

# SENECA: Change Detection in Optical Imagery using Siamese Networks with Active-Transfer Learning

Giuseppina Andresini<sup>a,b,\*</sup>, Annalisa Appice<sup>a,b</sup>, Dino Ienco<sup>c</sup>, Donato Malerba<sup>a,b</sup>

<sup>a</sup>*Department of Computer Science, University of Bari Aldo Moro, Bari, Italy*

<sup>b</sup>*Consorzio Interuniversitario Nazionale per l'Informatica - CINI, Bari, Italy*

<sup>c</sup>*INRAE, UMR TETIS, University of Montpellier, Montpellier, France*

---

## Abstract

<sup>1</sup> Change Detection (CD) aims to distinguish surface changes based on bi-temporal remote sensing images. In the recent years, deep neural models have made a breakthrough in CD processes. However, training a deep neural model requires a large volume of labelled training samples that are time-consuming and labour-intensive to acquire. With the aim of learning an accurate CD model with limited labelled data, we propose **SENECA**: a method based on a CD Siamese network, which takes advantage of both Active Learning (AL) and Transfer Learning (TL) to handle the constraint of limited supervision. More precisely, we jointly use AL and TL to adapt a CD model trained on a labelled source domain to a (related) target domain featured by a restricted access to labelled data. We report results from an experimental evaluation involving five pairs of images acquired via Sentinel-2 satellites between 2015 and 2018 in various locations picked all over Asia and USA. The results show the beneficial effects of the proposed AL and TL strategies on

---

\*Corresponding author

*Email addresses:* [giuseppina.andresini@uniba.it](mailto:giuseppina.andresini@uniba.it) (Giuseppina Andresini), [annalisa.appice@uniba.it](mailto:annalisa.appice@uniba.it) (Annalisa Appice), [dino.ienco@inrae.fr](mailto:dino.ienco@inrae.fr) (Dino Ienco), [donato.malerba@uniba.it](mailto:donato.malerba@uniba.it) (Donato Malerba)

<sup>1</sup>This version of the contribution has been accepted for publication, after peer review but is not the Version of Record and does not reflect post-acceptance improvements, or any corrections. The Version of Record is available online at: <https://doi.org/10.1016/j.eswa.2022.119123>. Accepted Version is subject to the publisher's Accepted Manuscript terms of use <https://www.elsevier.com/about/policies-and-standards/copyright>

the accuracy of the decisions made by the CD Siamese network and depict the merit of the proposed approach over competing CD baselines.

*Keywords:* Active learning, Transfer learning, Siamese network, Change detection, Sentinel-2 data

---

## 1. Introduction

Change Detection (CD) is a central task in the field of computer vision since it has the objective to detect changes in multiple images of the same scene acquired at different period of time (Ru et al. 2021). Focusing on the analysis of optical remote sensing images depicting the same geographical area, the CD task is the process of detecting differences among various images of the same scene as a consequence of natural and/or human activities.

Due to the unprecedented availability of remote sensing imagery acquired by up-to-date Earth observation systems (a notable example is the Sentinel-2 mission belonging to the European Copernicus Programme <sup>2</sup>), it is becoming easier and easier to obtain images covering the same geographical area acquired with a regular revisit time. This technological revolution highlights the importance of conceiving and developing effective CD methods to fully exploit the amount of freely available remote sensing information.

Application-wise, CD methods are largely employed in remote sensing analysis (Lv et al. 2022) to cope with a diverse set of applications like land cover change detection (Shi, Zhong, Zhao, Lv, Liu & Zhang 2022), urban change detection (Hafner et al. 2022), disaster management (Sublime & Kalinicheva 2019) and environmental monitoring (Lewis et al. 2016), among the others.

Modern advances in CD methods mainly rely on deep learning (DL) approaches (Jiang et al. 2022) due to their ability to cope with imagery data through the automatic extraction of hierarchical multilevel features via representational learning (Bengio et al. 2013). Despite remarkable performances exhibited by neural network approaches in many different applications of image analysis, one of the main limitation of the use of DL methods is related to the label-hungry behaviour they exhibit. In fact, a large amount of labelled samples is, commonly, required for effective deep neural network

---

<sup>2</sup><https://www.copernicus.eu/en>

29 training, while facing this condition poses critical issues related to the deploy-  
30 ment of DL approaches in tasks characterised by limited amount of labelled  
31 data (Ouali et al. 2020). For the specific case of remote sensing CD, it can be  
32 highly labor-intensive and time-consuming to collect remote sensing image  
33 pairs with well-labelled change information each time a CD method must be  
34 deployed. This condition may affect the use of DL methods for CD, since  
35 a DL model may need to be transferred from the imagery data (pair of im-  
36 ages) on which it is learnt (source data) to new unseen imagery data (target  
37 data). Moreover, the crucial point is how to alleviate the dependence of DL  
38 models from large amount of labelled data, while meeting the requirement  
39 to transfer CD models from a source to a target scenario.

40 With the objective to reduce the dependence of remote sensing CD models  
41 from the necessity to access abundant amount of labelled data when trans-  
42 ferred on new scene, in this paper, we propose SENECA (Siamese nEtwork  
43 based chaNge detection in optical imagEry with aCtive trAnsfer learning): a  
44 method based on a Siamese network, which combines both Active Learning  
45 (AL) and Transfer Learning (TL) to effectively deal with remote sensing CD  
46 analysis in a scenario featured by limited supervision.

47 The Siamese network is a neural model especially tailored to compare  
48 together pairs of entities (Lu et al. 2017) with the aim to learn data em-  
49 beddings that satisfy pair-wise metric constraints. As a Siamese network  
50 is well-suited to deal with the class imbalance condition (Gautheron et al.  
51 2020), recent studies (Shi, Liu, Li, Liu, Wang & Zhang 2022, Ruzicka et al.  
52 2020) have started the investigation of Siamese networks in CD tasks, where  
53 the number of changed pixels is often much less than that of unchanged ones.

54 The proposed method firstly learns a CD Siamese network on source  
55 data, where change labelled data are available and, successively, it adapts  
56 the source Siamese network to the target data through a Transfer Learning  
57 (TL) strategy. TL is performed with fine tuning, that is one of the most  
58 widely used approach for TL when working with DL models. In particular,  
59 the fine tuning approach starts with a pre-trained deep neural model on the  
60 source data and trains it further on the target data. In this study, the fine  
61 tuning approach is performed with a limited supervision provided by means of  
62 an Active Learning (AL) strategy. More precisely, we adopt a segmentation-  
63 based AL strategy that allows the Siamese network pre-trained for CD in  
64 a source area to select samples that spatially span the target area. This  
65 may contribute to reduce possible issues exhibited to confidence-based AL  
66 strategy in remote sensing data analysis (Pasolli et al. 2014) that are more

67 prone to select redundant, in terms of spatial auto-correlation, samples.

68 The experimental evaluation, involving recent state of the art CD com-  
69 petitors, on five pair of images acquired via the Sentinel-2 satellite missions <sup>3</sup>  
70 between 2015 and 2018 in various locations picked all over Asia and USA  
71 has demonstrated the quality and the value of the proposed approach. More  
72 precisely, the results show the beneficial effects of combining AL and TL for  
73 all the downstream CD tasks.

74 In short, this paper provides the following contributions:

- 75 • The definition of a new CD method that is formulated combining both  
76 AL and TL, in order to reduce the necessity to access abundant amount  
77 of labelled samples when a CD neural network (Siamese network) is  
78 transferred from a source scene to a new target scene.
- 79 • The use of a segmentation-based AL strategy that allows us to effec-  
80 tively select active samples spanning the target area, in order to transfer  
81 the pre-trained CD Siamese network from the source to the target pair  
82 of images.
- 83 • An in-depth and extensive evaluation of the proposed method **SENECA**  
84 w.r.t. recent competing CD methods on five co-registered, bi-temporal  
85 multispectral images acquired with Sentinel-2 satellites in locations  
86 picked all over both Asia and USA.

87 The rest of this manuscript is organised as follows. Section 2 presents  
88 the recent related literature in remote sensing CD analysis. Section 3 in-  
89 troduces the background and the CD problem definition we adopt in our  
90 work. Section 4 describes the proposed Active-Transfer Learning (ATL) CD  
91 framework. Section 5 introduces the experimental settings, the performances  
92 evaluation as well as the discussion related to the results while Section 6 con-  
93 cludes and pave the way to possible future works.

## 94 **2. Related Work**

95 A general overview of CD approaches for land cover dynamics is presented  
96 in (Lv et al. 2022). Here, the authors review the main issues in terms of  
97 methods, applications and available benchmarks related to remote sensing

---

<sup>3</sup><https://sentinel.esa.int/nl/web/sentinel/missions/sentinel-2>

98 CD with a particular focus on Very High spatial Resolution (VHR) imagery.  
99 Recently, Jiang et al. (2022) have provided a review of CD methods for  
100 remote sensing imagery under the lens of DL-based techniques underlying  
101 the fact that the community still lacks of a comprehensive review of the  
102 recent progress concerning neural network methods in remote sensing CD.

103 With a focus on the advances on unsupervised CD methods, Celik (2009)  
104 defined a PCA-based method for CD in multitemporal satellite images. This  
105 method partitions a difference image into non-overlapping blocks and per-  
106 forms the PCA, in order to extract the orthonormal eigenvectors of the set  
107 of non-overlapping blocks and build an eigenvector space. Subsequently, it  
108 represents each pixel of the difference image with a new feature vector that  
109 is the projection of the block-based difference image samples onto the gen-  
110 erated eigenvector space. Finally, the CD map is built by partitioning the  
111 feature vector space into two clusters using k-means clustering with  $k = 2$   
112 and then assigning each pixel to the one of the two clusters by using the  
113 minimum Euclidean distance between the pixel’s feature vector and mean  
114 feature vector of clusters.

115 Appice et al. (2020) introduced an unsupervised learning method for CD.  
116 This method combines clustering, PCA and classification with the aim to  
117 separate changed areas from unchanged background. More in detail, firstly  
118 a clustering stage is performed on the bi-temporal images with the aim to  
119 identify an initial set of labelled samples and, successively, the extracted  
120 labelled samples are used to feed a supervised binary classification stage. The  
121 classification stage trains a Random Forest from the principal components  
122 of the fusion (concatenation) of bi-temporal pixel vectors using the labels  
123 produced in the clustering stage.

124 López-Fandiño et al. (2019) experimented a change vector analysis (CVA)  
125 method in the field of imagery CD. The proposed CVA method computes the  
126 difference between two optical images of a scene with the spectral angle dis-  
127 tance and uses the Otsu’s thresholding to separate the changed areas from  
128 the unchanged areas. Andresini et al. (2022) investigated the use of autoen-  
129 coder neural networks for CVA in a pair of optical images. More precisely,  
130 given a pair of optical images, the method learns an autoencoder model on  
131 the first image of the bi-temporal image pair then, the model is employed to  
132 reconstruct both the first and second images. Successively, the spectral angle  
133 distance is computed pixel-wise and, finally a threshold approach is adopted  
134 to separate changed from non-changed pixels in a totally unsupervised way.

135 Ma et al. (2019) illustrated a matrix factorisation method for CD in syn-

136 thetic aperture radar images. In this method, the factorisation model of  
137 the low-rank and sparse matrix is used to extract both (unchanged) back-  
138 ground and (changed) foreground information from images. More in detail,  
139 mean and variance matrices related to both unchanged and changed areas are  
140 summarised through statistical features that are, subsequently, used to learn  
141 a naive Bayes classifier. At the end, the classification model is employed to  
142 derive a CD map that distinguishes between changed and unchanged areas.

143 Wu et al. (2020) described an unsupervised method for CD in optical  
144 images based on generative adversarial networks. The proposed method uses  
145 CVA to build an initial CD map. Subsequently it applies a training sample  
146 selection method to select training samples that are processed to train the  
147 generative adversarial network. The generator of the generative adversarial  
148 network is used to build the final CD map.

149 Wu et al. (2021) illustrated an unsupervised method defined for CD in  
150 heterogeneous images. This method takes a pair of images acquired with  
151 different sensors (e.g., optical images and synthetic aperture radar images)  
152 as input. It combines together convolutional autoencoder and commonality  
153 autoencoder with the aim to firstly extract a vector-based representation of  
154 the input images and, successively, extract the common features by means of  
155 a reconstruction process. Finally, it deploys an unsupervised segmentation  
156 approach on the difference map to extract the changed areas.

