NACALA-ROOF-MATERIAL: DRONE IMAGERY FOR ROOF DETECTION, CLASSIFICATION, AND SEGMEN TATION TO SUPPORT MOSQUITO-BORNE DISEASE RISK ASSESSMENT

Anonymous authors

Paper under double-blind review

ABSTRACT

As low-quality housing and in particular certain roof characteristics are associated with an increased risk of malaria, classification of roof types based on remote sensing imagery can support the assessment of malaria risk and thereby help prevent the disease. To support research in this area, we release the Nacala-Roof-Material dataset, which contains high-resolution drone images from Mozambique with corresponding labels delineating houses and specifying their roof types. The dataset defines a multi-task computer vision problem, comprising object detection, classification, and segmentation. In addition, we benchmarked various state-of-theart approaches on the dataset. Canonical U-Nets, YOLOv8, and a custom decoder on pretrained DINOv2 served as baselines. We show that each of the methods has its advantages but none is superior on all tasks, which highlights the potential of our dataset for future research in multi-task learning. While the tasks are closely related, accurate segmentation of objects does not necessarily imply accurate instance separation, and vice versa. We address this general issue by introducing a variant of the deep ordinal watershed (DOW) approach that additionally separates the interior of objects, allowing for improved object delineation and separation. We show that our DOW variant is a generic approach that improves the performance of both U-Net and DINOv2 backbones, leading to a better trade-off between semantic segmentation and instance segmentation.

032 033 034

008

009

010 011 012

013

015

016

017

018

019

021

024

025

026

027

028

029

031

1 INTRODUCTION

Mosquito-borne diseases refer to a group of infectious illnesses transmitted by the bite of mosquitoes. Malaria is a mosquito-borne disease caused by single-celled parasites of the Plasmodium group spread 037 through bites of infected female Anopheles mosquitoes. It ranks among the world's most severe public health problems and is a leading cause of mortality and disease in many developing countries. It is therefore crucial to improve prevention, control, and surveillance measures of malaria, particularly in 040 sub-Saharan Africa (Venkatesan, 2024; WHO, 2023). Low-quality housing built of natural materials, 041 for example, having a thatched roof of grass or palm and having cane, grass, shrub, or mud as internal 042 and external walls, is associated with an increased risk of malaria infection (Dlamini et al., 2017). 043 Sub-standard housing has more mosquito entry points and most malaria transmissions in sub-Saharan 044 Africa occur inside dwellings while the inhabitants are asleep (Tusting et al., 2020; 2017; Jatta et al., 2018; Tusting et al., 2019). Houses with metal roofs are hotter in the daytime than houses with thatched roofs. This may reduce mosquito survival and inhibit parasite development within the 046 mosquito in metal roof houses. On this basis, the proliferation of modern construction materials in 047 sub-Saharan Africa may have contributed decisively to the reduction of malaria cases (Tusting et al., 048 2019). Classification of roof characteristics thus holds potential to support malaria surveillance and control programs. Roof characteristics, such as geometry, material, and condition can be monitored using remote sensing imagery to advance risk assessment of mosquito-borne diseases and guide 051 mitigation strategies, especially when detailed health and socioeconomic data are scarce. 052

Here, we introduce the Nacala-Roof-Material drone-imagery dataset to support the development of machine learning algorithms for automated building *and* roof type mapping in low-income areas

prone to malaria risk. Our dataset is based on high-resolution drone imagery (≈ 4.4 cm) of peri-urban and rural settlements in Nacala, Mozambique. An NGO (anonymous during peer review) delineated 17954 buildings and categorized them according to five roof types, and the authors again carefully verified all annotations.

058 We define three tasks on the Nacala-Roof-Material dataset, building detection, multi-class roof type classification, and pixel-level building segmentation. While these tasks are related, closer 060 inspection reveals a misalignment between their objectives. Accurate segmentation as measured by 061 the intersection over union (IoU) does not necessarily imply accurate object separation, and vice 062 versa. For accurate detection and classification, it would be sufficient to only detect the interior of 063 an object as long as the segmented area allows to correctly classify the type. If the roofs of two 064 buildings are (almost) touching, then some segmentation may have a high IoU but could make it difficult to separate buildings for counting. This is also a common issue in other applications, e.g., 065 when studying cells in medical images (Ronneberger et al., 2015) or trees from satellite images 066 (Brandt et al., 2020; Mugabowindekwe et al., 2022)). 067

068 We benchmark three conceptually different state-of-the-art approaches on our multi-task dataset. 069 First, we evaluate YOLOv8 (Jocher et al., 2023) developed for object detection, classification, and instance segmentation. Second, we build a segmentation model based on DINOv2 (Oquab et al., 071 2024), a state-of-the-art pretrained vision transformer. Lastly, we evaluate U-Net (Ronneberger et al., 2015) a fully-convolutional encoder-decoder architecture, designed for semantic segmentation. To 072 address the potential conflicts between pixel-level segmentation and correct object separation as 073 outlined above, we propose a simple approach based on the recent work by Cheng et al. (2024), which 074 we refer to as the Deep Ordinal Watershed (DOW) method. We extend both U-Net and DINOv2 to 075 produce an additional output map that predicts the interior of objects. While the original exterior 076 segmentation map maximizes the IoU, we show that the interior map supports object separation. 077

- The main contributions of our work are the following:
- We provide the Nacala-Roof-Material dataset containing drone imagery from peri-urban and rural areas in a sub-Saharan African region. The dataset contains accurate segmentation labels for buildings, categorized into five roof types.
- Based on the dataset, we define a multi-task machine learning benchmark for binary and multiclass object detection and semantic segmentation. We implemented and benchmarked different carefully adopted baseline methods, reflecting three different approaches to address these tasks.
- 3. We propose a general and simple approach to extend models for semantic segmentation to yield good segmentation *and* object separation results.

 The data and code for reproducing the experiments are available through this anonymous url: https://osf.io/us628/?view_only=3c25a48d420f4ec7a43cb76e66e92b26. A
 project page, with link to data and code, will be setup upon acceptance.

The next section presents the Nacala-Roof-Material data, provides some background about roof types
 and risk of vector-borne diseases, and briefly discussed related datasets. Section 3 describes the deep
 learning models we evaluated with an emphasis on deep watershed methods. Experimental results
 are presented in Section 4 before we conclude.

094 095

2 NACALA-ROOF-MATERIAL DATA

096 097

099

098 2.1 BACKGROUND: HOUSING CONDITIONS AND RISK OF MOSQUITO-BORNE DISEASES

In sub-Saharan Africa, housing conditions, health outcomes, and socioeconomic status of the residents are interrelated (Gram-Hansen et al., 2019; Degarege et al., 2019; Tusting et al., 2020). As poverty is widespread, diseases are more prevalent, and data are scarce in this region, automatic profiling of housing conditions based on analysis of satellite imagery holds the potential to estimate the socioeconomic status of the inhabitants and assess the risk of disease. This may in turn support targeted public health interventions.

Mosquitoes are vectors for diseases such as malaria, dengue, Zika, West Nile fever, Chikungunya, and Yellow fever. In 2022, more than 600 000 deaths occurred due to malaria globally and out of the approximately 249 million documented cases, around 233 million occurred within the WHO African

Region, accounting for roughly around 94% of the total documented cases. The economic impact of
malaria in Sub-Saharan Africa not only impedes progress towards achieving Sustainable Development
Goal 3 (Good Health and Well-being) but also undermines efforts to attain SDG 1 (No Poverty) and
SDG 8 (Decent Work and Economic Growth) by compromising economic productivity. Extreme
weather conditions caused by climate change will likely exacerbate problems with mosquito-borne
diseases in sub-Saharan Africa, as floods are expected to increase in frequency and have been linked
to outbreaks of malaria in Africa (Githeko et al., 2000).

115 Low-quality housing increases the risk of transmission of diseases by mosquitoes, as sub-standard 116 houses have more mosquito entry points and thereby increase human exposure to infection in the 117 home (Tusting et al., 2015; Dlamini et al., 2017). Mosquito survival is lower in metal-roof houses 118 compared to thatched-roof houses due to higher daytime temperatures (Tusting et al., 2015). Most malaria transmissions in sub-Saharan Africa occur indoors at night, and poor climatic performance of 119 housing has been linked to increased malaria risk (Jatta et al., 2018). This is because elevated indoor 120 temperatures can cause discomfort for inhabitants, which may result in decreased use of mosquito 121 nets during the night. Roof materials, geometry, and conditions are critical for indoor climate, as 122 roofs comprise the primary surface exposed to the sun. Automatic classification of roof characteristics 123 thus holds potential for informing risk assessment of malaria and support targeted interventions. 124

- 125
- 126 2.2 THE NACALA-ROOF-MATERIAL DATASET

127 We gathered drone imagery of the Nacala region in Mozambique. The burden of malaria in Mozam-128 bique is approximately 10-fold the world average (number of documented cases compared to the 129 total population, Venkatesan, 2024). The data covers three informal settlements of Nacala, a city of 130 350 000 inhabitants on the northern coast of Mozambique. Aerial imagery was collected using a DJI 131 Phantom 4 Pro drone and processed using AgiSoft Metashape software. The flight height was 120 m, 132 the total flight duration was 504 minutes (the drone flight protocols are available in the supplementary 133 material). All data was recorded between October and December 2021, under a development project 134 led by an NGO and supported by the Nacala Municipal Council.

The total number of buildings in the study areas is 17 954. We distinguished five major types of roof materials in Nacala, namely metal sheet, thatch, asbestos, concrete, and no-roof, and their counts are 9776, 6428, 566, 174, and 1010, respectively. The region is mostly dominated by metal sheets and thatch roofs.

From the three informal settlements, see Figure 1, the first two areas were split into training $\mathcal{D}_{\text{train}}$, 140 validation \mathcal{D}_{val} , and test \mathcal{D}_{test} using stratified sampling. We created a square grid of 225 meters and 141 counted the roof types in these cells. Then we partitioned the cells into three sets based on the class 142 counts to achieve a similar class distribution in each set, where we prioritized the distribution of 143 minority classes (i.e., concrete and asbestos). We defined that a building only belongs to a specific 144 grid cell if its centroid falls into the cell. If a building area falls into two grid cells and those two cells 145 belong to two different sets (e.g., training and test set), we choose to have data pixels in the set where 146 the centroid of the building is placed. The remaining part of the building in the other set was masked to avoid data leaking between sets. 147

Although objects in training, validation, and test sets are from different cells, they stem from the same two areas. To evaluate the generalization to a new area without adjacent training data, we hold out the third settlement as a second test set referred to as \mathcal{D}_{ext} .

151 152

153

2.3 ANNOTATION PROCESS AND QUALITY

154 Building boundary and roof type annotations were collected and corrected in three steps. First, local 155 university students from the Nacala region and members of the NGO mapping community manually 156 traced building boundaries, and collected roof types on-site, as part of a wider survey campaign. 157 Local mapping teams used field papers and GPS tracking apps on smartphones. Secondly, the field 158 data was then digitized, corroborated by observation of the drone orthomosaics. Finally, all building 159 boundaries and roof types were verified by new observation of the drone orthomosaics conducted by the authors, and corrections made whenever necessary (almost all of the annotations were manually 160 adjusted slightly in this second step, missing annotations were added). The authors created a grid over 161 all annotations and verified all buildings in each grid cell to ensure double-checking every building.



Figure 1: (a) Visualisation of the training, validation and test sets with reference to longitude and latitude; (b) Drone imagery with labels; (c) Instance counts for each class in all sets.

For the building boundaries, the quality can be accurately measured by direct observation of the drone orthomosaics. In rare cases, under very specific lighting and imagery conditions, some uncertainties can arise between two similarly looking roof types, for instance asbestos and concrete. These cases are, however, rare and do not compromise the overall quality of the annotated data. The three-steps process – having at least two people independently looking at the images and the labeling – ensured a high label quality. An estimated 120 person-hours were required for the first of these steps and around 40 person-hours for the second.

193 194 195

196

184

185

187

188

189

190

191

192

2.4 RELATED DATASETS

197 The project "Mapping Informal Settlements in Developing Countries using Machine Learning with Noisy Annotations and Multi-resolution Multi-spectral Data" (Helber et al., 2018; Gram-Hansen et al., 2019) is most closely related to our work. They used freely available 10m/pixel resolution 199 imagery from the Sentinel-2 satellite and obtained labels for three roof types (metal, shingles, thatch) 200 from geo-located survey data provided by Afrobarometer¹. These labels are very noisy in space 201 and time. The labels are often not aligned with buildings because the geo-located coordinates were 202 distorted for privacy reasons. Furthermore, the survey questions and satellite image observations may 203 not be aligned in time. While the low spatial resolution of the Sentinel-2 imagery might allow to 204 cover large geographic regions, it makes roof type classification challenging (Helber et al., 2018). 205

There are many datasets that contain remote sensing imagery with building labels, which, however, typically do not distinguish roof types. In particular, *Open Buildings* is a freely available continentalscale building dataset covering the whole of Africa (Sirko et al., 2021). In comparison, Nacala-Roof-Material is much more focused, providing significantly higher resolution images, more accurate delineations, and in particular roof type classifications.

Alidoost and Arefi (2018) distinguish between roof types in aerial images. However, they map a rather high-income town in Germany, where they distinguish between three roof shapes common in that region (flat, gable, and hip). Another dataset for classifying roof geometry is provided by Persello et al. (2023), who distinguish 12 fine-grained details of roof geometry.

