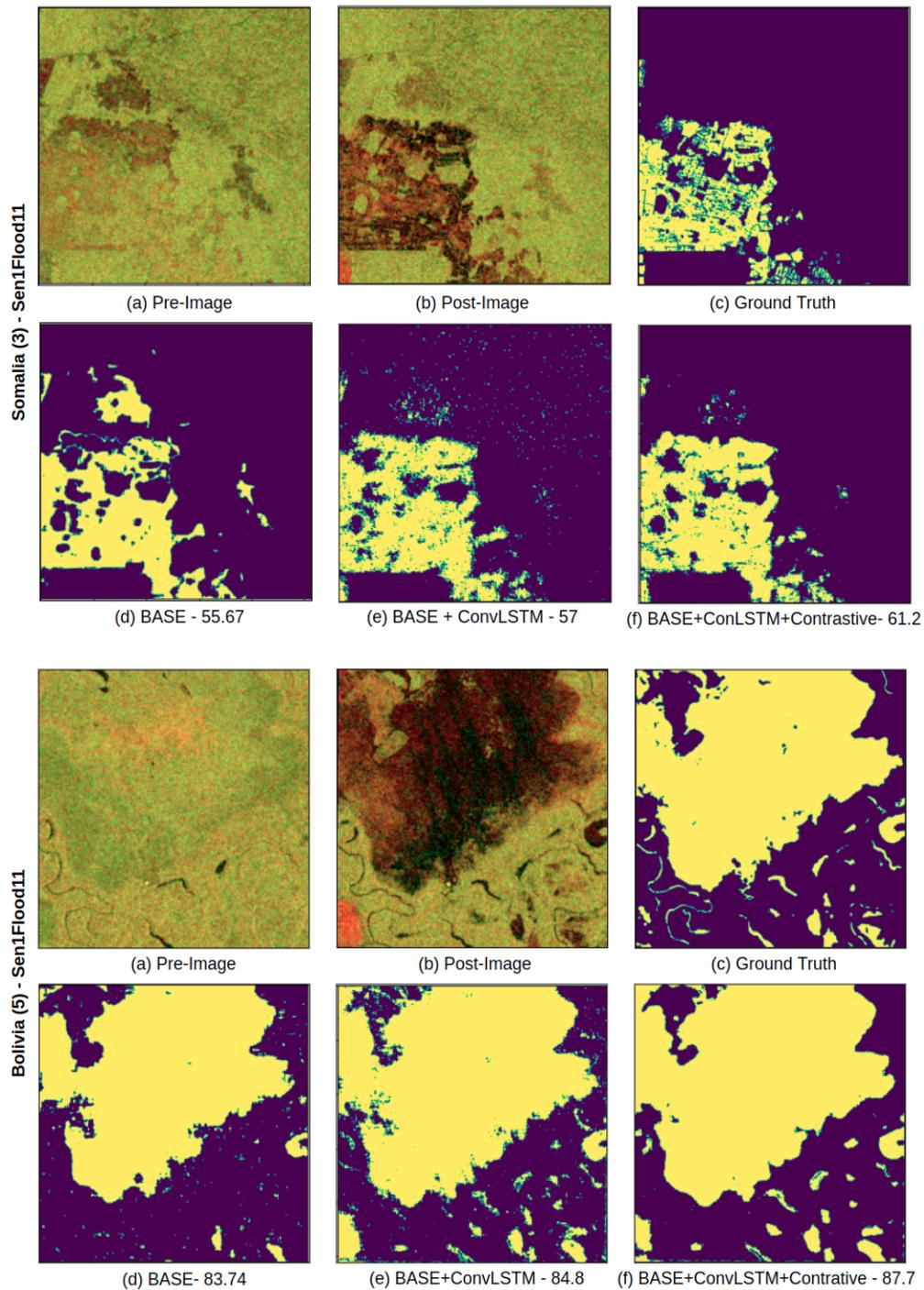


Fig. 11. Qualitative Comparison for Ablation Study: In first row (a), (b) and (c) represents latest pre-flood image, post-flood image and ground truth. In second row (d), (e) and (f) represents detection results from BASE network(encoder-decoder based VAE), BASE with Convolutional LSTM, and proposed CLVAE CD method(BASE with Convolutional LSTM trained using contrastive learning). The number below each change map is corresponding percentage IoU score.