157 Regarding recent supervised CD methods, Daudt et al. (2018) proposed  
158 a CD framework based on a fully convolutional Siamese network. In the  
159 proposed method, firstly the image pairs are encoded via a Siamese network  
160 with the aim to extract new data representation from each of the images  
161 and, then, the extracted bi-temporal representations are combined with the  
162 aim to produce a CD map in a fully supervised fashion. Here, the method  
163 is developed to make inference on the same data on which it is learnt with-  
164 out taking into account possible shifts in the underlying data distribution  
165 between training and test data.

166 Yang et al. (2019) proposed a DL-based CD method especially tailored to  
167 transfer a CD model from a source to a target domain. The proposed method  
168 involves two stages: i) a pre-training step in which the model is trained  
169 on the label abundant source domain and ii) a refinement step in which  
170 the model is fine-tuned according to pseudo-labels generated on the target  
171 domain in a self-training manner. The refinement stage exploits pseudo CD  
172 maps generated on the target data on which spatial reasoning, at region- and  
173 boundary-scale, is deployed to select target samples with associated pseudo-

174 labels.

175 Shi, Liu, Li, Liu, Wang & Zhang (2022) designed and evaluated a deeply  
176 supervised attention metric-based network. CD maps are learnt by means of  
177 Siamese networks, while convolutional attention blocks are integrated with  
178 the aim to provide highly discriminative features. In addition, supervision  
179 is employed to enhance the feature extractor’s learning ability and gener-  
180 ate more useful features that are subsequently used to discriminate between  
181 changed and unchanged areas.

182 Finally, Ruzicka et al. (2020) explored AL in the context of neural network  
183 based remote sensing CD. The proposed work evaluates AL in a scenario char-  
184 acterised by a reduced amount of labelled source data to train the CD model.  
185 The method leverages as backbone model a Siamese network with an encoder  
186 pre-trained on the Imagenet dataset. To implement the AL process, the un-  
187 certainty related to an ensemble of Siamese network models is exploited as a  
188 criterion to sample new labelled data thus enriching the training set. While  
189 the inference is performed at pixel level, the method selects new samples at  
190 tile level thus, possibly introducing noisy information in the training data.  
191 The findings of this study highlight that the AL process permits to automati-  
192 cally balance the training distribution reaching out similar performances as a  
193 model supervised with a large pre-annotated training set. While this method  
194 shares with our proposal the idea to use AL, it differs from **SENECA** on two  
195 main aspects: firstly, the AL sampling strategy is purely based on model  
196 uncertainty without taking into account the spatial dimension that strongly  
197 characterizes remote sensing data and, secondly, it integrates new samples  
198 at tile level (patch of  $256 \times 256$  pixels) thus introducing possible noisy labels  
199 conversely to our method in which pixel-level samples are integrated.

200 To sum up, the majority of CD methods only use the information from the  
201 current images themselves without taking into account possible distribution  
202 shifts between the training data (here referred as source domain) and the test  
203 data (here referred as target domain) that can negatively impact the final  
204 detection performances. When methods are proposed to explicitly manage  
205 such a data distribution shift, they mainly rely on heuristics (Shi, Liu, Li,  
206 Liu, Wang & Zhang 2022), sample selection based on uncertainty derived by  
207 the model output (Ruzicka et al. 2020) or self-training approaches that can  
208 introduce issues related to confirmation bias (Tarvainen & Valpola 2017) as  
209 well as mistakes due to large gaps between the source and the target domains.

210 **3. Basics**

211 Let us consider a MultiSpectral (MS) sensor technology (e.g., Sentinel-2)  
 212 to observe the Earth’s surface over  $K$  spectral bands. Every spectral band is a  
 213 numeric feature proportional to the ultraviolet and short wavelength infrared  
 214 for a given band. Let **scene** be a geographic scene spanned over  $m_{\text{scene}} \times$   
 215  $n_{\text{scene}}$  pixels, where a pixel denotes an area of around a few square meters  
 216 of the Earth’s surface (i.e., it is a function of the sensor’s spatial resolution),  
 217 which is unequivocally referenced with spatial coordinates  $(i, j)$ , with  $1 \leq i \leq$   
 218  $m_{\text{scene}}$  and  $1 \leq j \leq n_{\text{scene}}$ , according to the usual matrix representation. A  
 219 bi-temporal MS dataset  $\mathbf{D}_{\text{scene}}$  is composed of two co-registered MS images,  
 220 i.e.,  $\mathbf{D}_{\text{scene}} = (\mathbf{X}_{\text{scene}}^1, \mathbf{X}_{\text{scene}}^2)$ , which describe MS data of **scene** acquired  
 221 by using the Sentinel-2 MS sensor technology. Note that  $\mathbf{X}_{\text{scene}}^1$  and  $\mathbf{X}_{\text{scene}}^2$   
 222 are acquired in two distinct time periods, denoted as  $t^1$  and  $t^2$ , respectively,  
 223 with  $t^1 < t^2$ . Every MS image of  $\mathbf{D}_{\text{scene}}$  is represented as a tensor of  $m_{\text{scene}} \times$   
 224  $n_{\text{scene}}$  pixels and  $K$  spectral bands. For each dataset, the pixel indexed by  
 225 row  $i$  and column  $j$  contains a vector of data sensed on that resolution cell  
 226 over  $K$  spectral bands (MS vector). The pair  $(\mathbf{X}_{\text{scene}}^1(i, j), \mathbf{X}_{\text{scene}}^2(i, j))$   
 227 denotes the bi-temporal MS vectors of  $\mathbf{D}_{\text{scene}}$  associated with pixel  $(i, j)$ .  
 228 Finally, the CD map  $\mathbf{Y}_{\text{scene}}$  of a bi-temporal dataset  $\mathbf{D}_{\text{scene}}$  is a matrix of  
 229  $m_{\text{scene}} \times n_{\text{scene}}$  binary labels with  $\mathbf{Y}(i, j) = 1$  if a change occurred in the  
 230 surface covered by pixel  $(i, j)$  from  $t^1$  to  $t^2$ ; 0 otherwise.

231 **4. Proposed Method**

232 We assume that a MS sensor technology with  $K$  MS bands is used to  
 233 monitor both a source scene  $\mathbf{S}$  and a target scene  $\mathbf{T}$ , respectively. Both  $\mathbf{S}$   
 234 and  $\mathbf{T}$  covering different geographical areas. Each MS image of scene  $\mathbf{S}$   
 235 is represented as  $m_{\mathbf{S}} \times n_{\mathbf{S}}$  pixels and  $K$  spectral bands. Each MS image of scene  
 236  $\mathbf{T}$  is represented as  $m_{\mathbf{T}} \times n_{\mathbf{T}}$  pixels and  $K$  spectral bands. Let us consider:  
 237 (1) a bi-temporal MS dataset  $\mathbf{D}_{\mathbf{S}} = (\mathbf{X}_{\mathbf{S}}^1, \mathbf{X}_{\mathbf{S}}^2)$  of  $\mathbf{S}$ ; (2) the ground truth  
 238 CD map  $\mathbf{Y}_{\mathbf{S}}$  of  $\mathbf{D}_{\mathbf{S}}$ ; and (3) a bi-temporal MS dataset  $\mathbf{D}_{\mathbf{T}} = (\mathbf{X}_{\mathbf{T}}^1, \mathbf{X}_{\mathbf{T}}^2)$   
 239 of  $\mathbf{T}$ . The CD methodology of SENECA, schematised in Figure 1, is mainly  
 240 based on four components:

- 241 • The training of a CD model (pre-trained source CD model) from the  
 242 labelled source, bi-temporal MS dataset.
- 243 • The use of a segmentation-based AL strategy to divide the target scene  
 244 in super-pixel objects, select the medoids of the super-pixel objects and



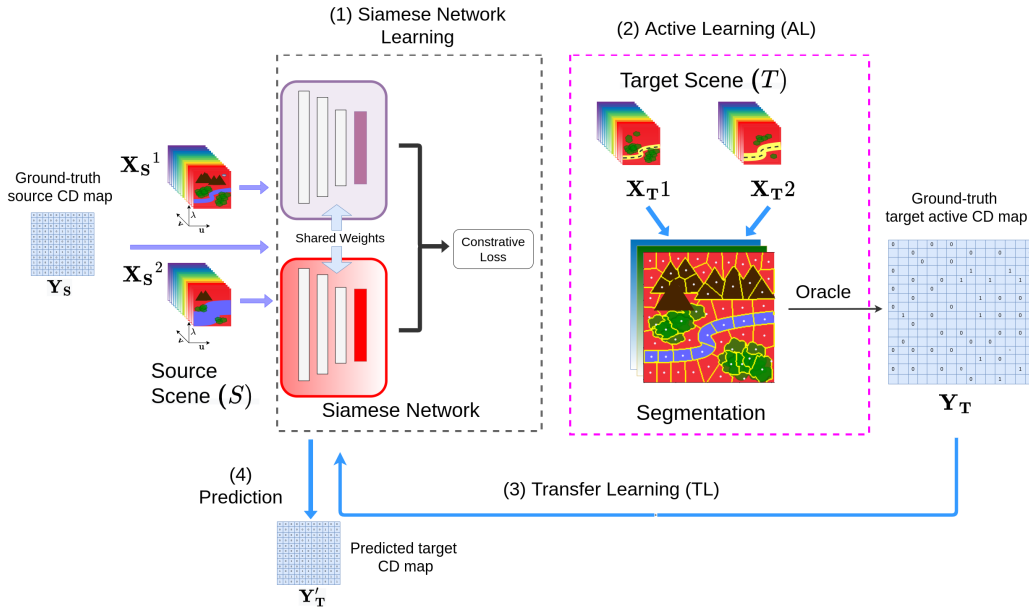


Figure 1: Schema of SENECA: (1) A CD Siamese network is trained from the pair of bi-temporal images  $\mathbf{X}_S^1$  and  $\mathbf{X}_S^2$  of a source scene  $\mathbf{S}$  and the ground truth CD map  $\mathbf{Y}_S$ . (2) A segmentation-based AL strategy is used to select samples of bi-temporal images  $\mathbf{X}_T^1$  and  $\mathbf{X}_T^2$  of a target scene  $\mathbf{T}$  and acquire their CD labels  $\mathbf{Y}_T$ . (3) A fine tuning-based TL strategy is used to update the parameters of the source Siamese network model with the limited samples of the target bi-temporal images  $\mathbf{X}_T^1$  and  $\mathbf{X}_T^2$  for which the CD labels  $\mathbf{Y}_T$  have been acquired through the AL strategy. (4) The fine-tuned target CD Siamese network is used to predict the still unknown labels of the CD map for the target images and build the complete CD map  $\mathbf{Y}'_T$ .

245        acquire the labels of the pixel medoids selected through the segmenta-  
 246        tion step.

247        • The use of a fine tuning-based TL strategy to update the parameters of  
 248        the source CD model with limited MS data of the target, bi-temporal  
 249        MS dataset for which the labels have been acquired through the AL  
 250        strategy (target CD model).

251        • The use of the target CD model, updated with fine tuning, to predict  
 252        the still unknown labels of the CD map of the target, bi-temporal MS  
 253        dataset.