²¹⁵

¹www.afrobarometer.org



Figure 2: Baseline (top) and DOW (bottom) variants of our systems using either ResNet34 (in the case of the U-Net architectures) or DINOv2 as encoders. When using DOW, The watershed algorithm takes two segmentation masks as input, the predicted objects (level 1) and their interiors (level 2). In the two-stage approach, the classifier shown in Figure 4 is using the binary building segmentation (left). In the end-to-end setting, the roof material is predicted directly with a multi-class segmentation approach (right).

BENCHMARKED METHODS

This section presents the approaches we benchmarked on the Nacala-Roof-Material dataset. The goal is to accurately segment the buildings (as assessed by metrics based on the IoU), separate individual buildings, and classify the roof materials. As baselines, we considered U-Net (Ronneberger et al., 2015), YOLOv8 (Jocher et al., 2023), and a model performing segmentation based on DINOv2 (Oquab et al., 2024). Furthermore, we extend the U-Net and the DINOv2 based systems with the deep ordinal watershed method recently proposed by Cheng et al. (2024). These approaches are compared in two settings. In the *two-stage* setting, we first solved the building segmentation and separation tasks and afterwards classified the roof material for each detected building. In the end-to-end setting, segmentation and classification were done in parallel.

3.1 BASELINE MODELS

U-Net. The U-Net is arguably the most common architecture for semantic segmentation (Ron-neberger et al., 2015). We utilized a ResNet34 (He et al., 2016) encoder pretrained on ImageNet and a decoder similar to the original U-Net, except that we used nearest-neighbor upsampling instead of transposed convolutions (Odena et al., 2016), see Figure A.6 in the Appendix.

To identify individual instances in the semantic segmentation output map, the connected components in the map were determined (Brandt et al., 2020). To better separate individual buildings, we used a pixel-wise weight map during training that puts more emphasis on the space between buildings as already suggested by Ronneberger et al. (see Appendix A.1 for details) and commonly used in remote sensing (e.g. Brandt et al., 2020). However, this is not sufficient to separate buildings that are very close to each other or touch each other. Thus, we modified the target segmentation masks during training: Some border pixels were relabeled as background to ensure that there is a minimum gap of $n_{gap} = 7$ pixels between roofs. This modification of the target masks was only applied during training, before computing the weight map but not when calculating any performance metrics.

- YOLOv8. We trained YOLOv8 (Jocher et al., 2023), which is among the state-of-the-art methods for instance segmentation. We fine-tuned a model pretrained on the COCO dataset. While the original

YOLO architecture was designed for object detection, YOLOv8 allows for instance segmentation by integrating concepts from YOLACT (Bolya et al., 2019).

DINOv2. We benchmarked an approach based on DINOv2 (Oquab et al., 2024), a state-of-the-art pretrained vision transformer. It uses the DINOv2 *Base* model as an encoder, which is extended by a convolutional decoder. The DINOv2 output, a patch embedding with the shape of $\mathbb{R}^{1024 \times 768}$, is reshaped into feature maps of size $\mathbb{R}^{32 \times 32 \times 768}$. Then convolutional and linear upsampling layers are used on top of these feature maps as a decoder (see Appendix A.3). We used the same loss function, weighting function, training label adjustment, and training strategy as for U-Net. We froze the encoder weights and only the convolutional decoder was trained.



Figure 3: The U-Net_{DOW} architecture creates two output maps that segment objects and their interiors, respectively. The architecture differs from the baseline U-Net only in the definition of their output heads.

309 310 3.2 DEEP ORDINAL WATERSHED

306

307

308

U-Nets and the DINOv2 based method described above try to classify each pixel as accurately as
possible. However, for proper separation of objects it is sufficient – and typically preferable – if only
the interior of an object is segmented. If the border of a building can be classified as background,
even touching buildings can be separated. This reasoning leads to the deep ordinal watershed (DOW)
model introduced by Cheng et al. (2024).

316 In the watershed approach, each pixel is assigned a height and the image is viewed as a topological 317 map (Soille and Ansoult, 1990). A DOW architecture does not only predict a single segmentation 318 mask but n_{lev} feature maps for $n_{\text{lev}} + 1$ discrete height levels, $\{0, 1, \dots, n_{\text{lev}}\}$, where 0 corresponds 319 to the highest and n_{lev} to the lowest elevation. Background pixels are assumed to have level 0. The 320 Euclidean distance transformation is computed for each object, and the distances are discretized into 321 the remaining n_{lev} height levels. Target feature map $m \in \{1, \dots, n_{\text{lev}}\}$ marks all pixel with a distance level of m or higher. That is, the objects in the target feature maps get smaller with increasing m (if 322 $n_{\text{lev}} = 1$ we recover the standard U-Net). Learning the discrete height levels of pixels this way solves 323 an ordinal regression task (Frank and Hall, 2001; Cheng et al., 2008). Given the pixel heights, the

324 watershed algorithm can be applied as a post-processing step for instance segmentation (Soille and 325 Ansoult, 1990). Local minima in the elevation map define basins, each of which defines a distinct 326 object. Adopting a flooding metaphor, the watershed algorithm now floods the basins until basins 327 attributed to different starting points meet on watershed lines. Pixels attributed to the same basin 328 belong to the same object.

329 Cheng et al. (2024) employ a DOW U-Net for individual tree segmentation, however, without 330 a comparison with a standard U-Net or exploring different numbers of levels. For our task, we 331 hypothesize that a minimal number of $n_{\text{lev}} = 2$ different non-background heights is sufficient. In this 332 setting, the system outputs two masks representing the full object and its interior, respectively. Let 333 $n_{\rm pix}$ denote the difference in distance between two levels. The smallest building in our dataset has 334 size 1.463 m^2 . Thus, for the given image resolution, the number of pixels per side is approximately 335 $\sqrt{1.46}/0.044$. This suggests to define the levels such that $n_{\text{pix}} < 13$, and we picked $n_{\text{pix}} = 10$.

336 We empirically evaluated DOW variants of both our U-Net and DINOv2 based systems, see Figure 2. 337 We describe the U-Net extension in more detail in Appendix A.2, the DINOv2 based systems 338 were modified analogously. The DOW U-Net network architecture U-Net_{DOW} used in our study is 339 illustrated in Figure 3. For a comparison with a DOW U-Net with $n_{lev} = 6$ we refer to Appendix A.2 340 and Appendix B. 341

Although the approaches are related, we would like to stress the DOW method is conceptually 342 different from *deep level sets*, where deep neural networks learn a (continuous) level set function, the 343 zero-set of which defines object boundaries (Hu et al., 2017; Hatamizadeh et al., 2020), as well as 344 from predicting interior and border of an object as, for instance, done by Girard et al. (2021). 345

346 347

348

3.3 TWO-STAGE VS. END-TO-END

349 All the neural network architectures described above can directly classify the roof types of detected 350 buildings by predicting multi-class segmentation masks. However, encouraged by good classification 351 results using DINOv2 features, we also studied an alternative two-stage approach: First we seg-352 mented and separated the buildings using the algorithms described above ignoring the roof material 353 information. That is, we reduced the multi-class problem to a binary task. After that, we predicted the roof material of each detected building. We used DINOv2 to processes a 448×448 patch centered 354 355 around each building, see Figure 4. The output of DINOv2, a patch embedding with the shape of $\mathbb{R}^{1024 \times 768}$ was reshaped into feature maps of $\mathbb{R}^{32 \times 32 \times 768}$. These feature maps were then upsampled 356 to the input patch size, masked with a target binary building mask, and average pooling was applied to 357 obtain the final feature vector for the building. Standard machine learning classifiers were applied to 358 this embedding to predict the roof material, where linear probing gave the best results (see Appendix 359 B.2 for a comparison of different classifiers). 360

The two-stage methods are referred to as U-Net, DINOv2, U-Net_{DOW}, and DINOv2_{DOW}, and the 361 corresponding end-to-end methods are denoted by U-Net_{Multi}, DINOv2_{DOW-Multi}, U-Net_{DOW-Multi}, and 362 DINOv2_{DOW-Multi}, see Figure 2 for an overview. 363



Figure 4: The architecture of the DINOv2 based roof material classifier used in the two-stage setting. A classifier (e.g., logistic regression) is applied to the resulting feature vector.

376 377

375

4 EXPERIMENTS AND RESULTS 379

4.1 EXPERIMENTAL SETUP

380

381

382 All models, except for YOLOv8 where we followed its original training protocol, were trained using cross-entropy loss with pixel-wise weighting. We employed the AdamW optimizer (Loshchilov and Hutter, 2019) with an initial learning rate of 0.0003. All models were trained for 300 epochs, utilizing 384 a learning rate scheduler that decreased the learning rate by a factor of 10 every 50 epochs. The final 385 weight configuration and hyperparameters for each model were selected based on the highest IoU 386 score achieved on the validation dataset. The hyperparameters of the U-Net were chosen by observing 387 results on the validation dataset in an iterative process. The high training speed of YOLOv8 allowed 388 for more systematic model selection: We applied the genetic algorithm that comes as part of the 389 YOLOv8 framework for hyperparameter optimization (Jocher et al., 2023). The input patch sizes for 390 the U-Net variants, YOLOv8, and DINOv2 models were 512, 640, and 448, respectively. 391

3924.2 EVALUATION METRICS

The semantic segmentation performance was evaluated by the IoU. We considered both the IoU of the binary building segmentation and the mean IoU for class-specific roof segmentation. The roof materials concrete and asbestos are very rare. While \mathcal{D}_{train} , \mathcal{D}_{val} , and \mathcal{D}_{test} are stratified samples containing all classes, the spatially distinct data \mathcal{D}_{ext} does not contain any example of the two roof types, see Figure 1. To allow for a better comparison between the two test sets and to see the effect of the rare classes on the macro-averaged mean IoU, we provide the mean IoU of the three main classes (mIoU³) alongside with the mean IoU of all five classes (mIoU⁵).

401 Instance segmentation was assessed using the AP_{50} score, that is, the average precision evaluated 402 at an IoU threshold of 0.5 (Everingham et al., 2010; Lin et al., 2014). We evaluated the AP for both the predictions of building instances and the predictions of multi-class roof type instances (i.e., 403 in the latter case an object is only detected if the roof material is correctly identified). Similar to 404 IoU, mAP₅₀ and mAP₅₀ denote the mean AP₅₀ over three and five classes. To estimate the average 405 precision, a confidence score is required for each building segment. The confidence score of binary 406 and multi-class segmentation models was obtained by interpreting the neural networks' outputs 407 as probability distributions over classes and calculating the mean probability of belonging to the 408 predicted class over all pixel within a predicted segment. The exception was YOLOv8, which 409 provides its own confidence score. When a classifier using DINOv2 features was used on top of 410 binary segmentation models, the confidence score was derived from the canonical probability score 411 of the classifier. Additional metrics, AP₅₀₋₉₅ and TP_s, are shown in Appendix B. Information on the 412 computer resources is provided in Appendix A.4.

413 414

415

4.3 **RESULTS AND DISCUSSION**

Our experimental results on \mathcal{D}_{test} and \mathcal{D}_{ext} are presented in Table 1, additional details can be found in 416 Appendix B. All metrics on the test sets were computed on raw images instead of patches to avoid 417 artifacts when splitting images. We report averages over five trials on the corresponding standard 418 deviations. The methods reached AP₅₀ and IoU values on the spatially separated test set of up to 419 0.963 and 0.880, respectively. Thus the tasks can be solved with an accuracies high enough for 420 subsequent analysis while still leaving room for improvement. Detecting thatch roofs is particularly 421 relevant, as they are associated with an increased malaria risk (Tusting et al., 2019), and these roofs 422 can be identified particularly well, see Table B.4 in the Appendix. 423

423

Comparison of methods. When comparing the different approaches, we find that there is no 425 method that was better than the others across all metrics. The U-Nets and YOLOv8 did well on their 426 home grounds: YOLOv8 gave good object detection results (e.g., the best AP_{50} scores), while the 427 U-Nets performed well for semantic segmentation as measured by IoU. DINOv2 combined with a 428 simple decoder was also competitive. Exemplary results are shown in Figure 5. As could be expected, classifying the minority roof types asbestos and especially concrete (which resembles concreted 429 background areas) was most difficult, in particular for end-to-end YOLOv8, see Table B.4. YOLOv8 430 had the tendency to produce artefacts when applied to the larger images. This is one of the reasons 431 for its lower IoU score.

