



**HAL**  
open science

# Domain Adaptation for Mapping LCZs in Sub-Saharan Africa with Remote Sensing: A Comprehensive Approach to Health Data Analysis

Basile Rouse, Sylvain Lobry, Géraldine Duthé, Valérie Golaz, Laurent Wendling

## ► To cite this version:

Basile Rouse, Sylvain Lobry, Géraldine Duthé, Valérie Golaz, Laurent Wendling. Domain Adaptation for Mapping LCZs in Sub-Saharan Africa with Remote Sensing: A Comprehensive Approach to Health Data Analysis. IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, 2024, 17, pp.13016-13029. 10.1109/JSTARS.2024.3421284 . hal-04750599

**HAL Id: hal-04750599**

**<https://hal.science/hal-04750599v1>**

Submitted on 23 Oct 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

# Domain Adaptation for Mapping LCZs in Sub-Saharan Africa With Remote Sensing: A Comprehensive Approach to Health Data Analysis

Basile Rouse<sup>1</sup>, Sylvain Lobry<sup>2</sup>, *Member, IEEE*, Géraldine Duthé, Valérie Golaz, and Laurent Wendling

**Abstract**—Environment and population are closely linked, but their interactions remain challenging to assess. To fill this gap, modeling the environment at a fine resolution brings a significant value, if combined with population-based studies. This is particularly challenging in regions where the availability of both population and environmental data are limited. In low- and middle-income countries, many demographic and health data are from nationally representative household surveys, which now provide approximate geolocations of the sampled households. In parallel, freely available remote sensing data, due to their high spatial and temporal resolution, make it possible to capture the local environment at any time. This study aims to correlate standard demographic and health information with a high-resolution environment characterization derived from satellite data, encompassing both rural and urban areas in Sub-Saharan Africa. We use the malaria indicator survey conducted in 2017–2018 in Burkina Faso. We first present a deep semisupervised domain adaptation strategy based on the intertropical climatic characteristics of the country for precisely mapping local climate zones (LCZs). This strategy models seasonal variations through contrastive learning to extract useful information for the mapping process. We then use this high-resolution LCZ map to characterize, in four groups, the immediate environment of the sampled households. We find a significant association between these local environments and malaria among households' children. Going beyond the traditional dichotomous urban/rural characterization, our results provide interesting insights for public health. This innovative method offers new avenues for exploring population and environment interactions, especially in the growing climate change concern.

**Index Terms**—Deep learning, demography, domain adaptation (DA), land cover, remote sensing.

## I. INTRODUCTION

RECENT studies highlight the interactions between population and environmental characteristics, especially in

Manuscript received 27 December 2023; revised 23 May 2024 and 24 June 2024; accepted 26 June 2024. Date of publication 1 July 2024; date of current version 24 July 2024. This work was supported in part by the Data Intelligence Institute of Paris (DiIP), in part by IdEx Université Paris Cité under Grant ANR-18-IDEX-0001, and in part by HPC resources from GENCI-IDRIS under Grant 2021-AD011013527. (*Corresponding author: Basile Rouse.*)

Basile Rouse is with the LIPADE, Université Paris Cité, F-75006 Paris, France, and also with the French Institute for Demographic Studies (INED), 93300 Aubervilliers, France (e-mail: basile.rousse@u-paris.fr).

Sylvain Lobry and Laurent Wendling are with the LIPADE, Université Paris Cité, F-75006 Paris, France (e-mail: sylvain.lobry@u-paris.fr; laurent.wendling@u-paris.fr).

Géraldine Duthé and Valérie Golaz are with the French Institute for Demographic Studies, French Institute for Demographic Studies (INED), 93300 Aubervilliers, France.

Digital Object Identifier 10.1109/JSTARS.2024.3421284

countries where population, strongly depends on the local environment for resources and economic activities [1], [2], [3]. With regard to health in many intertropical countries, malaria, which is a mosquito vector-borne disease, has an epidemiological facies that directly depends on the climate and its variability over the year (e.g., temperature, rainfalls), as well as the local environment (e.g., urbanity, presence of stagnated water) [4]. In Sub-Saharan Africa, most countries are located in the intertropical zone (apart from the southern part), and malaria is still a major health issue [5]. High-frequency environmental monitoring would be beneficial for public policy. However, studying the link between population and environment is challenging as it requires matching data from both fields. In Sub-Saharan Africa, where most of the low-income countries are located, the challenge is even higher because both population and environmental data are limited. In these countries, most of the demographic and health data are from nationally representative household surveys as those conducted by the demographic and health survey (DHS) program. To address the rising concern of the environment, the DHS program provides approximate coordinates of the areas where households are located, as well as some associated environmental indicators (e.g., normalized difference vegetation index, rainfall). However, these indicators are rarely frequently updated or spatially precise. This gap can be filled using very frequent and freely available remote sensing images, such as those from the Copernicus program. Such images allow us to create complex environmental indicators (e.g., land cover maps) for specific areas and periods. This ability to picture the environment, when linked to population and health data, can provide valuable insights about environmental risk factors [6], [7], [8]. Gibb et al. [7] used remote sensing to study interactions between climate and dengue in Vietnam. These studies often use basic indicators, such as vegetation or precipitation, and would gain from more detailed environmental characterization.

Many products and land use/land cover classification systems have been developed in the last years to offer resources for environmental characterization: the global human settlement layer produced global spatial data about human presence until 2014 [9], the global human footprint estimates the human influence on the environment [10], and the ESA world cover project generated worldwide 10 m resolution land cover maps for 2020 and 2021 [11], [12]. However, these classification schemes have some limitations, such as the ESA World Cover Project's

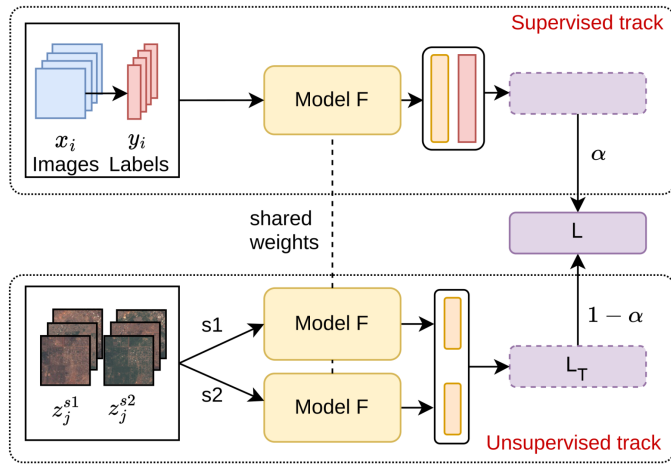


Fig. 1. Training process using s-SSDA, including a supervised and an unsupervised track. These two tracks are used simultaneously during training. First, a batch of labeled  $32 \times 32 \times 10$  images is fed into the model to compute the cross-entropy loss  $L_S$ . Then, the contrastive loss  $L_T$  is computed from prediction vectors from positive and negative pairs. The two losses are combined into one total loss  $L$  for backpropagation.

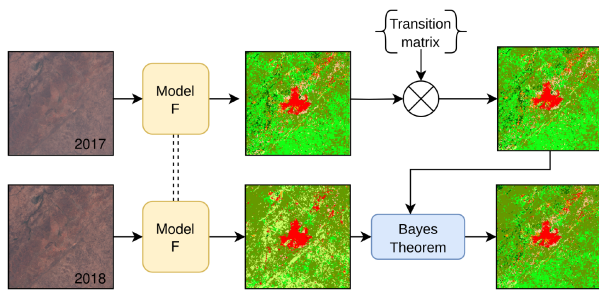


Fig. 2. Mapping process. We perform a temporal regularization process using a Markov Chain. Sentinel-2 images from early 2017 and 2018 are selected to match the survey period. A prior LCZ map for early 2018 is computed from a Markov chain and a map for early 2017. This prior is linked to the 2018 map using the Bayes theorem. Color legend can be found in Fig. 1.

definition of urban areas, which is closely tied to cultural aspects, limiting its transferability [13].

To address these limitations and provide a more globally applicable and detailed characterization, Stewart et al. [14] introduced the local climate zones (LCZs) classification scheme. It is made of 10 urban classes and 7 rural classes based on their surface physical properties and human activities, which allows a better transferability than ESA world cover and a more comprehensive description level than the Global Human Settlement Layer. Classes and their reference colors are represented in Fig. 4. LCZs were first developed to support Urban Heat Island research [15] but have been used to support various fields, such as energy usage [16], climate [17], or geoscience modeling [18]. They also offer a great opportunity to analyze population data better, thanks to their comprehensive and culturally independent classification. Indeed, LCZs have proven to be a good representation for linking the environment to many different health issues, as depression [19], cardiovascular diseases [20], or urban health issues in Sub-Saharan cities [21].

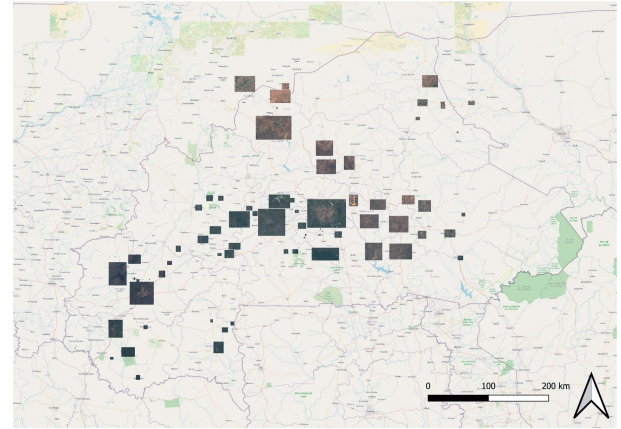


Fig. 3. Regions of interests used to create the seasonal dataset over Burkina Faso.

