



**HAL**  
open science

## An image-based Deep Learning workflow for 3D Heritage point cloud semantic segmentation

Eugénio Pellis, Arnadi Murtiyoso, Andrea Masiero, Grazia Tucci, Michele Betti, Pierre Grussenmeyer

### ► To cite this version:

Eugénio Pellis, Arnadi Murtiyoso, Andrea Masiero, Grazia Tucci, Michele Betti, et al.. An image-based Deep Learning workflow for 3D Heritage point cloud semantic segmentation. 9th International Workshop 3D-ARCH "3D Virtual Reconstruction and Visualization of Complex Architectures". 2–4 March 2022, Mantua, Italy, Mar 2022, Mantua, France. pp.426-434, 10.5194/isprs-archives-XLVI-2-W1-2022-429-2022 . hal-03976192

**HAL Id: hal-03976192**

**<https://hal.science/hal-03976192>**

Submitted on 6 Feb 2023

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

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

# AN IMAGE-BASED DEEP LEARNING WORKFLOW FOR 3D HERITAGE POINT CLOUD SEMANTIC SEGMENTATION

E. Pellis<sup>1,3</sup>, A. Murtiyoso<sup>2</sup>, A. Masiero<sup>1</sup>, G. Tucci<sup>1</sup>, M. Betti<sup>1</sup>, P. Grussenmeyer<sup>3</sup>

<sup>1</sup> Department of Civil and Environmental Engineering (DICEA), University of Florence, 50139 Florence, Italy - (eugenio.pellis, andrea.masiero, grazia.tucci, michele.betti)@unifi.it

<sup>2</sup> Forest Resources Management Group, Institute of Terrestrial Ecosystems, Department of Environmental Systems Science, ETH Zürich, Switzerland - arnadidhstaratri.murtiyoso@usys.eth.ch

<sup>3</sup> Université de Strasbourg, INSA Strasbourg, CNRS, ICube Laboratory UMR 7357, Photogrammetry and Geomatics Group, 67000 Strasbourg, France – pierre.grussenmeyer@insa-strasbourg.fr

## Commission II

**KEY WORDS:** 3D Point Cloud, Deep Learning, Semantic Segmentation, Cultural Heritage

### ABSTRACT:

The interest in high-resolution semantic 3D models of historical buildings continuously increased during the last decade, thanks to their utility in protection, conservation and restoration of cultural heritage sites. The current generation of surveying tools allows the quick collection of large and detailed amount of data: such data ensure accurate spatial representations of the buildings, but their employment in the creation of informative semantic 3D models is still a challenging task, and it currently still requires manual time-consuming intervention by expert operators. Hence, increasing the level of automation, for instance developing an automatic semantic segmentation procedure enabling machine scene understanding and comprehension, can represent a dramatic improvement in the overall processing procedure. In accordance with this observation, this paper aims at presenting a new workflow for the automatic semantic segmentation of 3D point clouds based on a multi-view approach. Two steps compose this workflow: first, neural network-based semantic segmentation is performed on building images. Then, image labelling is back-projected, through the use of masked images, on the 3D space by exploiting photogrammetry and dense image matching principles. The obtained results are quite promising, with a good performance in the image segmentation, and a remarkable potential in the 3D reconstruction procedure.

## 1. INTRODUCTION

In recent years, the high level of automation achieved by the latest 3D acquisition technologies, like laser scanner or photogrammetry, allows to collect a large amount of data in a short time, and to obtain a high level of details in shape, geometry and semantic information (Shan and Toth, 2018). Thanks to these continuous technological developments, the use of such surveying methods had a significant increase, in particular for the study of heritage buildings, with the main aim of documentation, interpretation, protection, conservation and restoration of cultural heritage (Yastikli, N., 2007). Certain of the most advanced applications also deals with virtual reality and Building Information Modelling (H-BIM) (López et al., 2018).

Given the huge amount of 3D data that are usually collected nowadays to properly describe the sites of interests, automatic processing procedures shall be preferable in order to reduce the processing times. Nevertheless, automatic raw data processing for the creation of as-built 3D models still faces several significant challenges, often causing the need of time-consuming manual intervention, in particular to deal with the complexity of heritage buildings.

One of the key points to enable automation in 3D data processing is the development of a semantic segmentation procedure. This procedure involves the task of classifying each point of a 3D point cloud into classes or categories according to its semantic meaning, for example according to the constructive element typology (e.g., wall, columns, vaults, etc.) (Matrone et al., 2020). Certain of the main issues related to working with 3D point clouds are: the large data size, which implies long computing time, the unordered structure of point clouds, and the low availability of shared datasets and tools, which may be the basis of a common processing strategy, e.g. based on the use of properly developed artificial intelligence tools (Malinverni et al., 2019).