-101	
433	Table 1: Benchmarking results on the Nacala-Roof-Material dataset. The table reports averages over
434	five trials \pm standard deviations. The upper five models were trained in the two-stage setting. The
435	lower half of the models was trained in the end-to-end setting, where multi-class classification is
436	performed together with the segmentation as indicated by the subscript Multi. Models that used the
/37	DOW extension are indicated by the subscript DOW . IoU and AP_{50} were computed on the binary
400	output, where the predictions of multi-class models were binarized. mIoU and mAP ₅₀ are macro
430	averages, the superscipts indicate whether the averaging was done over all five classes or over the
439	three frequent roof types. Results for individual roof types can be found in Appendix B.
440	

			$\mathcal D$	test		$\mathcal{D}_{\mathrm{ext}}$					
		pixel level			object level			pixel level		object level	
Model Name	IoU	mIoU ³	mIoU ⁵	AP ₅₀	mAP ₅₀ ³	mAP ₅₀ ⁵	IoU	mIoU ³	AP ₅₀	mAP ₅₀ ³	
YOLOv8	$\begin{array}{c} 0.866 \\ \pm \ 0.012 \end{array}$	$\begin{array}{c} 0.713 \\ \pm \ 0.019 \end{array}$	$\begin{array}{c} 0.568 \\ \pm \ 0.015 \end{array}$	0.941 ± 0.003	$\begin{array}{c} 0.815 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} 0.698 \\ \pm \ 0.018 \end{array}$	$\begin{array}{c} 0.896 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.761 \\ \pm \ 0.006 \end{array}$	0.963 ± 0.005	$\begin{array}{c} 0.846 \\ \pm 0.008 \end{array}$	
U-Net	$\begin{array}{c} \textbf{0.895} \\ \pm \ \textbf{0.003} \end{array}$	$\begin{array}{c} 0.757 \\ \pm \ 0.024 \end{array}$	$\begin{array}{c} 0.570 \\ \pm \ 0.016 \end{array}$	$\begin{array}{c} 0.910 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.810 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.688 \\ \pm \ 0.014 \end{array}$	$\begin{array}{c} 0.909 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.748 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} 0.929 \\ \pm \ 0.000 \end{array}$	$\begin{array}{c} 0.787 \\ \pm \ 0.011 \end{array}$	
U-Net _{DOW}	$\begin{array}{c} \textbf{0.895} \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.775 \\ \pm \ 0.013 \end{array}$	$\begin{array}{c} 0.577 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} 0.935 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.836 \\ \pm \ 0.005 \end{array}$	0.730 ± 0.011	$\begin{array}{c} \textbf{0.911} \\ \pm \ \textbf{0.002} \end{array}$	$\begin{array}{c} 0.764 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.947 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.812 \\ \pm \ 0.008 \end{array}$	
DINOv2	$\begin{array}{c} 0.880 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.741 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.549 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.881 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.783 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.673 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} 0.904 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.718 \\ \pm \ 0.015 \end{array}$	$\begin{array}{c} 0.922 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.804 \\ \pm \ 0.008 \end{array}$	
DINOv2 _{DOW}	$\begin{array}{c} 0.881 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.761 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.566 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.931 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.836 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.718 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} 0.904 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.767 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.961 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.861 \\ \pm \ 0.009 \end{array}$	
YOLOv8 _{Multi}	$\begin{array}{c} 0.824 \\ \pm \ 0.023 \end{array}$	$\begin{array}{c} 0.708 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.550 \\ \pm \ 0.017 \end{array}$	$\begin{array}{c} 0.910 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.816 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} 0.597 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} 0.885 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.785 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.948 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.849 \\ \pm \ 0.015 \end{array}$	
U-Net _{Multi}	$\begin{array}{c} 0.879 \\ \pm \ 0.012 \end{array}$	$\begin{array}{c} 0.783 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.634 \\ \pm \ 0.024 \end{array}$	$\begin{array}{c} 0.924 \\ \pm \ 0.004 \end{array}$	0.850 ± 0.011	$\begin{array}{c} 0.716 \\ \pm \ 0.018 \end{array}$	$\begin{array}{c} 0.903 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.805 \\ \pm \ 0.020 \end{array}$	$\begin{array}{c} 0.943 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.844 \\ \pm \ 0.039 \end{array}$	
U-Net _{DOW-Multi}	$\begin{array}{c} 0.893 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.779 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} 0.674 \\ \pm \ 0.041 \end{array}$	$\begin{array}{c} 0.933 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.838 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.710 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.906 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.798 \\ \pm \ 0.012 \end{array}$	$\begin{array}{c} 0.944 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.830 \\ \pm \ 0.017 \end{array}$	
DINOv2 _{Multi}	$\begin{array}{c} 0.879 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.768 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.694 \\ \pm \ 0.013 \end{array}$	$\begin{array}{c} 0.894 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.810 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.680 \\ \pm \ 0.019 \end{array}$	$\begin{array}{c} 0.899 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.819 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.940 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.871 \\ \pm \ 0.015 \end{array}$	
$DINOv2_{DOW-Multi}$	$\begin{array}{c} 0.884 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.783 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.732 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.918 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.810 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} 0.701 \\ \pm \ 0.027 \end{array}$	$\begin{array}{c} 0.901 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.823 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} 0.944 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.843 \\ \pm \ 0.012 \end{array}$	

463

464

467

468

469

432

462 In general, the DOW extension improved both U-Nets and DINOv2 based architectures. Comparing DINOv2 with DINOv2_{DOW} and U-Net with U-Net_{DOW}, the DOW variants were better in all ten performance indices (except for IoU on \mathcal{D}_{test} where U-Net and U-Net_{DOW} gave the same result). Comparing DINOv2_{Multi} with DINOv2_{DOW-Multi}, the latter was better in all indicators except mAP $_{50}^3$ 465 on \mathcal{D}_{ext} . Only for U-Net_{DOW-Multi} the results were mixed, using DOW gave lower values for five 466 indices and higher values for the other half. Overall, the DOW extension had a statistically significant positive effect on the object separation as intended. If we pool all 20 DOW trials and compare with the corresponding trials predicting a single mask, then the AP50 improved significantly (two-sided Wilcoxon rank sum test, p < 0.001) while the difference in IoU was not significant (p > 0.05). 470

471 **Computational requirements.** Since compute resources might be limited for researchers interested 472 in this application, we analyze the runtime for deployment of these models (see Table A.2). While all 473 methods run in a reasonable time, the end-to-end approach is faster than the two-stage approach, with 474 U-Net_{Multi} being the fastest and DINOv2_{DOW} the slowest model. For example, mapping the entire 475 city of Nacala (31910 ha), U-Net_{Multi} would take approximately 6.36 GPU hours on a single AMD 476 MI250X GPU with 64 GB VRAM.

478 **Dataset size.** We ran the experiments with an 80% stratified subset of the data, see Table B.8 in the 479 supplementary material. The results changed only very slightly, indicating that our training dataset is 480 representative and large enough for the defined task given the regional constraints.

481

477

482 **Limitations.** The Nacala-Roof-Material dataset is not a large-scale dataset by current standards 483 and it is restricted to a single region. However, considering the proliferation of low-cost drone technologies, high-resolution geospatial surveying is becoming increasingly affordable and common 484 in sub-Saharan Africa. Accordingly, similar but unlabelled data will likely become available in the 485 coming years at large scale, which makes it important to develop methods to make good use of these



Figure 5: Exemplary predictions on \mathcal{D}_{test} by different models. The predictions are polygonized and colored by class. The roof types with few training examples, asbestos and concrete, are particularly difficult, see bottom row.

data now. The Nacala-Roof-Material dataset covering informal settlements is a good example for the target areas of our risk disease monitoring and prevention research. In this context, Mozambique is particularly relevant because the country suffers from a high malaria incidence rate (Venkatesan, 2024). The second test set allows for testing generalization in an area geographically separated from the main training/test/validation data. In general, we would argue that there is a need for medium size benchmark datasets such as the Nacala-Roof-Material data to support equity in machine learning research, as we need benchmarks that can be utilized by researchers with limited compute resources.

5 CONCLUSIONS

508

509

510

511 512 513

514

515

516

517

518

519 520 521

522

523 The Nacala-Roof-Material dataset contains high-resolution drone imagery from informal settlements 524 in Mozambique, where buildings and their roof material were carefully annotated. We curated 525 the dataset as part of an intercontinental and interdisciplinary research project on risk assessment of mosquito-borne diseases, especially malaria, with the goal to predict risk maps and to develop 526 and support measures for risk reduction. From a methodological perspective, the dataset defines 527 a multi-task problem. We are interested in accurate semantic segmentation to determine the roof 528 areas and also in identifying the individual buildings and classifying their roof types. Thus, the 529 dataset adds to the landscape of computer vision benchmarks by providing a relevant resource for the 530 development and evaluation of frameworks that strive at solving semantic segmentation as well as 531 object detection and classification simultaneously with a high accuracy. For example, working on the 532 Nacala-Roof-Material data has led us to the proposed deep ordinal watershed (DOW) approach, a 533 reduced variant of the method described by Cheng et al. (2024). This variant method first segments 534 objects along with their interiors into two elevation levels and then performs a watershed segmentation to separate objects. The DOW idea is applicable beyond the Nacala-Roof-Material data, on which it 536 improved both the standard U-Net architectures as well as a system based on DINOv2 features for segmentation. Implementations of all algorithms will be publicly available together with the data. With the Nacala-Roof-Material dataset, we invite the machine learning community to develop new 538 approaches for interpreting high-resolution drone images that can ultimately support risk assessments of vector-borne diseases.

6 REPRODUCIBILITY STATEMENT

The data and code for reproducing the experiments are available anonymously through https: //osf.io/us628/?view_only=3c25a48d420f4ec7a43cb76e66e92b26. All material will be made freely available on an official project homepage upon acceptance.

7 ETHICS STATEMENT

We did not identify any ethical issues; we refer to the data sheet in Appendix C.

549 550

552

553

554 555

556

561

562

563 564

565

566 567

569

574

575

576

577

578

579

580

582

583

584

585

586

588

589

590

591

540

541 542

543

544

545 546

547 548

551 **REFERENCES**

- F. Alidoost and H. A. Arefi. CNN-based approach for automatic building detection and recognition of roof types using a single aerial image. *PFG Journal of Photogrammetry, Remote Sensing and Geoinformation Science*, 86:235–248, 2018.
- D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee. YOLACT: Real-time instance segmentation. In *International Conference on Computer Vision (ICCV)*, pages 9157–9166, 2019.
- M. Brandt, C. J. Tucker, A. Kariryaa, K. Rasmussen, C. Abel, J. Small, J. Chave, L. V. Rasmussen, P. Hiernaux, A. A. Diouf, L. Kergoat, O. Mertz, C. Igel, F. Gieseke, J. Schöning, S. Li, K. Melocik, J. Meyer, S. Sinno, E. Romero, E. Glennie, A. Montagu, M. Dendoncker, and R. Fensholt. An unexpectedly large count of trees in the western Sahara and Sahel. *Nature*, 587:78–82, 2020.
 - J. Cheng, Z. Wang, and G. Pollastri. A neural network approach to ordinal regression. In *International Joint Conference on Neural Networks (IJCNN)*, pages 1279–1284. IEEE, 2008.
 - Y. Cheng, S. Oehmcke, M. Brandt, L. Rosenthal, A. Das, A. Vrieling, S. Saatchi, F. Wagner, M. Mugabowindekwe, W. Verbruggen, C. Beier, and S. Horion. Scattered tree death contributes to substantial forest loss in California. *Nature Communications*, 15:641, 2024.
 - 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*, 14 (1), 2019.
- N. Dlamini, M. S. Hsiang, N. Ntshalintshali, D. Pindolia, R. Allen, N. Nhlabathi, J. Novotny, M.-S. Kang Dufour, A. Midekisa, R. Gosling, A. LeMenach, J. Cohen, G. Dorsey, B. Greenhouse, and S. Kunene. Low-quality housing is associated with increased risk of malaria infection: A national population-based study from the low transmission setting of Swaziland. *Open Forum Infectious Diseases*, 4(2):ofx071, 2017.
 - M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The Pascal visual object classes (VOC) challenge. *International Journal of Computer Vision*, 88:303–338, 2010.
 - E. Frank and M. Hall. A simple approach to ordinal classification. In *European Conference on Machine Learning* (*ECML*), pages 145–156. Springer, 2001.
 - N. Girard, D. Smirnov, J. Solomon, and Y. Tarabalka. Polygonal building extraction by frame field learning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5891–5900, 2021.
 - A. K. Githeko, S. W. Lindsay, U. E. Confalonieri, and J. A. Patz. Climate change and vector-borne diseases: a regional analysis. *Bulletin of the World Health Organization*, 78(9):1136 1147, 2000.
 - B. Gram-Hansen, P. Helber, I. Varatharajan, F. Azam, A. Coca-Castro, V. Kopackova, and P. Bilinski. Mapping informal settlements in developing countries using machine learning and low resolution multi-spectral data. In AAAI/ACM Conference on AI, Ethics, and Society, 2019.
 - A. Hatamizadeh, D. Sengupta, and D. Terzopoulos. End-to-end trainable deep active contour models for automated image segmentation: Delineating buildings in aerial imagery. In *European Conference on Computer Vision (ECCV)*, page 730–746. Springer, 2020. ISBN 978-3-030-58609-6.
 - K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- 592 P. Helber, B. Gram-Hansen, I. Varatharajan, F. Azam, A. Coca-Castro, V. Kopackova, and P. Bilinski. Generating
 593 material maps to map informal settlements. In *NeurIPS 2018 Workshop on Machine Learning for the Developing World*, 2018.

