

A CNN REGRESSION MODEL TO ESTIMATE BUILDINGS HEIGHT MAPS USING SENTINEL-1 SAR AND SENTINEL-2 MSI TIME SERIES

Ritu Yadav¹, Andrea Nascetti^{1,2}, Yifang Ban¹

¹KTH Royal Institute of Technology (Sweden), ²University of Liège (Belgium)

ABSTRACT

Accurate estimation of building heights is essential for urban planning, infrastructure management, and environmental analysis. In this study, we propose a supervised Multimodal Building Height Regression Network (MBHR-Net) for estimating building heights at 10m spatial resolution using Sentinel-1 (S1) and Sentinel-2 (S2) satellite time series. S1 provides Synthetic Aperture Radar (SAR) data that offers valuable information on building structures, while S2 provides multispectral data that is sensitive to different land cover types, vegetation phenology, and building shadows. Our MBHR-Net aims to extract meaningful features from the S1 and S2 images to learn complex spatio-temporal relationships between image patterns and building heights. The model is trained and tested in 10 cities in the Netherlands. Root Mean Squared Error (RMSE), Intersection over Union (IOU), and R-squared (R^2) score metrics are used to evaluate the performance of the model. The preliminary results (3.73m RMSE, 0.95 IoU, 0.61 R^2) demonstrate the effectiveness of our deep learning model in accurately estimating building heights, showcasing its potential for urban planning, environmental impact analysis, and other related applications.

Index Terms— Building Height Estimation, Sentinel, Deep Learning, Fusion, Regression, Time Series.

1. INTRODUCTION

More than half of the world’s population currently lives in cities. By 2050, an estimated 7 out of 10 people will likely live in urban areas. While cities contribute more than 80% of global GDP they are also accountable for major energy consumption and carbon emission [1]. Urbanization monitoring is essential to assess its impact on the environment and support sustainable development. Accurate building height estimation plays an important role in urban planning, as it is an indicator of urban heat islands effect, population, energy consumption, and urban climate.

Earth Observation (EO) has been highlighted as an effective tool for mapping large-scale human settlements. Several

methodologies have been developed to extract building footprints in the last decades. Various large-scale global urban footprint data sets are available and widely used by the scientific community [2, 3]. However, these data sets are intrinsically two-dimensional and do not provide information on building height. In recent years, some studies have tried to fill this gap and estimate building heights from satellite imagery. For example, [4], proposed a Random Forest (RF) regressor for continental-scale height mapping at 1 km spatial resolution. The authors used Landsat-8 OLI, Sentinel-1 SAR and various handcrafted spatial features along with auxiliary data. Reference data was derived from a combination of open street maps, government websites and commercial maps. The authors of [5] developed a VVH building height indicator from Sentinel-1 SAR data and estimated building heights at 500m resolution. The indicator was evaluated in major cities in the US with ICESat data as reference and achieved an RMSE of 1.5m. [6], extended the World Settlement Footprint (WFS) [2], including the building heights derived by the DSM collected by the TanDEM-X mission, and generated the WFS 3D data set at 90m resolution. The estimated building heights have been validated showing a promising accuracy with 6.01m RMSE score. However, it relies on a commercial DSM that is not easy to update frequently. [7] presented a Support Vector Machine (SVM) regression model to derive building heights at 10m resolution with RMSE of 3.2m to 4.2m; the authors used a set of handcrafted spatial and temporal features from Sentinel-1 and Sentinel-2 time series as input to the model. The approach is tested in several cities in Germany using available ALS (Airborne Laser Scanner) data as a reference. [8] estimated building height for China at 10m resolution and achieved 6.1m RMSE. The authors used a combined approach from [4] and [7] with additional ALOS PALSAR, WFS footprints and DEM data. The reference data is derived from Baidu map services with an assumption of each floor height to be 3m.

The aim of this study is to investigate a supervised Convolutional Neural Network (CNN) based regression model for estimating building height using only freely available Sentinel-1 SAR and Sentinel-2 MSI time series data. We frame the task of generating building height maps as a pixel-wise regression task, assuming the following: (1) zero-pixel values represent no buildings (as usual in urban footprint

The research is part of the project ‘EO-AI4GlobalChange’ funded by Digital Futures.