depict that convLSTM helped the network to learn better representation and lead to better recall, precision, F1 score, and IoU. Contrastive learning, on the other hand, improved the precision of the results at the risk of lower recall. Contrastive learning also provided better F1 and IoU score, indicating that training the model with contrastive learning reduce false detections (FN and FP). Figure 10 gives further insight into the ablation study in form of box plots depicting IoU score on y-axis. From left to right we see an increase in median value or increase in length of the upper quartile, which indicates an increase in IoU score for at least 50% of the sites.

For the qualitative analysis of the effect of ConvLSTM and contrastive learning can be seen in the two samples shown in Figure 11. In the first sample from Bolivia, ConvLSTM (e) detected more changes but also added speckle noise in some regions. The speckle-noise is then removed by contrastive learning (f). In the second sample from 'Somalia' site, ConvLSTM (e) eliminates false detections and gave more true positives at the risk of small speckle noise. The speckle noise is then removed with the help of contrastive learning(see (f)), leading to a better IoU score.

ConvLSTM helps in learning better feature representation and can detect the changes more efficiently. But at the same time, it introduces speckle noise to the result, which is handled by training the model with contrastive learning. With contrastive learning, our model learns speckle noise and avoids that to be part of a change. This in turn helps the framework identify meaningful differences between pre and post-flood images.

D. Performance Variation With Respect to Time Series Length

In this section, we show how our proposed CD network performed with smaller and bigger time series data. Before selecting the settings of the proposed CLVAE model we experimented with different time series lengths. The average quantitative results corresponding to 2, 4, and 8 time series lengths are given in Table VII. Our network's performance improved as we increased the time series length from two to four. We see a small decline when the network is supplied with eight pre-images. The observed reason is an increase in seasonal changes and frequent partial flooding. Our network gave the best detection results with pre-images from the same season. It is noteworthy that, at the cost of a small decline in performance of the network is still reliable to use for a longer time series.

TABLE VII
PERFORMANCE VARIATION WITH RESPECT TO LENGTH(2, 4 AND 8) OF TIME SERIES.

	Mean R	Mean P	Mean F1	Mean IoU
2	72.20	68.40	70.24	54.14
4	77.01	74.90	75.43	64.53
8	77.19	71.68	73.65	62.40

E. Performance Variation With Respect to Patch Size.

We also experimented with the network's input patch size. The average metric values of the results are shown in Table

VIII. The network performed best with patch size 16x16. Our CD network uses patch-wise distribution differences to generate the final change map. As we increase the patch size, the network fails to capture small changes efficiently through the distribution difference. Therefore smaller patch gave better results shown below.

TABLE VIII
PERFORMANCE VARIATION WITH RESPECT TO INPUT PATCH SIZE.

	Mean R	Mean P	Mean F1	Mean IoU
16x16	77.01	74.90	75.43	64.53
32x32	76.08	69.45	71.56	59.22

F. Generalizability

The generalizability of a network is its capability to generate good results on unseen sites. Our proposed CLVAE is an unsupervised CD network and hence requires no label for training. This means that we do not need to rely on the generalizability of our model. Rather we can train it on unlabeled SAR data from any area of interest covered by SAR satellite (Sentinel-1) and use it for CD. However, we conducted experiments to investigate the generalizability of our proposed unsupervised CD method compared to the supervised methods.