- 594 P. Hu, B. Shuai, J. Liu, and G. Wang. Deep level sets for salient object detection. In Computer Vision and Pattern Recognition (CVPR), pages 2300-2309, 2017.
 - E. Jatta, M. Jawara, J. Bradley, D. Jeffries, B. Kandeh, J. B. Knudsen, A. L. Wilson, M. Pinder, U. D'Alessandro, and S. W. Lindsay. How house design affects malaria mosquito density, temperature, and relative humidity: an experimental study in rural Gambia. The Lancet Planetary Health, 2(11):e498-e508, 2018.
 - G. Jocher, A. Chaurasia, and J. Qiu. Ultralytics YOLO. https://github.com/ultralytics/ ultralytics, 2023. Accessed 20/03/2024.
 - T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In European Conference on Computer Vision (ECCV), pages 740–755. Springer, 2014.
 - I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In International Conference on Leaning Representations (ICLR), 2019.
- M. Mugabowindekwe, M. Brandt, J. Chave, F. Reiner, D. L. Skole, A. Kariryaa, C. Igel, P. Hiernaux, P. Ciais, 608 O. Mertz, X. Tong, S. Li, G. Rwanyiziri, T. Dushimiyimana, A. Ndoli, V. Uwizeyimana, J.-P. Barnekow 609 Lillesø, F. Gieseke, C. J. Tucker, S. Saatchi, and R. Fensholt. Nation-wide mapping of tree-level aboveground 610 carbon stocks in Rwanda. Nature Climate Change, pages 91-97, 2022.
- A. Odena, V. Dumoulin, and C. Olah. Deconvolution and checkerboard artifacts. Distill, 2016. 612
- 613 M. Oquab, T. Darcet, T. Moutakanni, H. V. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. HAZIZA, F. Massa, 614 A. El-Nouby, M. Assran, N. Ballas, W. Galuba, R. Howes, P.-Y. Huang, S.-W. Li, I. Misra, M. Rabbat, V. Sharma, G. Synnaeve, H. Xu, H. Jegou, J. Mairal, P. Labatut, A. Joulin, and P. Bojanowski. DINOv2: 615 Learning robust visual features without supervision. Transactions on Machine Learning Research, 2024. 616
- 617 C. Persello, R. Hänsch, G. Vivone, K. Chen, Z. Yan, D. Tang, H. Huang, M. Schmitt, and X. Sun. 2023 IEEE 618 GRSS data fusion contest: Large-scale fine-grained building classification for semantic urban reconstruction [technical committees]. IEEE Geoscience and Remote Sensing Magazine, 11(1):94–97, 2023. 619
- 620 O. Ronneberger, P. Fischer, and T. Brox. U-Net: Convolutional networks for biomedical image segmentation. In International Conference on Medical Image Computing and Computer-assisted Intervention (MICCAI), 622 pages 234-241. Springer, 2015.
 - W. Sirko, S. Kashubin, M. Ritter, A. Annkah, Y. S. E. Bouchareb, Y. Dauphin, D. Keysers, M. Neumann, M. Cisse, and J. Quinn. Continental-scale building detection from high resolution satellite imagery. arXiv preprint arXiv:2107.12283, 2021.
 - P. J. Soille and M. M. Ansoult. Automated basin delineation from digital elevation models using mathematical morphology. Signal Processing, 20(2):171-182, 1990.
 - L. S. Tusting, M. Ippolito, B. Willey, I. Kleinschmidt, G. Dorsey, R. Gosling, and S. Lindsay. The evidence for improving housing to reduce malaria: A systematic review and meta-analysis. *Malaria Journal*, 14, 12 2015.
- 631 L. S. Tusting, C. Bottomley, H. Gibson, I. Kleinschmidt, A. J. Tatem, S. W. Lindsay, and P. W. Gething. Housing 632 improvements and malaria risk in sub-Saharan Africa: a multi-country analysis of survey data. PLoS Medicine, 633 14(2):e1002234, 2017.
 - L. S. Tusting, D. Bisanzio, G. Alabaster, E. Cameron, R. E. Cibulskis, M. Davies, S. Flaxman, H. S. Gibson, J. B. T. Knudsen, C. M. Mbogo, F. O. Okumu, L. von Seidlein, D. J. Weiss, S. W. Lindsay, P. W. Gething, and S. Bhatt. Mapping changes in housing in sub-Saharan Africa from 2000 to 2015. Nature, 568(7752):391 -394, 2019.
 - L. S. Tusting, P. Gething, H. Gibson, B. Greenwood, J. Knudsen, S. Lindsay, and S. Bhatt. Housing and child health in sub-Saharan Africa: A cross-sectional analysis. PLoS Medicine, 17:e1003055, 03 2020.
- P. Venkatesan. The 2023 WHO world malaria report. The Lancet Microbe, 5(3):e214, 2024. 641
- 642 WHO. World malaria report 2023. https://www.who.int/publications/i/item/ 643 9789240086173, 2023. (Accessed on 05/16/2024).
- 644

597

598

600

601 602

603

604 605

606

607

611

621

623

624

625

626

627

628 629

630

634

635

636

637

638

639

- 645
- 646 647

649

682

683 684

685 686

687 688

689

690

694

А

650 A.1 U-Net 651 652 653 654 655 1x1 Conv, K, /1 656 657 Decoder block Map C=16 658 Input Image 659 ple, SF=2 7x7 Conv, 64, /2, +3 ReLU Decoder block Previous input 661 C=32 х 3x3 Maxpool, /2 662 3x3 Conv, C, S le. SF=2 idual 663 Residual block х3 ReLU C=64, S=1 Res 3x3 Conv, C, S den Decoder block 665 C=64 F(x) 666 Residual block + x 1 C=128, S=2 667 Upsample, SF=2 Residual block 668 y=F(x)+xх3 C=128, S=1 Out Channels (C) = ? 669 Decoder block Stride (S) = ?C=128 670 **Residual Block** Residual block 671 x 1 C=256, S=2 672 Upsample, SF=2 ReLU Residual block 673 x 5 3x3 Conv, C, /1 C=256, S=1 ReLU 674 Decoder block C=256 3x3 Conv, C, /1 675 Residual block 676 x 1 C=512, S=2 Concatenate mple, SF=2 677 Residual block x 2 Out Channels (C) = ? 678 C=512, S=1 Decoder Block 679 680 Figure A.6: Basic U-Net architecture 681

DETAILS ON MODELS AND TRAINING PROCEDURE

The basic U-Net architecture we used is shown in Figure A.6.

During training, the loss of each background pixel x is multiplicatively weighted by w(x) defined as

$$w(x) = w_0 \cdot \exp\left(-\frac{(d_1(x) + d_2(x))^2}{2\sigma^2}\right)$$
 (A.1)

following Ronneberger et al. (2015). Here, $d_1(x)$ denotes the distance to the border of the nearest segment, and $d_2(x)$ is the distance to the border of the second nearest segment. We set $w_0 = 10$ and $\sigma = 5$ according to Ronneberger et al. (2015).

⁶⁹¹ During training, we modified the target masks to ensure that $d_1(x) + d_2(x) \ge n_{gap} = 7$ for each background pixel x before we computed the weights w(x).

A.2 DEEP ORDINAL WATERSHED U-NETS

We considered a stripped down version of the DOW U-Net proposed by Cheng et al. (2024) and set the number of elevation levels to $n_{lev} = 2$. The architecture of the resulting DOW network is depicted in main text Figure 3, which extends the basic U-Net architecture shown in Figure A.6. In contrast to the original U-Net, the DOW model has two heads. One is predicting an object's area, while the other predicts its interior. The interior is defined by removing pixels within a 10-pixel distance from the border of the building segment. Each head comprises a convolutional layer, batch normalization, ReLU activation, and finally a pointwise convolutional layer with outputs equal to the number of classes. While the first head had filters of size 3×3 in its first convolutional layer, the second head for the interior used 64 filters. The class label of an object was derived from the second head. If no interior was predicted, which can happen in the case of small objects, the output from the first head defined the class.

We compared this DOW variant, referred to as U-Net_{DOW}, to the original DOW with several elevation levels, in which the levels are added to the standard U-Net architecture (Figure A.6) simply by increasing the number of output masks. We considered $n_{\text{lev}} = 6$ discrete height levels and accordingly refer to the model as U-Net_{DOW-6}. The pixel margin n_{pix} for each height level was determined experimentally by testing $n_{\text{pix}} \in \{1, 3, 5, 7, 9, 11, 13, 15\}$ on validation data, leading to $n_{\text{pix}} = 5$ for U-Net_{DOW-6}. An experimental comparison of U-Net_{DOW} and U-Net_{DOW-6} can be found in the extended results in Section B in the appendix.

713 714

715

A.3 SEGMENTATION AND CLASSIFICATION USING DINOV2

The segmentation architecture based on DINOv2 is illustrated in Figure A.7. We refer to it simply as
 DINOv2. From this architecture, we derived DINOv2_{DOW} in the same way as we extended U-Net to
 U-Net_{DOW}.



A.4 COMPUTE RESOURCES

All experiments were conducted on AMD MI250X GPUs provided by LUMI². A total of 8550 GPU hours were used for the project, including preliminary experiments not included in the paper. The computation time for training semantic segmentation model was approximately 20 hours for 300 epochs when the entire data were loaded to GPU memory.

744 745 746

747

748

739

740

741

742

743

B ADDITIONAL RESULTS

B.1 DETAILED RESULTS FOR DIFFERENT ROOF MATERIALS

Additional results on \mathcal{D}_{test} are presented in Table B.3 and Table B.4. The tables report the IoU scores for the individual roof material classes. They also show the true positive rates TP_s in addition to the AP₅₀ the AP₅₀₋₉₅. The AP₅₀₋₉₅ is defined as the mean AP over IoU thresholds from 50% to 95% with an interval of 5%. The mean of AP₅₀₋₉₅ over all classes is mAP₅₀₋₉₅. TP_s are the number of segments that overlap with ground truth segments with a minimum IoU of 0.5, we used this metric to assess the counting of buildings.

²https://lumi-supercomputer.eu

Table A.2: Computation time deploying the models using a single AMD MI250X GPU with 64 GB
 VRAM. The prediction time includes segmentation, polygonisation, and post processing in case of
 the DOW method. The DOW method takes longer since its post-processing step currently runs on
 the CPU. We used our research code without any additional optimization for speed. We used our
 research code without optimization for speed.

101		
762	Model	Prediction time on \mathcal{D}_{test} (minutes)
763		Two-stage approach
764		3.84
765	I OLOVO U Net	3.64
766	U-Netrow	5.55
767	DINOv2	6.81
768	DINOv2 _{DOW}	11 10
769		
770		End-to-end approach
771	YOLOv8 _{Multi}	1.59
772	U-Net _{Multi}	1.52
773	U-Net _{DOW-Mult}	ti 4.62
77/	DINOv2 _{Multi}	5.36
775	DINOv2 _{DOW-N}	Aulti 6.52
115		

Beyond the performance metrics already discussed, we have included the results for U-Net_{DOW-6} as
 described in Section A.2 in the appendix, showing that the two DOW architectures perform on par.

The corresponding results on \mathcal{D}_{ext} are given in Table B.5 and Table B.6 The mean IoU in Table B.5, and mAP₅₀ and mAP₅₀₋₉₅ in Table B.6 estimated on only four classes as there are no asbestos roofs in \mathcal{D}_{ext} . Also, there are only two buildings of concrete found in \mathcal{D}_{ext} and these two buildings were not identified from any of the experimental models, so results for the concrete class were not added to both tables.

789 790

- 791
- 792

793

794 795

B.2 PERFORMANCE OF DIFFERENT CLASSIFIERS

- 796 797
- 798

799 In the two-stage approach, we used a classifier based on DINOv2 features, as described in Section 3.3 800 and illustrated in Figure 4. The input representation was fixed and was processed by standard classification algorithms. We compared linear probing based on logistic regression with L_2 -regularization 801 and k-nearest neighbours (kNN) classification trained on our data. For evaluating the classifiers 802 and tuning their hyperparameters, we combined the training and validation data and performed 803 10-fold cross-validation (CV) with F1-score as performance metric. The best CV results gave logistic 804 regression with L_2 -regularization, and this model was used for all subsequent two-stage experiments, 805 see Table B.2. 806

We also performed an ablation study to show the importance of the masking and the upsampling in our architecture shown in Figure 4. The results are also depicted in Table B.2. When we omitted the masking and considered all features, the results got considerably worse. If we omitted the upsampling of the DINOv2 output and downsampled the masks instead, the performance also slightly dropped.

811 Тя

Table B.3: Pixel-level accuracies on \mathcal{D}_{test} . IoU refers to the IoU computed on the binary outputs, where the predictions of multi-class models were binarized. mIoU⁵ refers to the macro average of the IoUs for the individual classes. The subscript *Multi* indicates the end-to-end setting.