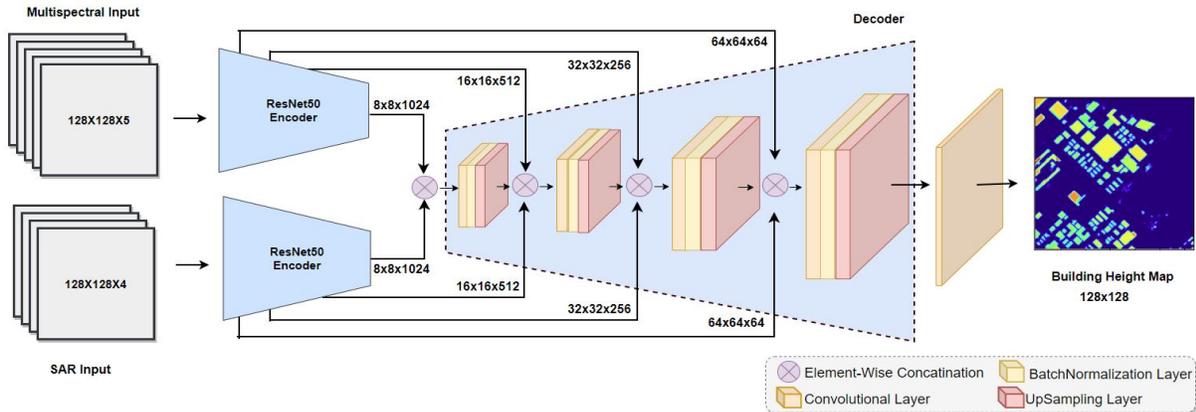


Fig. 1: The proposed MBHR-Net for building height estimation.

data), and (2) pixel values greater than 1 directly correspond to estimated building height. We developed a CNN regression model based on the U-Net architecture that takes Sentinel-1 and Sentinel-2 multi-temporal data and estimates building height at a 10m spatial resolution.

2. DATA DESCRIPTION

The data used in this study include the Sentinel-1 Ground Range Detected (GRD) and Sentinel-2 MSI Level-2A. We collected data on ten largest cities in The Netherlands namely Amsterdam, Rotterdam, The Hague, Utrecht, Eindhoven, Groningen, Breda, Tilburg, Nijmegen and Almere. For reference, we used 3D Bag data developed by the 3D geoinformation Group of the Technical university of Delft. The database is completely open source and automatically generated from the Buildings and Addresses Register (BAG) providing 2D footprints and from the National Altimetric Model (AHN) derived from ALS data. The database contains multiple Levels Of Detail (LOD) building models (LOD1.2, LOD1.3 and LOD 2.2). We selected LOD1.3 data for this study. We created non-overlapping tiles of size 128x128 pixels within the administrative boundaries of the 10 metropolitan areas provided by the European Environment Agency. Corresponding to each tile we collected S1 and S2 time series using Google Earth Engine’s python API [9] and collected LOD1.3 data from 3D BAG database. For S2 MSI data, we generated monthly cloud-free composites and downloaded 5 bands (Red, Green, Blue, NIR and SWIR). For S1 SAR data we first computed the monthly average to reduce the speckle for both ascending and descending orbits, then downloaded 4 bands (VV, VH polarizations for both orbits). The S2 data contains 12 images per tile (one image for each month) and the year is matched with the acquisition of AHN data. The input data is downloaded at 10m spatial resolution, and the reference data are rasterized and resampled to match 10m. We divided the dataset into training and test set using an 80-20 ratio, resulting in 1,737 training samples and 434 test samples.

3. METHODOLOGY

3.1. MBHR-Net Architecture

Figure 1 shows the architecture of the proposed MBHR-Net, consisting one branch for learning multispectral features of S2 images and another branch for learning SAR backscatter features of S1 images. The S2 branch takes a five-channel input (red, green, blue, NIR and SWIR bands) and the S1 branch takes a four-channel input (VV, VH for both descending and ascending orbits).