254        A detailed description of the four components is reported in the following.

Table 1: List of used symbols

Symbol	Meaning
$\mathbf{S}$	Source scene of $m_{\mathbf{S}} \times n_{\mathbf{S}}$ pixels
$\mathbf{T}$	Target scene of $m_{\mathbf{T}} \times n_{\mathbf{T}}$ pixels
$\mathbf{D}_{\mathbf{S}}$	Source bi-temporal MS dataset composed of co-registered MS images $\mathbf{X}_{\mathbf{S}}^1$ and $\mathbf{X}_{\mathbf{S}}^2$ of the scene $\mathbf{S}$
$\mathbf{D}_{\mathbf{T}}$	Target bi-temporal MS dataset composed of co-registered MS images $\mathbf{X}_{\mathbf{T}}^1$ and $\mathbf{X}_{\mathbf{T}}^2$ of the scene $\mathbf{T}$
$\mathbf{Y}_{\mathbf{S}}$	Ground-truth source CD map
$\mathbf{A}_{\mathbf{T}}$	Target active scene
$\mathbf{Y}_{\mathbf{T}}$	Ground-truth target active CD map
$\mathbf{Y}_{\mathbf{T}}'$	Predicted target CD map
$f_{\mathbf{S}}(\cdot)$	Embedding learned with the source Siamese network
$f_{\mathbf{T}}(\cdot \dots)$	Embedding learned with the target Siamese network
$\theta$	Otsu’s threshold

255 The list of used symbols is reported in Table 1.

256 *4.1. Source Siamese network*

257 A Siamese network is pre-trained as a source CD deep neural model. In  
 258 particular, the source CD Siamese network is trained by minimising a loss  
 259 function computed on the sample distance of all the  $m_{\mathbf{S}} \times n_{\mathbf{S}}$  bi-temporal  
 260 MS vectors  $(\mathbf{X}_{\mathbf{S}}^1(i, j), \mathbf{X}_{\mathbf{S}}^2(i, j)) \in \mathbf{D}_{\mathbf{S}}$  having labels  $\mathbf{Y}_{\mathbf{S}}(i, j) \in \mathbf{Y}_{\mathbf{S}}$ .

The Siamese network architecture consists of two identical supervised neural networks with shared weights, in order to learn the hidden representation (embedding) of the bi-temporal, MS vectors recorded in  $\mathbf{D}_{\mathbf{S}}$ . The two neural networks are both feed-forward multi-layer perceptrons, and employ error back-propagation during training. They work in parallel comparing the embedding outputs at the end through Euclidean distance. With the aim to learn the model weights, we use the contrastive loss function that was originally proposed by Hadsell et al. (2006) to minimise the Euclidean distance, in the embedding space, between two samples that belong to the same class label and maximises the distance between two samples with different labels. In SENECA, all the bi-temporal MS vectors  $(\mathbf{X}_{\mathbf{S}}^1(i, j), \mathbf{X}_{\mathbf{S}}^2(i, j)) \in \mathbf{D}_{\mathbf{S}}$  that

are labelled with the class  $\mathbf{Y}_{\mathbf{S}}(i, j) = 1$  (*change*) are handled as pairs of samples with different land cover, while all the bi-temporal MS vectors that are labelled with the class 0 (*non-change*) are handled as pairs of samples labelled with the same land cover. Hence, the contrastive loss is defined as follows:

$$\mathcal{L}_c = \sum_{i,j \in \mathbf{S}} ((1 - \mathbf{Y}_{\mathbf{S}}(i, j))d_{\mathbf{S}}(i, j)^2 + \mathbf{Y}_{\mathbf{S}}(i, j) \max(\alpha - d_{\mathbf{S}}(i, j), 0)^2), \quad (1)$$

261 where  $f_{\mathbf{S}}(\cdot)$  is the embedding learned with the source Siamese network,  
 262  $d_{\mathbf{S}}(i, j) = \|f_{\mathbf{S}}(\mathbf{X}_{\mathbf{S}}^1(i, j)) - f_{\mathbf{S}}(\mathbf{X}_{\mathbf{S}}^2(i, j))\|_2$  and  $\alpha$  is the margin. Notice  
 263 that, during this training stage, the desired source embedding  $f_{\mathbf{S}}(\cdot)$  is learned  
 264 achieving that the distance between the bi-temporal MS vectors of the changed  
 265 pixels of  $(\mathbf{D}_{\mathbf{S}}, \mathbf{Y}_{\mathbf{S}})$  get larger than the unchanged pixel distances of  $(\mathbf{D}_{\mathbf{S}}, \mathbf{Y}_{\mathbf{S}})$   
 266 by a margin of  $\alpha$ .

#### 267 4.2. Active learning

268 An AL strategy is used to identify a portion of the target scene  $\mathbf{A}_{\mathbf{T}} \subseteq \mathbf{T}$   
 269 (active target scene) that covers few relevant pixels of  $\mathbf{T}$  (active pixels) for  
 270 which it is suitable to acquire the unknown CD labels associated with the  
 271 bi-temporal MS vectors contained in  $\mathbf{D}_{\mathbf{T}}$ . To this aim, a segmentation algo-  
 272 rithm is used, in order to group together similar adjacent pixels in visually  
 273 meaningful spatial regions – super-pixel objects – that can be used to re-  
 274 duce the number of primitives for the AL analysis. In this study, we use  
 275 the Simple Linear Iterative Clustering algorithm, referred as SLIC (Achanta  
 276 et al. 2012), as segmentation approach. SLIC is inspired by the standard  
 277 k-means clustering algorithm, in order to generate super-pixel object. The  
 278 complexity of SLIC is linear in the number of pixels and independent of the  
 279 number of super-pixels. It adopts a weighted distance measure combines  
 280 colour and spatial proximity while simultaneously providing control over the  
 281 size and compactness of super-pixel objects. In *SENECA*, the segmentation  
 282 is performed to divide the target scene  $\mathbf{T}$  into super-pixel objects and, suc-  
 283 cessively, sample an active pixel to label for each super-pixel object. This  
 284 segmentation step is expected to allow the selection of active pixels for both  
 285 classes (*change*=1 and *non-change*=0) in  $\mathbf{T}$  based on the MS information  
 286 enclosed in  $\mathbf{D}_{\mathbf{T}}$ .

To perform the segmentation step, first the tensor  $\mathbf{X}_{\mathbf{T}} = \mathbf{X}_{\mathbf{T}}^1 \bullet \mathbf{X}_{\mathbf{T}}^2$  is built by applying pixel-wise the concatenation operator  $\bullet$  through the MS

bands of both  $\mathbf{X}_T^1$  and  $\mathbf{X}_T^2$ . More precisely,  $\mathbf{X}_T$  is a tensor of  $m_T \times n_T$  pixels and  $2K$  spectral bands. The spectral dimensionality of  $\mathbf{X}_T$  is then reduced from  $2K$  bands to 2 principal components –  $PC_1$  and  $PC_2$ . This pre-processing step is based on a previous study (Deng et al. 2008) that used the Principal Component Analysis (PCA) as transformation to better highlight the difference between two images. Based upon the theory reported in this previous study, the change can be identified in the second component, while the first component is assumed to be the sum of the common information. Subsequently, SLIC is used to segment  $\mathbf{T}$  into  $\kappa$  super-pixel objects based on the information enclosed in  $\mathbf{X}^T$ . The user-defined parameter  $\kappa$  allows us to control the number of super-pixel objects and, therefore, the number of active exemplars sampled through the super-pixel objects. In particular, for each super-pixel object  $\mathbf{o}$ , the medoid pixel of  $\mathbf{o}$ , i.e., the pixel of  $\mathbf{o}$  that is the closest in space to the centre of  $\mathbf{o}$ , is identified. Formally,

$$\text{medoid}(\mathbf{o}) = \underset{(i,j) \in \mathbf{o}}{\operatorname{argmin}} \left( (i - i_c)^2 + (j - j_c)^2 \right), \quad (2)$$

where  $(i_c, j_c)$  is the centre of  $\mathbf{o}$  having coordinates  $i_c = \frac{\sum_{(i,j) \in \mathbf{o}} i}{|\mathbf{o}|}$  and  $j_c = \frac{\sum_{(i,j) \in \mathbf{o}} j}{|\mathbf{o}|}$ . Finally, the active target scene  $\mathbf{A}_T$  is populated with the medoid pixels of the super-pixel objects:

$$\mathbf{A}_T = \{\text{medoid}(\mathbf{o}) | \mathbf{o} \in \text{SLIC}(\mathbf{X}^T)\}. \quad (3)$$

287 Notice that the active target scene  $\mathbf{A}_T$  defines the AL-based set of rel-  
 288 evant pixel exemplars of  $\mathbf{T}$  whose ground truth CD labels  $\mathbf{Y}_T$  are acquired  
 289 with respect to the bi-temporal MS vectors of  $\mathbf{D}_T$ .  $\mathbf{A}_T$  is used to complete  
 290 the limited supervision of the target Siamese network with the TL strat-  
 291 egy. Henceforth, we rely on  $\mathbf{Y}_T(i, j) = 0/1$  for each  $(i, j) \in \mathbf{A}_T$ , *unknown*  
 292 otherwise.

### 293 4.3. Transfer learning

294 A TL strategy is used to adapt the embedding  $f_S(\cdot)$  pre-trained on  $\mathbf{D}_S$   
 295 to  $\mathbf{D}_T$ . This adaptation is completed using the limited supervision provided  
 296 by the labels acquired in  $\mathbf{Y}_T$  in correspondence of active pixels of  $\mathbf{A}_T$ . In  
 297 particular, the fine tuning strategy is applied. This is an application of the

298 transfer learning principle in deep learning (Tan et al. 2018) that allows  
 299 us to train a deep neural model using limited labelled samples of a target  
 300 distribution. Instead of weights being randomly initialised, they are those  
 301 pre-trained on samples from a different – but related – source distribution. In  
 302 this study, the fine tuning strategy starts with the weights of the pre-trained  
 303 Siamese network that has learned the source embedding  $f_{\mathbf{S}}$ . Subsequently,  
 304 it updates these weights to minimise the contrastive loss formulated in Eq.  
 305 1 and right now evaluated on the bi-temporal MS vectors of  $\mathbf{D}_{\mathbf{T}}$  and the  
 306 labels of  $\mathbf{Y}_{\mathbf{T}}$ , which are associated with active pixels of  $\mathbf{A}_{\mathbf{T}}$ . This allows  
 307 us to adapt the pre-trained Siamese network to new changes in the target  
 308 bi-temporal MS dataset without retraining from scratch with limited class  
 309 estimates only, which would incur in significant overhead and cause artefacts.  
 310 We denote  $f_{\mathbf{T}}(\cdot)$  the target embedding trained with the fine tuning strategy.

#### 311 4.4. Target CD map

Finally,  $f_{\mathbf{T}}(\cdot)$  is used to predict the unknown CD map  $\mathbf{Y}_{\mathbf{T}}'$  associated with  $\mathbf{D}_{\mathbf{T}}$ . For each pixel  $(i, j) \in \mathbf{A}_{\mathbf{T}}$ ,  $\mathbf{Y}_{\mathbf{T}}'(i, j) = \mathbf{Y}_{\mathbf{T}}(i, j)$ , where  $\mathbf{Y}_{\mathbf{T}}(i, j)$  is the CD label acquired in the AL step. For each pixel  $(i, j) \in \mathbf{T} - \mathbf{A}_{\mathbf{T}}$ ,  $\mathbf{Y}_{\mathbf{T}}'(i, j)$  is predicted as follows:

$$\mathbf{Y}_{\mathbf{T}}'(i, j) = \begin{cases} 1 & \|f_{\mathbf{T}}(\mathbf{X}_{\mathbf{T}}^1(i, j)) - f_{\mathbf{T}}(\mathbf{X}_{\mathbf{T}}^2(i, j))\|_2 \geq \theta \\ 0 & \text{otherwise} \end{cases}. \quad (4)$$

In Eq. 4, the threshold  $\theta$  is automatically identified with the Otsu’s algorithm (Otsu 1972). This is an adaptive threshold algorithm that is commonly used in image binarization problems to turn a single intensity threshold that separates pixels into two classes. Using the Otsu’s algorithm, the threshold is determined by minimising the intra-class intensity variance defined as a weighted sum of variances of the two classes<sup>4</sup>. To this aim, we assume that the MS bi-temporal vector distances, computed pixel-wise in  $\mathbf{D}_{\mathbf{T}}$ , are represented in an histogram with  $L$  equal-width bins (levels) denoted as  $[1, \dots, L]$ .