		IoU-	Score of eac	ch class			
Model Name	Metal Sheet	Thatch	Asbestos	Concrete	No Roof	mIoU ⁵	IoU
YOLOv8	$\begin{array}{c} 0.807 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.852 \\ \pm \ 0.038 \end{array}$	$\begin{array}{c} 0.450 \\ \pm \ 0.023 \end{array}$	$\begin{array}{c} 0.250 \\ \pm \ 0.027 \end{array}$	$\begin{array}{c} 0.480 \\ \pm \ 0.034 \end{array}$	$\begin{array}{c} 0.568 \\ \pm \ 0.015 \end{array}$	0.866 ± 0.012
DINOv2	$\begin{array}{c} 0.796 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.855 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.335 \\ \pm \ 0.019 \end{array}$	$\begin{array}{c} 0.189 \\ \pm \ 0.017 \end{array}$	$\begin{array}{c} 0.571 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.549 \\ \pm \ 0.005 \end{array}$	0.880 ± 0.00
DINOv2 _{DOW}	$\begin{array}{c} 0.814 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.868 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.351 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.195 \\ \pm \ 0.019 \end{array}$	$\begin{array}{c} 0.602 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.566 \\ \pm \ 0.004 \end{array}$	0.881 ± 0.00
U-Net	$\begin{array}{c} 0.813 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} 0.881 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.408 \\ \pm \ 0.012 \end{array}$	$\begin{array}{c} 0.171 \\ \pm \ 0.021 \end{array}$	$\begin{array}{c} 0.577 \\ \pm \ 0.073 \end{array}$	$\begin{array}{c} 0.570 \\ \pm \ 0.016 \end{array}$	0.89 ± 0.00
U-Net _{DOW}	$\begin{array}{c} 0.824 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.879 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.384 \\ \pm \ 0.042 \end{array}$	$\begin{array}{c} 0.174 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.623 \\ \pm \ 0.028 \end{array}$	$\begin{array}{c} 0.577 \\ \pm \ 0.009 \end{array}$	0.89 ± 0.00
U-Net _{DOW-6}	$\begin{array}{c} 0.824 \\ \pm \ 0.006 \end{array}$	0.887 ± 0.002	$\begin{array}{c} 0.424 \\ \pm \ 0.055 \end{array}$	$\begin{array}{c} 0.160 \\ \pm \ 0.026 \end{array}$	$\begin{array}{c} 0.591 \\ \pm \ 0.057 \end{array}$	$\begin{array}{c} 0.577 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} 0.888 \\ \pm \ 0.00 \end{array}$
YOLOv8 _{Multi}	$\begin{array}{c} 0.750 \\ \pm \ 0.030 \end{array}$	$\begin{array}{c} 0.824 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.405 \\ \pm \ 0.021 \end{array}$	$\begin{array}{c} 0.223 \\ \pm \ 0.059 \end{array}$	$\begin{array}{c} 0.549 \\ \pm \ 0.026 \end{array}$	$\begin{array}{c} 0.550 \\ \pm \ 0.017 \end{array}$	0.824 ± 0.02
DINOv2 _{Multi}	$\begin{array}{c} 0.824 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.866 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.491 \\ \pm \ 0.022 \end{array}$	$\begin{array}{c} 0.675 \\ \pm \ 0.044 \end{array}$	$\begin{array}{c} 0.614 \\ \pm \ 0.015 \end{array}$	$\begin{array}{c} 0.694 \\ \pm \ 0.013 \end{array}$	0.879 ± 0.00
DINOv2 _{DOW-Multi}	$\begin{array}{c} \textbf{0.840} \\ \pm \ \textbf{0.003} \end{array}$	$\begin{array}{c} 0.874 \\ \pm \ 0.002 \end{array}$	0.545 ± 0.015	0.767 ± 0.021	$\begin{array}{c} 0.634 \\ \pm \ 0.014 \end{array}$	$\begin{array}{c}\textbf{0.732}\\ \pm \ 0.008\end{array}$	0.884 ± 0.00
U-Net _{Multi}	$\begin{array}{c} 0.819 \\ \pm \ 0.012 \end{array}$	$\begin{array}{c} 0.880 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.514 \\ \pm \ 0.025 \end{array}$	$\begin{array}{c} 0.306 \\ \pm \ 0.091 \end{array}$	$\begin{array}{c} \textbf{0.650} \\ \pm \ 0.029 \end{array}$	$\begin{array}{c} 0.634 \\ \pm \ 0.024 \end{array}$	$\begin{array}{c} 0.879 \\ \pm \ 0.01 \end{array}$
U-Net _{DOW-Multi}	$\begin{array}{c} 0.830 \\ \pm \ 0.022 \end{array}$	$\begin{array}{c} 0.884 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.502 \\ \pm \ 0.039 \end{array}$	$\begin{array}{c} 0.533 \\ \pm \ 0.194 \end{array}$	$\begin{array}{c} 0.623 \\ \pm \ 0.017 \end{array}$	$\begin{array}{c} 0.674 \\ \pm \ 0.041 \end{array}$	$\begin{array}{c} 0.893 \\ \pm 0.00 \end{array}$

Table B.4: Object-level accuracy on $\mathcal{D}_{\text{test}}$. We report the AP for each roof type, and mAP_{50} and $mAP_{50.95}$ are macro averages over the roof types. The rightmost three columns give the results when we discard the roof type information and just consider building detection. The TP_s columns count true positives, where TP_s are the number of objects that overlap with ground truth objects with a minimum IoU of 0.5. The total number of ground truth objects in the $\mathcal{D}_{\text{test}}$ is 2527.

		A	P ₅₀ of each o	class		aver	rage over cla	sses	igno	ring roof t	уре
Model Name	Metal Sheet	Thatch	Asbestos	Concrete	No Roof	mAP ₅₀	mAP ₅₀₋₉₅	TPs	AP ₅₀	AP ₅₀₋₉₅	TPs
YOLOv8	$\begin{array}{c} 0.841 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.945 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.505 \\ \pm \ 0.032 \end{array}$	$\begin{array}{c} 0.542 \\ \pm \ 0.055 \end{array}$	$\begin{array}{c} 0.661 \\ \pm \ 0.026 \end{array}$	$\begin{array}{c} \textbf{0.698} \\ \pm \ \textbf{0.018} \end{array}$	$\begin{array}{c} 0.548 \\ \pm \ 0.010 \end{array}$	2262.2 ± 7.386	0.941 ± 0.003	$\begin{array}{c} 0.798 \\ \pm \ 0.002 \end{array}$	2405.0 ± 5.514
DINOv2	$\begin{array}{c} 0.799 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.892 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.455 \\ \pm \ 0.023 \end{array}$	$\begin{array}{c} 0.565 \\ \pm \ 0.034 \end{array}$	$\begin{array}{c} 0.657 \\ \pm \ 0.016 \end{array}$	$\begin{array}{c} 0.673 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} 0.526 \\ \pm \ 0.005 \end{array}$	$\underset{\pm \ 6.841}{2131.0}$	$\begin{array}{c} 0.881 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.732 \\ \pm \ 0.004 \end{array}$	$\underset{\pm 8.089}{2257.4}$
DINOv2 _{DOW}	$\begin{array}{c} 0.849 \\ \pm 0.005 \end{array}$	0.944 ± 0.004	0.478 ± 0.005	0.601 ± 0.036	0.717 ± 0.009	0.718 ± 0.009	$\begin{array}{c} \textbf{0.568} \\ \pm \ \textbf{0.004} \end{array}$	2236.2 ± 5.636	0.931 ± 0.004	$\begin{array}{c} 0.780 \\ \pm 0.002 \end{array}$	2375.8 ± 5.115
U-Net	0.826 ± 0.005	0.924 ± 0.006	$0.499 \\ \pm 0.016$	0.511 ± 0.042	0.679 ± 0.015	0.688 ± 0.014	0.578 ± 0.014	2191.2 ± 11.25	0.910 ± 0.005	0.797 ± 0.003	2323.0 ± 6.033
U-Net _{DOW}	0.855 + 0.005	0.946 + 0.005	0.545 + 0.019	$0.596 \\ \pm 0.049$	0.707 + 0.012	0.730 + 0.011	0.614 + 0.007	2249.4 + 4.128	0.935 + 0.001	0.819 + 0.003	2383.6 + 5.314
U-Net _{DOW-6}	$\begin{array}{c} 0.851 \\ \pm 0.006 \end{array}$	$\begin{array}{c} 0.003\\ 0.943\\ \pm \ 0.004\end{array}$	0.551 ± 0.011	0.587 ± 0.049	$\begin{array}{c} 0.687 \\ \pm \ 0.022 \end{array}$	$\begin{array}{c} 0.724 \\ \pm 0.007 \end{array}$	0.606 ± 0.005	2243.2 ± 3.487	$0.929 \\ \pm 0.004$	$\begin{array}{c} 0.818 \\ \pm 0.002 \end{array}$	2374.4 ± 5.783
YOLOv8 _{Multi}	0.849	0.923	0.467	0.070	0.676	0.597	0.481	2195.6	0.910	0.751	2328.2
DINOv2 _{Multi}	${ \pm 0.008 \atop \pm 0.004 }$	0.927 ± 0.007	0.445 ± 0.043	0.525 ± 0.076	± 0.020 0.637 ± 0.028	0.680 ± 0.019	0.509 ± 0.008	± 13.383 2230.8 ± 5.344		$\begin{array}{c} \pm 0.007 \\ 0.732 \\ \pm 0.002 \end{array}$	$ \pm 9.988 2305.6 \pm 5.678 $
DINOv2 _{DOW-Multi}	$\begin{array}{c} \textbf{0.885} \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.942 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.542 \\ \pm \ 0.026 \end{array}$	$\begin{array}{c} 0.529 \\ \pm \ 0.122 \end{array}$	$\begin{array}{c} 0.605 \\ \pm \ 0.019 \end{array}$	$\begin{array}{c} 0.701 \\ \pm \ 0.027 \end{array}$	$\begin{array}{c} 0.557 \\ \pm \ 0.027 \end{array}$	$\begin{array}{c} 2270.2 \\ \pm 8.376 \end{array}$	$\begin{array}{c} 0.918 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.769 \\ \pm \ 0.004 \end{array}$	$\underset{\pm 7.467}{2360.8}$
U-Net _{Multi}	$\begin{array}{c} 0.883 \\ \pm 0.009 \end{array}$	$\begin{array}{c} 0.940 \\ \pm 0.010 \end{array}$	0.531 ± 0.022	0.498 ± 0.051	0.728 ± 0.040	0.716 ± 0.018	0.603 ± 0.011	2262.4 ± 6.119	0.924 ± 0.004	$\begin{array}{c} 0.797 \\ \pm 0.007 \end{array}$	2358.4 ± 9.091
U-Net _{DOW-Multi}	$\begin{array}{c} 0.894 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.943 \\ \pm 0.000 \end{array}$	$\begin{array}{c} 0.516 \\ \pm \ 0.038 \end{array}$	$\begin{array}{c} 0.517 \\ \pm \ 0.038 \end{array}$	$\begin{array}{c} 0.678 \\ \pm \ 0.018 \end{array}$	$\begin{array}{c} 0.710 \\ \pm 0.006 \end{array}$	$\begin{array}{c} 0.606 \\ \pm 0.007 \end{array}$	2275.2 ± 5.879	$\begin{array}{c} 0.933 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.813 \\ \pm 0.003 \end{array}$	2382.0 ± 9.055

Table B.5: Pixel-level accuracies on \mathcal{D}_{ext} . IoU refers to the IoU computed on the binary outputs, where the predictions of multi-class models were binarized. mIoU⁵ refers to the macro average of the IoUs for the individual classes. The subscript *Multi* indicates the end-to-end setting.

	IoU-Sc	core of eac	ch class		
Model Name	Metal Sheet	Thatch	No Roof	IoU (Mean)	IoU (Binary)
YOLOv8	$\begin{array}{c} 0.888 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.879 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.516 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} 0.761 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.896 \\ \pm \ 0.002 \end{array}$
DINOv2	$\begin{array}{c} 0.844 \\ \pm \ 0.012 \end{array}$	$\begin{array}{c} 0.854 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.456 \\ \pm \ 0.037 \end{array}$	$\begin{array}{c} 0.718 \\ \pm \ 0.015 \end{array}$	$\begin{array}{c} 0.904 \\ \pm \ 0.001 \end{array}$
DINOv2 _{DOW}	$\begin{array}{c} 0.886 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.875 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.541 \\ \pm \ 0.013 \end{array}$	$\begin{array}{c} 0.767 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.904 \\ \pm \ 0.001 \end{array}$
U-Net	$\begin{array}{c} 0.896 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.883 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.463 \\ \pm \ 0.017 \end{array}$	$\begin{array}{c} 0.748 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} 0.909 \\ \pm \ 0.001 \end{array}$
U-Net _{DOW}	$\begin{array}{c} 0.905 \\ \pm \ 0.002 \end{array}$	0.895 ± 0.003	$\begin{array}{c} 0.493 \\ \pm \ 0.018 \end{array}$	$\begin{array}{c} 0.764 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} \textbf{0.911} \\ \pm \ 0.002 \end{array}$
U-Net _{DOW-6}	$\begin{array}{c} 0.900 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.889 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.452 \\ \pm \ 0.031 \end{array}$	$\begin{array}{c} 0.747 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} 0.902 \\ \pm \ 0.003 \end{array}$
YOLOv8 _{Multi}	$\begin{array}{c} 0.890 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.860 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.606 \\ \pm \ 0.019 \end{array}$	$\begin{array}{c} 0.785 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.885 \\ \pm \ 0.002 \end{array}$
DINOv2 _{Multi}	$\begin{array}{c} 0.907 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.874 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} \textbf{0.676} \\ \pm \ 0.022 \end{array}$	$\begin{array}{c} 0.819 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.899 \\ \pm \ 0.001 \end{array}$
DINOv2 _{DOW-Multi}	$\begin{array}{c} 0.912 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.881 \\ \pm \ 0.003 \end{array}$	0.676 ± 0.020	$\begin{array}{c} \textbf{0.823} \\ \pm \ \textbf{0.007} \end{array}$	$\begin{array}{c} 0.901 \\ \pm \ 0.001 \end{array}$
U-Net _{Multi}	$\begin{array}{c} \textbf{0.913} \\ \pm \ \textbf{0.005} \end{array}$	$\begin{array}{c} 0.884 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.617 \\ \pm \ 0.061 \end{array}$	$\begin{array}{c} 0.805 \\ \pm \ 0.020 \end{array}$	$\begin{array}{c} 0.903 \\ \pm \ 0.002 \end{array}$
U-Net _{DOW-Multi}	$\begin{array}{c} \textbf{0.913} \\ \pm \ \textbf{0.005} \end{array}$	$\begin{array}{c} 0.884 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.617 \\ \pm \ 0.061 \end{array}$	$\begin{array}{c} 0.805 \\ \pm \ 0.020 \end{array}$	$\begin{array}{c} 0.903 \\ \pm \ 0.002 \end{array}$

010	
919	Table B.6: Object-level accuracies on \mathcal{D}_{ext} . We report the AP for each roof type, and mAP_{50} and
920	
520	mAP_{50-95} are macro averages over the classes. The rightmost three columns give the results when we
921	discard the roof type information and just consider building detection. TPs are the number of objects
922	that overlap with ground truth objects with a minimum IoU of 0.5. The total number of ground truth
923	objects in the \mathcal{D}_{ext} is 1541.