We adopted U-Net architecture, a widely used encoder-decoder based segmentation network. An encoder compresses the salient features (feature maps) of the input images and a decoder upsamples the compressed features to predict output of same size as input. The encoder block has a number of repetitions (level) of the sequence; convolutional layer, maxpooling layer and batch normalization layer. The decoder has a sequence of convolutional layers, upsampling layer and back normalization layer to output a segmentation map.

Our MBHR-Net contains two encoders one in each branch to learn different modality features separately. We adopted ResNet50 with four levels as encoder. The output feature maps of each level is of different size capturing different semantics. From each encoder, four feature maps are extracted. The size of the feature maps are (64x64), (32x32), (16x16) and (8x8). These multiscale features are fused using element-wise concatenation operation. Via skip connection, the fused features maps are combined with the same size decoder layers. These skip connections help the decoder network to condition not only on the latent representation but also on intermediate representations of the encoder, which lead to fine-grained details in predictions[10, 11]. The combined feature maps are upsampled and processed through decoder layers. At the end of the decoder network we use a convolutional later with (1x1) kernel and *ReLU* activation function making it a regression layer.

3.2. Augmentation strategy and Training

For one reference patch we have 24 input images i.e. 12 time series images for each modality. These 12 images capture surrounding features in different months. We assume that in a 12-month period, there are negligible changes in building heights but the season changes surrounding conditions. Therefore, the 12-month images can be treated as augmented images, creating 12 augmented pairs. With augmentation, the size of training samples increased from 1734 to 1734x12. As training losses we used a weighted combination of two regression losses, Mean Squared Error loss (MSE) and Cosine Similarity (CS) loss given in equation Eq. 1 – 2. We used 0.8 weight for CS loss and 0.2 for MSE loss.

$$MSE = (y_{true} - y_{pred})^2 \quad (1)$$

$$CS = - \sum \|y_{true}\|^2 * \|y_{pred}\|^2 \quad (2)$$

We trained the model for 100 epochs with batch size 4, adam optimizer and 0.0001 as initial learning rate. For better convergence, the learning rate is decayed until 0.00001. The decay steps are controlled with the "reduce on plateau" method. The code is implemented in Keras and the model is trained for 6 hours on Google colab GPU.

4. RESULTS AND EVALUATION

The predicted building height maps MBHR-Net are first filtered using building footprints obtained by binarizing the reference height maps. For binarization, a pixel is set to be building pixel (1.0) if the pixel value is > 1.0 otherwise set to no building (0.0 value). The filtered building height maps are evaluated using two metrics Root Mean Square Error (RMSE) and R^2 score given in Eq. 3, 4, where n is the number of validation samples, $BH_{est,i}$ is the estimated value of the height of the building and $BH_{ref,i}$ is reference building height. RMSE indicates the accuracy of predicted heights with respect to reference and R^2 score estimates the model effectiveness in learning variance in the building heights.

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (BH_{ref,i} - BH_{pred,i})^2}{n}} \quad (3)$$

$$R^2 = 1 - \frac{(n-1) \sum_{i=1}^n (BH_{ref,i} - BH_{pred,i})^2}{(n-2) \sum_{i=1}^n (BH_{ref,i} - BH_{pred,i})^2} \quad (4)$$

For accurate building height estimation, it is important to ensure the alignment of predicted buildings with the reference. We evaluated building alignment using well known metric Intersection over Union (IoU). RMSE and R^2 scores are calculated using reference and network predictions directly whereas IoU is calculated on binarized (building and no building) references and predictions.

A high-quality building height estimation model is characterized by a low RMSE, a high IoU, and a high R^2 score.