We trained our CLVAE network only on pre-images from 'Spain' flood site and took inference on all CEMS sites. The generated average results are given in Table IX. This experiment shows that our unsupervised method is still performing better than the supervised Sen1Floods11 benchmark method. On average CLVAE gave better precision, F1 score and IoU score. It is also worth noting that, the supervised method DAUSAR gave a high recall but really low precision. Big difference between recall(high) and precision(low) indicates a lot of false positives, which we can see in the qualitative results presented in Figure 12. ADS-Net on the other hand gave lower recall as well as lower precision compared to our proposed unsupervised change detection method,

TABLE IX
GENERALIZABILITY COMPARISON OF PROPOSED UNSUPERVISED CLVAE METHOD WITH TWO SUPERVISED METHODS ADS-NET AND DAUSAR TRAINED ON SEN1FLOODS11 BENCHMARK.

	Mean R	Mean P	Mean F1	Mean IoU
CLVAE	78.49	71.52	74.8	59.79
ADS-Net	62.21	69.74	60.36	47.39
DAUSAR	88.49	53.86	63.70	51.26

For qualitative comparison, change maps for two flood sites 'Australia' and 'Bosnia' are shown in Figure 12. In both samples, the supervised method detected a large portion of the non-flooded area as flooded. Therefore generate false positives. Whereas our CLVAE gave significantly better CD results. On 'Australia' site CLVAE outperformed the supervised method by 3 to 13% and on 'Bosnia' site by 6 to 14%. Even though we see a drop in CLVAE performance when it is trained on one site and inference is taken on others, the drop is not that high. The CD by CLVAE is still generalizable to other geographically different and unseen sites.

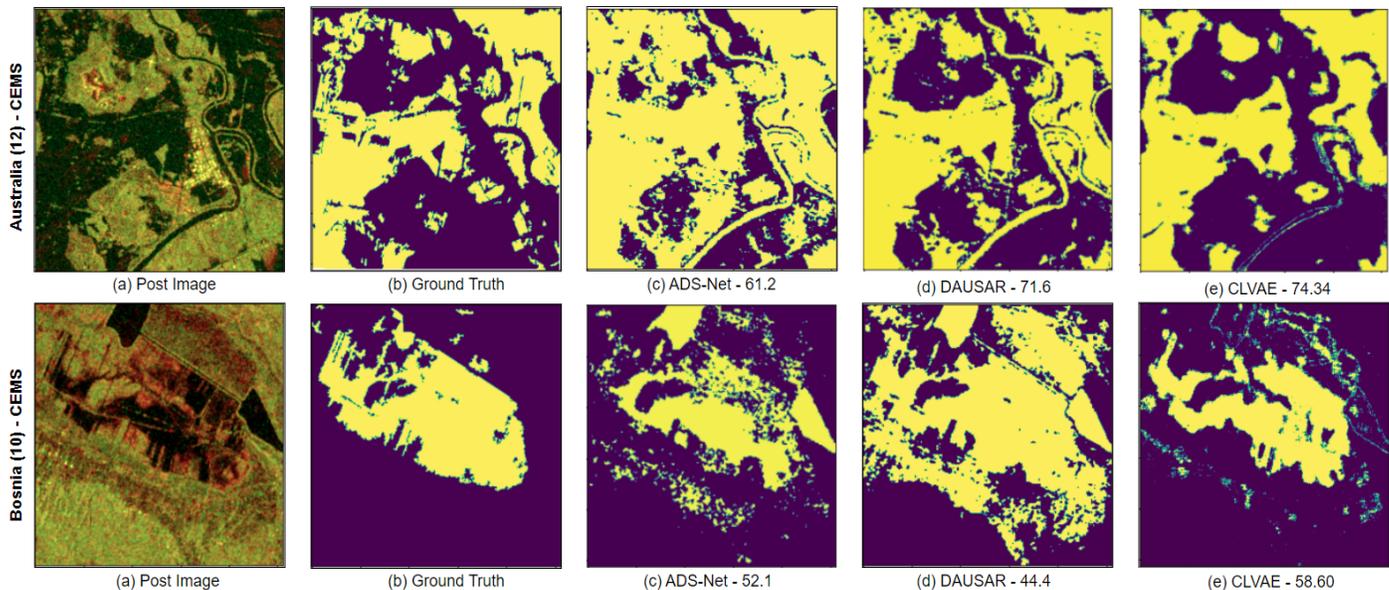


Fig. 12. Qualitative Comparison for Generalizability Test: (a) represents post flood image, (b) represents ground truth, (c) represents change map from supervised SenIFloods11 benchmark method, and (d) represents change map from our proposed CLVAE method. The number below each change map is corresponding percentage IoU score.