	AP	50 of each c	lass	Obj	ects with Cla	isses	Only I	Building O	bjects
Model Name	Metal Sheet	Thatch	No Roof	mAP ₅₀	mAP ₅₀₋₉₅	TPs	AP ₅₀	AP ₅₀₋₉₅	TPs
YOLOv8	$\begin{array}{c} 0.928 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} \textbf{0.947} \\ \pm \ \textbf{0.000} \end{array}$	$\begin{array}{c} 0.661 \\ \pm \ 0.023 \end{array}$	$\begin{array}{c} 0.846 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.428 \\ \pm \ 0.002 \end{array}$	1447.2 ± 4.534	0.963 ± 0.005	$\begin{array}{c} 0.838 \\ \pm \ 0.002 \end{array}$	$\begin{array}{r}1493.8\\\pm3.826\end{array}$
DINOv2	$\begin{array}{c} 0.896 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.883 \\ \pm \ 0.012 \end{array}$	$\begin{array}{c} 0.634 \\ \pm \ 0.019 \end{array}$	$\begin{array}{c} 0.804 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.390 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c}1382.8\\\pm 8.818\end{array}$	$\begin{array}{c} 0.922 \\ \pm \hspace{0.1cm} 0.005 \end{array}$	$\begin{array}{c} 0.785 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c}1434.2\\\pm \ 6.524\end{array}$
DINOv2 _{DOW}	$\begin{array}{c} 0.932 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} 0.944 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.709 \\ \pm \ 0.016 \end{array}$	$\begin{array}{c} 0.861 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} 0.426 \\ \pm \ 0.004 \end{array}$	$\begin{array}{r}1444.8\\\pm 5.455\end{array}$	$\begin{array}{c} 0.961 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.830 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} \textbf{1494.4} \\ \pm \text{ 4.363} \end{array}$
U-Net	$\begin{array}{c} 0.915 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.921 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.520 \\ \pm \ 0.027 \end{array}$	$\begin{array}{c} 0.590 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.407 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c}1399.4\\\pm 4.758\end{array}$	$\begin{array}{c} 0.929 \\ \pm \ 0.000 \end{array}$	$\begin{array}{c} 0.836 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c}1438.4\\\pm \textbf{4.499}\end{array}$
U-Net _{DOW}	$\begin{array}{c} 0.932 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.946 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.559 \\ \pm \ 0.027 \end{array}$	$\begin{array}{c} 0.812 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.528 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c}1429.0\\\pm \ 6.229\end{array}$	$\begin{array}{c} 0.947 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} \textbf{0.858} \\ \pm \ 0.004 \end{array}$	$\begin{array}{c}1468.6\\\pm \ 6.499\end{array}$
U-Net _{DOW-6}	$\begin{array}{c} 0.935 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.940 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.509 \\ \pm \ 0.022 \end{array}$	$\begin{array}{c} 0.795 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.518 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c}1421.0\\\pm3.688\end{array}$	$\begin{array}{c} 0.939 \\ \pm \ 0.000 \end{array}$	$\begin{array}{c} 0.851 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c}1458.8\\\pm 4.118\end{array}$
YOLOv8 _{Multi}	$\begin{array}{c} 0.949 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.934 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} 0.664 \\ \pm \ 0.044 \end{array}$	$\begin{array}{c} 0.849 \\ \pm \ 0.015 \end{array}$	$\begin{array}{c} 0.423 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c}1446.0\\\pm 4.899\end{array}$	$\begin{array}{c} 0.948 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.808 \\ \pm \ 0.005 \end{array}$	1477.2 ± 3.655
DINOv2 _{Multi}	$\begin{array}{c} 0.951 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.929 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} \textbf{0.732} \\ \pm \text{ 0.043} \end{array}$	$\begin{array}{c} \textbf{0.871} \\ \pm \ \textbf{0.015} \end{array}$	$\begin{array}{c} 0.424 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} \textbf{1453.2} \\ \pm 5.154 \end{array}$	$\begin{array}{c} 0.940 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.796 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c}1467.8\\\pm5.154\end{array}$
DINOv2 _{DOW-Multi}	$\begin{array}{c} 0.955 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.942 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.633 \\ \pm \ 0.028 \end{array}$	$\begin{array}{c} 0.843 \\ \pm \ 0.012 \end{array}$	$\begin{array}{c} 0.503 \\ \pm \ 0.038 \end{array}$	$\begin{array}{r}1462.2\\\pm 5.192\end{array}$	$\begin{array}{c} 0.944 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.810 \\ \pm \ 0.002 \end{array}$	$\begin{array}{r}1481.4\\\pm \textbf{4.499}\end{array}$
U-Net _{Multi}	$\begin{array}{c} \textbf{0.956} \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.926 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.651 \\ \pm \ 0.107 \end{array}$	$\begin{array}{c} 0.844 \\ \pm \ 0.039 \end{array}$	$\begin{array}{c} \textbf{0.548} \\ \pm \ 0.017 \end{array}$	$\begin{array}{c}1439.0\\\pm 16.358\end{array}$	$\begin{array}{c} 0.943 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.838 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c}1463.6\\\pm 13.063\end{array}$
U-Net _{DOW-Multi}	$\begin{array}{c} 0.956 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.926 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.651 \\ \pm \ 0.107 \end{array}$	$\begin{array}{c} 0.844 \\ \pm \ 0.039 \end{array}$	$\begin{array}{c} 0.438 \\ \pm \ 0.017 \end{array}$	$\begin{array}{c}1439.0\\\pm 16.358\end{array}$	$\begin{array}{c} 0.943 \\ \pm \hspace{0.1cm} 0.010 \end{array}$	$\begin{array}{c} 0.838 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c}1463.6\\\pm 13.063\end{array}$

⁹⁴⁷ 948

950

951

952

953

954

955

Table B.7: Multi-class prediction using DOW-Multi models. There are several ways to derive multiclass predictions in the DOW models. The approach in the main part of this study derives class labels by taking the majority vote of all inner segmentations combined with the border pixels from the outer mask that do not overlap with the inner one. Alternatively, one could solely consider only the inner mask or outer mask, indicated by the suffixes *-inner* and *-outer*, respectively, in the subscripts of the model suffixes. This table adds the results for these alternative methods, reporting again averages over five trials \pm standard deviations. The upper five models were added to compare DINOv2 models and the lower half of the models were added using U-Net-based methods.

				\mathcal{D}	test			$\mathcal{D}_{\mathrm{ext}}$			
			pixel leve	1		object leve	el	pixe	l level	object	t level
N	Iodel Name	IoU	mIoU ³	mIoU ⁵	AP ₅₀	mAP ₅₀ ³	mAP ₅₀	IoU	mIoU ³	AP ₅₀	mAP ₅₀
D	DINOv2 _{Multi}	0.879	0.768	0.694	0.894	0.810	0.680	0.899	0.819	0.940	0.871
D	DINOv2 _{DOW-Multi}	± 0.001 0.884 ± 0.002	± 0.005 0.783 ± 0.005	± 0.013 0.732 ± 0.009	± 0.004 0.918 ± 0.002	± 0.008 0.810 ± 0.009	± 0.019 0.701 ± 0.027	± 0.001 0.901 ± 0.001	± 0.006 0.822 ± 0.007	± 0.003 0.944 ± 0.005	± 0.015 0.843 ± 0.012
D	DINOv2 _{DOW-Multi-outer}	$\begin{array}{c} 0.884 \\ \pm 0.002 \end{array}$	$\begin{array}{c} 0.783 \\ \pm 0.005 \end{array}$	$\begin{array}{c} 0.732 \\ \pm 0.008 \end{array}$	$\begin{array}{c} 0.918 \\ \pm 0.002 \end{array}$	$\begin{array}{c} 0.810 \\ \pm 0.009 \end{array}$	0.701 ± 0.027	$\begin{array}{c} 0.901 \\ \pm 0.001 \end{array}$	$\begin{array}{c} 0.823 \\ \pm 0.007 \end{array}$	$\begin{array}{c} 0.944 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} -0.843 \\ \pm 0.012 \end{array}$
D	DINOv2 _{DOW-Multi-inner}	$\begin{array}{c} 0.884 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.782 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.730 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} 0.918 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.810 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} 0.699 \\ \pm \ 0.027 \end{array}$	$\begin{array}{c} 0.901 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.821 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} 0.944 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.842 \\ \pm \ 0.012 \end{array}$
U	J-Net _{Multi}	$\begin{array}{c} 0.879 \\ \pm \ 0.012 \end{array}$	$\begin{array}{c} 0.783 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.634 \\ \pm \ 0.024 \end{array}$	$\begin{array}{c} 0.924 \\ \pm \ 0.004 \end{array}$	0.850 ± 0.011	$\begin{array}{c} 0.716 \\ \pm \ 0.018 \end{array}$	$\begin{array}{c} 0.903 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.805 \\ \pm \ 0.020 \end{array}$	$\begin{array}{c} 0.943 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.844 \\ \pm \ 0.039 \end{array}$
U	J-Net _{DOW-Multi}	$\begin{array}{c} 0.893 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.779 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} 0.674 \\ \pm \ 0.041 \end{array}$	$\begin{array}{c} 0.933 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.838 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.709 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.906 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.798 \\ \pm \ 0.012 \end{array}$	$\begin{array}{c} 0.944 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.830 \\ \pm \ 0.017 \end{array}$
U	J-Net _{DOW-Multi-outer}	0.893 ± 0.002	0.779	0.674 ± 0.041	$0.933 \\ \pm 0.003$	0.838 ± 0.005	0.710 ± 0.006	0.906 ± 0.002	0.798 ± 0.012	0.944	0.830 ± 0.017
U	J-Net _{DOW-Multi-inner}	$\begin{array}{c} 0.893 \\ \pm 0.002 \end{array}$	$\begin{array}{c} 0.011\\ 0.779\\ \pm \ 0.011\end{array}$	0.673 ± 0.040	$\begin{array}{c} 0.933 \\ \pm 0.003 \end{array}$	$\begin{array}{c} 0.838 \\ \pm 0.006 \end{array}$	$\begin{array}{c} \pm 0.000\\ 0.709\\ \pm 0.006\end{array}$	$\begin{array}{c} 0.906 \\ \pm 0.002 \end{array}$	0.797 ± 0.012	$\begin{array}{c} \pm 0.005\\ 0.944\\ \pm 0.005\end{array}$	$\begin{array}{c} \pm 0.017\\ 0.830\\ \pm 0.017\end{array}$

972	٥	7	0
070	J	1	-
	~	_	~

Table B.8: We ran the experiments with an 80% stratified subset of the data. The results changed only very slightly, indicating that our training dataset is representative and large enough for the defined

			\mathcal{D}	test			\mathcal{D}_{ext}			
	pixel level				object level			l level	object level	
Model Name	IoU	mIoU ³	mIoU ⁵	AP ₅₀	mAP_{50}^3	mAP ₅₀ ⁵	IoU	mIoU ³	AP ₅₀	mAP_{50}^3
YOLOv8	$\begin{array}{c} 0.855 \\ \pm \ 0.020 \end{array}$	$\begin{array}{c} 0.715 \\ \pm \ 0.018 \end{array}$	$\begin{array}{c} 0.560 \\ \pm \ 0.014 \end{array}$	0.937 ± 0.004	$\begin{array}{c} 0.823 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.679 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} 0.894 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.773 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.960 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.854 \\ \pm \ 0.008 \end{array}$
U-Net	$\begin{array}{c} 0.893 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} 0.758 \\ \pm \ 0.021 \end{array}$	$\begin{array}{c} 0.571 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.914 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.824 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.683 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} \textbf{0.912} \\ \pm \ \textbf{0.002} \end{array}$	$\begin{array}{c} 0.738 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} 0.941 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.818 \\ \pm \ 0.017 \end{array}$
U-Net _{DOW}	$\begin{array}{c} \textbf{0.896} \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} \textbf{0.792} \\ \pm \ \textbf{0.006} \end{array}$	$\begin{array}{c} 0.582 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.933 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} \textbf{0.846} \\ \pm \ \textbf{0.005} \end{array}$	$\begin{array}{c} 0.695 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} 0.911 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.743 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} 0.953 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.820 \\ \pm \ 0.014 \end{array}$
DINOv2	$\begin{array}{c} 0.880 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.721 \\ \pm \ 0.027 \end{array}$	$\begin{array}{c} 0.535 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.883 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.793 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.657 \\ \pm \ 0.014 \end{array}$	$\begin{array}{c} 0.904 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.729 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.928 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.825 \\ \pm \ 0.006 \end{array}$
DINOv2 _{DOW}	$\begin{array}{c} 0.879 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.758 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.554 \\ \pm \ 0.006 \end{array}$	$\begin{array}{c} 0.931 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.841 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.697 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} 0.903 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.755 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} \textbf{0.961} \\ \pm \ \textbf{0.004} \end{array}$	$\begin{array}{c} 0.863 \\ \pm \ 0.006 \end{array}$
YOLOv8 _{Multi}	$\begin{array}{c} 0.806 \\ \pm \ 0.028 \end{array}$	$\begin{array}{c} 0.697 \\ \pm \ 0.018 \end{array}$	$\begin{array}{c} 0.536 \\ \pm \ 0.029 \end{array}$	$\begin{array}{c} 0.901 \\ \pm \ 0.002 \end{array}$	$\begin{array}{c} 0.798 \\ \pm \ 0.010 \end{array}$	$\begin{array}{c} 0.583 \\ \pm \ 0.009 \end{array}$	$\begin{array}{c} 0.881 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.777 \\ \pm \ 0.007 \end{array}$	$\begin{array}{c} 0.942 \\ \pm \ 0.003 \end{array}$	$\begin{array}{c} 0.832 \\ \pm \ 0.004 \end{array}$
U-Net _{Multi}	$\begin{array}{c} 0.886 \\ \pm \ 0.005 \end{array}$	$\begin{array}{c} 0.786 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.666 \\ \pm \ 0.023 \end{array}$	$\begin{array}{c} 0.924 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.839 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} \textbf{0.712} \\ \pm \ \textbf{0.023} \end{array}$	$\begin{array}{c} 0.906 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.799 \\ \pm \ 0.011 \end{array}$	$\begin{array}{c} 0.948 \\ \pm \ 0.000 \end{array}$	$\begin{array}{c} 0.839 \\ \pm \ 0.012 \end{array}$
U-Net _{DOW-Multi}	$\begin{array}{c} 0.887 \\ \pm 0.003 \end{array}$	0.782 ± 0.009	$\begin{array}{c} 0.670 \\ \pm 0.044 \end{array}$	$0.929 \\ \pm 0.005$	$\begin{array}{c} 0.837 \\ \pm \ 0.011 \end{array}$	0.709 ± 0.013	$0.905 \\ \pm 0.003$	0.812 ± 0.012	$\begin{array}{c} 0.951 \\ \pm \ 0.004 \end{array}$	0.855 ± 0.020
DINOv2 _{Multi}	$\begin{array}{c} 0.870 \\ \pm \ 0.004 \end{array}$	0.771 ± 0.004	$0.675 \\ \pm 0.004$	$\begin{array}{c} 0.901 \\ \pm \ 0.005 \end{array}$	0.824 ± 0.013	0.690 ± 0.023	$\begin{array}{c} 0.895 \\ \pm \ 0.003 \end{array}$	$0.809 \\ \pm 0.013$	$0.945 \\ \pm 0.007$	0.870 ± 0.022
DINOv2 _{DOW-Multi}	$\begin{array}{c} 0.880 \\ \pm \ 0.001 \end{array}$	$\begin{array}{c} 0.785 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} \textbf{0.738} \\ \pm \ \textbf{0.008} \end{array}$	$\begin{array}{c} 0.911 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.823 \\ \pm \ 0.008 \end{array}$	$\begin{array}{c} 0.703 \\ \pm \ 0.036 \end{array}$	$\begin{array}{c} 0.899 \\ \pm \ 0.002 \end{array}$	0.819 ± 0.003	$\begin{array}{c} 0.955 \\ \pm \ 0.004 \end{array}$	$\begin{array}{c} 0.868 \\ \pm \ 0.009 \end{array}$