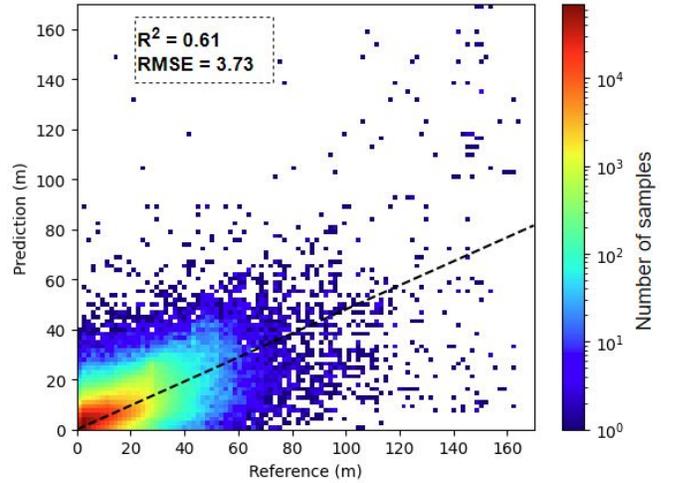
These metrics provide insight into different aspects of the model's performance. The RMSE, IoU and R^2 scores of the proposed MBHR-Net are 3.73 m, 0.95 and 0.61. The RMSE score of 3.73 m suggests that, on average, the model's accuracy is approximately one floor showing a remarkable accuracy considering the spatial resolution of S1 and S2 imagery (10-20 m). The IoU score of 0.95 indicates that the predicted height regions overlap significantly with the ground truth regions, demonstrating the model's ability to precisely identify buildings. The R^2 score is 0.61, suggesting that the model can explain a substantial portion of the variance in building heights but the model may have limitations in accurately capturing certain factors or complexities affecting building heights.

Table 1: Performance metrics on test set.

	RMSE (m)	R^2	IoU
mean \pm std	3.73 \pm 2.01	0.61 \pm 0.12	0.95 \pm 0.10

Figure 1 shows a scatter plot between predicted heights and reference heights of all pixels with values greater than 1 in the test set. In the scatterplot we reported that the model shows a fair correlation between reference heights and predictions. But there is an overall underestimation of heights, which we aim to improve in our future works.

Fig. 2: Scatter plot of the predicted and reference height values. The evaluation is on test set and the height values are in meters (m).



We also present a few test samples for qualitative evaluation (see Figure 3). The first two rows show good height estimation examples where the predicted heights are accurate with few underestimations. The third sample shows comparatively less accurate height estimations with error ranges of 0 to 4 meters. This is possibly an example of the type of complexity which is not accurately learned by the model. The high IOU score of the model can be verified from all samples showing a good alignment between the predicted buildings