Let  $\eta_i$  be the number of pixels at level  $i$ , so that  $\sum_{i=1}^L \eta_i$  corresponds to the

---

<sup>4</sup>Minimising the intra-class variance is equivalent to maximising the inter-class variance, since the total variance (the sum of the intra-class variance and the inter-class variance) is constant for different partitions.

total number of pixels in the target scene  $T$ , i.e.,  $\sum_{i=1}^L \eta_i = n_{\mathbf{T}} m_{\mathbf{T}}$ . According to this, the probability of each level  $i$  is computed as  $p_i = \frac{\eta_i}{n_{\mathbf{T}} m_{\mathbf{T}}}$ . The Otsu’s algorithm identifies the optimal threshold level  $\theta$ , in order to divide the pixels of the target scene into the class 0 (no-change), spanned over the distance levels  $[1, 2, \dots, \theta]$ , and the class 1 (change), spanned over the distance levels  $[\theta + 1, \dots, L]$ , respectively. The optimal  $\theta$  is chosen with the goal to minimize the intra-class variance that is defined as a weighted sum of variances of the two classes:

$$\theta = \operatorname{argmin}_{1 \leq \theta \leq L} (w_0(\theta) \sigma_1^2(\theta) + w_1(\theta) \sigma_2^2(\theta)), \quad (5)$$

where  $\sigma_1^2(\theta)$  and  $\sigma_2^2(\theta)$  are the variance computed on the two classes separated by  $\theta$ . Finally, the weights  $w_0(\theta)$  and  $w_1(\theta)$  are the probabilities of the two classes, which are computed as follows:

$$w_0(\theta) = \sum_{i=1}^{\theta} p_i \text{ and } w_1(\theta) = \sum_{i=\theta+1}^L p_i. \quad (6)$$

312 Final considerations concern the fact that the predicted CD map can contain  
 313 errors or mistakes. To cope with these issues, we may apply the principle of  
 314 local auto-correlation of objects, according to which detected clusters, com-  
 315 prising changed objects, generally expand across contiguous regions (Appice  
 316 & Malerba 2019). Based on this principle, we may decide to change the as-  
 317 signment of pixels that strongly disagree with surrounding assignments. It  
 318 mainly corresponds to performing a spatial-aware correction of the change  
 319 assignment defined with Otsu’s threshold. This correction, also used in (Ap-  
 320 pice et al. 2020, Andresini et al. 2022), assigns each pixel to the label that  
 321 originally groups the majority of its neighbouring pixels reached within a  
 322 fixed radius, in order to ensure spatial smoothness reducing salt and pepper  
 323 errors.

#### 324 4.5. Time complexity

325 The time complexity of **SENECA** is the sum of the time costs of training a  
 326 Siamese Network, selecting active samples with the segmentation-based AL  
 327 strategy and performing the fine-tuning of the Siamese Network on the active  
 328 samples. The time cost of both training and fine-tuning a Siamese Network  
 329 depends on the cost of training a deep neural network. This mainly de-  
 330 pends on the cost of computing the gradient descent (in the back-propagation

Table 2: Characteristics (acquisition time points, scene size and percentage of changed pixels) of the bi-temporal MS images gathered with Sentinel-2 satellites in five scenes (**Abu Dhabi**, **Beihai**, **Beirut**, **Cupertino** and **Las Vegas**)

Scene	Timestamp 1	Timestamp 2	Scene size	%Change
<b>Abu Dhabi</b>	Jan 20, 2016	Mar 28, 2018	785 × 799	0.037
<b>Beihai</b>	Dec 09, 2016	Mar 09, 2018	772 × 902	0.024
<b>Beirut</b>	Aug 20, 2015	Oct 03, 2017	1070 × 1180	0.026
<b>Cupertino</b>	Sep 08, 2015	Mar 26, 2018	788 × 1015	0.023
<b>Las Vegas</b>	Aug 20, 2015	Feb 05, 2018	824 × 716	0.076

stage), that is,  $\mathbf{O}(lwde)$ , where  $l$  is the number of layers in the network,  $w = \mathbf{O}(r^2)$  is the number of weights per layer,  $r$  is the maximum number of neurons per layer,  $d$  is the number of samples and  $e$  is the number of epochs. The time cost of the segmentation-based AL step mainly depends on the complexity of SLIC that is  $\mathbf{O}(N)$  with  $N$  the number of segmented pixels.

## 5. Experimental Evaluation and Discussion

We evaluated the effectiveness of the CD methodology implemented by **SENECA** on five co-registered, bi-temporal MS images (see Section 5.1) that were acquired with Sentinel-2 satellites in locations picked all over both Asia and USA. The implementation of **SENECA** used in this evaluation is illustrated in Section 5.2. The measured performance metrics are described in Section 5.3, while the results are discussed in Section 5.4.

### 5.1. Imagery data description

We considered five co-registered, bi-temporal MS images<sup>5</sup> with various levels of visible urbanisation. Image were picked over Asia (**Abu Dhabi**, **Beihai** and **Beirut**) and USA (**Cupertino** and **Las Vegas**), respectively (Caye Daudt et al. 2019). Each image was gathered by the Sentinel-2 satellites of the Copernicus program, in 13 spectral bands between visible and short wavelength infrared in the period between 2015 and 2018. All bands

<sup>5</sup><https://rcdaudt.github.io/oscd/>

350 are resampled at a spatial resolution of 10m. The pixel-level change ground  
351 truth was provided for each bi-temporal image with the annotated changes  
352 focused on urban land cover (e.g., new buildings or new roads). In this  
353 study, each bi-temporal MS imagery dataset was used for both learning a  
354 CD Siamese network with labelled data, as well as for fine tuning a pre-  
355 trained CD Siamese network with limited labelled samples. A summary of  
356 the characteristics of the bi-temporal MS images is reported in Table 2.

### 357 *5.2. Implementation details*

358 SENECA was implemented in Python 3.8, using Keras 2.4— a high-level  
359 neural network API with TensorFlow as the backend (Abadi et al. 2015). In  
360 the pre-processing step, the spectral data were scaled in the range  $[0, 1]$  using  
361 the Min-Max normalization <sup>6</sup>.

362 The Siamese network was implemented with two base feed-forward net-  
363 works with shared weights. Each base network is a deep neural network with  
364 three layers with  $256 \times 128 \times 64$  neurons and two dropout layers. The Rectified  
365 Linear Unit (ReLU) activation function was used as activation to each hid-  
366 den layer and the contrastive function (Hadsell et al. 2006) was used as loss  
367 function. In the supervised initialization step, the weights were initialised  
368 following the Xavier scheme, while, in the fine tuning step, the weights saved  
369 from the previous network were used as a starting point. For each dataset,  
370 we optimized the hyper-parameter using the tree-structured Parzen estima-  
371 tor algorithm as implemented in the Hyperopt library (Bergstra et al. 2013).  
372 This hyper-parameter optimization was performed by using 20% of the entire  
373 training set as a validation set according to the Pareto Principle. We selected  
374 the hyper-parameter configuration that achieved the lowest validation loss.  
375 The hyper-parameters and their corresponding possible values are reported in  
376 Table 3. We trained the network with mini-batches using back-propagation,  
377 and the gradient-based optimization was performed using the Adam update  
378 rule (Kingma & Ba 2014).

379 For the AL strategy, we performed the segmentation step using the SLIC  
380 algorithm as implemented in Scikit-image library <sup>7</sup>. The number  $\kappa$  of seg-  
381 ments to detect in the target scene through SLIC was set as a percentage

---

<sup>6</sup><https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.OneHotEncoder.html>

<sup>7</sup><https://scikit-image.org/docs/dev/api/skimage.segmentation.html>



Table 3: Hyperparameter search space for the Siamese model.

Hyper-parameter	Values
batch size	$\{2^5, 2^6, 2^7, 2^8, 2^9\}$
learning rate	$[0.0001, 0.01]$
dropout	$[0, 0.5]$

382  $\kappa\%$  of the target scene size, where  $\kappa\%$  is a user-defined parameter. By de-  
 383 fault  $\kappa\% = 1\%$ . Since SLIC algorithm processes RGB images, we scaled the  
 384 two principal components extracted from the bi-temporal target dataset for  
 385 the segmentation step in a range 0-255. In addition, we added a dummy  
 386 component, set equal to 0, in order to create the third channel of the RGB  
 387 representation.

388 Threshold  $\theta$  used to predict the target CD map was estimated using the  
 389 implementation of Otsu’s algorithm from scikit-image library <sup>8</sup>. Finally, the  
 390 radius of the kernel used for spatial-aware correction  $R$  was set equal to 5 for  
 391 all the target scenes.

### 392 5.3. Performance metrics

393 In this Section we introduce the metrics measured to evaluate the accu-  
 394 racy of the predicted CD maps, the homogeneity of super-pixel segmentation  
 395 and the efficiency of the learning process.

396 We measured the accuracy of the predicted CD maps with the Fscore  
 397 (F1)(Tan et al. 2005), the Area Under the ROC curve (AUCROC)(Tan et al.  
 398 2005) and the Geometric mean (G-mean) (Kubat & Matwin 1997). These  
 399 metrics are commonly considered in the remote sensing field for the eval-  
 400 uation of CD methods. Let us consider:  $tp$  – the number of pixels of the  
 401 scene with the class *change* that are correctly predicted as belonging to that  
 402 class type;  $fp$  – the number of pixels not belonging to the class *change* that  
 403 are wrongly predicted as belonging to the class *change*;  $tn$  – the number  
 404 of pixels not belonging to class *change* that are predicted as not belong-  
 405 ing to class *change*;  $fn$  – the number of pixels of the class *change* that are  
 406 wrongly predicted as not belonging to that class type;  $n$  is the total num-  
 407 ber of pixels in the scene. The F1 measures the harmonic mean of precision

---

<sup>8</sup>[https://scikit-image.org/docs/dev/api/skimage.filters.html#skimage.filters.threshold\\_otsu](https://scikit-image.org/docs/dev/api/skimage.filters.html#skimage.filters.threshold_otsu)

408 and recall, i.e.,  $Fscore = 2 \frac{precision \times recall}{precision + recall}$ . The higher the F1, the better the  
409 balance between precision and recall achieved by the evaluated method. In  
410 particular, the precision measures how many pixels are correctly classified  
411 for the class *change*, given all predictions of that class type in the scene,  
412 i.e.,  $precision = \frac{tp}{tp+fp}$ . The recall measures how many pixels are correctly  
413 predicted for the class *change* given all occurrences of that class type in the  
414 scene, i.e.,  $recall = \frac{tp}{tp+fn}$ . The AUCROC measures the Area Under the ROC  
415 curve as it was defined with the False Positive Rate (FPR) on the x-axis and  
416 the True Positive Rate (TPR) on the y-axis. The FPR measures how many  
417 pixels are wrongly classified in the class *change* given all the occurrences of  
418 negative samples of that class type, i.e.,  $FPR = \frac{fp}{fp+tn}$ . The TPR measures  
419 how many pixels are correctly predicted for the class *change* given all occur-  
420 rences of that class, i.e.,  $TPR = \frac{tp}{tp+fn}$ . Hence, the AUCROC value expresses  
421 the probability that a given method will rank a positive sample of the class  
422 *change* higher than a negative sample of the considered class. The G-mean  
423 measures the geometric mean of *specificity* and *recall* by equally consider-  
424 ing the errors on both classes, i.e.,  $G - mean = \sqrt{specificity \times recall}$ . In  
425 particular, the *specificity* measures how many pixels are correctly predicted  
426 for the class *unchange* given all occurrences of that class in the scene, i.e.,  
427  $specificity = \frac{tn}{tn+fp}$ .

428 In addition, we measured the homogeneity of the super-pixel objects with  
429 the Purity and F1. The Purity is a simple evaluation criterion of cluster  
430 quality. To compute Purity, each super-pixel object  $\mathbf{o}_i$  identified through the  
431 segmentation step is assigned to the class  $c_j$  (*change* vs *non-change*) that is  
432 the most frequent in the object. The accuracy of this assignment is measured  
433 by counting the number of correctly assigned target pixels and dividing by the

434 total number of target pixels  $m_{\mathbf{T}} \times n_{\mathbf{T}}$ , i.e.,  $Purity = \frac{1}{m_{\mathbf{T}} \times n_{\mathbf{T}}} \sum_{i=1}^{\kappa} \max_j |\mathbf{o}_i \cap c_j|$ ,

435 where  $\kappa$  is the number of super-pixel objects detected in the segmentation  
436 step, while  $|\cdot|$  denotes the cardinality operator. Similarly, the F1 of the  
437 segmentation output is measured by assuming the most frequent CD class  
438 observed in a super-pixel object as the CD class assigned by the segmentation  
439 to each pixel grouped in the super-pixel object. We measured F1 per class  
440 considering firstly the *change* class (F1 (*change*)) and, successively, the *non-*  
441 *change* class (F1 (*non-change*)). These two scores allow us to monitor the  
442 ability of the segmentation step of depicting super-pixel objects covering  
443 both homogeneous changed regions and homogeneous non-changed regions,

444 respectively.