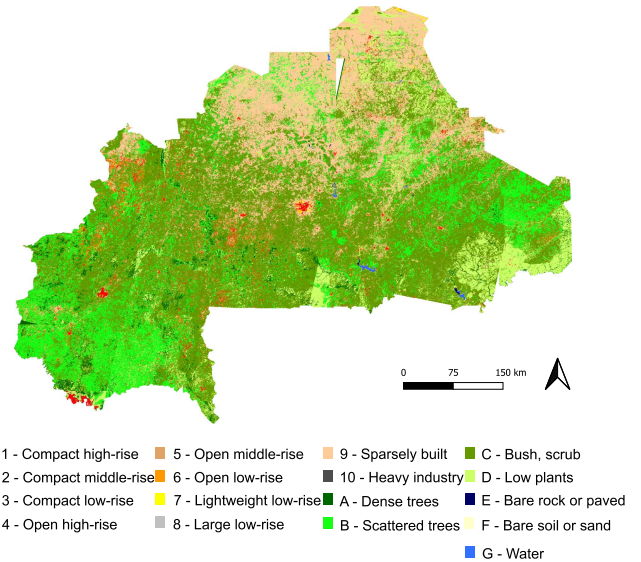


Fig. 4. LCZ map of Burkina Faso for early 2018.

However, mapping LCZs is a challenging research topic. Several works on LCZs have focused on the generation of high-quality maps using administrative data, vector data, or remotely sensed data [13], [22], [23], [24]. In particular, the community-based project “World Urban Database and Access Portal Tools” provides a worldwide database on urban morphology [25]. It has been used alongside other available LCZ products to produce continental and worldwide LCZ maps with a 100 m resolution [26], [27]. The So2Sat LCZ42 dataset was introduced to support research on automatic classification of LCZ using deep neural networks [28]. This dataset provides labeled  $32 \times 32$  pixel Sentinel-1 and Sentinel-2 patches from 42 agglomerations worldwide. It has been used for the generation of 1642 LCZ maps using deep neural networks, to further support urban research [29].

Deep neural networks have substantially contributed to computer vision [30], [31]. Due to their improved transferability compared to conventional methods (e.g., random forests)

and their ability to incorporate contextual information, these techniques are well suited for remote sensing tasks [32], [33]. Moreover, deep networks have been successfully applied to LCZ classification using convolutional neural networks [13], [34], [35], as well as with recurrent neural networks [23], [36], taking advantage of the large-scale So2Sat dataset. However, these models remain significantly reliant on training data, resulting in spatial domain adaptation (DA) issues. For instance, only 3 of the 42 training cities are located in Africa, but only 1 in the intertropical zone of the continent (Nairobi). The intertropical African climate characteristics, with its large seasonal changes between wet and dry seasons, cannot be learned by a neural network fully supervised on this data. More data are required to include such information in the training process. Although Sentinel-2 data are freely accessible, its annotation is time-consuming and requires expertise. Taking advantage of the large amount of available Sentinel-2 data using DA methods is required to generate accurate LCZ maps of Sub-Saharan Africa.

In machine learning, DA involves transferring a model from a labeled dataset, referred to as the source dataset, to a target dataset that has a different underlying distribution. The objective of DA methods is to enable models to generalize better and perform well in the target domain, even with limited or no supervision (i.e., labeled data) from the target domain. It is suitable for many real-life tasks where few labeled data are available [37]. DA techniques have been explored for remote sensing [38], [39] for transferring knowledge from a sensor to another [40], [41], [42] and from a region to another [43], [44], [45], [46]. This enables transferring the model on large-scale datasets, as SEN12MS [47] or So2Sat [28], on other regions of the world. Common DA techniques include generative networks to modify the data distributions [48], [49], latent space alignment with adversarial training [50], [51], [52] or semisupervised learning (SSL) that leverages unlabeled data to improve model's performances [53], [54]. Most DA strategies, when applied to remote sensing, are focused on small DA problems, such as cross-city or regional adaptation. Tasar et al. [43] performed DA using 4 European cities, whereas Zhu et al. used ISPRS Vaihingen and Potsdam datasets. These strategies are then not adapted for land cover mapping at a country level, as they do not include the environmental complexities of a whole country.

This article is built upon a preliminary work [55] conducted for large scale LCZs mapping in Sub-Saharan Africa. In the following, we go further by investigating the potential links between malaria and local environment, as described by the generated high-resolution map. The focus of our study is on Burkina Faso, using a nationally representative household survey related to malaria conducted in 2017–2018. We take advantage of the high-temporal frequency of Sentinel-2 images to capture the temporal climate specifics of Burkina Faso, as it is crucial for linking maps to data collected over defined periods.

Our contributions are as follows.

1) We introduce a novel semisupervised DA (SSDA) method that includes seasonal data and contrastive learning into training for mapping LCZs in Sub-Saharan Africa.

- 2) We propose a method for linking the resulting LCZ map to the approximate buffer areas where households can be found in DHS-like surveys.
- 3) We show that local environments bring interesting insights for public health issues related to malaria exposure. More specifically, we highlight that there are intrarural and intraurban factors impacting malaria propagation, which goes beyond the traditional dichotomous urban/rural characterization.
- 4) Our approach opens new perspectives for population studies: as we use freely available data, this method is replicable and offers new avenues to explore the interaction between population and environment.

Section II describes the method used for mapping LCZs in Burkina Faso. The application of this method on Burkina Faso is presented in Section III. Section IV addresses the creation and assessment of Burkina Faso LCZ maps and their integration into a demographic analysis procedure. The impact of our proposed method for the mapping of LCZs and its association with demographic data is discussed in Section V. Finally, Section VI concludes this article.

## II. METHOD

Although LCZs have been proposed as a universal tool, their mapping from remote sensing imagery depends on the background landscape and seasonal changes [26]. To tackle this problem, we propose an SSDA method that takes advantage of the seasonal changes of countries in Sub-Saharan Africa to extract spatial and temporal information about the target areas. Our objective is to define a mapping process taking a set of images of spatial size  $32 \times 32$  with 10 channels and returning a series of LCZ labels. The resulting map will be created from this set of images. To this effect, we define a neural network  $F(\cdot)$  taking  $32 \times 32 \times 10$  images as input and yielding a 1-D vector  $s = [s_1, s_2, \dots, s_{17}]^t$  of 17 elements. Each element is a prediction score of the 17 LCZ classes for a given input. To start with, we provide a description of the source and target datasets used in this process. Then, we describe our method, named seasonal-SSDA (s-SSDA), to train  $F(\cdot)$  on the spatial DA problem.

### A. Source and Target Datasets

DA methods aim to learn relevant features from a labeled source dataset and transfer them to the target dataset. In this work, the target dataset is not labeled. Both datasets are made of Sentinel-2 images with 10 spectral bands at resolutions of 10 m and 20 m, up-sampled at 10 m using a bicubic interpolation. Let us define the source dataset  $D_S = (x_i, y_i)_{i \in \llbracket 1, n_S \rrbracket}$  where  $x_i$  is a  $32 \times 32 \times 10$  image,  $y_i$  its associated ground truth LCZ class, and  $n_S$  is the number of samples in the dataset. In this work, the source dataset is the So2Sat dataset [28]. We supplement  $D_S$  with a target dataset that contains the areas covered by the demographic survey of interest, Burkina Faso. We denote this dataset as  $D_T = (z_i^{s1}, z_i^{s2})_{i \in \llbracket 1, n_T \rrbracket}$  made of  $n_T$  pairs of

$32 \times 32 \times 10$  images ( $z_i^{s_1}, z_i^{s_2}$ ) from the same area at seasons  $s_1$  and  $s_2$ .

### B. Seasonal Semisupervised DA

In order to extract seasonal features from the target images without requiring labels, we use the contrastive loss and teach the model to be invariant to the seasons. This loss has been widely used for self-supervised learning [56], [57] and aims to train the model to increase a similarity measure between positive (similar) pairs while decreasing this similarity measure within negative (dissimilar) pairs

$$L_{i,j} = -\log \frac{\exp(\text{sim}(F(z_i), F(z_j))/\tau)}{\sum_{k=1, k \neq i}^{2B} \exp(\text{sim}(F(z_i), F(z_k))/\tau)} \quad (1)$$

where  $B$  is the number of pair of images,  $\text{sim}(\cdot, \cdot)$  is the cosine similarity,  $\tau$  is the temperature,  $(i, j) \in \llbracket 1, B \rrbracket^2$ ,  $(z_l)_{l \in \llbracket 1, 2B \rrbracket}$  samples from a batch of  $B$  samples, and  $(z_i, z_j)$  is a positive pair. In our case, the positive pair is made of two images of the same area in different seasons and the negative pair is made of other two images of different areas.

This loss is integrated into a two-tracks SSDA process involving the labeled source dataset  $D_S$  and the unlabeled target dataset  $D_T$ . This process is shown in Fig. 1. The first track is a regular supervised learning process using  $D_S$  to minimize a supervised cross-entropy loss  $L_S$ : the model  $F(\cdot)$  is trained to classify LCZs on various regions of the world. The second process, using the unlabeled samples from  $D_T$ , is done simultaneously.  $F(\cdot)$  is used as a Siamese neural network similarly as shown in [56] and [58] to classify two different images from the exact same spatial area but at two different seasons  $s_1$  and  $s_2$ . The image from the second season  $z_i^{s_2}$  can be seen as a seasonal augmentation of  $z_i^{s_1}$ . To keep the model training simple, we do not consider possible new settlements that are built between the recording dates of the two images, which could change the LCZ class of a defined area. For example, we assume that an area classified as ‘‘low plants’’ in the first season will remain ‘‘low plants in the second season, even if houses have been built. Then, the contrastive loss  $L_T$  is computed between the outputs of the positive pair and the negative pair. It is worth noting that rural LCZ labels can change throughout the year due to seasonal variations. However, urban areas should remain unaffected by seasonal changes, despite having a different visual aspect. We introduce this prior by adding weights to the contrastive loss to penalize inconsistency in urban area predictions. This second track aims to enforce robustness to the seasons as well as transferring its knowledge to unseen areas, which are not present in  $D_S$ . The loss used for the training of the model  $F(\cdot)$  is a combination of the results of the two tracks with a regularization coefficient  $\alpha \in [0, 1]$

$$L = \alpha \times L_S + (1 - \alpha) \times L_T. \quad (2)$$

This regularization term is determined empirically.