As a basic step for scene understanding and comprehension, several semantic segmentation techniques have already been tested in a wide range of different applications. Nowadays machine learning (ML) and deep learning (DL) techniques are the most extensively investigated and the most promising (Zhang et al., 2019). The main approaches to face semantic segmentation of 3D point clouds can be divided in two main groups: (i) projection-based methods and (ii) point-based methods.

Projection based methods leverage on an intermediate representation of the cloud in order to face with the unordered structure of 3D points, then they apply standard 2D approaches to perform the segmentation (Badrinarayanan et al., 2015), (Chen et al., 2017), (Cordts et al., 2016), and finally they re-project the extracted features on the starting shape or point cloud (Su et al., 2015), (Boulch et al., 2018).

Differently, point based methods work directly with the 3D points and they leverage of the full use of the characteristic of the raw point cloud data, considering all the spatial and geometrical information (Qi et al., 2017a).

This paper focuses on the development of a workflow for the classification of 3D point clouds based on a multi-view approach, exploring and testing the use of state-of-the-art neural networks on heritage scenarios, with the main aim of improving automation in 3D heritage model generation.

## 2. RELATED WORK

Semantic analysis has become a central topic in several applications, such as computer vision, robotics, or remote sensing. Numerous approaches have been developed to automate this task, including algorithmic, machine learning and deep learning approaches. Despite the increasing demand in heritage survey data processing, few works focus on dealing with the classification of 3D heritage point clouds.

Malinverni et al. (2019) described a method to label and cluster automatically a point cloud based on a supervised Deep Learning approach, using a state-of-the-art neural network called PointNet++ (Qi et al., 2017b). The results are promising but they still face several challenges due to the complexity of the training scene. Further developments of this research should include the association to the Neural Network of a structured ontology for the semantic parsing, the future specification of more detailed classes and the creation of a synthetic dataset to overcome the problem of very poor-annotated training data.

Matrone et al. (2020) make a comparison between machine learning and deep learning methods working directly with point cloud, and then, they proposed a new architecture named DGCNN-Mod+3DFeat that combined the positive aspects and advantages of Machine Learning and Deep Learning techniques, and they test it on the ARCH dataset benchmark with promising results. Both ML and DL proved to be approximately equally valuable: each of the two techniques alternatively outperformed the other one.

In their work, Grilli et al. (2019) presented a research for the classification of 3D heritage data (point clouds or polygonal mesh models) using different supervised learning approaches, including Machine Learning and Deep Learning. Their method works on 2D data ("texture-based approach) or directly the 3D data ("geometric-based approach) depending on the needs and scope of classification. The method was applied and validated on four different archaeological scenarios, proving its reliability and replicability, being effective in providing metric information as well.

Murtiyoso et al. (2021) developed an approach to semantically segment building façades in 6 categories using Deep Learning on rectified images by deploying pre-existing and pre-trained networks. The network outcome is back-projected into the 3D space by exploiting a depth map to generate a semantically segmented point cloud. The obtained overall accuracy is quite promising, yielding to a value of 79,8%. A mask-based approach has also been presented in Murtiyoso et al. (2022), investigating the use of back-projection into 3D space to more complex scenarios.

This work proposes the combined use of deep learning based semantic image segmentation with photogrammetric 3D reconstruction. The main aim is to introduce the semantic classification at the beginning of the classical photogrammetric workflow. The initial results obtained in our tests are quite promising, experimentally proving the potential effectiveness of using projection-based techniques for the generation of semantically segmented building point clouds.

### 3. DEVELOPED METHODOLOGY

In this paper we propose a multi-view based methodology to perform semantic segmentation of the heritage 3D point clouds. This approach is based on the segmentation of a 2D intermediate representation of the cloud, and then on the re-projection of the extracted features on the initial cloud. Despite working directly with 3D data provides an opportunity for a better understanding of the spatial and geometrical information, dealing with 2D images is usually a quite effective strategy. On one hand, it allows to exploit the tried-and-tested results in 2D image processing, in particular with Convolutional Neural Networks (CNNs), and, on the other hand, it allows to develop an automatic procedure for the creation of a directly-segmented cloud starting from photogrammetric images. In addition, at this time, a multi-view approach for heritage buildings semantic segmentation has never been tested.

Our methodology is composed of two main steps: (i) at first, the images are processed by a neural network to perform the

semantic segmentation directly on the images, and (ii) secondly, a back-projection procedure allows to transfer the image features on the 3D point cloud.