445 Finally, we evaluate the time performance (**TIME**) spent both learning the  
446 pre-trained CD model from the source scene and fine tuning a pre-trained  
447 CD model to the target scene. They were collected on a Linux machine with  
448 an Intel(R) Core(TM) i9-10900K CPU @ 3.70GHz and 64GB RAM. All the  
449 experiments are executed on a single GeForce RTX 3070. In this study, the  
450 training **TIME** was measured in minutes.

#### 451 5.4. Results

452 The empirical validation was done with the Siamese network as CD model,  
453 in order to answer the following questions:

454 Q1 To what extent the number of active pixels selected in the target scene  
455 by the AL labelling strategy has an effect on the performance of the  
456 CD model adapted with TL? (Sensitivity analysis in Section 5.4.1)

457 Q2 Is the CD model adapted to a target scene with the proposed ATL  
458 strategy more powerful in labelling the target scene than the CD model  
459 pre-trained on the source scene? (Ablation study in Section 5.4.2)

460 Q3 How does the performance of a CD model pre-trained on a source scene  
461 and adapted to a target scene through the proposed ATL strategy  
462 change with either the source scene or the target scene? (Source/target  
463 scene study in Section 5.4.3)

464 Q4 Does the defined CD method outperform recent, state-of-the-art CD  
465 systems? (Competitor study in Section 5.4.4)

466 Experiments were performed by considering 20 configurations of source-  
467 target scenes. More precisely, for each target scene we generated four config-  
468 urations by varying the source scene among the left-out scenes. For example,  
469 given the target scene **Abu Dhabi**, four configurations were generated by  
470 selecting the source scene among: **Beihai**, **Beirut**, **Cupertino** and **Las**  
471 **Vegas**, respectively.

##### 472 5.4.1. Sensitivity analysis (Q1)

473 The sensitivity analysis was performed, in order to assess the influence  
474 of  $\kappa$ , i.e., the number of active pixels selected through the segmentation  
475 step on the behaviour of **SENECA**. As in the implementation of **SENECA**,  
476  $\kappa = \kappa\% \times n_{\mathbf{T}} \times m_{\mathbf{T}}$ , we analysed the performance of **SENECA** by varying  $\kappa\%$

Table 4: F1, AUCROC, G-mean and TIME (in mins) of SENECA by varying  $\kappa\%$  among = 0.1%, 1% and 5%. We report the mean  $\pm$  standard deviation of performances measured on all the target scenes with every CD model pre-trained with each left-out source scene.

$\kappa\%$	F1	AUCROC	G-mean	TIME
0.1%	0.40 ( $\pm 0.20$ )	0.73 ( $\pm 0.09$ )	0.68 ( $\pm 0.10$ )	13.94 ( $\pm 3.35$ )
1%	0.53 ( $\pm 0.18$ )	0.76 ( $\pm 0.05$ )	0.73 ( $\pm 0.07$ )	92.92 ( $\pm 48.97$ )
5%	0.57 ( $\pm 0.17$ )	0.79 ( $\pm 0.06$ )	0.76 ( $\pm 0.09$ )	61040.86 ( $\pm 2209.95$ )

477 among 0.1%, 1% and 5%. The mean and standard deviation of F1, AUCROC,  
 478 G-mean and TIME measured for SENECA in all the tested configurations are  
 479 reported in Table 4. Figure 2 reports the F1 computed for each target scene  
 480 by varying both the source scene and  $\kappa\%$ . These results show that the higher  
 481 the value of  $\kappa\%$  (and, consequently, the higher the number  $\kappa$  of active pixels),  
 482 the higher the accuracy of SENECA. On the other hand, this gain in accuracy  
 483 is at the cost of the extra time spent completing the learning stage, as well  
 484 as the higher effort and cost spent by experts acquiring the ground truth CD  
 485 labels for the active pixels.

486 Additional conclusions can be drawn by analysing the homogeneity of  
 487 super-pixel objects extracted through the segmentation step and considered  
 488 to sample the active pixels of each scene. Figure 3 shows the segmentation’s  
 489 output of each considered scene as it was detected with  $\kappa\% = 1\%$ . Figure 4  
 490 reports the Purity, F1 for the class *change* and F1 for the class *non-change* as  
 491 they were measured on the output of the segmentation step by varying  $\kappa\%$   
 492 among 0.1%, 1% and 5%. These results reveal that the higher the value of  
 493  $\kappa\%$ , the finer-grained the segmentation of each scene in super-pixel objects  
 494 and the higher the homogeneity of CD labels grouped in each super-pixel  
 495 object. Detecting finer-grained super-pixel objects allows us to better depict  
 496 homogeneous segments that mainly contain either changed pixels or non-  
 497 changed pixels. Notably, the gain in the homogeneity of super-pixel objects  
 498 is greater with respect to the class *change* than with respect to the class  
 499 *non-change*.

500 In general, we note that  $\kappa\% = 1\%$  allows us to achieve a good trade-  
 501 off among homogeneity of segmentation, accuracy of final CD predictions,  
 502 computation time spent completing the learning process, as well as effort

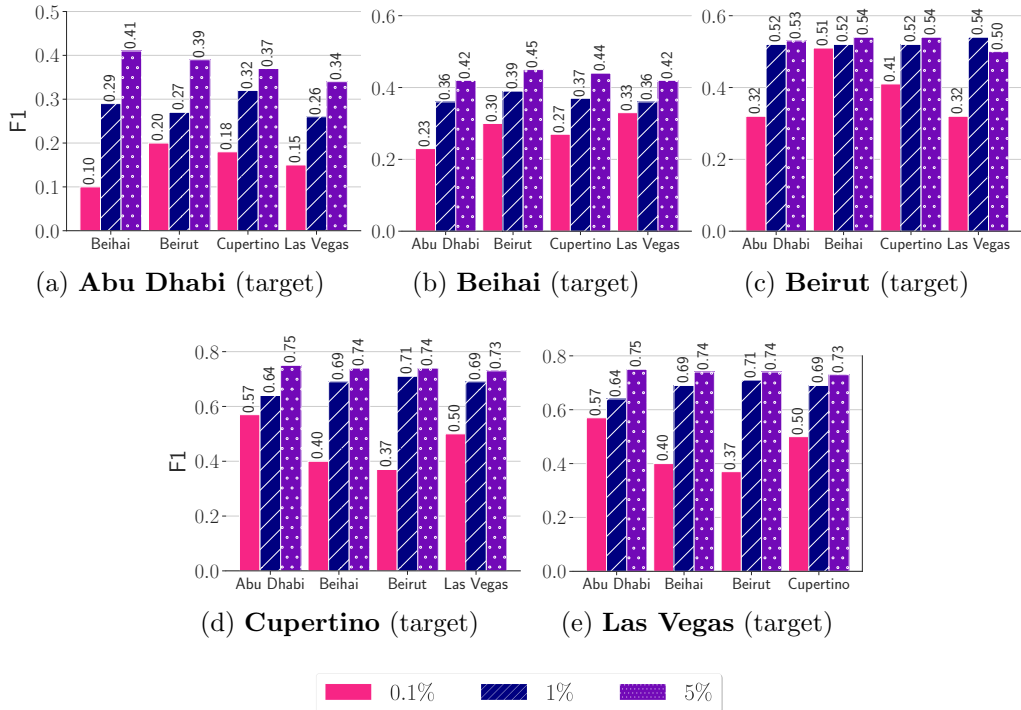


Figure 2: F1 of SENECA by varying  $\kappa\%$  among 0.1%, 1% and 5%. For each target scene (Figures 2a-2e), we compare the F1 of the CD maps predicted by SENECA by varying the source scene.

503 and human cost spent acquiring CD labels. Due to these reasons, we report  
 504 results achieved with  $\kappa\% = 1\%$  in the rest of the experimental evaluation.

#### 505 5.4.2. Ablation analysis (Q2)

506 The ablation analysis of SENECA was conducted, in order to explore how  
 507 the ATL strategy can impact the performance of the CD model pre-trained on  
 508 a source scene by adapting it to each left-out target scene. To this purpose, we  
 509 ran the ATL strategy of SENECA with  $\kappa\% = 1\%$  and measured the accuracy  
 510 of the changes detected in each target scene by varying the source scene  
 511 considered to learn the pre-trained CD model. For the ablation study, we  
 512 also report the performance of Siamese that is the configuration that discards  
 513 the ATL strategy. Specifically, Siamese used the CD model pre-trained on  
 514 a source scene to detect changes of a target scene without performing any  
 515 adaptation of the pre-trained CD model. The mean and standard deviations  
 516 of F1, AUCROC, G-mean and TIME of both SENECA and Siamese are reported

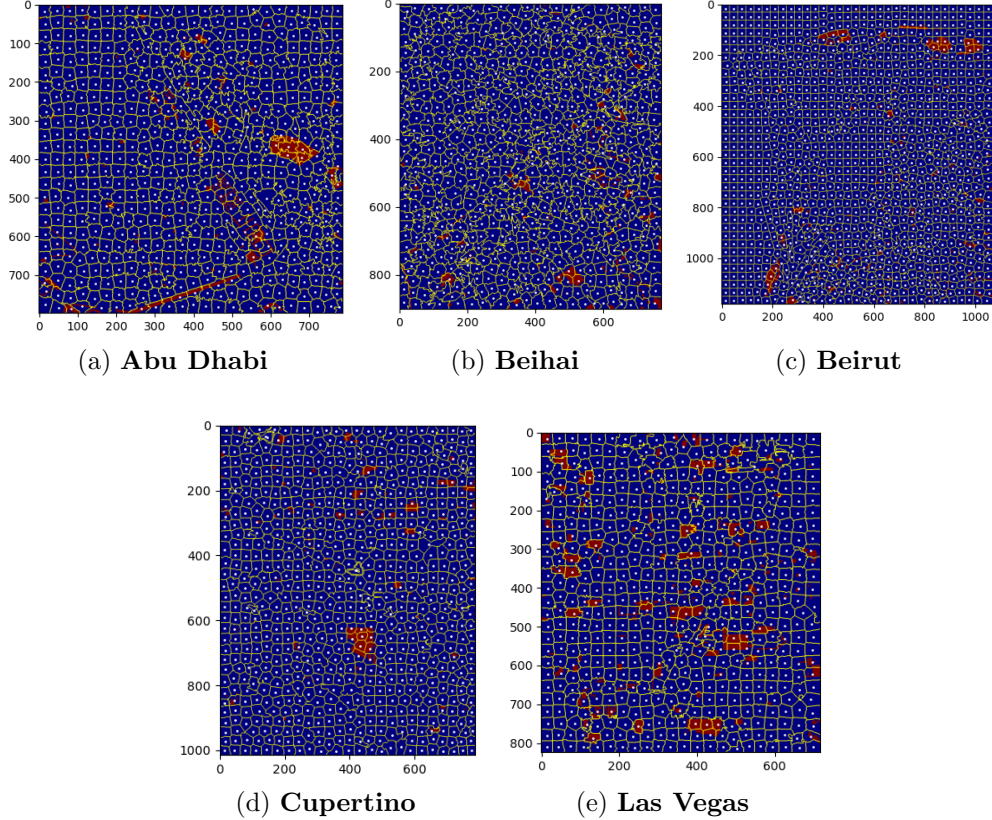


Figure 3: Super-pixel objects detected with the segmentation step performed with  $\kappa\% = 0.1\%$ . The red areas denote the changed regions, while the blue areas denote the unchanged regions in the corresponding scenes. White circles denote the active pixels sampled throughout the segmentation step.