### C. Temporal Regularization

Sentinel-2 products are available at a very high frequency (maximum five days). LCZ maps of the target areas can be generated not only at the designated time but also during the

same month in years preceding the year of interest. We take advantage of Sentinel-2 temporal data to ensure temporal continuity between years using a Markov chain. We, thus, define the following.

- 1)  $I_N \in \mathbb{R}^{32 \times 32}$  an array built upon be the output of  $F(\cdot)$  at time  $N$ , where each coefficient of  $I_N$  is the class with the highest score given by  $F(\cdot)$ .
- 2)  $\text{LCZ}_N$  the LCZ class ( $c_N \in \llbracket 1, 17 \rrbracket$ ) of the patch at time  $N$ .
- 3)  $M \in \mathbb{R}^{17 \times 17}$  a matrix where  $m_{r,c} \in \llbracket 1, 17 \rrbracket^2$  is the coefficient of  $M$  at row  $r$  and column  $c$ .

$m_{i,j}$  is the probability in the first Markov process to go from  $\text{LCZ}_{N-1} = i$  to  $\text{LCZ}_N = j$ ,  $(i, j) \in \llbracket 1, 17 \rrbracket^2$ . These probabilities are dependent upon the environmental and political context of the target area. The definition of these weights is discussed in Section III.

If the LCZ classification of a specific area at time  $N$  (i.e.,  $\text{LCZ}_N$ ) follows a first-order Markov process (from two consecutive years at the same season), for all  $N$

$$P(\text{LCZ}_N = c_N) = m_{c_{N-1}, c_N} \times P(\text{LCZ}_{N-1} = c_{N-1}) \quad (3)$$

then, according to the Bayes theorem

$$\begin{aligned} P(\text{LCZ}_N = c_N | I_N) &= \frac{P(I_N | \text{LCZ}_N = c_N)}{P(I_N)} \times m_{c_{N-1}, c_N} \\ &\times P(\text{LCZ}_{N-1} = c_{N-1}). \end{aligned} \quad (4)$$

$P(I_N | \text{LCZ}_N = c_N)$  is the prediction vector of the model. After predicting mono-temporal LCZ scores with  $F$ , the Markov chain can be applied to obtain the final LCZ maps. This regularization process is shown in Fig. 2. In the following section, we explain how this was used for generating an LCZ map of Burkina Faso in early 2018, during the survey period.

## III. GENERATION OF AN LCZ MAP OF BURKINA FASO

### A. Target Dataset Creation

This section describes the procedure employed to generate the target dataset for the SSL part of the training process. To reduce the domain gap between available training data and Burkina Faso, we supplement the source dataset So2Sat by Sentinel-2 images over Burkina Faso at the end of the dry and of the wet seasons. We use level L1C images in order to match the So2Sat template. This target dataset has been created using the following procedure.

- 1) *Downloading of L1C sentinel-2 tiles*: Linked to each of Burkina Faso’s region’s capital cities at the end of the dry and of the wet seasons to maximize variances between the two tiles. The two tiles of the same region at different times were selected to have under 5% of cloud coverage to reduce errors caused by clouds.
- 2) *Selection of regions of interest*: Where areas of interest can be found: cities, villages, industries, natural parks, forests, lakes, or rivers. Tiles are cropped into rectangular shapes centered on areas of interest and large enough to

TABLE I  
TRANSITION WEIGHTS FOR THE MARKOV PROCESS DURING THE LCZ MAP GENERATION

	1	2	3	4	5	6	7	8	9	10	A	B	C	D	E	F	G
1) Compact high-rise	1.00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2) Compact mid-rise	0.05	0.95	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3) Compact low-rise	0	0.05	0.95	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4) Open high-rise	0.05	0	0	0.95	0	0	0	0	0	0	0	0	0	0	0	0	0
5) Open mid-rise	0	0.05	0	0	0.95	0	0	0	0	0	0	0	0	0	0	0	0
6) Open low-rise	0	0	0.1	0	0	0.9	0	0	0	0	0	0	0	0	0	0	0
7) Lightweight low-rise	0	0	0.1	0	0	0	0.9	0	0	0	0	0	0	0	0	0	0
8) Large low-rise	0	0	0	0	0	0	0	1.0	0	0	0	0	0	0	0	0	0
9) Sparsely built	0	0	0	0	0	0.1	0	0	0.8	0	0	0.1	0	0	0	0	0
10) Heavy industry	0	0	0	0	0	0	0	0	0	1.0	0	0	0	0	0	0	0
A) Dense Trees	0	0	0	0	0	0	0	0	0	0	0.95	0.05	0	0	0	0	0
B) Scattered Trees	0	0	0	0	0	0	0	0	0.30	0	0.01	0.69	0	0	0	0	0
C) Bush, scrubs	0	0	0	0	0	0	0	0	0.10	0	0	0	0.9	0	0	0	0
D) Low plants	0	0	0	0	0	0	0	0	0.10	0	0	0	0	0.9	0	0	0
E) Bare rock or paved	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1.0	0	0
F) Bare soil or sand	0	0	0	0	0	0	0	0	0.01	0	0	0	0	0	0	0.99	0
G) Water	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1.0

Correspondence between identifiers and names is available in Fig. 4.

include nearby environments. The same regions of interest are selected for both tiles from the dry and wet seasons.

- 3) *Splitting of the regions of interest*: Into patches of  $32 \times 32$  pixels to match the patch size of the So2Sat dataset. Patches pairs (same patch region, different seasons) are created to feed the neural network during training.

This procedure results in 225K patch pairs distributed throughout all of Burkina Faso, as shown in Fig. 3. As it can be seen on the spatial distribution of the samples, all regions and climates are included in the training dataset.

### B. Training Settings

We use ResNet50 architecture [30], pretrained on the complete So2Sat dataset. This pretraining stage is performed to initialize the model weights for the subsequent semisupervised training step. For this second step, the Adam optimizer is used with a learning rate of 0.001. The batch size for both the supervised and unsupervised phases is set to 256. The temperature parameter  $\tau$  is set to 0.5. Based on our experiments, we set the parameter  $\alpha$  from (2) to 0.9 in this work.

### C. Transition Weights

The assumption that LCZs follow a Markov chain requires the definition of transition weights, which represent the probabilities of moving from one state, i.e., LCZ class, to another. The context of the country highly influences these probabilities, as they may result from urbanization (e.g., the transition from “open low-rise” to “compact low-rise”), forest management (e.g., prohibition of deforestation), or the geographical situation of the country. For instance, in Burkina Faso, urbanization is proceeding rapidly and the country’s terrain is predominantly flat. Thus, cities are likely to expand horizontally rather than vertically, suggesting that the transition weights to “compact high-rise” should be very low. Transition weights have been set empirically regarding Burkina Faso’s spatial and political features, such as the urbanization plan. Transition weights can be found in Table I.

### D. Resulting Map

Malaria indicator survey (MIS) data were collected in Burkina Faso from November 2017 to March 2018—from the end of the wet season to the dry season. However, we generate the map for early 2018 as the majority of interviews were conducted in January and February. We gather Sentinel-2 images from early 2017 and 2018 to perform the Markov process. Only images with a cloud cover of less than 5% are selected. To match the model’s input size of  $32 \times 32$  pixels, each Sentinel tile is gridded into  $320 \text{ m} \times 320 \text{ m}$  image patches. Patches from 2017 and 2018 are classified and used for the Markov process. The map is produced by concatenating all the generated LCZ patches and is shown in Fig. 4.<sup>1</sup> One map of Burkina Faso, halfway through the survey period (early 2018) when most of the interviews have been done, has been generated to link to the demographic survey. The initial map has a resolution of  $320 \text{ m} \times 320 \text{ m}$  and is upsampled to the input Sentinel image resolution,  $10 \text{ m} \times 10 \text{ m}$ , within the subdivision of the initial pixels ( $320 \text{ m} \times 320 \text{ m}$ ) into smaller pixels ( $10 \text{ m} \times 10 \text{ m}$ ). As expected, the map is divided into 3 main parts (4 if we consider cities/urban areas), which correspond to the climate profiles of Burkina Faso. The southern wetter part is mostly covered by “scattered trees” areas, and the more temperate part is mostly covered by “bush/scrub” areas, except for the nature reserves in the eastern part of the country. The LCZ classification of northern dryer areas presents a greater challenge for classification. The precision of this classification, through the use of Sentinel-2, is restricted by its 10 m resolution and the up-sampling of the 20 m resolution bands to 10 m. This may result in the inability to detect small settlements or houses smaller than the Sentinel-2 resolutions. The absence of this detection leads to a misinterpretation of areas with a very low rate of construction, e.g., the LCZ class “sparsely built,” and areas without any buildings, e.g., “bare soil or sand” and “low vegetation,” which can be found in the desert-like part in the north. This confusion is reinforced by the use of SSL that struggles when input samples are difficult, such as “sparsely built” samples as indicated by Bechtel et al. [59].

