

# MULTISPECTRAL REMOTE SENSING

Master's Degree in Geospatial Engineering

## LABORATORY PROJECT 2

DEEP LEARNING FOR BUILDING DETECTION USING HIGH-RESOLUTION IMAGES	
<b>OBJECTIVE</b>	Development of a complete workflow for building segmentation in Greater Lisbon using U-Net (PyTorch) and high-resolution orthophotos
<b>DURATION</b>	5 May 2026 – 26 May 2026
<b>EVALUATION</b>	Oral presentation – 2 June 2026
<b>BIBLIOGRAPHY</b>	<p><a href="#">The U-Net : A Complete Guide   Medium</a></p> <p><a href="#">Building a U-Net Architecture for Image Segmentation with Python and Keras   by Musti Öztemiz   Medium</a></p> <p><a href="#">Cook your First U-Net in PyTorch. A magic recipe to empower your image...   by Mostafa Wael   TDS Archive   Medium</a></p> <p><a href="#">UNet for Building Segmentation (PyTorch)</a></p> <p><a href="https://www.youtube.com/watch?v=IHq1t7NxS8k">https://www.youtube.com/watch?v=IHq1t7NxS8k</a></p> <p><a href="https://campus.datacamp.com/courses/deep-learning-for-images-with-pytorch/image-segmentation?ex=9">https://campus.datacamp.com/courses/deep-learning-for-images-with-pytorch/image-segmentation?ex=9</a></p>

## Building detection with U-Net workflow

<p style="text-align: center;"><b>STEP 1</b></p> <p style="text-align: center;"><b>Area of interest (AOI) definition and data preparation in QGIS</b></p>	<p><b>1.1 Define the AOI</b></p> <p>Select a sub-region inside the high-resolution orthophoto, mainly covered by buildings</p> <p><b>1.2 Download building data from OSM<sup>1</sup></b></p> <p><b>1.3 Load the high-resolution orthophoto</b></p> <p><b>1.4 Check CRS consistency</b></p> <p>Both layers must use the same CRS; if different, reproject both to EPSG: 32629</p> <p><b>1.5 Clip data using AOI</b></p> <p>Clip both the raster and vector data using AOI</p> <p><b>Output: AOI orthophoto and AOI building shapefile</b></p>
<p style="text-align: center;"><b>STEP 2</b></p> <p style="text-align: center;"><b>Data preprocessing in PYTHON</b></p>	<p><b>2.1 Load input data (orthophoto and building shapefile)</b></p> <p><b>2.2 CRS check</b></p> <p><b>2.2 Convert the buildings vector file into a raster mask</b></p> <p><b>2.4 Save ground truth mask</b></p> <p><b>Output: binary building mask (0= background, 1= building)</b></p>
<p style="text-align: center;"><b>STEP 3</b></p> <p style="text-align: center;"><b>Tile generation for U-NET in PYTHON</b></p>	<p><b>3.1. Generate training patches</b></p> <p>Split orthophoto and mask into tiles (256 x 256 pixels)</p> <p><b>Output: image tiles (X) and mask tiles (y)</b></p>
<p style="text-align: center;"><b>STEP 4</b></p> <p style="text-align: center;"><b>U-Net model training in PYTHON</b></p>	<p><b>4.1 Define the model architecture</b></p> <p>Encoder: ResNet34 (pretrained on ImageNet)</p> <p>Decoder: U-Net with skip connections</p> <p><b>4.2 Training configuration</b></p> <p>Loss function: BCE + Dice</p> <p>Optimizer: Adam (lr= 1e-3)</p> <p>Epochs: 10 to 30</p> <p><b>Output: trained U-Net model for building segmentation</b></p>