Table B.9: Cross-validation accuracies on combined training and validation data of k-nearest neighbour classification (kNN) and logistic regression applied to the DINOv2 features. The baseline is the architecture depicted in Figure 4, *w/o mask* refers to omitting the masking and averaging the DINOv2 features across the whole input patch, and *w/o upsampling* did not upsample the DINOv2 features but downsampled the building mask instead.

			F1	-Score		
	base	eline	w/o	mask	w/o upsampling	
Classifier	Mean	Std	Mean	Std	Mean	Std
Logistic Regression	0.770	0.063	0.573	0.077	0.768	0.067
kNN	0.734	0.045	0.389	0.029	0.733	0.051

C DATASHEET

C.1 MOTIVATION

C.1.1 FOR WHAT PURPOSE WAS THE DATASET CREATED? WAS THERE A SPECIFIC TASK IN MIND? WAS THERE A SPECIFIC GAP THAT NEEDED TO BE FILLED?

The dataset was created to support research on multi-task computer vision problems and to support mosquito-borne disease risk assessment in African cities. The list of tasks include classification, semantic segmentation, and instance segmentation of roofs and their material. While these tasks are closely related, each serves a different purpose and accurate segmentation of objects need not imply accurate object separation, and vice versa. The dataset is ideal for bench-marking methods for the above mentioned tasks.

1026 1027 C.1.2 WHO CREATED THE DATASET (E.G., WHICH TEAM, RESEARCH GROUP) AND ON BEHALF OF WHICH ENTITY (E.G., COMPANY, INSTITUTION, ORGANIZATION)?

The dataset was created by authors. The drone imagery and building footprints were captured by an NGO. The imagery and the building footprints were fused, re-registered, cleaned, verified, and split into given datasets by the authors.

1032
1033C.1.3Who funded the creation of the dataset? If there is an associated grant,
please provide the name of the grant or and the grant name and number

¹⁰³⁵ We will provide the project and grant details later.

1037 C.1.4 ANY OTHER COMMENT?

1038 1039 None.

1036

1046

1052 1053 1054

1055

1056

1057

1058

1061 1062

1063 1064

1069

1040 1041 С.2 Сомрозітіон

1042
1043
1044
1044
1044
1045
C.2.1 WHAT DO THE INSTANCES THAT COMPRISE THE DATASET REPRESENT (E.G., DOCUMENTS, PHOTOS, PEOPLE, COUNTRIES)? ARE THERE MULTIPLE TYPES OF INSTANCES (E.G., MOVIES, USERS, AND RATINGS; PEOPLE AND INTERACTIONS BETWEEN THEM; NODES AND EDGES)? PLEASE PROVIDE A DESCRIPTION

The dataset comprises of very-high resolution orthophotos captured through a drone and expert drawn polygons for all buildings with annotation of their roof material. The dataset covers three informal settlements in Nacala. Five classes of roof material are identified: metal sheet, thatch, asbestos, concrete, and no-roof. An example of a portion of the orthophoto and roof labels is shown in Fig. C.8.



Figure C.8: Drone Imagery with RGB (Red, Green, and Blue channels) and annotations

C.2.2 HOW MANY INSTANCES ARE THERE IN TOTAL (OF EACH TYPE, IF APPROPRIATE)?

The total number of building polygons in the data is 17954. The distribution of roof material classes
is imbalanced. The number of buildings belonging to metal sheet, thatch, asbestos, concrete, and
no-roof classes are 9776, 6428, 566, 174, and 1010, respectively.

C.2.3 DOES THE DATASET CONTAIN ALL POSSIBLE INSTANCES OR IS IT A SAMPLE (NOT NECESSARILY RANDOM) OF INSTANCES FROM A LARGER SET? IF THE DATASET IS A SAMPLE, THEN WHAT IS THE LARGER SET? IS THE SAMPLE REPRESENTATIVE OF THE LARGER SET (E.G., GEOGRAPHIC COVERAGE)? IF SO, PLEASE DESCRIBE HOW THIS REPRESENTATIVENESS WAS VALIDATED/VERIFIED. IF IT IS NOT REPRESENTATIVE OF THE LARGER SET, PLEASE DESCRIBE WHY NOT (E.G., TO COVER A MORE DIVERSE RANGE OF INSTANCES, BECAUSE INSTANCES WERE WITHHELD OR UNAVAILABLE)

1077 The dataset contains all available instances. All informal settlements in Nacala that have drone
1078 orthophotos available are prepared as a dataset. Furthermore, all buildings visible in the orthophotos
1079 are included, and the five identified building classes cover all possible roof materials in the area, and
the most predominant roof materials present in the wider Nacala region.

1080 WHAT DATA DOES EACH INSTANCE CONSIST OF? "RAW" DATA (E.G., UNPROCESSED C.2.4 1081 TEXT OR IMAGES)OR FEATURES? IN EITHER CASE, PLEASE PROVIDE A DESCRIPTION 1082 The data consists of aerial images and corresponding labels. Labels are building footprints with the 1084 attribute of roof class. The raw images are GeoTiff images tagged with a spatial reference system. The raw labels are GeoJSON files with the same spatial reference system as images. 1086 1087 1088 C.2.5 IS THERE A LABEL OR TARGET ASSOCIATED WITH EACH INSTANCE? IF SO, PLEASE 1089 PROVIDE A DESCRIPTION. 1090 1091 The labels on the image are polygons describing the geometry of the building footprints and their 1092 associated roof material classes, as described above. In the raw data, the material class is saved 1093 under the attribute name of **mater id** in GeoJSON files. The values of metal sheet, thatch, asbestos, concrete, and no-roof in the attribute are 1, 2, 3, 4, and 5, respectively. The same values are assigned 1094 to the patch labels. 1095 1096 C.2.6 IS ANY INFORMATION MISSING FROM INDIVIDUAL INSTANCES? IF SO, PLEASE PROVIDE A DESCRIPTION, EXPLAINING WHY THIS INFORMATION IS MISSING (E.G., 1099 BECAUSE IT WAS UNAVAILABLE). THIS DOES NOT INCLUDE INTENTIONALLY REMOVED 1100 INFORMATION BUT MIGHT INCLUDE, E.G., REDACTED TEXT. 1101 1102 Everything is included. No data is missing. 1103 1104 1105 ARE RELATIONSHIPS BETWEEN INDIVIDUAL INSTANCES MADE EXPLICIT (E.G., USERS' C.2.7 1106 MOVIE RATINGS, SOCIAL NETWORK LINKS)? IF SO, PLEASE DESCRIBE HOW THESE 1107 RELATIONSHIPS ARE MADE EXPLICIT. 1108 1109 No, the geometry and material attribute of each building footprint is independently recorded. 1110 1111 1112 C.2.8 ARE THERE RECOMMENDED DATA SPLITS (E.G., TRAINING, 1113 DEVELOPMENT/VALIDATION, TESTING)? IF SO, PLEASE PROVIDE A DESCRIPTION OF 1114 THESE SPLITS, EXPLAINING THE RATIONALE BEHIND THEM. 1115 1116 The roof material classes are not balanced and are not geographically distributed uniformly. The 1117 data is split into training, validation, and test sets using stratified random sampling to account for 1118 the class imbalance. We created a square grid of 225 meters and counted the roof types in these 1119 cells. Then we partitioned the cells into three sets based on the class counts to achieve a similar class 1120 distribution in each dataset, where we prioritized the distribution of minority classes (i.e., concrete 1121 and asbestos). We defined that a building only belongs to a specific grid cell if its centroid falls into 1122 the cell. These grid cells separate the images and labels into training, validation and test sets. See 1123 Fig. C.9 as an example. Initially, only two informal settlements were labelled and therefore, only 1124 these two settlements are divided into the 3 sets. The third informal settlement was labelled later and treated as a second test set. 1125 1126 1127 ARE THERE ANY ERRORS, SOURCES OF NOISE, OR REDUNDANCIES IN THE DATASET? IF C.2.9 1128 SO, PLEASE PROVIDE A DESCRIPTION. 1129 1130 The images exhibit a high level of details and the building footprint geometry and material attributes 1131 are meticulous noted by experts. The dataset is free from errors, noise, or redundancies to the greatest 1132 extent possible but we acknowledge that even with expert craftsmanship, there is always a chance of 1133

human error.



1188	C.2.14	Is it possible to identify individuals (i.e., one or more natural persons),
1189		EITHER DIRECTLY OR INDIRECTLY (I.E., IN COMBINATION WITH OTHER DATA) FROM
1190		THE DATASET? IF SO, PLEASE DESCRIBE HOW.
1191		
1192	It is not	possible to identify individuals in the drone imagery. In any publicly available and geo-coded
1193	1mage, 1	it is possible to identify individual houses and reverse geocode into a human-readable address.
1194		
1195	C.2.15	DOES THE DATASET CONTAIN DATA THAT MIGHT BE CONSIDERED SENSITIVE IN ANY
1196		WAY (E.G., DATA THAT REVEALS RACE OR ETHNIC ORIGINS, SEXUAL ORIENTATIONS,
1197		RELIGIOUS BELIEFS, POLITICAL OPINIONS OR UNION MEMBERSHIPS, OR LOCATIONS;
1198		FINANCIAL OR HEALTH DATA; BIOMETRIC OR GENETIC DATA; FORMS OF
1199		GOVERNMENT IDENTIFICATION, SUCH AS SOCIAL SECURITY NUMBERS; CRIMINAL
1200		HISTORY)? IF SO, PLEASE PROVIDE A DESCRIPTION
1201		
1202	No.	
1203		
1204	C.2.16	Any other comments?
1205	0.2.110	
1206	None.	
1207		
1208		
1209	C.3 C	Collection Process
1210	~ • •	
1211	C.3.1	HOW WAS THE DATA ASSOCIATED WITH EACH INSTANCE ACQUIRED? WAS THE DATA
1212		DIRECTLY OBSERVABLE (E.G., RAW TEXT, MOVIE RATINGS), REPORTED BY SUBJECTS
1213		(E.G., SURVEY RESPONSES), OR INDIRECTLY INFERRED/DERIVED FROM OTHER DATA
1214		(E.G., PARI-OF-SPEECH TAGS, MODEL-BASED GUESSES FOR AGE OR LANGUAGE): IF
1215		DATA WAS THE DATA VALIDATED/VERIFIED? IF SO PLEASE DESCRIBE HOW
1216		
1217	The dat	ta is observable as images. The ground sampling distance of pixel or spatial resolution of
1218	the ima	gery is ≈ 4.4 cm/pixel. QGIS was used for the visualization of images and re-registration,
1219	cleaning	g, and verification of building footprints and their attributes.
1220		
1221	C 3 2	WHAT MECHANISMS OF PROCEDURES WERE USED TO COLLECT THE DATA (E.C.
1222	C.J.2	W HAT MECHANISMS ON PROCEDURES WERE USED TO COLLECT THE DATA (E.G.,
1223		PROGRAM SOFTWARE API)? HOW WERE THESE MECHANISMS OF PROCEDURES
1224		VALIDATED?
1225		
1226	The dro	one imagery was captured using a DJI Phantom 4 Pro drone and processed using AgiSoft
1227	Metash	ape software. All building footprints are annotated using Open Street Map (JOSM) that use
1228	OpenSt	reetMap in the backend. Before splitting into train, validation, and test sets, the missing labels
1229	and all	geometric and attribute errors were corrected in QGIS software.
1230		
1231		
1232		Table C.10: Drone flight information summary
1233		Number of flighter 4
1234		Drone: DII Phantom / Dro
1235		Camera Brand [.] DII
1236		Camera Model: FC6310
1237		Image Resolution: 4864×3648 (~18MP)
1238		Flight Altitude: 120 m
1239		Flight dates: 27-30, October 2021
1240		Total flight duration: 504 minutes
1241		