<sup>1</sup>The LCZ map can be downloaded from the following link: <https://drive.google.com/file/d/1KkpZfNj0DGOZuuVb-6nvOCWCvCp2C1Qk/view?usp=sharing>

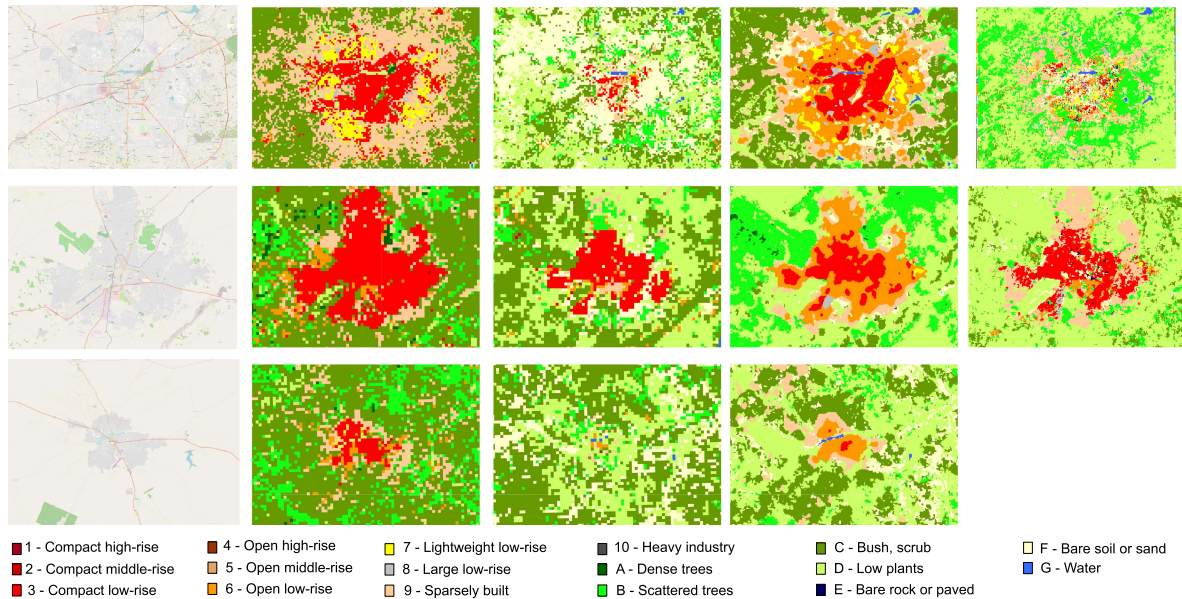


Fig. 5. Visual comparison of LCZ maps of 3 cities (Ouagadougou, Bobo-Dioulasso, and Fada-Ngourma). In addition to a map from OpenStreetMap, different methods are shown: the proposed s-SSDA, supervised training on So2Sat (Baseline), Global LCZ map [27], and GUL [29], from left to right.

### E. Comparison With Other LCZ Products

Several LCZ products have been introduced in recent years. In particular, So2Sat Global Urban LCZ (GUL) [29] and global LCZ map [27] enabled a global overview of the morphology of the Earth's surface. So2Sat GUL is a set of 1642 cities worldwide, based on So2Sat LCZ42. A two-track model has been trained in a supervised manner. The first track is a residual network taking as inputs Sentinel-2 images of an area. Several seasons can be used as input and the model is asked to predict the LCZ class for each season. The second track is a different residual network taking as input a Sentinel-1 image of the same area. The mean of the prediction vectors is then computed to obtain the final prediction vector. The model is used to generate LCZ maps of unseen cities all around the world. Two maps for Burkina Faso are available: The capital city, Ouagadougou, in the center of the country, and Bobo-Dioulasso, the second biggest city in the country, in the southwest.

The Global LCZ map covers all the planet surface [27] with a resolution of 100 m. It is the result of pixel-based Random Forest classifiers trained on 46 spatial features on training areas previously made by urban experts or in crowd-sourcing platforms using guidelines [60]. The whole surface of Burkina Faso was mapped.

The validation of land cover maps in intertropical Sub-Saharan Africa presents significant challenges. First, the region's diverse and complex ecosystems contribute to a wide range of land cover types. This variability makes it difficult to develop generalized validation approaches, especially for LCZs. Second, few references are available in Sub-Saharan Africa due to the lack of resources. To validate the effectiveness of our training approach, we gathered images over large areas over four cities in the different climate zones of Burkina Faso: Ouagadougou, Bobo-Dioulasso, Fada-Ngourma, and Ouahigouya [55]. These

images are split into patches and manually labeled from VHR remote sensing images.

Fig. 5 is a visual comparison of the LCZ maps of Ouagadougou, Bobo-Dioulasso, and Fada N'Gourma generated using different models. Maps from the baseline and the proposed s-SSDA model have been generated from the same Sentinel-2 tiles for early January 2018. As expected, the baseline model struggles to predict urban areas in Ouagadougou and Fada N'Gourma, as such information may not be included in the training data. Similarly, So2Sat GUL does not predict the whole of Ouagadougou as an urban area. Interestingly, the two previous models yield good visual results over Bobo-Dioulasso. Moreover, we quantitatively validated these products using the validation set used in the ablation study in [55]. Four Sentinel-2 images from early 2018 were collected and cropped over Ouahigouya, Bobo-Dioulasso, Fada-Ngourma, and Ouagadougou. The resulting images were gridded into  $32 \times 32$  pixel areas and labeled using very high-resolution (VHR) satellite images. As So2Sat GUL is only publicly available for areas over Ouagadougou and Bobo-Dioulasso, metrics have been computed on the 186 overlapping patches with the validation set. As GUL and the Global LCZ map are available at a lower resolution than the validation patches, we compute the validation metrics for two different aggregation procedures.

- 1) We consider that the area is well classified if the majority of the pixels image area has the same label as our validation set (MR).
- 2) We consider that the area is well classified if at least one of the subpixels has the same label as our validation set (IS). This second aggregation method is, therefore, the most favorable for GUL and Global LCZ map.

Quantitative comparison results are given in Table II. Both s-SSDA methods outperform So2sat GUL in all cases and for all metrics. In the best-case scenario, the temporally regularized

TABLE II  
COMPARISON RESULTS

Method	OA	F1	IoU
GUL - MR	0.140	0.122	0.070
GUL - IS	0.194	0.203	0.119
Global LCZ - MR	0.360	0.397	0.258
Global LCZ - IS	0.447	0.481	0.329
s-SSDA	0.427	0.402	0.278
s-SSDA + Markov	<b>0.561</b>	<b>0.538</b>	<b>0.389</b>

We compare our map to other LCZ products on 494 manually labeled patches. OA is the overall accuracy, F1 is the f1 score, and IoU is the intersection over the union. Best results are indicated in bold.

s-SSDA outperforms GUL by more than 30%. These results highlight the necessity to perform DA during training in unseen areas of the world, even after training on a global dataset. The global LCZ map, when taking the best scenario, achieved further improvements, obtaining an overall accuracy of 0.447, an F1-score of 0.481, and an IoU of 0.329, which are similar but slightly better than s-SSDA without regularization. The temporal regularization enables better mapping results, with a close to 10% improvement for each metric. The results on our validation set highlight the potential of taking advantage of seasonal features to perform DA.

#### IV. APPLICATION: MALARIA IN BURKINA FASO

In this section, we explore the integration of high-resolution LCZ characterization of the environment into a demographic analysis procedure. Specifically, we investigate malaria in Burkina Faso using the MIS 2017–2018. This section is structured as follows: In Section IV-A, we present the MIS survey: the collected data and the approximated geolocalization of the households. In Section IV-B, we provide a novel characterization of the households' environment *at the finest spatial granularity* based on the LCZ map produced. In Section IV-C, we estimate the presence of malaria among households according to four categories of environment. In Section IV-D, using univariate and multivariate analysis, we estimate the correlation between malaria and the environment, considering socio-economic characteristics that could interfere in the interaction.

##### A. MIS 2017–2018

In Burkina Faso, malaria is one of the most important cause of death (the first one in 2019 according to IHME<sup>2</sup>). A MIS from the DHS Program was conducted in 2018–2018 in the country [61]. In the survey, the sample was built so that results on malaria prevalence for 6–59 month-old children are representative of each of the 17 study areas (administrative regions and big cities, e.g., Ouagadougou). In total, 252 enumeration zones (EZs, i.e., geographic areas for conducting structured population surveys) were selected from the national sampling frame. EZs are the smallest spatial units where geolocations are available. While some of the selected EZs could not be visited

for safety reasons, 245 EZs were visited, and a total of 6322 households were interviewed. Among households, information about socio-demographic characteristics, preventive behaviors, health care related to fever occurrence, and knowledge on malaria were collected. Among each household, several tests were done on all eligible children aged 6–59 months with their legal representatives' consent, starting with a malaria rapid test, providing the first results in 15 min. A hemoglobin test was also performed in order to detect anemia, which is highly correlated to malaria. When the rapid test proved positive, a treatment for malaria was freely provided to the parents/legal representatives, and another blood sample was collected for a laboratory test to confirm the rapid test result and characterize the malaria parasite. In order to preserve household confidentiality, their spatial coordinates are not public. For each EZ, the average geolocation of households is computed and randomly displaced in an area of radius  $R$  depending on the type of EZ ( $R = 2$  km for urban EZs,  $R = 5$  km for rural EZs, except for 1% of the latter for which  $R = 10$  km).

##### B. Characterization of Households' Local Environments

To account for the displacement of the EZs geolocations, we model buffer areas around each EZ centroid by disk  $C_e$  (with  $e$  being the EZ identifier) of radius of two or five kms depending on the urban or rural type of the EZ. As only 1% of rural EZ centroids are displaced within a circle of 10 kms, we make the assumption that reducing the buffer area to five km disks will not alter the results. Different processes can be used to model the EZ environment. DHS experts suggest averaging the value of interest on the buffer area  $C_e$  [62], i.e., taking the mean LCZ distribution in the context of this study. This method implies considering pixels that contradict survey data on whether the area is urban or rural, as urban areas can be included in rural 5 km disks. Furthermore, it yields artificially heterogeneous environment indicators as the large area of  $C_e$  is taken into account. Grace et al. [63] suggested using VHR images or existing global population database as from World-Pop<sup>3</sup> to select interesting pixels. In this study, we rely on the urban/rural characterization of the EZ provided in MIS data. For each EZ  $e$ , we semirandomly sample  $n_{\text{random}}$  squared areas of size  $A = 10 \times 10$  (100 m  $\times$  100 m) inside  $C_e$  to model the potential true geolocations of interviewed households. To ensure sampling consistency with MIS information on the urban/rural type of  $e$ , these  $n_{\text{random}}$  areas have been selected to contain at least  $\delta = 10\%$  of pixels, which LCZ classifications belong to the urban or rural area. This semirandom sampling procedure results in a total of  $n_{\text{random}}$  areas that cover all EZs. These  $n_{\text{sampled}}$  areas give a global view on the local environment in which the households were interviewed, excluding impossible values yet preserving borderline cases as city borders and urban parks. A visual example of such a semirandom area selection is given in Fig. 6.