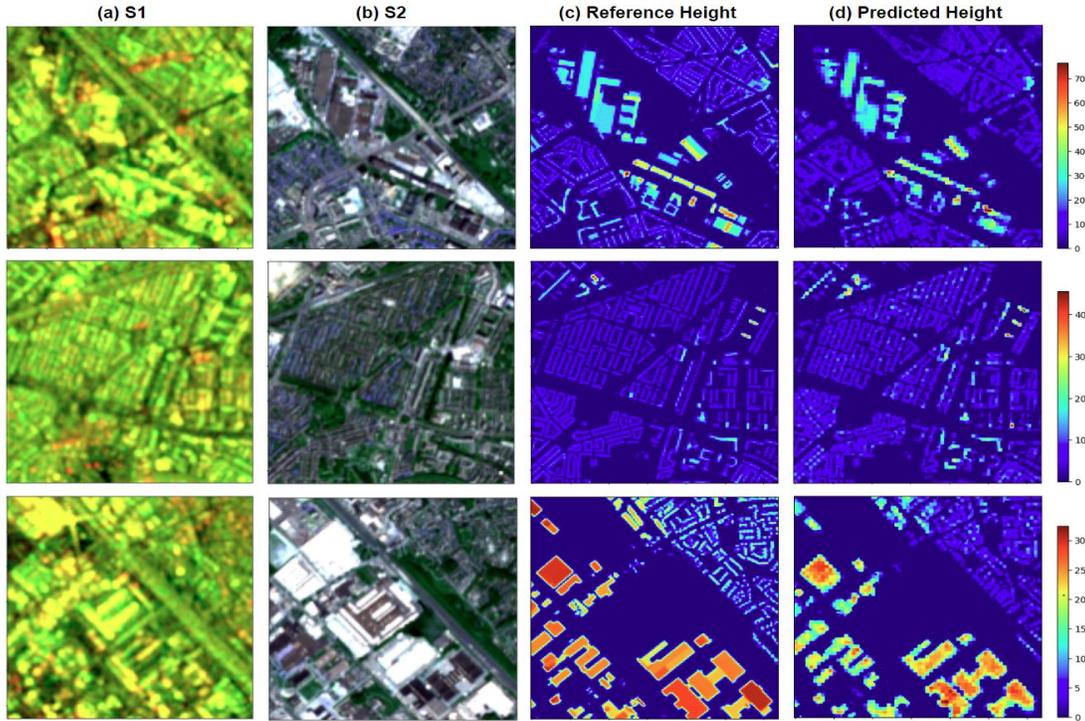


Fig. 3: Result Samples. Three sample results are visualized in 3 rows. From left to right, Sentinel-1 SAR, Sentinel-2 MSI, Reference Building Height and Predicted Building Height (in meters) from our MBHR-Net.

and the ground truth buildings, including the boundary areas.

5. CONCLUSION

In this study, we developed a deep learning model for building height estimation using combined Sentinel-1 SAR and Sentinel-2 MSI time series. The performance evaluation of MBHR-Net demonstrated promising accuracy in both height estimation, with an RMSE of 3.73m, and building footprint delineation, with a 95% IoU. These results indicate the potential of MBHR-Net for estimating building heights with accurate building footprint delineation. The potential future investigation directions are to expand the data set to include different geographic regions, building types, incorporate additional data sources, and develop a more advanced deep learning model to handle the complexity of such large data to build a more generalized approach.

6. REFERENCES

- [1] UN, “The sustainable development goals report 2022,” 2022.
- [2] Mattia Marconcini, Annkatrin Metz-Marconcini, Thomas Esch, and Noel Gorelick, “Understanding current trends in global urbanisation—the world settlement footprint suite,” *GI Forum*, 2021.
- [3] Sebastian Hafner, Yifang Ban, and Andrea Nascetti, “Unsupervised domain adaptation for global urban extraction using sentinel-1 sar and sentinel-2 msi data,” *Remote Sensing of Environment*, 2022.
- [4] Mengmeng Li, Elco Koks, Hannes Taubenböck, and Jasper van Vliet, “Continental-scale mapping and analysis of 3d building structure,” *Remote Sensing of Environment*, 2020.
- [5] Xuecao Li, Yuyu Zhou, Peng Gong, Karen C Seto, and Nicholas Clinton, “Developing a method to estimate building height from sentinel-1 data,” *Remote Sensing of Environment*, 2020.
- [6] Thomas Esch, Elisabeth Brzoska, Stefan Dech, Benjamin Leutner, Daniela Palacios-Lopez, Annkatrin Metz-Marconcini, Mattia Marconcini, Achim Roth, and Julian Zeidler, “World settlement footprint 3d—a first three-dimensional survey of the global building stock,” *Remote sensing of environment*, 2022.
- [7] David Frantz, Franz Schug, Akpona Okujeni, Claudio Navacchi, Wolfgang Wagner, Sebastian van der Linden, and Patrick Hostert, “National-scale mapping of building height using sentinel-1 and sentinel-2 time series,” *Remote Sensing of Environment*, 2021.

- [8] Wan-Ben Wu, Jun Ma, Ellen Banzhaf, Michael E Meadows, Zhao-Wu Yu, Feng-Xiang Guo, Dhritiraj Sen-gupta, Xing-Xing Cai, and Bin Zhao, “A first chinese building height estimate at 10 m resolution (cnbh-10 m) using multi-source earth observations and machine learning,” *Remote Sensing of Environment*, 2023.
- [9] Noel Gorelick, Matt Hancher, Mike Dixon, Simon Ilyushchenko, David Thau, and Rebecca Moore, “Google earth engine: Planetary-scale geospatial analysis for everyone,” *Remote Sensing of Environment*, 2017.
- [10] Ritu Yadav, Andrea Nascetti, Hossein Azizpour, and Yi-fang Ban, “Unsupervised flood detection on sar time series,” *arXiv preprint arXiv:2212.03675*, 2022.
- [11] Lei Ma, Yu Liu, Xueliang Zhang, Yuanxin Ye, Gaofei Yin, and Brian Alan Johnson, “Deep learning in re-mote sensing applications: A meta-analysis and review,” *ISPRS journal of photogrammetry and remote sensing*, 2019.