VII. CONCLUSION

In this paper, we proposed a novel unsupervised remote sensing CD method based on a probabilistic model. Our method CLVAE cumulatively benefits from the reconstruction approach, latent parameters learning of probabilistic auto-encoder, distribution difference method, convolutional LSTM, and contrastive learning techniques. Our model strongly learns the spatiotemporal correlation between time series SAR data. The extensive experimental results on SenIFloods11 data and CEMS data display the potential of the proposed CD method. Our method yield 64.53% average IoU and 75.43% average F1 score. Our results have surpassed the performance of existing unsupervised non-DL methods i.e. log-ratio, CVA, and unsupervised DL method RaVAEn. On average our CLVAE has a lead of 10-21% in terms of F1 score and 7-27% in terms of IoU score. As 6 out of 9 sites in our evaluation data are from SenIFloods11 test data, we also compared our results with the supervised methods ADS-Net and DAUSAR trained on SenIFloods11 dataset. Our unsupervised CD method CLVAE shows a lead of 6% F1 score and 7% IoU over compared supervised method. Further to this, we also presented a change point detection framework based on our CD method. The framework can detect changes at early stages which in turn can save lives through timely evacuation, alerts, and other disaster management activities. In light of new satellite missions, a better temporal frequency of one or two images per day can be immensely helpful in monitoring the change point. The proposed method and framework are light on memory, have low computation time (faster training and inference), and are also inexpensive in terms of data preparation as no annotation is required.

In this study, we proposed a CD method on Sentinel-1 SAR data and presented its efficiency in detecting floods. In future

we will test and extend our change detection method in other applications areas as well. Different remote sensing sensors capture specific features of the scene. In our future works we would like to train our network with complimenting features from multi-sensors. Another potential direction is urban flood detection using high-resolution data, where our model can be trained to detect flooded buildings and roads at a better resolution. This can help local transport agencies to reroute the traffic saving lives of pedestrians and drivers.

REFERENCES

- [1] "2021 disasters in numbers," <https://reliefweb.int/report/world/2021-disasters-numbers#:~:text=In%202021%2C%20a%20total%20of,across%20the%202001%2D2020%20period,> accessed: 2022-06-30.
- [2] N. Anusha and B. Bharathi, "Flood detection and flood mapping using multi-temporal synthetic aperture radar and optical data," *The Egyptian Journal of Remote Sensing and Space Science*, vol. 23, no. 2, pp. 207–219, 2020.
- [3] G. Di Baldassarre, G. Schumann, L. Brandimarte, and P. Bates, "Timely low resolution sar imagery to support floodplain modelling: a case study review," *Surveys in geophysics*, vol. 32, no. 3, pp. 255–269, 2011.
- [4] D. Lu, P. Mausel, E. Brondizio, and E. Moran, "Change detection techniques," *International journal of remote sensing*, vol. 25, no. 12, pp. 2365–2401, 2004.
- [5] L. T. Luppino, F. M. Bianchi, G. Moser, and S. N. Anfinsen, "Unsupervised image regression for heterogeneous change detection," *arXiv preprint arXiv:1909.05948*, 2019.
- [6] W. A. Malila, "Change vector analysis: An approach for detecting forest changes with landsat," in *LARS symposia*, 1980, p. 385.
- [7] F. Bovolo, "A multilevel parcel-based approach to change detection in very high resolution multitemporal images," *IEEE Geoscience and Remote Sensing Letters*, vol. 6, no. 1, pp. 33–37, 2008.
- [8] F. Thonfeld, H. Feilhauer, M. Braun, and G. Menz, "Robust change vector analysis (rcva) for multi-sensor very high resolution optical satellite data," *International Journal of Applied Earth Observation and Geoinformation*, vol. 50, pp. 131–140, 2016.
- [9] J. Deng, K. Wang, Y. Deng, and G. Qi, "Pca-based land-use change detection and analysis using multitemporal and multisensor satellite data," *International Journal of Remote Sensing*, vol. 29, no. 16, pp. 4823–4838, 2008.