<sup>2</sup>[Online]. Available: <https://www.healthdata.org/research-analysis/health-by-location/profiles/burkina-faso>

<sup>3</sup>[Online]. Available: <https://www.worldpop.org/datacatalog/>



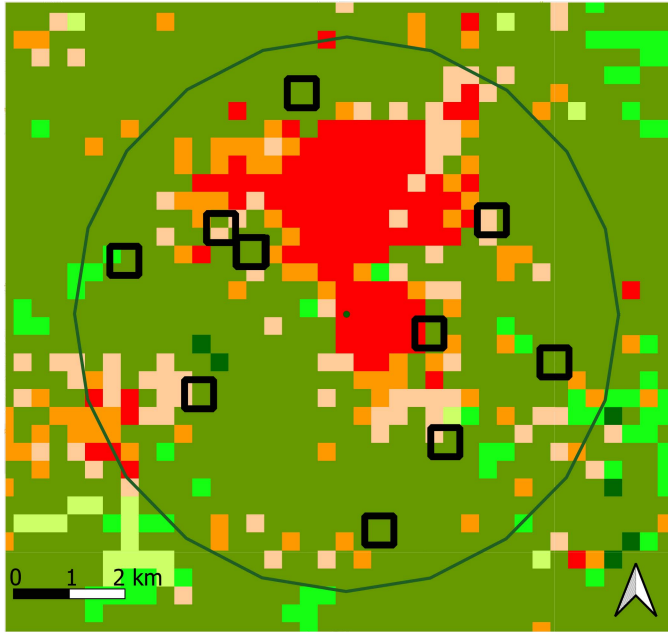


Fig. 6. Semirandom selection of potential location of households. Example of a rural EZ near Dédougou, Burkina Faso. The locations are randomly sampled in the rural part of  $C_k$ .

### C. Environment and Malaria Rates at the EZ Level

For a first visualization of the effects of the environment on the presence of malaria with our map, we find correlations between types of environment and malaria rate. We do not attempt to directly predict malaria prevalence from LCZ distributions as it would obscure socio-economic disparities potentially treating two households—one wealthy and one poor—as equals, despite socio-economic factors such as mosquito nets, types of sanitation facilities, and others potentially influencing the outcomes. Then, we clustered the EZs’ LCZ distribution into  $n_E = 4$  types of environments using the fuzzy C-means algorithm [64]. The LCZ distributions of the centers of each cluster, or types of environment, are depicted in Fig. 7. The definition of clusters using the LCZ distribution goes beyond the urban/rural dichotomy and takes advantage of the variety of classes in the classification scheme. Each cluster is highly polarized by a single LCZ class. Clusters 1 and 2 tend to be urban, and clusters 3 and 4 tend to be rural. Cluster 1 is highly urban and mostly includes “compact low-rise” LCZs. This cluster can be associated with towns and cities. Cluster 2 is less urban, mostly made of sparsely built areas, as can be found in the outer part of towns and cities and in the northern part of the country. Clusters 3 and 4 are two rural clusters dominated by scattered trees and bush/scrub, respectively.

MIS data provide rapid test results for each child tested in all EZs. In this work, we define *malaria rate* as the number of positive cases among children aged 6–59 months divided by the total number of children aged 6–59 months. We plot in Fig. 8, the distribution of malaria rates for each type of environment. Visually, malaria rates seem associated to the type of environment.

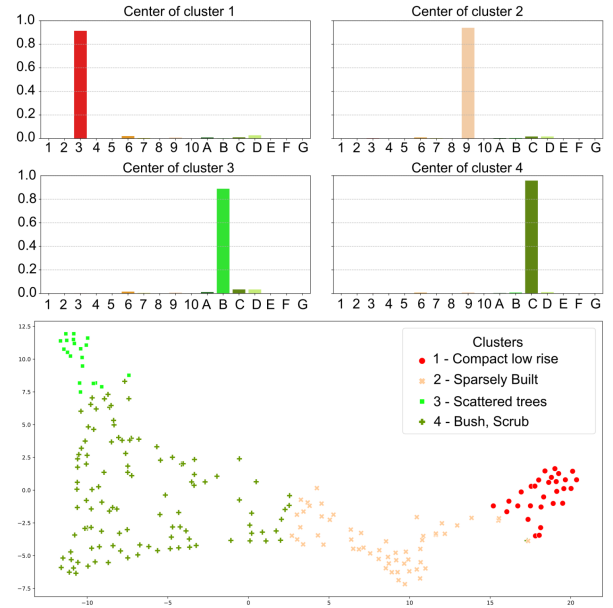


Fig. 7. LCZ distributions of clusters centers (up) and 2-D representation using t-SNE (down) [65].

TABLE III  
P-VALUES IN STUDENT’S T-TEST RESULTS

P-values	Cluster 1	Cluster 2	Cluster 3	Cluster 4
Cluster 1	1	< 0.001	< 0.001	< 0.001
Cluster 2	x	1	0.004	0.033
Cluster 3	x	x	1	0.160
Cluster 4	x	x	x	1

On average, households located in clusters 1 and 2 show lower malaria rates than rural clusters 3 and 4. Within these urban clusters, the most urban one presents lower rates. Cluster 3 shows slightly higher rates than cluster 4, but the difference is not significant according to the Student’s t-test shown in Table III. That nonsignificance may be explained by the over-representation of the bush/scrub class in the map presented in the mapping section.

### D. Environmental and Socio-Economic Influences on Malaria at the Household Level

This section introduces the method to link the local environments of households and their socio-economic data to malaria. First, we describe the dependant covariate that we want to explain: the presence of malaria cases in households. Then, we present the households’ characteristics, which will be used to explain the dependent covariate. Finally, we show that the influence of the environment on malaria is significant, even when considering household socio-economic data.

1) *Dependent Covariate*: Malaria rapid diagnostic test for all tested children are provided in MIS data. These malaria tests are blood-based tests detecting specific antigens (proteins) produced by malaria parasites in the blood of infected individuals. To study malaria at the household level, we now define our dependent covariate as the presence of at least one malaria-positive case, according to these tests, among the children aged 6–59

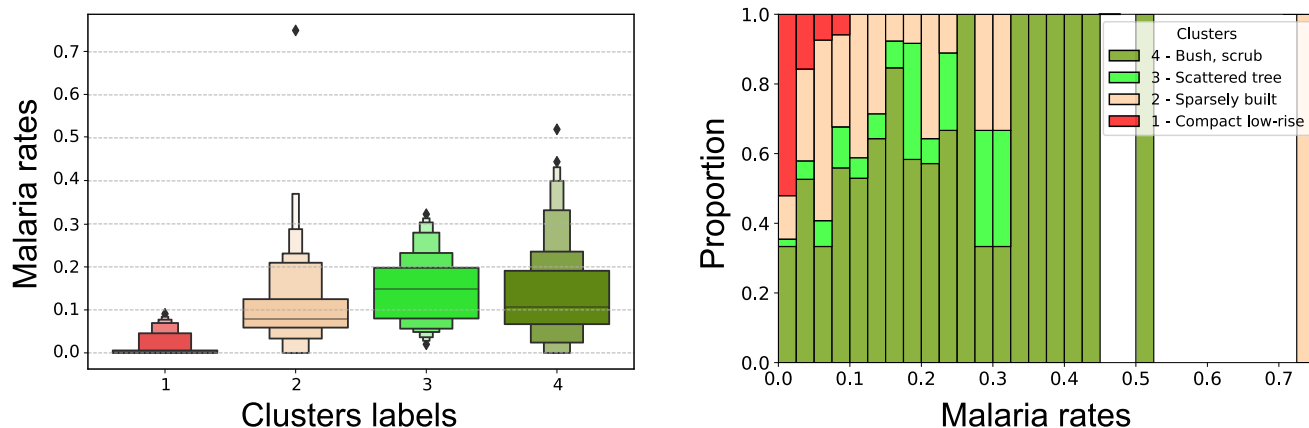


Fig. 8. Distribution of EZs malaria rates (left) and proportion of malaria rates by intervals, grouped by cluster.

months in the household. Thereafter, we use the expression *positive household* to refer to such households. Thus, this covariate is binary and is explained using a logistic regression and household-related data. The logistic regression aims to estimate the probability of an event occurring based on a given set of independent covariates. In this work, we estimate the probability of having at least one positive case in a household according to the covariates described as follows. That model allows us to analyze the interactions of the environment with other covariates at the household level.

2) *Households' Characteristics*: As mentioned in the survey description, MIS data provide socio-economic information about households in addition to malaria test results. Our explanatory covariates, at the household level, include information about wealth, type of toilets used, source of drinking water, the children's mother's level of education and the total number of children between 6 and 59 months in the household. The wealth index represents financial well-being. In this study, it is split into quintiles (poorest, poor, middle, rich, richest). According to the MIS data report [61], the types of toilets fall into three categories: no facility, nonimproved (e.g., bucket toilet, hanging toilet/latrine), and improved (e.g., flush toilet to piped sewer system, composting toilet). Sources of drinking water were also split into similar categories: nonimproved (unprotected spring or well, river), improved sources (e.g., piped water), and other sources (e.g., tanker truck, protected well). The mother's level of education cannot be indicated trivially. In the context of this article, one household may be composed of several children with different mothers. To be as close as possible to reality, we define the level of education of the mothers in a household as the most represented level (no education, primary, secondary, and higher) within the mothers present in the household. If there are two mothers in the household or if there is equality, we choose the highest level of education. The number of children aged 6–59 months is also computed and used as an explanatory covariate.

We use the clusters described in the previous section as our first environmental variable. The type of environment for households is determined by the cluster in which their corresponding EZs are classified. Our generated map has a label resolution of

320 m  $\times$  320 m, which may be too high for classifying small water bodies that may increase the population of mosquitoes in the area. To compensate for that lack, we add the presence of water in the buffer area using OpenStreetMap data. Table IV gives a summary of the considered covariates. In this table, "Rate+" is the proportion of households in each category where at least one child between 6 and 59 months have been tested positive to malaria.

3) *Results*: After deleting invalid values, i.e., households with missing data, there are 4357 households with at least one child between 6 and 59 months in the survey dataset. Their socio-economic backgrounds are given in Table IV. All these households are given a binary label according to their malaria positivity: 0 being negative and 1 being positive. Univariate association results are shown in Table V and multivariate analysis results, using all the explanatory covariates are given in Table VI. No weight was used to balance the households depending on the number of children under 5 years old. The effect of the number of children is controlled in the final regression model.