<sup>1</sup> <https://download.geofabrik.de/europe/portugal.html>

<p><b>STEP 5</b></p> <p><b>Prediction and post-processing in PYTHON</b></p>	<p><b>5.1 Sliding window prediction</b></p> <p>Split the image into overlapping patches (e.g., 256 x 256); run prediction for each patch; store outputs in a full-size probability map</p> <p><b>Output: raw probability map (continuous values 0-1)</b></p> <p><b>5.2 Seamless reconstruction</b></p> <p>Combine overlapping predictions by averaging overlapping regions for a smooth output (removes tile boundary artifacts)</p> <p><b>Output: seamless probability raster (no tile borders)</b></p> <p><b>5.3 Thresholding (binary building map)</b></p> <p>Convert probabilities to classes by applying a threshold (commonly 0.5); if <math>\geq 0.5</math>, building (1), else background (0)</p> <p><b>Output: binary building map</b></p> <p><b>5.4 Export the prediction as GeoTIFF</b></p> <p><b>Output: Geotiff file fully compatible with QGIS and GIS workflows</b></p> <p><b>5.5 Evaluate the model performance</b></p> <p>Compare the predicted map versus the ground truth mask</p> <ul style="list-style-type: none"> <li>• IoU (spatial overlap quality)</li> <li>• Dice (segmentation similarity, especially effective for buildings)</li> <li>• Accuracy (pixel-wise correctness)</li> </ul> <p><b>Output: quantitative performance report</b></p>
-----------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

## U-NET HYPERPARAMETERS

Hyperparameter	Meaning	Standard Values
<b>Epochs (EPOCHS)</b>	The number of times the model sees the entire training dataset during training	50-200 epochs are common for U-Net training <i>Too few epochs: underfitting (model does not learn enough patterns)</i> <i>Too many: overfitting (model memorizes training data and performs poorly on new data)</i>
<b>Patch size (PATCH_SIZE)</b>	Dimensions of the image tiles (patches) used as input to the network	128 x 128, 256 x 256, or 512 x 512 <i>Small patches: less spatial context but lower GPU memory</i> <i>Large patches: more contextual information and better object continuity, but require more GPU memory and longer training time</i>
<b>Prediction stride (stride)</b>	Distance (in pixels) that the sliding window moves when extracting patches; controls patch overlap	128 with a patch size of 256 produces 50% overlap <i>Smaller stride: more overlap, smoother predictions, better border continuity, but more computation</i> <i>Larger stride: faster processing, but may create visible tile boundaries</i>
<b>Batch size (batch_size)</b>	The number of image patches processed together before updating weights	2-16 are typical values for segmentation (8 is a common compromise between stability and memory usage) <i>Small batches (2-8): lower memory consumption but noisier gradient updates</i> <i>Large batches (32-64): more stable and faster training, but require substantially more GPU memory</i>
<b>Number of filters</b>	Number of feature maps learned at each convolution level; determines model capacity and feature extraction complexity	64 → 128 → 256 → 512 → 1024 (typical U-Net architecture) More filters: improved feature extraction and potentially higher accuracy; however, this also increases computational cost, training time, and GPU memory requirements
<b>Optimizer (optimizer)</b>	Algorithm responsible for updating model weights during training	Common optimizers: Adam: fast convergence and stable training; most used in U-Net SGD: slower training, but may provide better generalization RMSprop: common in medical image segmentation
<b>Learning rate (lr)</b>	Controls how much the network weights are adjusted after each training step	0.001 (1e-3) or 0.0001 (1e-4) (a suitable learning rate produces a smooth and progressive decrease in training loss) <i>Too high: unstable training and oscillating loss</i> <i>Too low: very slow convergence</i>
<b>Loss function (BCE+Dice)</b>	Measures the segmentation error by comparing predictions with reference labels	BCE + Dice loss (balances pixel accuracy and object shape quality) <i>BCE (Binary Cross Entropy): evaluates pixel-wise classification accuracy</i> <i>Dice loss: evaluates spatial overlap between prediction and ground truth</i>
<b>Activation function (ReLU)</b>	Introduces non-linearity into the network, enabling the model to learn complex spatial patterns	ReLU (Rectified Linear Unit): converts negative values to zero; computationally efficient Sigmoid: converts outputs to probabilities between 0 and 1 (commonly used in binary segmentation output layers) Softmax: used for multiclass segmentation problems
<b>Threshold (threshold)</b>	Converts predicted probabilities into final binary segmentation classes; Pixels above the threshold are assigned to the positive class (in our case, buildings)	0.5 (threshold selection strongly influences the final appearance and accuracy of the segmentation map) <i>Lower threshold (0.3-0.4): detects more objects and reduces false negatives (FN), but increases false positives (FP) and noise</i> <i>High threshold (0.7-0.8): cleaner segmentation with fewer false positives (FP), but may miss small or uncertain objects and produce fragmented results</i>