- [10] M. Dharani and G. Sreenivasulu, "Land use and land cover change detection by using principal component analysis and morphological operations in remote sensing applications," *International Journal of Computers and Applications*, vol. 43, no. 5, pp. 462–471, 2021.
- [11] A. Asokan and J. Anitha, "Change detection techniques for remote sensing applications: a survey," *Earth Science Informatics*, vol. 12, no. 2, pp. 143–160, 2019.
- [12] R. C. Daudt, B. Le Saux, and A. Boulch, "Fully convolutional siamese networks for change detection," in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 4063–4067.
- [13] W. G. C. Bandara and V. M. Patel, "Revisiting consistency regularization for semi-supervised change detection in remote sensing images," *arXiv preprint arXiv:2204.08454*, 2022.
- [14] C. Zhang, P. Yue, D. Tapete, L. Jiang, B. Shangguan, L. Huang, and G. Liu, "A deeply supervised image fusion network for change detection in high resolution bi-temporal remote sensing images," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 166, pp. 183–200, 2020.
- [15] J. Chen, Z. Yuan, J. Peng, L. Chen, H. Huang, J. Zhu, Y. Liu, and H. Li, "Dasnet: Dual attentive fully convolutional siamese networks for change detection in high-resolution satellite images," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 14, pp. 1194–1206, 2020.
- [16] D. Wang, X. Chen, M. Jiang, S. Du, B. Xu, and J. Wang, "Ads-net: An attention-based deeply supervised network for remote sensing image change detection," *International Journal of Applied Earth Observation and Geoinformation*, vol. 101, p. 102348, 2021.
- [17] R. Yadav, A. Nascetti, and Y. Ban, "Attentive dual stream siamese u-net for flood detection on multi-temporal sentinel-1 data," *arXiv preprint arXiv:2204.09387*, 2022.
- [18] H. Chen, Z. Qi, and Z. Shi, "Remote sensing image change detection with transformers," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–14, 2021.
- [19] D. F. Muñoz, P. Muñoz, H. Mofatkhari, and H. Moradkhani, "From local to regional compound flood mapping with deep learning and data fusion techniques," *Science of the Total Environment*, vol. 782, p. 146927, 2021.
- [20] P. Manjusree, L. Prasanna Kumar, C. M. Bhatt, G. S. Rao, and V. Bhanumurthy, "Optimization of threshold ranges for rapid flood inundation mapping by evaluating backscatter profiles of high incidence angle sar images," *International Journal of Disaster Risk Science*, vol. 3, no. 2, pp. 113–122, 2012.
- [21] Z. Zhong, J. Li, Z. Luo, and M. Chapman, "Spectral–spatial residual network for hyperspectral image classification: A 3-d deep learning framework," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 2, pp. 847–858, 2017.
- [22] H. Noh, J. Ju, M. Seo, J. Park, and D.-G. Choi, "Unsupervised change detection based on image reconstruction loss," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 1352–1361.
- [23] W. Jing, S. Zhu, P. Kang, J. Wang, S. Cui, G. Chen, and H. Song, "Remote sensing change detection based on unsupervised multi-attention slow feature analysis," *Remote Sensing*, vol. 14, no. 12, p. 2834, 2022.
- [24] V. Ržička, A. Vaughan, D. De Martini, J. Fulton, V. Salvatelli, C. Bridges, G. Mateo-Garcia, and V. Zantedeschi, "Unsupervised change detection of extreme events using ml on-board," *arXiv preprint arXiv:2111.02995*, 2021.
- [25] T. Zhan, M. Gong, X. Jiang, and S. Li, "Log-based transformation feature learning for change detection in heterogeneous images," *IEEE Geoscience and Remote Sensing Letters*, vol. 15, no. 9, pp. 1352–1356, 2018.
- [26] J. Liu, M. Gong, K. Qin, and P. Zhang, "A deep convolutional coupling network for change detection based on heterogeneous optical and radar images," *IEEE transactions on neural networks and learning systems*, vol. 29, no. 3, pp. 545–559, 2016.
- [27] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [28] X. Chen, H. Fan, R. Girshick, and K. He, "Improved baselines with momentum contrastive learning," *arXiv preprint arXiv:2003.04297*, 2020.
- [29] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar *et al.*, "Bootstrap your own latent—a new approach to self-supervised learning," *Advances in neural information processing systems*, vol. 33, pp. 21 271–21 284, 2020.
- [30] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 132–149.