*Association between malaria and socio-economic covariates*: Taken independently, most of the socio-economic covariates have a significant association with the positive households. The level of education of the mother has a negative association with the presence of malaria: households with more educated women is less likely to be positive households compared to less educated mothers. The source of drinking water also shows a similar association. Interestingly, only improved toilets are significantly associated with positive households whereas having nonimproved toilets does not seem to reduce malaria rates. Similarly, only the richest quintile of the population is associated with lower malaria presence. Contrary to what could be expected, differences between the other quintiles are not significant. This can be explained by the construction of the EZs where the prevalence is similar for the 4 poorest quintiles as shown in Table IV.

*Association between malaria and environmental covariates*: The presence of water in the buffer area, i.e., presence of water near the households' locations, is not significant. The types of environments defined above are all significant when the reference covariate is the type "sparsely built." "compacted low-rise"

TABLE IV  
SOCIO-ECONOMIC AND ENVIRONMENTAL BACKGROUND OF HOUSEHOLDS

Variables	Categories	Count	Proportion (%)	Rate+ (%)	Prevalence (%)
Wealth index	Poorest	932	21.39	27.0	19.0
	<b>Poorer</b>	939	21.55	25.0	18.0
	Middle	897	20.59	23.0	17.0
	Richer	817	18.75	26.0	18.0
	Richest	772	17.72	8.0	6.0
Type of toilets	No facility	1878	43.10	28.0	20.0
	Not improved	323	7.41	24.0	18.0
	Improved	2156	49.48	17.0	12.0
Source of drinking water	<b>No facility</b>	842	19.33	28.0	20.0
	Not improved	392	9.00	21.0	14.0
	Improved	3123	71.68	21.0	15.0
Level of education	<b>No education</b>	2934	75.68	26.0	18.0
	Primary	512	13.21	21.0	16.0
	Secondary	431	11.12	10.0	8.0
Number of children	<b>1</b>	2176	49.94	14.0	14.0
	2	1546	35.48	27.0	16.0
	3+	635	14.57	40.0	20.0
Type of environment	Bush/scrub	2506	57.52	26.0	19.0
	Scattered Trees	386	8.86	26.0	18.0
	Sparsely built	991	22.75	20.0	14.0
	Compact low-rise	474	10.88	4.0	3.0
Water in buffer	<b>0</b>	1774	44.74	22.0	16.0
	1	621	15.66	24.0	18.0
	2+	1570	39.60	22.0	15.0

Categories of the reference household are in bold.

TABLE V  
RESULTS OF UNIVARIATE LOGISTIC REGRESSIONS

		5%	95%	Odds Ratio	P-values	Significant
<b>Number of children</b> (ref. 1)	2	0.364	0.507	0.429	0.000	True
	3+	1.521	2.243	1.847	0.000	True
<b>Type of toilets</b> (ref. no facility)	Improved	0.437	0.592	0.508	0.000	True
	Not improved	0.625	1.079	0.821	0.157	False
<b>Source of drinking water</b> (ref. no facility)	Improved	0.556	0.785	0.661	0.000	True
	Not improved	0.510	0.901	0.678	0.007	True
<b>Level of education</b> (ref. no education)	Primary	0.583	0.921	0.732	0.008	True
	Secondary +	0.238	0.450	0.327	0.000	True
<b>Type of environment</b> (ref. sparsely built)	Compact low-rise	0.079	0.223	0.133	0.000	True
	Scattered trees	1.091	1.860	1.424	0.009	True
	Bush/scrub	1.169	1.669	1.397	0.000	True
<b>Wealth index</b> (ref. poorer)	Poorest	0.888	1.343	1.092	0.404	False
	Middle	0.740	1.134	0.916	0.420	False
	Richer	0.852	1.310	1.056	0.617	False
	Richest	0.198	0.358	0.266	0.000	True
<b>Water in buffer area</b> (ref. 0)	1	0.914	1.404	1.133	0.256	False
	2+	0.836	1.161	0.986	0.862	False

It explains the positivity of households in Burkina Faso, 2017–2018 according to MIS data, considering each covariate separately.

TABLE VI  
RESULTS OF A MULTIVARIATE LOGISTIC REGRESSION

		5%	95%	Odds ratio	P-values	Significant
Intercept		0.313	0.573	0.423	0.000	True
<b>Number of children</b> (ref. 1)	2	0.433	0.629	0.522	0.000	True
	3+	1.421	2.180	1.760	0.000	True
<b>Type of toilets</b> (ref. no facility)	Improved	0.620	0.918	0.754	0.005	True
	Not improved	0.622	1.152	0.846	0.288	False
<b>Source of drinking water</b> (ref. no facility)	Improved	0.650	0.965	0.792	0.021	True
	Not improved	0.567	1.061	0.776	0.112	False
<b>Level of education</b> (ref. no education)	Primary	0.863	1.427	1.110	0.416	False
	Secondary +	0.501	1.026	0.717	0.069	False
<b>Type of environment</b> (ref. sparsely built)	Compact low-rise	0.162	0.513	0.288	0.000	True
	Scattered trees	1.191	2.128	1.592	0.002	True
	Bush/scrub	1.187	1.777	1.453	0.000	True
<b>Wealth index</b> (ref. poorer)	Poorest	0.859	1.379	1.088	0.484	False
	Middle	0.788	1.279	1.004	0.975	False
	Richer	0.863	1.420	1.107	0.422	False
	Richest	0.490	1.048	0.716	0.086	False
<b>Water in buffer area</b> (ref. 0)	1	0.842	1.345	1.064	0.603	False
	2+	0.815	1.172	0.977	0.806	False

The model explains the malaria positivity of households in Burkina Faso, 2017–2018 according to MIS data taking into account all the covariates.

is associated with lower rates, “bush/scrub” with higher rates and “scattered trees” with the highest rates.

## V. DISCUSSION

This study aims to estimate the impact of the local environment on a major health issue in a Sub-Saharan African country, by the following:

- 1) characterizing local environmental data from freely available satellite data and using LCZs classification;
- 2) linking this information with demographic and health data collected from household-based surveys that provide approximated geolocalization of surveyed households;
- 3) assessing interaction between environment and health, taking into account socioeconomic factors.

In Section V-A, we discuss the LCZ mapping method using DA presented in Sections II and III, and in Section V-B, the impact of the environment on malaria depicted in Section IV.

### A. LCZ Mapping

The framework presented in this article aims to map LCZs in a country where no ground truth is available using contrastive deep learning methods applied to Sentinel-2 images. It makes use of seasonal variations in the target country to extract useful information for LCZ mapping. To focus the training step on that very country, a specific training dataset is created. For this study, we chose to work on Burkina Faso where a recent DHS was conducted and collected data as well as geolocalization ones were available for such analysis. It is worth noting that the resulting model is significantly focused on the target country and loses its adaptive capacity in other areas. For example, its classification performance on the So2Sat dataset drops by 40%. The training step has, therefore, reversed the DA problem: the model is now focused on the target country and cannot be generalized to other areas. However, creating training datasets on other countries (areas with similar climate characteristics) should allow models to accurately classify these new areas.

In addition, the distribution unbalance between the source and target LCZ classes can lead to data bias. The So2Sat dataset was indeed built on cities and around, which may limit the characterization of more rural environments. Rural areas such as “bare areas” may be under-represented in this data. However, our seasonal dataset includes data from all over Burkina Faso from very urban areas to very rural ones. This difference in the number of classes in each dataset may limit the performance of the model after DA. Moreover, Burkina Faso has specific areas that are not considered in So2Sat (e.g., Saharan desert areas in the north) and, thus, are more challenging to classify accurately. Maps generated using this data may alter the result of the multivariate analysis in rural areas, as the population in Burkina Faso is mostly rural.

The OA being 56% suggests that the accuracy of the LCZ maps could be improved. Having 17 classes, the results suggest a fine level of classification and above other state-of-the-art methods. Moreover, it still provides valuable insight into land cover and use patterns. Classification errors are smoothed by the

random selection of pixels performed in Section IV-B. Nonetheless, results on the link between environment and malaria are in line with our expectations.

### B. LCZ and Malaria

A significant association between the type of environments as defined above and malaria presence has been found for this survey even after controlling for demographic and socioeconomic factors as the number of children in the household, the type of drinking water source, the types of toilets, the level of education of the mother, and the presence of water bodies. Some covariates included in this study provide results that differ from our expectations. The poorer populations are often associated with higher malaria rates due to the lack of facilities (medical center, improved toilets) and knowledge [66], [67] but no significant association between low wealth and malaria rates was found in this study. The sampling of the survey resulted in the 4 lowest wealth quintiles’ populations (poorest, poorer, middle, and richer) having similar malaria prevalence. The only significant association has been found for the richest quintile, with lower presence of malaria in comparison to the other quintiles. The presence of water bodies in the buffer area also yields conflicting results. In general, mosquitoes are more common in regions where the humidity is high and where water can be found. In this study, due to the voluntary spatial displacement of household locations, it is not possible to compute the distance of households to small water bodies, indicated on OSM data, for example, which may indicate a greater presence of mosquitoes. As mosquitoes are unlikely to travel very high distances, the modified geolocations of EZs are not precise enough to conclude on this covariate. MIS data also provide a “proximity to water” covariate for each EZ, based on international databases such as the lakes dataset and the shoreline dataset. This covariate does not include the proximity to small water bodies as can be found in OSM data and does not give further information about the proximity to water bodies.

The associations between positive households and types of toilets or sources of drinking water lead to similar results. The associations with positive households are only significant for the “improved” covariates taking as reference the absence of facility. Indeed, improved sources of water are protected and associated with a reduced malaria presence. For the case of sources of drinking water, improved facilities lead to less stagnant water, reducing the presence of mosquitoes. However, unimproved (not protected) sources or the absence of sources leads to a greater presence of mosquitoes. The same reasoning can be applied to the types of toilets used, as improved facilities be more protected, and reduce the presence of mosquitoes. In this study, there is no evidence, however, that having unimproved facilities contributes to the reduction of malaria presence in the households, compared with households with no facilities at all.