IF THE DATASET IS A SAMPLE FROM A LARGER SET, WHAT WAS THE SAMPLING STRATEGY (E.G., DETERMINISTIC, PROBABILISTIC WITH SPECIFIC SAMPLING PROBABILITIES)?	
The data was prepared based on its availability. All data from the project was made available and	
this dataset.	
WHO WAS INVOLVED IN THE DATA COLLECTION PROCESS (E.G., STUDENTS, CROWDWORKERS, CONTRACTORS) AND HOW WERE THEY COMPENSATED (E.G., HOW MUCH WERE CROWDWORKERS PAID)?	
ne imagery was captured by an NGO. Nacala residents and local university students performed d data collection, receiving stipends and data bundles. The polygons and attributes of building nts were corrected by authors.	
Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please describe the timeframe in which the data associated with the instances was created.	
ne imagery was captured between October and December 2021. All labels are manually ed on the imagery beginning January 2022 until May 2024.	
WERE ANY ETHICAL REVIEW PROCESSES CONDUCTED (E.G., BY AN INSTITUTIONAL REVIEW BOARD)? IF SO, PLEASE PROVIDE A DESCRIPTION OF THESE REVIEW PROCESSES, INCLUDING THE OUTCOMES, AS WELL AS A LINK OR OTHER ACCESS POINT TO ANY SUPPORTING DOCUMENTATION.	
No ethical review was conducted. All data collection, including drone flights and on-site mapping, was approved and supported by the Nacala Municipal Council and facilitated on the ground by neighbourhood-level authorities.	
DID YOU COLLECT THE DATA FROM THE INDIVIDUALS IN QUESTION DIRECTLY, OR OBTAIN IT VIA THIRD PARTIES OR OTHER SOURCES (E.G., WEBSITES)?	
a was not collected from individuals.	
Were the individuals in question notified about the data collection? If so, please describe (or show with screenshots or other information) how notice was provided, and provide a link or other access point to, or otherwise reproduce, the exact language of the notification itself.	
DID THE INDIVIDUALS IN QUESTION CONSENT TO THE COLLECTION AND USE OF THEIR DATA? IF SO, PLEASE DESCRIBE (OR SHOW WITH SCREENSHOTS OR OTHER INFORMATION) HOW CONSENT WAS REQUESTED AND PROVIDED, AND PROVIDE A LINK OR OTHER ACCESS POINT TO, OR OTHERWISE REPRODUCE, THE EXACT LANGUAGE TO WHICH THE INDIVIDUALS CONSENTED.	
IF CONSENT WAS OBTAINED, WERE THE CONSENTING INDIVIDUALS PROVIDED WITH A MECHANISM TO REVOKE THEIR CONSENT IN THE FUTURE OR FOR CERTAIN USES? IF SO, PLEASE PROVIDE A DESCRIPTION, AS WELL AS A LINK OR OTHER ACCESS POINT TO THE MECHANISM (IF APPROPRIATE).	

N/A.

1296 C.3.11 HAS AN ANALYSIS OF THE POTENTIAL IMPACT OF THE DATASET AND ITS USE ON 1297 DATA SUBJECTS (E.G., A DATA PROTECTION IMPACT ANALYSIS) BEEN CONDUCTED? IF 1298 SO, PLEASE PROVIDE A DESCRIPTION OF THIS ANALYSIS, INCLUDING THE OUTCOMES, 1299 AS WELL AS A LINK OR OTHER ACCESS POINT TO ANY SUPPORTING DOCUMENTATION. 1300 N/A. 1301 1302 C.3.12 ANY OTHER COMMENTS? 1303 1304 None. 1305 1306 1307 C.4 PREPROCESSING/CLEANING/LABELING 1308 WAS ANY PREPROCESSING/CLEANING/LABELING OF THE DATA DONE (E.G., C.4.1 1309 DISCRETIZATION OR BUCKETING, TOKENIZATION, PART-OF-SPEECH TAGGING, SIFT 1310 FEATURE EXTRACTION, REMOVAL OF INSTANCES, PROCESSING OF MISSING VALUES)? 1311 IF SO, PLEASE PROVIDE A DESCRIPTION. IF NOT, YOU MAY SKIP THE REMAINDER OF 1312 THE QUESTIONS IN THIS SECTION. 1313 1314 Because of the large size of raw aerial imagery, the images of training and validation sets were 1315 cropped to 512×512 pixels of **patches**. The data processing optimized is usefulness in training 1316 deep learning models. The total number of patches in the training and validation sets are 8366 and 1799, respectively, after cropping them without any overlap. The test sets were provided without 1317 cropping into patches, so these images are provided in different sizes. The images in test sets are not 1318 cropped because dividing them into patches may lead to under- or over-estimation of instances and 1319 may influence counting accuracy over large areas. The test-1 set consists of 22 images and the test-2 1320 set consists of a single image. There are 10930, 2956, 2527 and 1541 buildings in train, validation, 1321 test-1 and test-2 sets, respectively. 1322 For the classification of buildings into different roof classes, the DINOv2 features were extracted 1323 for train, validation and test-1 sets and made available with the dataset. These features of the train, 1324 validation, and test-1 buildings are saved in train.npy, test.npy and valid.npy files, respectively. Each 1325 row in the NumPy file is a DINOv2 feature of a single building along with a label in its last column. 1326 The five roof classes are there in the data: 1-Metal Sheet, 2-Thatch, 3-Asbestos, 4-Concrete and 5-No 1327 Roof. The feature extraction is further explained in our research paper. 1328 1329 C.4.2 WAS THE "RAW" DATA SAVED IN ADDITION TO THE 1330 PREPROCESSED/CLEANED/LABELED DATA (E.G., TO SUPPORT UNANTICIPATED FUTURE 1331 USES)? IF SO, PLEASE PROVIDE A LINK OR OTHER ACCESS POINT TO THE "RAW" DATA. 1332 1333 Yes. Along with patches and their labels, the dataset contains the raw data. 1334 1335 IS THE SOFTWARE USED TO PREPROCESS/CLEAN/LABEL THE INSTANCES AVAILABLE? C.4.3 1336 IF SO, PLEASE PROVIDE A LINK OR OTHER ACCESS POINT. 1337 1338 The Python script is used to prepare patches and label different deep-learning models (e.g., UNet and 1339 YOLOv8). The Python script is available in the provided zip files. 1340 1341 C.4.4 ANY OTHER COMMENTS? 1342 1343 None. 1344 1345 C.5 USES 1346 1347 HAS THE DATASET BEEN USED FOR ANY TASKS ALREADY? IF SO, PLEASE PROVIDE A C.5.1 1348 DESCRIPTION 1349

At the time of preparing this datasheet, the dataset was only used for tasks performed in our paper.

1350 1351 1252	C.5.2	Is there a repository that links to any or all papers or systems that use the dataset? If so, please provide a link or other access point
1353 1354	No.	
1355 1356	C.5.3	What (other) tasks could the dataset be used for?
1357 1358 1359	There a etc.	re other objects in the images that can also be mapped, for example, trees, roads, water bodies,
1360 1361 1362 1363 1364 1365 1366 1367	C.5.4	Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses? For example, is there anything that a dataset consumer might need to know to avoid uses that could result in unfair treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other risks or harms (e.g., legal risks, financial harms)? If so, please provide a description. Is there anything a dataset consumer could do to mitigate these risks or harms?
1368 1369	There is	s no risk of using this dataset.
1370 1371 1372	C.5.5	ARE THERE TASKS FOR WHICH THE DATASET SHOULD NOT BE USED? IF SO, PLEASE PROVIDE A DESCRIPTION.
1373 1374	None.	
1375 1376 1377	C.5.6	ANY OTHER COMMENTS
1378	None.	
1380	C.6 I	DISTRIBUTION
1381 1382 1383 1384	C.6.1	WILL THE DATASET BE DISTRIBUTED TO THIRD PARTIES OUTSIDE OF THE ENTITY (E.G., COMPANY, INSTITUTION, ORGANIZATION) ON BEHALF OF WHICH THE DATASET WAS CREATED? IF SO, PLEASE PROVIDE A DESCRIPTION.
1385 1386 1387 1388	Current ?view our web	ly, the dataset is made available through this anonymous url: https://osf.io/us628/ _only=3c25a48d420f4ec7a43cb76e66e92b26. But later, we will publish through opage.
1389 1390 1391	C.6.2	How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)?
1392 1393	Same a	s above answer.
1394 1395	C.6.3	WHEN WILL THE DATASET BE DISTRIBUTED?
1396 1397	We will	distribute our dataset through other sources as soon as our paper is accepted.
1398 1399 1400 1401 1402 1403	C.6.4	WILL THE DATASET BE DISTRIBUTED UNDER A COPYRIGHT OR OTHER INTELLECTUAL PROPERTY (IP) LICENSE, AND/OR UNDER APPLICABLE TERMS OF USE (TOU)? IF SO, PLEASE DESCRIBE THIS LICENSE AND/OR TOU, AND PROVIDE A LINK OR OTHER ACCESS POINT TO, OR OTHERWISE REPRODUCE, ANY RELEVANT LICENSING TERMS OR TOU, AS WELL AS ANY FEES ASSOCIATED WITH THESE RESTRICTIONS

The dataset will be released under Open Data Commons Open Database License (ODbL) v1.0 licence.

1404 1405 1406 1407 1408	C.6.5	Have any third parties imposed IP-based or other restrictions on the data associated with the instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions.	
1409 1410 1411	No		
1412 1413 1414 1415	C.6.6	DO ANY EXPORT CONTROLS OR OTHER REGULATORY RESTRICTIONS APPLY TO THE DATASET OR TO INDIVIDUAL INSTANCES? IF SO, PLEASE DESCRIBE THESE RESTRICTIONS, AND PROVIDE A LINK OR OTHER ACCESS POINT TO, OR OTHERWISE REPRODUCE, ANY SUPPORTING DOCUMENTATION	
1416 1417	No		
1418 1419	C.6.7	Any other comments?	
1420 1421	None.		
1422 1423	C.7 I	DATASET MAINTENANCE	
1424	C.7.1	WHO IS SUPPORTING/HOSTING/MAINTAINING THE DATASET?	
1425 1426 1427	Two of details	our authors will be responsible for hosting and maintaining our dataset. We will add more here as soon as our paper accepted.	
1428 1429 1430	C.7.2	How can the owner/curator/manager of the dataset be contacted (e.g., email address)?	
1431 1432 1433	Curator details	Curators of this dataset can be communicated through our email addresses. We will provide more details as soon as our paper is accepted.	
1434 1435	C.7.3	IS THERE AN ERRATUM? IF SO, PLEASE PROVIDE A LINK OR OTHER ACCESS POINT.	
1436 1437	No, this	s is the initial release.	
1438 1439 1440 1441	C.7.4	WILL THE DATASET BE UPDATED (E.G., TO CORRECT LABELING ERRORS, ADD NEW INSTANCES, DELETE INSTANCES)? IF SO, PLEASE DESCRIBE HOW OFTEN, BY WHOM, AND HOW UPDATES WILL BE COMMUNICATED TO DATASET CONSUMERS (E.G., MAILING LIST, GITHUB)?	
1443	In case	of any updates, we will communicate through our webpage (we will provide the link later).	
1444 1445 1446 1447 1448 1449	C.7.5	IF THE DATASET RELATES TO PEOPLE, ARE THERE APPLICABLE LIMITS ON THE RETENTION OF THE DATA ASSOCIATED WITH THE INSTANCES (E.G., WERE THE INDIVIDUALS IN QUESTION TOLD THAT THEIR DATA WOULD BE RETAINED FOR A FIXED PERIOD OF TIME AND THEN DELETED)? IF SO, PLEASE DESCRIBE THESE LIMITS AND EXPLAIN HOW THEY WILL BE ENFORCED.	
1450 1451	N/A.		
1452 1453 1454 1455 1456 1457	C.7.6	WILL OLDER VERSIONS OF THE DATASET CONTINUE TO BE SUPPORTED/HOSTED/MAINTAINED? IF SO, PLEASE DESCRIBE HOW. IF NOT, PLEASE DESCRIBE HOW ITS OBSOLESCENCE WILL BE COMMUNICATED TO DATASET CONSUMERS. THE DATASET HAS ALREADY BEEN UPDATED; OLDER VERSIONS ARE KEPT AROUND FOR CONSISTENCY	

1458	C.7.7	IF OTHERS WANT TO EXTEND/AUGMENT/BUILD ON THIS DATASET. IS THERE A
1459		MECHANISM FOR THEM TO DO SO? IF SO, IS THERE A PROCESS FOR
1460		TRACKING/ASSESSING THE QUALITY OF THOSE CONTRIBUTIONS. WHAT IS THE
1461		PROCESS FOR COMMUNICATING/DISTRIBUTING THESE CONTRIBUTIONS TO USERS?
1462	NT/ A	
1463	N/A	
1464	C 7 8	ANY OTHER COMMENTS?
1465	C.7.0	ANY OTHER COMMENTS?
1466	None.	
1467		
1468		
1469		
1470		
1471		
1472		
1473		
1474		
1475		
1470		
1478		
1479		
1480		
1481		
1482		
1483		
1484		
1485		
1486		
1487		
1488		
1489		
1490		
1491		
1492		
1493		
1494		
1495		
1496		
1497		
1498		
1499		
1500		
1501		
1502		
1503		
1504		
1505		
1500		
1507		
1500		
1510		
1511		
1011		