- [31] H. Dong, W. Ma, L. Jiao, F. Liu, and L. Li, "A multiscale self-attention deep clustering for change detection in sar images," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–16, 2021.
- [32] S. Saha, P. Ebel, and X. X. Zhu, "Self-supervised multisensor change detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–10, 2021.
- [33] L. T. Luppino, M. A. Hansen, M. Kampffmeyer, F. M. Bianchi, G. Moser, R. Jenssen, and S. N. Anfinsen, "Code-aligned autoencoders for unsupervised change detection in multimodal remote sensing images," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [34] S. Deldari, D. V. Smith, H. Xue, and F. D. Salim, "Time series change point detection with self-supervised contrastive predictive coding," in *Proceedings of the Web Conference 2021*, 2021, pp. 3124–3135.
- [35] D. Bonafilia, B. Tellman, T. Anderson, and E. Issenberg, "Sen1floods11: a georeferenced dataset to train and test deep learning flood algorithms for sentinel-1," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 210–211.
- [36] "Copernicus emergency management services," <https://emergency.copernicus.eu/mapping/copernicus-emergency-management-service#zoom=2&lat=31.47858&lon=7.20923&layers=0BT00>, accessed: 2020-09-30.
- [37] N. Gorelick, M. Hancher, M. Dixon, S. Ilyushchenko, D. Thau, and R. Moore, "Google earth engine: Planetary-scale geospatial analysis for everyone," *Remote Sensing of Environment*, vol. 202, pp. 18–27, 2017, big Remotely Sensed Data: tools, applications and experiences. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0034425717302900>
- [38] "Copernicus list of ems.risk and recovery mapping activations," <https://emergency.copernicus.eu/mapping/list-of-activations-risk-and-recovery>, accessed: 2020-09-30.
- [39] "Cems bosnia flood," https://emergency.copernicus.eu/mapping/system/files/components/EMSR572_AOI01_DEL_PRODUCT_r1_RTP02_v1.pdf, accessed: 2022-09-30.
- [40] "Cems australia flood," https://emergency.copernicus.eu/mapping/system/files/components/EMSR570_AOI01_DEL_MONIT01_r1_RTP01_v2.pdf, accessed: 2022-09-30.
- [41] "Cems scotland flood," https://emergency.copernicus.eu/mapping/system/files/components/EMSR640_AOI02_DEL_MONIT01_r1_RTP02_v3.pdf, accessed: 2022-09-30.
- [42] X. Liu, G. Hu, Y. Chen, X. Li, X. Xu, S. Li, F. Pei, and S. Wang, "High-resolution multi-temporal mapping of global urban land using Landsat images based on the Google Earth Engine platform," *Remote sensing of environment*, vol. 209, pp. 227–239, 2018.
- [43] P. Gong, X. Li, J. Wang, Y. Bai, B. Chen, T. Hu, X. Liu, B. Xu, J. Yang, W. Zhang *et al.*, "Annual maps of global artificial impervious area (GAIA) between 1985 and 2018," *Remote Sensing of Environment*, vol. 236, p. 111510, 2020.
- [44] X. Zhang, L. Liu, C. Wu, X. Chen, Y. Gao, S. Xie, and B. Zhang, "Development of a global 30 m impervious surface map using multisource and multitemporal remote sensing datasets with the Google Earth Engine platform," *Earth System Science Data*, vol. 12, no. 3, pp. 1625–1648, 2020.
- [45] R. Goldblatt, M. F. Stuhlmacher, B. Tellman, N. Clinton, G. Hanson, M. Georgescu, C. Wang, F. Serrano-Candela, A. K. Khandelwal, W.-H. Cheng *et al.*, "Using Landsat and nighttime lights for supervised pixel-based image classification of urban land cover," *Remote Sensing of Environment*, vol. 205, pp. 253–275, 2018.
- [46] R. Ravanelli, A. Nascetti, R. V. Cirigliano, C. Di Rico, G. Leuzzi, P. Monti, and M. Crespi, "Monitoring the impact of land cover change on surface urban heat island through Google Earth Engine: Proposal of a global methodology, first applications and problems," *Remote Sensing*, vol. 10, no. 9, 2018.
- [47] D. P. Kingma and M. Welling, "Auto-encoding variational bayes. corr abs/1312.6114 (2013)," *arXiv preprint arXiv:1312.6114*, vol. 482, 2013.
- [48] X. Shi, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo, "Convolutional lstm network: A machine learning approach for precipitation nowcasting," *Advances in neural information processing systems*, vol. 28, 2015.
- [49] S. Sun, L. Mu, L. Wang, and P. Liu, "L-unet: An lstm network for remote sensing image change detection," *IEEE Geoscience and Remote Sensing Letters*, 2020.
- [50] D. J. Rezende, S. Mohamed, and D. Wierstra, "Stochastic backpropagation and approximate inference in deep generative models," in *International conference on machine learning*. PMLR, 2014, pp. 1278–1286.