The association results are in line with the data given by MIS and current knowledge about malaria prevalence in Burkina Faso. Malaria risk clusters can be found in the southern part of the country, mostly covered by “scattered tree” areas. However, we also found that using the climate regions (south, center, north)

does not lead to significant associations. Unlike other covariates, the characteristics of households' local environment remain invariant when including other explanatory covariates. Interestingly, the association between types of environments and malaria positivity is significant in both univariate and multivariate setups. This consistent significance of the environment characterization in both univariate and multivariate logistic regression analyses strengthens the evidence for its association with the presence of malaria, indicating a robust and persistent relationship even after controlling for other household covariates. This finding suggests that it possesses a genuine and independent effect on the outcome covariate. This significance is lost between the two rural types of environment (scattered trees and bush/scrub) when considering one of them as a reference. As indicated in Table III, those two clusters are not significantly different, it is, therefore, not surprising. Both LCZ classes can be found in similar regions of the country, and bush/scrub areas are substantially represented on the map. Likewise, another classification of the north of Burkina Faso than the unexpected sparsely built areas may lead to different clusters. In particular, more "bare soil or sand" or "low plants" areas were expected. A more accurate classification of some areas, together with the current one may refine the logistic regression results. The conclusions of this work should remain the same, as northern areas do not suffer from high malaria rates.

This study suggests that the use of the LCZ classification system is suitable for population data analysis, following the generation of a map with a 320 m resolution. Moving to a finer resolution, resulting in a more precise map (as for Global LCZ map [27]), may alter the results. LCZ classes definitions are based on the built and impervious surface fractions and then are directly linked to the total surface area covered by the input image of the model. Further studies are required to analyze the effect of the map resolution of the results on population data.

## VI. CONCLUSION

This work introduces a new strategy to link population data to up-to-date remote sensing images in Sub-Saharan Africa. First, we propose an SSDA strategy that takes advantage of the large amount of Sentinel-2 data as well as the seasonal variations in the target country. The deep neural network is taught to transfer its knowledge from the So2Sat dataset to a specific country by extracting useful features thanks to contrastive learning. Finally, we demonstrate that LCZs can be successfully linked to the population at the country level. Our work also highlights the necessity to consider local characteristics, such as seasonal variations, directly into the training of our models to overcome DA challenges. This innovative method offers new avenues for exploring population and environment interactions, especially in the growing concern of climate change.

## REFERENCES

[1] M. Bakhtsiyarava, K. Grace, and R. J. Nawrotzki, "Climate, birth weight, and agricultural livelihoods in Kenya and Mali," *Amer. J. Public Health*, vol. 108, no. S2, pp. S144–S150, Apr. 2018.

[2] S. Eissler, B. C. Thiede, and J. Strube, "Climatic variability and changing reproductive goals in Sub-Saharan Africa," *Glob. Environ. Change*, vol. 57, Jul. 2019, Art. no. 101912.

[3] T. A. Kugler et al., "People and Pixels 20 years later: The current data landscape and research trends blending population and environmental data," *Popul. Environ.*, vol. 41, no. 2, pp. 209–234, Dec. 2019.

[4] J. Mouchet and P. Carnevale, "Malaria endemicity in the various phyto-geographic and climatic areas of Africa, South of Sahara," *Southeast Asian J. Trop. Med. Public Health*, vol. 12, pp. 439–440, 1981.

[5] WHO, "World malaria report 2022," 2022. [Online]. Available: <https://www.who.int/publications/i/item/9789240064898>

[6] T. V. Ha et al., "Spatial distribution of Culex mosquito abundance and associated risk factors in Hanoi, Vietnam," *PLoS Neglected Trop. Dis.*, vol. 15, no. 6, Jun. 2021, Art. no.e0009497. [Online]. Available: <https://journals.plos.org/plosntds/article?id=10.1371/journal.pntd.0009497>

[7] R. Gibb et al., "Interactions between climate change, urban infrastructure and mobility are driving dengue emergence in Vietnam," *Nature Commun.*, vol. 14, no. 1, Dec. 2023, Art. no. 8179. [Online]. Available: <https://www.nature.com/articles/s41467-023-43954-0>

[8] O. Mudele, A. C. Frery, L. F. Zandrez, A. E. Eiras, and P. Gamba, "Modeling dengue vector population with earth observation data and a generalized linear model," *Acta Tropica*, vol. 215, 2021, Art. no. 105809. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0001706X20317228>

[9] M. Pesaresi et al., "GHS built-up grid, derived from Landsat, multitemporal (1975, 1990, 2000, 2014)," *Eur. Commission, Joint Res. Centre*, Brussels, Belgium, 2015.

[10] T. Esch et al., "Urban footprint processor—fully automated processing chain generating settlement masks from global data of the tandem-x mission," *IEEE Geosci. Remote Sens. Lett.*, vol. 10, no. 6, pp. 1617–1621, Nov. 2013.

[11] D. Zanaga et al., "Esa worldcover 10 m 2020 v100," Oct. 2021.

[12] D. Zanaga et al., "Esa worldcover 10 m 2021 v200," Oct. 2022.

[13] J. Rosentreter, R. Hagenseker, and B. Waske, "Towards large-scale mapping of local climate zones using multitemporal sentinel 2 data and convolutional neural networks," *Remote Sens. Environ.*, vol. 237, 2020, Art. no. 111472.

[14] I. D. Stewart and T. R. Oke, "Local climate zones for urban temperature studies," *Bull. Amer. Meteorological Soc.*, vol. 93, pp. 1879–1900, 2012.

[15] A. Nassar, G. Blackburn, and D. Whyatt, "Dynamics and controls of urban heat sink and island phenomena in a desert city: Development of a local climate zone scheme using remotely-sensed inputs," *Int. J. Appl. Earth Observ. Geoinf.*, vol. 51, pp. 76–90, 2016.

[16] P. J. Alexander, G. Mills, and R. Fealy, "Using LCZ data to run an urban energy balance model," *Urban Climate*, vol. 13, pp. 14–37, 2015.

[17] J. Geletič, M. Lehnert, P. Dobrovolny, and M. Žuvela-Aloise, "Spatial modelling of summer climate indices based on local climate zones: Expected changes in the future climate of Brno, Czech Republic," *Climatic Change*, vol. 152, no. 3/4, pp. 487–502, 2019.

[18] H. Wouters et al., "The efficient urban canopy dependency parametrization (sury) v1. 0 for atmospheric modelling: Description and application with the cosmo-CLM model for a belgian summer," *Geoscientific Model Develop.*, vol. 9, no. 9, pp. 3027–3054, 2016.

[19] T.-H. K. Chen et al., "Higher depression risks in medium—than in high-density urban form across denmark," *Sci. Adv.*, vol. 9, no. 21, Art. no. eadf3760, 2023. [Online]. Available: <https://www.science.org/doi/abs/10.1126/sciadv.adf3760>

[20] P. I. D. Lin et al., "Associations of local climate zones with cardiovascular disease: Findings from the US-based nationwide nurses' health study from 2000 to 2016," *ISEE Conf. Abstr.*, vol. 2023, no. 1, 2023. [Online]. Available: <https://ehp.niehs.nih.gov/doi/abs/10.1289/isee.2023.FP-051>

[21] O. Brousse et al., "Can we use local climate zones for predicting malaria prevalence across Sub-Saharan African cities?," *Environ. Res. Lett.*, vol. 15, no. 12, 2020, Art. no. 124051.

[22] B. Bechtel et al., "Mapping local climate zones for a worldwide database of the form and function of cities," *ISPRS Int. J. Geo- Inf.*, vol. 4, no. 1, pp. 199–219, 2015.

[23] C. Qiu, L. Mou, M. Schmitt, and X. X. Zhu, "Local climate zone-based urban land cover classification from multi-seasonal sentinel-2 images with a recurrent residual network," *ISPRS J. Photogrammetry Remote Sens.*, vol. 154, pp. 151–162, 2019.