517 in Table 5. Figure 5 reports the F1 scores computed per each target scene by  
 518 varying the source scene. These results show that the use of the ATL strategy  
 519 allows **SENECA** to gain accuracy compared to the baseline **Siamese**. Notably,  
 520 this conclusion can be drawn independently of the source scene considered to  
 521 train the source CD model. As expected, the higher accuracy of **SENECA** is at  
 522 the cost of the more computation time spent performing the proposed ATL  
 523 strategy. Figure 6 shows the computation time spent completing the four  
 524 learning steps of **SENECA** in all the performed experiments. These results  
 525 reveal that **SENECA** spends the most of its computation time performing the  
 526 segmentation step in the AL component, while the time spent performing

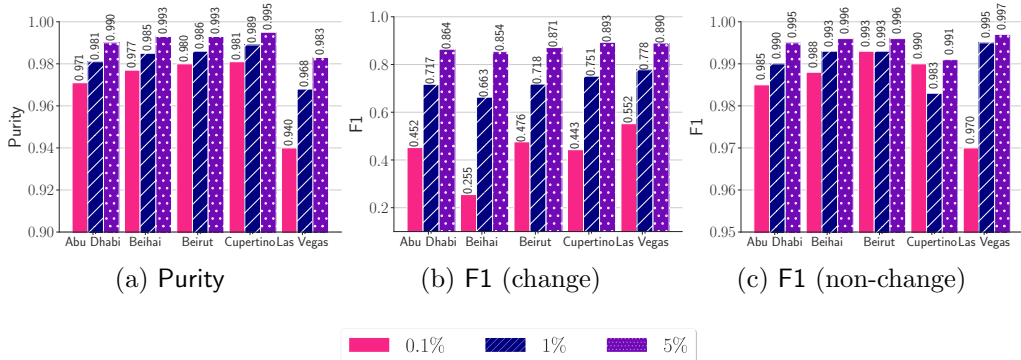


Figure 4: Purity, F1 of the class *change* and F1 of the class *non-change* measured on the output of the segmentation step performed by varying  $\kappa\%$  among 0.1%, 1% and 5%

Table 5: F1, AUCROC, G-mean and TIME (in mins) of SENECA with  $\kappa\% = 1\%$  and its baseline configuration Siamese. We report the mean  $\pm$  standard deviation of performances measured on all the target scenes with every CD model pre-trained with each left-out source scene.

Method	F1	AUCROC	G-mean	TIME
SENECA	0.53 ( $\pm 0.18$ )	0.76 ( $\pm 0.05$ )	0.73 ( $\pm 0.07$ )	92.92 ( $\pm 48.97$ )
Siamese	0.18 ( $\pm 0.09$ )	0.65 ( $\pm 0.09$ )	0.61 ( $\pm 0.11$ )	10.32 ( $\pm 2.39$ )

527 fine tuning in the TL component is generally small.

### 528 5.4.3. Source and target scenes (Q3)

529 This analysis was conducted to explore the effect of a specific source/target  
530 scene on the accuracy of SENECA. Figure 7 shows the F1 of SENECA by  
531 varying both the source scene and the target scene. Results show that the  
532 accuracy performance of SENECA changes significantly with the target scene.  
533 However, the differences in the F1 of SENECA are commonly negligible in each  
534 target scene by varying the source scene. The only exception is observed with  
535 the target scene **Beihai** where the F1 varies from 0.36 (with the source scene  
536 **Abu Dhabi**) to 0.54 (with the source scene **Las Vegas**). Interestingly, also  
537 the source CD Siamese network pre-trained on **Las Vegas** outperforms the  
538 source CD Siamese network pre-trained on **Abu Dhabi**, **Beirut** and **Cu-**  
539 **pertino** when they were used to predict the CD map of **Beihai** without the

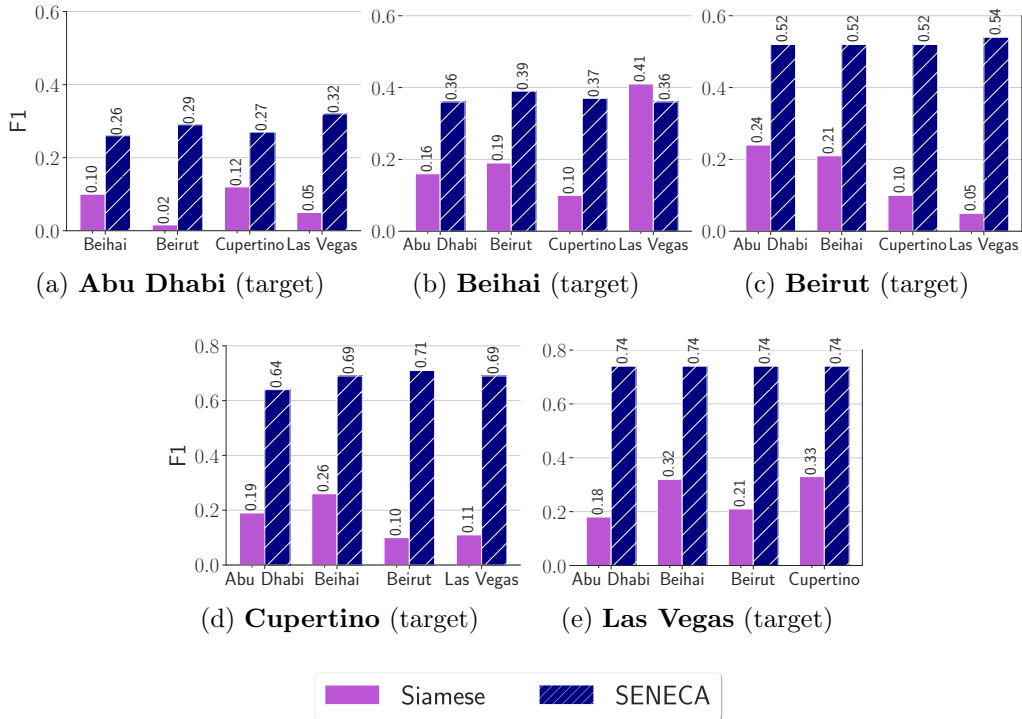


Figure 5: F1 of SENECA with  $\kappa = 1\%$  and its baseline configuration Siamese. For each target scene (Figures 5a-5e), we compare the F1 of the CD maps predicted by both SENECA and Siamese by varying the source scene.

540 ATL strategy (see result of Siamese in Figure 5b). This suggests that a future  
 541 research direction may focus on exploring which properties of the pre-trained  
 542 CD models may foster the better performance of the ATL strategy.

#### 543 5.4.4. Competitor analysis (Q4)

544 The comparative analysis is performed to assess the significance of accu-  
 545 racy and novelty of SENECA compared to several related methods, selected  
 546 from the state of the art in CD literature. Table 6 reports a summary of  
 547 the main characteristics of the considered competitors. We point out that  
 548 the competitors that integrate the Siamese network (Shi, Liu, Li, Liu, Wang  
 549 & Zhang 2022) and the ATL strategy (Ruzicka et al. 2020) are the closest  
 550 to SENECA. Specifically the method CBAM described in (Shi, Liu, Li, Liu,  
 551 Wang & Zhang 2022) introduces a convolutional attention block module in  
 552 the Siamese network, but neglects any TL mechanism to adapt a CD model



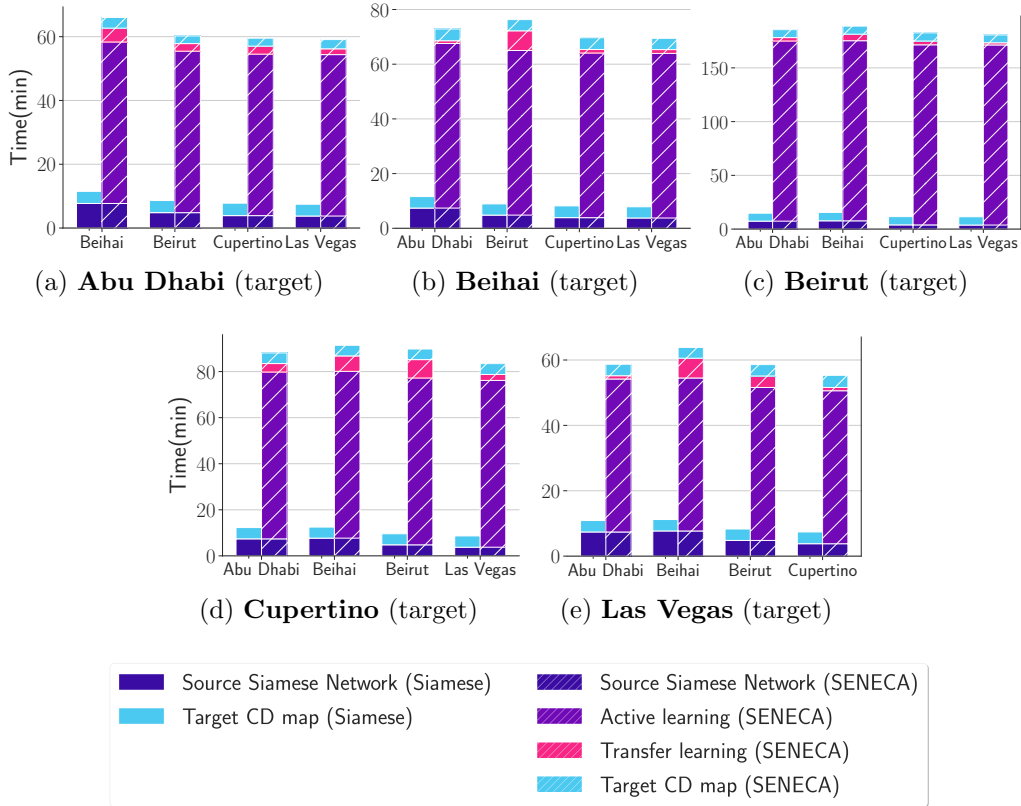


Figure 6: TIME of SENECA with  $\kappa = 1\%$  and its baseline configuration Siamese. For each target scene (Figures 6a-6e), we compare F1 of the CD map predicted by both SENECA and Siamese by varying the source scene.

553 trained in a source scene to a new target scene. The method SiameseU-Net  
 554 described in (Ruzicka et al. 2020) trains a Siamese network with ResNet-34  
 555 base networks from a target source and uses an AL strategy to fine tune  
 556 a source CD model to a target domain. In particular, it uses an ensemble  
 557 procedure to select the tiles of pixels for the active labelling. It extends  
 558 the source training set with the selected target active samples and re-trains  
 559 the Siamese network from scratch using the augmented training set. In this  
 560 comparative study, we experimented the AL strategy of both SiameseU-Net  
 561 and SENECA to acquire the labels of the 1% of target samples. The meth-  
 562 ods BIC<sup>2</sup>, ORCHESTRA, PCAK-Means and CVA perform an unsupervised  
 563 learning stage on the target scene by neglecting any information enclosed in  
 564 the source scene. Finally, the method CBAM performs a supervised learning

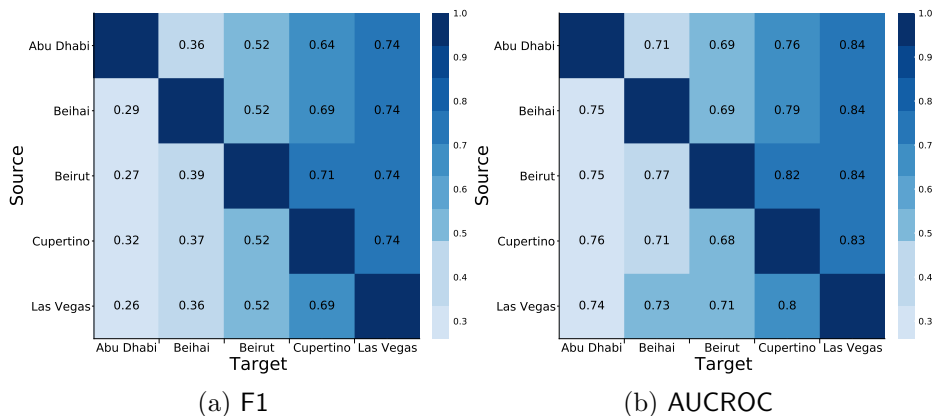


Figure 7: F1 of SENECA by varying both the target scene (axis X) and the source scene (axis Y)

Table 6: Compared algorithm description

Algorithm	Description
SENECA	Siamese network, ATL, Otsu’s method
BIC <sup>2</sup> (Appice et al. 2020)	GMM, PCA, Random Forest
ORCHESTRA(Andresini et al. 2022)	Autoencoder, CVA, spectral angle distance, Otsu’s method
CBAM(Shi, Liu, Li, Liu, Wang & Zhang 2022)	Siamese network, ResNet18, Attention
SiameseU-Net(Ruzicka et al. 2020)	Siamese network, ResNet34, Active learning
PCAK-means(Celik 2009)	PCA, k-Means
CVA(López-Fandiño et al. 2019)	CVA, spectral angle distance, Otsu’s algorithm

565 stage on the source scene and uses this pre-trained CD model on the target  
 566 scene.