- [51] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality reduction by learning an invariant mapping," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, vol. 2. IEEE, 2006, pp. 1735–1742.
- [52] N. Dong, M. Maggioni, Y. Yang, E. Pérez-Pellitero, A. Leonardis, and S. McDonagh, "Residual contrastive learning for joint demosaicking and denoising," *arXiv preprint arXiv:2106.10070*, 2021.
- [53] B. Huang, D. Reichman, L. M. Collins, K. Bradbury, and J. M. Malof, "Tiling and stitching segmentation output for remote sensing: Basic challenges and recommendations," *arXiv preprint arXiv:1805.12219*, 2018.
- [54] X. Yu, X. Wu, C. Luo, and P. Ren, "Deep learning in remote sensing scene classification: a data augmentation enhanced convolutional neural network framework," *GIScience & Remote Sensing*, vol. 54, no. 5, pp. 741–758, 2017.
- [55] H. Hu and Y. Ban, "Unsupervised change detection in multitemporal sar images over large urban areas," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 7, no. 8, pp. 3248–3261, 2014.
- [56] J.-S. Lee, "Speckle analysis and smoothing of synthetic aperture radar images," *Computer graphics and image processing*, vol. 17, no. 1, pp. 24–32, 1981.
- [57] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE transactions on systems, man, and cybernetics*, vol. 9, no. 1, pp. 62–66, 1979.
- [58] J.-C. Yen, F.-J. Chang, and S. Chang, "A new criterion for automatic multilevel thresholding," *IEEE Transactions on Image Processing*, vol. 4, no. 3, pp. 370–378, 1995.
- [59] "Cmarie-jean-pierre flourens," <https://www.britannica.com/biography/Marie-Jean-Pierre-Flourens>, accessed: 2020-04-11.