- [24] C. Yoo, D. Han, J. Im, and B. Bechtel, "Comparison between convolutional neural networks and random forest for local climate zone classification in mega urban areas using landsat images," *ISPRS J. Photogrammetry Remote Sens.*, vol. 157, pp. 155–170, 2019.
- [25] L. See et al., "Developing a community-based worldwide urban morphology and materials database (WUDAPT) using remote sensing and crowdsourcing for improved urban climate modelling," in *Proc. Joint Urban Remote Sens. Event*, 2015, pp. 1–4.
- [26] M. Demuzere, B. Bechtel, A. Middel, and G. Mills, "Mapping europe into local climate zones," *PLoS One*, vol. 14, no. 4, pp. 1–27, 2019. [Online]. Available: <https://doi.org/10.1371/journal.pone.0214474>
- [27] M. Demuzere et al., "A global map of local climate zones to support Earth system modelling and urban-scale environmental science," *Earth Syst. Sci. Data*, vol. 14, no. 8, pp. 3835–3873, 2022.
- [28] X. X. Zhu et al., "So2Sat LCZ42: A benchmark data set for the classification of global local climate zones [software and data sets]," *IEEE Geosci. Remote Sens. Mag.*, vol. 8, no. 3, pp. 76–89, Sep. 2020.
- [29] X. X. Zhu et al., "The urban morphology on our planet—global perspectives from space," *Remote Sens. Environ.*, vol. 269, 2022, Art. no. 112794.
- [30] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [31] J. Chai, H. Zeng, A. Li, and E. W. Ngai, "Deep learning in computer vision: A critical review of emerging techniques and application scenarios," *Mach. Learn. Appl.*, vol. 6, 2021, Art. no. 100134.
- [32] O. A. B. Penatti, K. Nogueira, and J. A. dos Santos, "Do deep features generalize from everyday objects to remote sensing and aerial scenes domains?," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, 2015, pp. 44–51.
- [33] E. Maggiori, Y. Tarabalka, G. Charpiat, and P. Alliez, "Fully convolutional neural networks for remote sensing image classification," in *Proc. IEEE Int. Geosci. Remote Sens. Symp.*, 2016, pp. 5071–5074.
- [34] N. Yokoya et al., "Open data for global multimodal land use classification: Outcome of the 2017 IEEE GRSS data fusion contest," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 11, no. 5, pp. 1363–1377, May 2018.
- [35] C. Qiu, X. Tong, M. Schmitt, B. Bechtel, and X. X. Zhu, "Multilevel feature fusion-based CNN for local climate zone classification from sentinel-2 images: Benchmark results on the So2Sat LCZ42 Dataset," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 13, pp. 2793–2806, 2020.
- [36] C. Qiu, M. Schmitt, L. Mou, P. Ghamisi, and X. X. Zhu, "Feature importance analysis for local climate zone classification using a residual convolutional neural network with multi-source datasets," *Remote Sens.*, vol. 10, no. 10, 2018, Art. no. 1572.
- [37] A. Oliver, A. Odena, C. A. Raffel, E. D. Cubuk, and I. Goodfellow, "Realistic evaluation of deep semi-supervised learning algorithms," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31, Scotland, U.K.: Curran Associates, Inc., 2018.
- [38] D. Tuia, C. Persello, and L. Bruzzone, "Domain adaptation for the classification of remote sensing data: An overview of recent advances," *IEEE Geosci. Remote Sens. Mag.*, vol. 4, no. 2, pp. 41–57, Jun. 2016.
- [39] J. Peng, Y. Huang, W. Sun, N. Chen, Y. Ning, and Q. Du, "Domain adaptation in remote sensing image classification: A survey," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 15, pp. 9842–9859, 2022. [Online]. Available: <https://ieeexplore.ieee.org/document/9944086/>
- [40] W. Huang, Y. Shi, Z. Xiong, Q. Wang, and X. X. Zhu, "Semi-supervised bidirectional alignment for remote sensing cross-domain scene classification," *ISPRS J. Photogrammetry Remote Sens.*, vol. 195, pp. 192–203, Jan. 2023. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S0924271622003069>
- [41] K. Gao, A. Yu, X. You, C. Qiu, and B. Liu, "Prototype and context-enhanced learning for unsupervised domain adaptation semantic segmentation of remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, 2023, Art. no. 5608316. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/10120939>
- [42] X. Chen, S. Pan, and Y. Chong, "Unsupervised domain adaptation for remote sensing image semantic segmentation using region and category adaptive domain discriminator," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 4412913. [Online]. Available: <https://ieeexplore.ieee.org/document/9864210>
- [43] O. Tasar, S. L. Happy, Y. Tarabalka, and P. Alliez, "ColorMapGAN: Unsupervised domain adaptation for semantic segmentation using color mapping generative adversarial networks," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 10, pp. 7178–7193, Oct. 2020. [Online]. Available: <http://arxiv.org/abs/1907.12859>
- [44] J. Hu, L. Mou, and X. X. Zhu, "Unsupervised domain adaptation using a teacher-student network for cross-city classification of sentinel-2 images," *Int. Arch. Photogrammetry, Remote Sens. Spatial Inf. Sci.*, vol. XLIII-B2-2020, pp. 1569–1574, Aug. 2020. [Online]. Available: <https://isprs-archives.copernicus.org/articles/XLIII-B2-2020/1569/2020/>
- [45] V. Marsocci, N. Gonthier, A. Garioud, S. Scardapane, and C. Mallet, "GeoMultiTaskNet: Remote sensing unsupervised domain adaptation using geographical coordinates," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops*, Jun. 2023, pp. 2075–2085. [Online]. Available: <https://ieeexplore.ieee.org/document/10208468/>
- [46] S. Hafner, Y. Ban, and A. Nascetti, "Unsupervised domain adaptation for global urban extraction using sentinel-1 SAR and sentinel-2 MSI data," *Remote Sens. Environ.*, vol. 280, 2022, Art. no. 113192.
- [47] M. Schmitt, L. Hughes, C. Qiu, and X. Zhu, "Sen12ms—A curated dataset of georeferenced multi-spectral sentinel-1/2 imagery for deep learning and data fusion," vol. IV-2/W7, pp. 153–160, 2019.
- [48] J. Hoffman et al., "CyCADA: Cycle-consistent adversarial domain adaptation," in *Proc. Int. Conf. Mach. Learn.* 2018, pp. 1989–1998.
- [49] S. Ettetdgui, S. Abu-Hussein, and R. Giryès, "ProCST: Boosting semantic segmentation using progressive cyclic style-transfer," Aug. 2022, *arXiv:2204.11891*.
- [50] Y.-H. Chen, W.-Y. Chen, Y.-T. Chen, B.-C. Tsai, Y.-C. F. Wang, and M. Sun, "No more discrimination: Cross city adaptation of road scene segmenters," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 2011–2020. [Online]. Available: <https://ieeexplore.ieee.org/document/8237482/>
- [51] T.-H. Vu, H. Jain, M. Bucher, M. Cord, and P. Perez, "ADVENT: Adversarial entropy minimization for domain adaptation in semantic segmentation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 2512–2521. [Online]. Available: <https://ieeexplore.ieee.org/document/8954439/>
- [52] Q. Zhao, S. Lyu, B. Liu, L. Chen, and H. Zhao, "Self-training guided disentangled adaptation for cross-domain remote sensing image semantic segmentation," May 2023, *arXiv:2301.05526*.
- [53] I. Alonso, A. Sabater, D. Ferstl, L. Montesano, and A. C. Murillo, "Semi-supervised semantic segmentation with pixel-level contrastive learning from a class-wise memory bank," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2021, pp. 8199–8208. [Online]. Available: <https://ieeexplore.ieee.org/document/9710242/>
- [54] K. Sohn et al., "FixMatch: Simplifying semi-supervised learning with consistency and confidence," in *Proc. Conf. Neural Inf. Process. Syst.*, 2020, pp. 596–608.
- [55] B. Rousse, S. Lobry, G. Duthé, V. Golaz, and L. Wendling, "Seasonal semi-supervised domain adaptation for linking population studies and Local Climate Zones," in *Proc. Joint Urban Remote Sens. Event*, 2023, pp. 1–4.
- [56] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *Proc. 37th Int. Conf. Mach. Learn.*, 2020, pp. 1597–1607.
- [57] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 9726–9735.
- [58] O. Mañas, A. Lacoste, X. Giró-i-Nieto, D. Vazquez, and P. Rodríguez, "Seasonal contrast: Unsupervised pre-training from uncurated remote sensing data," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2021, pp. 9394–9403.
- [59] B. Bechtel et al., "Generating WUDAPT Level 0 data—Current status of production and evaluation," *Urban Climate*, vol. 27, pp. 24–45, 2019.
- [60] M. Demuzere, S. Hankey, G. Mills, W. Zhang, T. Lu, and B. Bechtel, "Combining expert and crowd-sourced training data to map urban form and functions for the continental US," *Sci. Data*, vol. 7, 2020, Art. no. 264.
- [61] INSD, "Enquête sur les indicateurs du paludisme au burkina FASO," 2018. [Online]. Available: <https://dhsprogram.com/pubs/pdf/MIS32/MIS32.pdf>
- [62] C. Perez-Heydrich, J. Warren, C. Burgert, and M. Emch, "Influence of demographic and health survey point displacements on raster-based analyses," *Spatial Demography*, vol. 4, pp. 1–19, 2015.
- [63] K. Grace et al., "Integrating environmental context into DHS analysis while protecting participant confidentiality: A new remote sensing method," *Popul. Develop. Rev.*, vol. 45, 2019, Art. no. 1.
- [64] J. C. Bezdek, R. Ehrlich, and W. Full, "FCM: The fuzzy c-means clustering algorithm," *Comput. Geosci.*, vol. 10, no. 2, pp. 191–203, 1984.
- [65] L. van der Maaten and G. Hinton, "Visualizing data using t-SNE," *J. Mach. Learn. Res.*, vol. 9, pp. 2579–2605, 2008.
- [66] L. S. Tusting, "Why is malaria associated with poverty? Findings from a cohort study in rural Uganda," *Infect. Dis. Poverty*, vol. 5, no. 1, Aug. 2016, Art. no. 78.
- [67] A. Degarege, K. Fennie, D. Degarege, S. Chennupati, and P. Madhivanan, "Improving socioeconomic status may reduce the burden of malaria in sub saharan africa: A systematic review and meta-analysis," *PLoS One*, vol. 14, no. 1, pp. 1–26, 2019.



**Basile Rousse** received the Ms.Eng. degree in data science from IMT Atlantique, Nantes, France, in 2021. He is currently working toward the Ph.D. degree in image processing with Université Paris Cité, Paris, France, and with the French Institute for Demographic Studies (INED), Aubervilliers, France.

His research interests include image processing, domain adaptation, deep learning, and its application on population data.



**Valérie Golaz** received the Doctorate degree in demography from SciencePo Paris, Paris, France, in 2002.

She is currently a Research Director with the French Institute for Demographic Studies (INED), Aubervilliers, France. Her research focuses on family dynamics and the interactions between population and environment. On a broader scale, she is interested in the relationship between quantitative indicators and their political application, as well as the potential discrepancies between these metrics and the actual situation.



**Sylvain Lobry** (Member, IEEE) received the Ph.D. degree in signal and image processing from Télécom Paris, Paris, France, in 2017.

He is currently an Assistant Professor with the LIPADE Laboratory, Université de Paris, Paris. From 2017 to 2020, he was a Postdoctoral Researcher with Wageningen University, Wageningen, The Netherlands. His research interests include image processing and machine learning methods, for example visual question answering, using remote sensing data.



**Laurent Wendling** received the Ph.D. degree in computer science from Université Paul Sabatier, Toulouse, France, in 1997.

He is currently a Full Professor with the LIPADE Laboratory, Université Paris Cité, Paris, France. His research interests include pattern recognition and computer vision.



**Géraldine Duthé** received the Ph.D. degree in demography from the French Museum of Natural History in Paris, Paris, France, in 2006.

She is currently a Senior Researcher with the French Institute for Demographic Studies (INED), Aubervilliers, France. Her research mainly focuses on the health transition in high-mortality countries, especially in Sub-Saharan Africa where mortality is difficult to measure due to a lack or incompleteness of traditional data.