567 All the related methods were run using default parameters suggested  
 568 by the authors in the reference papers. In particular, BIC<sup>2</sup> was run with  
 569 the number of trees in the Random Forest set equal to 20, the number of  
 570 principal components set equal to 20, the threshold considered to select the  
 571 samples to train the Random Forests set equal to 0.85. Random Forests  
 572 were constructed with the number of random independent features to look  
 573 for the best split set equal to  $\sqrt{\#independentfeatures}$ , the bootstrap op-  
 574 tion was enabled with the bootstrap size set equal to the size of the training  
 575 set and the function to measure the quality of a split, was set equal to the  
 576 Gini index. The GMM was run with number of components set equal to  
 577 2, the covariance type set equal to diagonal (i.e. each component had its  
 578 own diagonal covariance matrix), the non-negative regularisation added to

579 the diagonal of covariance set equal to 0.00001. The image difference was  
580 computed with both the spectral angle distance and spectra-spatial cross  
581 correlation-based distance and the best results were considered for this com-  
582 parative study.<sup>9</sup> ORCHESTRA was run with the autoencoder architecture  
583 composed of 3 fully-connected (FC) layers of  $8 \times 4 \times 8$  neurons as proposed  
584 by the authors to process Sentinel-2 images. The learning rate and batch  
585 size were optimised with the tree-structured Parzen estimator in the range  
586  $[0.00001, 0.01]$  and the set  $\{32, 64, 128, 256, 512\}$ , respectively. The optimi-  
587 sation was done using 20% of the entire training set as a validation set. The  
588 dropout layer was used to prevent overfitting. The mean squared error was  
589 used as the loss function. The ReLu was selected as the activation function  
590 for each hidden layer, while *Linear* activation function was used for the last  
591 layer. The number of epochs was set equal to 150, retaining the best mod-  
592 els achieving the lowest loss on the validation set. CBAM was run with a  
593 Siamese Network implementing a ResNet18 (He et al. 2016) pre-trained on  
594 ImageNet. The ResNet18 was implemented with four basic blocks of depths  
595 equal to 64, 128, 256, and 512, respectively. Each basic block was com-  
596 posed by two convolutional layers with a kernel size of  $3 \times 3$  and two batch  
597 normalization layers. The model was fine-tuned for 150 epochs using Adam  
598 optimizer, ReLu activation was selected for each hidden layer. SiameseU-Net  
599 was run with a Siamese Network composed of two autoencoders with shared  
600 weights. The architecture of the encoder was implemented with a ResNet34  
601 (He et al. 2016) pre-trained on ImageNet with a kernel size of  $3 \times 3$ . The  
602 model was fine-tuned for 100 epochs using Adam optimizer and sigmoid as  
603 activation function for each hidden layer. The number of models used for the  
604 ensemble-based AL strategy was set equal to 5. PCAK-Means was run with  
605 the number of eigenvector equal to 3 and the block size equal to 4. CVA  
606 was run with the number of levels set equal to 256 in the Otsu’s algorithm.

607 The mean and standard deviations of F1, AUCROC, G-mean and TIME  
608 of both SENECA and related methods are reported in Table 7. These re-  
609 sults show that SENECA is able to outperform all the related methods in this  
610 study in terms of F1. On the other hand, PCAK-Means outperforms SENECA  
611 in terms of AUCROC and G-mean, where SENECA is the runner-up. This  
612 is a consequence of the fact that PCAK-Means discovers a higher number of

---

<sup>9</sup>The spectra-spatial cross correlation-based distance outperformed spectral angle distance in all scenes with the exception of **Las Vegas**.

Table 7: F1, AUCROC, G-mean and TIME (in mins) of SENECA with  $\kappa\% = 1\%$ , as well as the related methods. We report the mean  $\pm$  standard deviation of performances measured on all the target scenes with every CD model pre-trained with each left-out source scene.

Method	F1	AUCROC	G-mean	TIME
SENECA	0.53 ( $\pm 0.18$ )	0.76 ( $\pm 0.05$ )	0.73 ( $\pm 0.07$ )	92.92 ( $\pm 48.97$ )
BIC <sup>2</sup>	0.40 ( $\pm 0.29$ )	0.70 ( $\pm 0.14$ )	0.63 ( $\pm 0.17$ )	11.88 ( $\pm 5.36$ )
ORCHESTRA	0.23 ( $\pm 0.20$ )	0.66 ( $\pm 0.14$ )	0.65 ( $\pm 0.18$ )	57.94 ( $\pm 24.17$ )
CBAM	0.06 ( $\pm 0.04$ )	0.50 ( $\pm 0.07$ )	0.18 ( $\pm 0.13$ )	16.58 ( $\pm 3.46$ )
SiameseU-Net	0.33 ( $\pm 0.19$ )	0.65 ( $\pm 0.08$ )	0.65 ( $\pm 0.18$ )	114.04 ( $\pm 0.39$ )
PCAK-Means	0.40 ( $\pm 0.21$ )	0.79 ( $\pm 0.04$ )	0.78 ( $\pm 0.04$ )	1.22 ( $\pm 0.07$ )
CVA	0.24 ( $\pm 0.20$ )	0.70 ( $\pm 0.18$ )	0.64 ( $\pm 0.17$ )	1.40 ( $\pm 0.04$ )

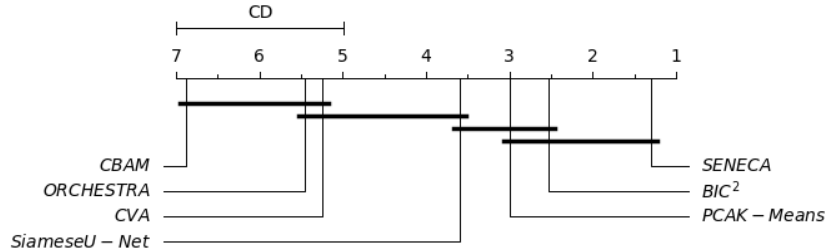


Figure 8: Nemenyi test of F1 of SENECA and related methods. Groups of methods that are not significantly different (at  $p \leq 0.05$ ) are connected.

613 *change* samples (and consequently a lower number of *non - change*  
614 samples) than SENECA. Hence, PCAK-Means performs a higher number of true  
615 positive samples, but also a higher number of false positive samples than  
616 SENECA. Therefore, SENECA outperforms PCAK-Means in terms of preci-  
617 sion ( $0.58 \pm 0.29$  in SENECA vs  $0.34 \pm 0.26$  in PCAK-Means) and specificity  
618 ( $0.97 \pm 0.04$  in SENECA vs  $0.91 \pm 0.08$  in PCAK-Means), while PCAK-Means  
619 outperforms SENECA in terms of recall ( $0.55 \pm 0.11$  in SENECA vs  $0.67 \pm 0.06$   
620 in PCAK-Means). The impact of recall is higher in the formulation of G-mean  
621 and AUCROC than in the formulation of F1. This motivates differences in  
622 the observed performances of the compared methods with respect to F1, AU-  
623 CROC and G-mean. In any case, a high number of false alarms (false positive)  
624 is not a desirable behaviour in imbalance classification problems such as CD

Table 8: F1 of SENECA ( $\kappa\% = 1\%$ ), as well as the related methods. The best results are in bold.

Source	Target	SENECA	BIC <sup>2</sup>	ORCHESTRA	CBAM	SiameseU-Net	PCAK-Means	CVA
Beihai	Abu Dhabi	<b>0.28</b>	0.19	0.15	0.07	0.24	0.19	0.17
Beirut		<b>0.27</b>	0.19	0.15	0.07	0.24	0.19	0.17
Cupertino		<b>0.31</b>	0.19	0.15	0.07	0.24	0.19	0.17
Las Vegas		<b>0.26</b>	0.19	0.15	0.07	0.24	0.19	0.17
Abu Dhabi	Beihai	<b>0.36</b>	0.09	0.05	0.09	0.34	0.41	0.05
Beirut		<b>0.39</b>	0.09	0.05	0.04	0.34	0.41	0.05
Cupertino		<b>0.36</b>	0.09	0.05	0.04	0.34	0.41	0.05
Las Vegas		<b>0.36</b>	0.09	0.05	0.04	0.34	0.41	0.05
Abu Dhabi	Beirut	<b>0.52</b>	0.35	0.06	0.05	0.08	0.20	0.06
Beihai		<b>0.51</b>	0.35	0.06	0.05	0.08	0.20	0.06
Cupertino		<b>0.51</b>	0.35	0.06	0.05	0.08	0.20	0.06
Las Vegas		<b>0.53</b>	0.35	0.06	0.05	0.08	0.20	0.06
Abu Dhabi	Cupertino	0.64	<b>0.68</b>	0.44	0.00	0.60	0.55	0.43
Beihai		<b>0.68</b>	<b>0.68</b>	0.44	0.04	0.60	0.55	0.43
Beirut		<b>0.70</b>	0.68	0.44	0.04	0.60	0.55	0.43
Las Vegas		<b>0.68</b>	<b>0.68</b>	0.44	0.03	0.60	0.55	0.43
Abu Dhabi	Las Vegas	<b>0.74</b>	0.72	0.46	0.00	0.41	0.67	0.45
Beihai		<b>0.73</b>	0.72	0.46	0.14	0.41	0.67	0.45
Beirut		<b>0.74</b>	0.72	0.46	0.15	0.41	0.67	0.45
Cupertino		<b>0.73</b>	0.72	0.46	0.04	0.41	0.67	0.45

625 tasks.

626 Further considerations concern the analysis of TIME. The DL-based meth-  
627 ods (i.e., SENECA, ORCHESTRA, CBAM and SiameseU-Net) spent more time  
628 than the remaining methods (BIC<sup>2</sup>, PCAK-Means and CVA). In any case, both  
629 SiameseU-Net and SENECA are the most time-consuming methods. Both  
630 methods train a Siamese network and integrate an AL-based strategy. How-  
631 ever, SENECA uses a segmentation-based AL strategy, while SiameseU-Net  
632 uses an ensemble-based AL strategy.

633 We proceed this comparative study by examining in depth the F1 re-  
634 sults per each scene. Results reported in Table 8 show that SENECA (with  
635  $\kappa\% = 1\%$ ) outperforms all the competitors of this study except for BIC<sup>2</sup>  
636 in the configuration with source **Abu Dhabi** and target **Cupertino**. However,  
637 SENECA outperforms (or performs equals to) BIC<sup>2</sup> on the target **Cupertino**  
638 when the source is **Beihai**, **Beirut** or **Las Vegas**. In addition, the high-  
639 est accuracy on the target **Beirut** is achieved by SENECA with the source  
640 **Beihai**. Finally, we ranked the compared methods by statistically testing  
641 whether the improvement of F1 of the computed CD maps is significant over  
642 the various experimental configurations. To this aim, we have used Fried-

643 man’s test (Demšar 2006). This is a non-parametric test that is commonly  
644 used to compare multiple methods over multiple experiments. It compares  
645 the average ranks of the methods, so that the best performing method gets  
646 the rank of 1. The second best gets rank 2. The null-hypothesis states that  
647 all the methods are equivalent. Under this hypothesis, the ranks of compared  
648 methods should be equal. In this study, we rejected the null hypothesis with  
649  $p\text{-value} \leq 0.05$ . As the null-hypothesis was rejected, that is, no method was  
650 singled out, we used a post-hoc test—the Nemenyi test—for pairwise compar-  
651 isons (Demšar 2006). The results of this test reported in Figure 8 shows  
652 that **SENECA** enables the production of the CD map that commonly achieve  
653 the highest F1 by having  $\text{BIC}^2$  as runner-up.

## 654 **6. Conclusions**

655 In this paper we have presented **SENECA**: an ATL methodology for CD  
656 in co-registered, bi-temporal MS images acquired with Sentinel-2 satellites in  
657 the same Earth’s scene, at different time points. The proposed methodology  
658 uses the TL strategy to adapt the Siamese network pre-trained from a source  
659 domain to a related target domain. The adaptation is performed with the  
660 limited supervision made available with the AL strategy. An experimental  
661 study was performed to show the effectiveness of the proposed CD method-  
662 ology, quantified in terms of CD accuracy. In particular, the results obtained  
663 have underlined that **SENECA** is able to produce decisions that outperform  
664 decisions produced with the baseline **Siamese** that is the configuration that  
665 discards the proposed ATL strategy. Furthermore, the experimental results  
666 clearly highlighted that **SENECA** achieves high quality performance with a  
667 limited amount of labels acquired via the AL process no matter the source  
668 data considered to learn the pre-trained Siamese network. Finally, the pro-  
669 posed ATL framework helps us to gain accuracy compared to various CD  
670 methods presented in the recent literature.