## ACCURACY METRICS (BINARY SEGMENTATION)

```
=====
Final full-image metrics
=====
Accuracy : 0.8691
Precision : 0.8855
Recall   : 0.8713
F1-score : 0.8783
IoU      : 0.7831

Confusion matrix
Rows = reference mask
Columns = prediction

          Pred 0   Pred 1
True 0   3443740  530487
True 1   605496  4100723
```

	Pred 0 (Background)	Pred 1 (Buildings)
True 0 (Background)	True Negative (TN)	False Positive (FP)
True 1 (Buildings)	False Negative (FN)	True Positive (TP)

TP = correctly detected building pixels

FP = background pixels wrongly classified as buildings

FN = building pixels missed by the model

TN = correctly detected background pixels

	Pred 0 (Background)	Pred 1 (Buildings)	Total	RECALL
True 0 (Background)	3443740	530487	3974227	
True 1 (Buildings)	605496	4100723	4706219	0.8713
Total	4049236	4631210	8680446	
PRECISION		0.8855		

$$\text{Recall} = \left( \frac{\text{TP}}{\text{TP} + \text{FN}} \right) = \left( \frac{4100723}{4706219} \right) = 0.8713$$

$$\text{Precision} = \left( \frac{\text{TP}}{\text{TP} + \text{FP}} \right) = \left( \frac{4100723}{4631210} \right) = 0.8855$$

$$\text{Accuracy} = \left( \frac{\text{TP} + \text{TN}}{\text{TN} + \text{FN} + \text{FP} + \text{TP}} \right) = \left( \frac{7544463}{8680446} \right) = 0.8691$$

$$\text{F1 - score} = \left( \frac{2 * \text{TP}}{2 * \text{TP} + \text{FP} + \text{FN}} \right) = \left( \frac{8201446}{9337429} \right) = 0.8783$$

For binary segmentation F1 - score = Dice coefficient

$$\text{IoU} = \left( \frac{\text{TP}}{\text{TP} + \text{FP} + \text{FN}} \right) = \left( \frac{4100723}{5236706} \right) = 0.7831$$

IoU (Intersection over Union) also called Jaccard Index  

$$= \left( \frac{\text{intersection of prediction and ground truth}}{\text{union of prediction and ground truth}} \right)$$



Reference Building Mask



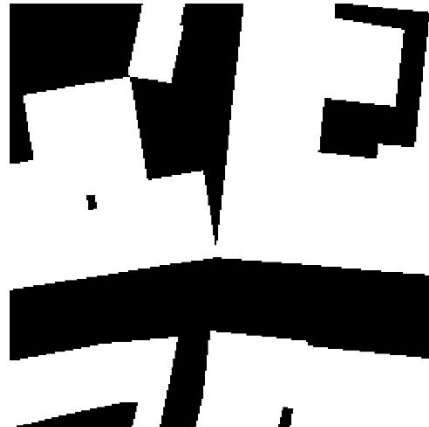
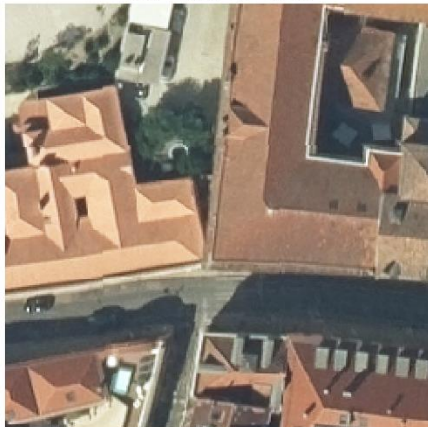
Building Detection - Final Prediction



Training pair: Orthophoto patch and corresponding mask  
Orthophoto patch 49  
row=1024, col=256  
Mask patch 49  
row=1024, col=256



Training pair: Orthophoto patch and corresponding mask  
Orthophoto patch 100  
row=2048, col=1024  
Mask patch 100  
row=2048, col=1024



Training pair: Orthophoto patch and corresponding mask  
Orthophoto patch 115  
row=2304, col=1792  
Mask patch 115  
row=2304, col=1792