671 One limitation of the proposed methodology is the absence of any expla-  
672 nation mechanism. A future research direction could be devoted to explore  
673 eXplainable Artificial Intelligence mechanism, e.g., attentions or transform-  
674 ers, possibly coupled with convolutions, to get insights about particular spa-  
675 tial characteristics that may help to better recognise specific changes through  
676 a CD model. Another limitation is that the proposed methodology does not  
677 discriminate among different change types. This may be explored as multi-  
678 class ATL problem where new change classes may appear or disappear in the

679 target domain with respect to the source domain. A further research direc-  
680 tion refers to the systematic investigation of expected properties of both the  
681 source scene and target scene, to better foster the performance of the ATL  
682 strategy. Finally, recent studies have explored the CD problem in time series  
683 of co-registered MS images that exhibit some temporal trend in the change  
684 phenomena. Temporal change patterns may be explored to extend the pro-  
685 posed ATL strategy from the bi-temporal to the multi-temporal setting.

## 686 **CRedit Authorship Contribution Statement**

687 **Giuseppina Andresini:** Conceptualization, Methodology, Software, Data  
688 curation, Investigation, Validation, Visualization, Writing - original draft,  
689 Writing - review & editing. **Annalisa Appice:** Conceptualization, Method-  
690 ology, Investigation, Validation, Supervision, Writing - original draft, Writing  
691 - review & editing. **Dino Ienco:** Conceptualization, Investigation, Writing  
692 - original draft, Writing - review & editing. **Donato Malerba:** Conceptual-  
693 ization, Writing - review & editing.

## 694 **References**

- 695 Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Cor-  
696 rado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., ... & Zheng,  
697 X. (2015), ‘TensorFlow: Large-scale machine learning on heterogeneous  
698 systems’.
- 699 Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P. & Süsstrunk, S. (2012),  
700 ‘Slic superpixels compared to state-of-the-art superpixel methods’, *IEEE*  
701 *Transactions on Pattern Analysis and Machine Intelligence* **34**(11), 2274–  
702 2282.
- 703 Andresini, G., Appice, A., Iaia, D., Malerba, D., Taggio, N. & Aiello, A.  
704 (2022), ‘Leveraging autoencoders in change vector analysis of optical satel-  
705 lite images’, *J. Intell. Inf. Syst.* **58**(3), 433–452.
- 706 Appice, A., Guccione, P., Acciaro, E. & Malerba, D. (2020), ‘Detecting  
707 salient regions in a bi-temporal hyperspectral scene by iterating clustering  
708 and classification’, *Applied Intelligence* **50**(10), 3179–3200.

- 709 Appice, A. & Malerba, D. (2019), ‘Segmentation-aided classification of hyper-  
710 spectral data using spatial dependency of spectral bands’, *ISPRS Journal*  
711 *of Photogrammetry and Remote Sensing* **147**, 215 – 231.
- 712 Bengio, Y., Courville, A. C. & Vincent, P. (2013), ‘Representation learning:  
713 A review and new perspectives’, *IEEE Trans. Pattern Anal. Mach. Intell.*  
714 **35**(8), 1798–1828.
- 715 Bergstra, J., Yamins, D. & Cox, D. D. (2013), Making a science of model  
716 search: Hyperparameter optimization in hundreds of dimensions for vision  
717 architectures, *in* ‘ICML’, pp. 115–123.
- 718 Caye Daudt, R., Le Saux, B., Boulch, A. & Gousseau, Y. (2019), ‘OSCD -  
719 Onera Satellite Change Detection’.
- 720 Celik, T. (2009), ‘Unsupervised change detection in satellite images using  
721 principal component analysis and  $k$ -means clustering’, *IEEE Geoscience*  
722 *and Remote Sensing Letters* **6**(4), 772–776.
- 723 Daudt, R. C., Saux, B. L. & Boulch, A. (2018), Fully convolutional siamese  
724 networks for change detection, *in* ‘2018 IEEE International Conference on  
725 Image Processing, ICIP 2018, Athens, Greece, October 7-10, 2018’, IEEE,  
726 pp. 4063–4067.
- 727 Demšar, J. (2006), ‘Statistical comparisons of classifiers over multiple data  
728 sets’, *J. Mach. Learn. Res.* **7**, 1–30.
- 729 Deng, J. S., Wang, K., Deng, Y. & Qi, G. J. (2008), ‘Pca-based land-use  
730 change detection and analysis using multitemporal and multisensor satel-  
731 lite data’, *International Journal of Remote Sensing* **29**(16), 4823–4838.
- 732 Gautheron, L., Habrard, A., Morvant, E. & Sebban, M. (2020), ‘Metric  
733 learning from imbalanced data with generalization guarantees’, *Pattern*  
734 *Recognition Letters* **133**, 298 – 304.
- 735 Hadsell, R., Chopra, S. & LeCun, Y. (2006), Dimensionality reduction by  
736 learning an invariant mapping, *in* ‘2006 IEEE Computer Society Confer-  
737 ence on Computer Vision and Pattern Recognition (CVPR’06)’, Vol. 2,  
738 pp. 1735–1742.



- 739 Hafner, S., Nascetti, A., Azizpour, H. & Ban, Y. (2022), ‘Sentinel-1 and  
740 sentinel-2 data fusion for urban change detection using a dual stream u-  
741 net’, *IEEE Geosci. Remote. Sens. Lett.* **19**, 1–5.
- 742 He, K., Zhang, X., Ren, S. & Sun, J. (2016), Deep residual learning for image  
743 recognition, *in* ‘2016 IEEE Conference on Computer Vision and Pattern  
744 Recognition (CVPR)’, pp. 770–778.
- 745 Jiang, H., Peng, M., Zhong, Y., Xie, H., Hao, Z., Lin, J., Ma, X. & Hu,  
746 X. (2022), ‘A survey on deep learning-based change detection from high-  
747 resolution remote sensing images’, *Remote. Sens.* **14**(7), 1552.
- 748 Kingma, D. P. & Ba, J. (2014), Adam: A method for stochastic optimization,  
749 *in* ‘ICLR’.
- 750 Kubat, M. & Matwin, S. (1997), Addressing the curse of imbalanced training  
751 sets: One-sided selection, *in* ‘In Proceedings of the Fourteenth Interna-  
752 tional Conference on Machine Learning’, Morgan Kaufmann, pp. 179–186.
- 753 Lewis, A., Lymburner, L., Purss, M. B. J., Brooke, B. P., Evans, B. J. K.,  
754 Ip, A., Dekker, A. G., Irons, J. R., Minchin, S., Mueller, N., Oliver, S.,  
755 Roberts, D. O., Ryan, B., Thankappan, M., Woodcock, R. & Wyborn,  
756 L. (2016), ‘Rapid, high-resolution detection of environmental change over  
757 continental scales from satellite data - the earth observation data cube’,  
758 *Int. J. Digit. Earth* **9**(1), 106–111.
- 759 Lu, J., Hu, J. & Zhou, J. (2017), ‘Deep metric learning for visual understand-  
760 ing: An overview of recent advances’, *IEEE Signal Processing Magazine*  
761 **34**(6), 76–84.
- 762 Lv, Z., Liu, T., Benediktsson, J. A. & Falco, N. (2022), ‘Land cover change  
763 detection techniques: Very-high-resolution optical images: A review’,  
764 *IEEE Geoscience and Remote Sensing Magazine* **10**(1), 44–63.
- 765 López-Fandiño, J., B. Heras, D., Argüello, F. & Dalla Mura, M. (2019), ‘Gpu  
766 framework for change detection in multitemporal hyperspectral images’,  
767 *Int J Parallel Prog* **47**, 272–292.
- 768 Ma, W., Wu, Y., Gong, M., Xiong, Y., Yang, H. & Hu, T. (2019), ‘Change  
769 detection in sar images based on matrix factorisation and a bayes classifier’,  
770 *International Journal of Remote Sensing* **40**(3), 1066–1091.

- 771 Otsu, N. (1972), ‘A threshold selection method from gray-level histograms’,  
772 *IEEE Trans. Geoscience and Remote Sensing* **9**(1), 62–66.
- 773 Ouali, Y., Hudelot, C. & Tami, M. (2020), ‘An overview of deep semi-  
774 supervised learning’, *CoRR* **abs/2006.05278**.
- 775 Pasolli, E., Melgani, F., Tuia, D., Pacifici, F. & Emery, W. J. (2014), ‘SVM  
776 active learning approach for image classification using spatial information’,  
777 *IEEE Trans. Geosci. Remote. Sens.* **52**(4), 2217–2233.
- 778 Ru, L., Du, B. & Wu, C. (2021), ‘Multi-temporal scene classification and  
779 scene change detection with correlation based fusion’, *IEEE Trans. Image  
780 Process.* **30**, 1382–1394.
- 781 Ruzicka, V., D’Aronco, S., Wegner, J. D. & Schindler, K. (2020), Deep active  
782 learning in remote sensing for data efficient change detection, *in* T. Cor-  
783 petti, D. Ienco, R. Interdonato, M. Pham & S. Lefèvre, eds, ‘Proceedings  
784 of MACLEAN: MACHine Learning for EArth ObservatioN Workshop co-  
785 located with the European Conference on Machine Learning and Principles  
786 and Practice of Knowledge Discovery in Databases (ECML/PKDD 2020),  
787 Virtual Conference, September 14-18, 2020’, Vol. 2766 of *CEUR Workshop  
788 Proceedings*, CEUR-WS.org.
- 789 Shi, Q., Liu, M., Li, S., Liu, X., Wang, F. & Zhang, L. (2022), ‘A deeply su-  
790 pervised attention metric-based network and an open aerial image dataset  
791 for remote sensing change detection’, *IEEE Transactions on Geoscience  
792 and Remote Sensing* **60**, 1–16.
- 793 Shi, S., Zhong, Y., Zhao, J., Lv, P., Liu, Y. & Zhang, L. (2022), ‘Land-  
794 use/land-cover change detection based on class-prior object-oriented con-  
795 ditional random field framework for high spatial resolution remote sensing  
796 imagery’, *IEEE Trans. Geosci. Remote. Sens.* **60**, 1–16.
- 797 Sublime, J. & Kalinicheva, E. (2019), ‘Automatic post-disaster damage map-  
798 ping using deep-learning techniques for change detection: Case study of  
799 the tohoku tsunami’, *Remote. Sens.* **11**(9), 1123.
- 800 Tan, C., Sun, F., Kong, T., Zhang, W., Yang, C. & Liu, C. (2018), A survey  
801 on deep transfer learning, *in* ‘Proc. of the International Conference on  
802 Artificial Neural Networks and Machine Learning (ICANN)’.

- 803 Tan, P.-N., Steinbach, M. & Kumar, V. (2005), *Introduction to Data Mining*,  
804 *(First Edition)*, Addison-Wesley Longman Publishing Co., Inc., Boston,  
805 MA, USA.
- 806 Tarvainen, A. & Valpola, H. (2017), Mean teachers are better role models:  
807 Weight-averaged consistency targets improve semi-supervised deep learn-  
808 ing results, *in* ‘ICLR’.
- 809 Wu, Y., Bai, Z., Miao, Q., Ma, W., Yang, Y. & Gong, M. (2020), ‘A clas-  
810 sified adversarial network for multi-spectral remote sensing image change  
811 detection’, *Remote Sensing* **12**(13).
- 812 Wu, Y., Li, J., Yuan, Y., Qin, A. K., Miao, Q.-G. & Gong, M.-G. (2021),  
813 ‘Commonality autoencoder: Learning common features for change detec-  
814 tion from heterogeneous images’, *IEEE Transactions on Neural Networks*  
815 *and Learning Systems* pp. 1–14.
- 816 Yang, M., Jiao, L., Liu, F., Hou, B. & Yang, S. (2019), ‘Transferred deep  
817 learning-based change detection in remote sensing images’, *IEEE Trans.*  
818 *Geosci. Remote. Sens.* **57**(9), 6960–6973.