#### 3.1 The dataset

To train the various networks we used an ongoing image-based dataset, proposed by Pellis et al. (2021), specifically designed for heritage semantic segmentation. Datasets play an indispensable role in the Neural Network training phase and few datasets are freely available in heritage scenarios. Two of the most important datasets in this context are Architectural Elements Dataset (AHE dataset), a collection of 10,000 images classified in 10 types of architectural elements for the task of classification, and ARCH dataset (Matrone et al., 2020b), a benchmark for large scale heritage point cloud semantic segmentation, composed of 17 annotated scenes in 10 categories.

The proposed dataset (Pellis et al., 2021) will be composed by 9 heritage buildings, built in different period and characterized by different architectural style. It is an ongoing project, and, at the current stage, the dataset is composed by 5 buildings: (1\_SC) Spedale del Ceppo, Pistoia, (2\_OSA) Ospedale Sant'Antonio, Lastra a Signa, (3\_SSA) Basilica della Santissima Annunziata, Firenze e (4\_CG) Certosa del Galluzzo, Firenze, (5\_CB) Cappella Buontalenti, Firenze. The segmentation categories considered in this dataset are structured following the guidelines defined in the ARCH dataset, which refers to the Industry Foundation Class (IFC) file format, to the CityGML (LOD3/4) and to the Art and Architecture Thesaurus (AAT). Hence, the dataset has been segmented in 10 classes, which correspond to certain BIM constructive elements, including arch, column, moulding, floor, door/window, wall, stair, vault, roof, other and background. Table 1 shows the percentage of pixels in each of the ten classes. For each building, a set of representative images from various perspective and angles, and their corresponding pixel-level ground truth are available. The total number of available images is currently about 4,700. All such images were collected during photogrammetric surveys, hence their alignment information is available as well, and it is possible to calculate the 3D point corresponding to each pixel.

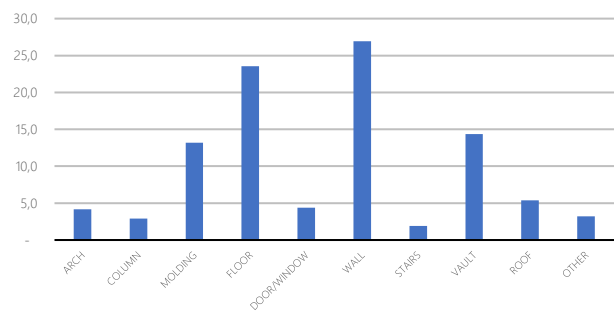


Table 1: General workflow of the developed methodology.

#### 3.2 Image Segmentation

The first step of the workflow consists in the segmentation of the intermediate images according with the conventions of the ARCHdataset. Image segmentation is a key topic in a lot of computer vision application and various algorithms have been developed in literature including thresholding methods, k-means clustering, region growing and other. However, over the past few years, deep learning networks have yielded a new generation of models, achieving the highest accuracy on the most popular benchmarks. The most performing architecture are Fully

Convolutional Networks (FCN), U-Net, SegNet and the DeepLab family.

In this paper for the task of image segmentation we used one of the most popular network architectures, DeepLabv3+. This network is an improvement of the previous DeepLab and it employs atrous convolution with upsampled filters to extract dense feature maps and to capture long range context. The network has been implemented with MATLAB using ResNet18 as a base classification architecture. Training a deep learning model from scratch is often unfeasible due to the dataset size required and due to long time for the initial tuning of the network. For these reasons we used transfer learning, a technique that allows to grab a model trained on a certain task and dataset and employ it on a different task. We used a pretrained version of ResNet18 trained on ImageNet dataset, one of the most popular and largest benchmarks for image semantic segmentation.

### 3.3 Back-projection procedure

The second step involves the back-projection of the classified pixels into the 3D point cloud. This is the most challenging step of the procedure because feature propagation on 3D space involves inevitably a loss of geometrical and spatial information, a loss of detail of shapes and the propagation of errors. In our proposed methodology we leverage on the photogrammetry principles using a masking method for the creation of a directly labelled point cloud starting from the images. The back-projection workflow starts from the initial segmentation of the survey images performed by a trained neural network. Secondly, according with the segmentation results, for each image and for each segmentation class a binary mask is created. The masked images were then employed during dense image matching in order to constrain the point cloud creation process. The results of this process are a separate dense 3D point clouds for each class. In the next section we are going to show some results of the developed procedure on our ongoing dataset.

## 4. RESULTS

### 4.1 Image Semantic Segmentation Results

This section reports some results on the semantic segmentation of the 2D images of historical buildings.

It is worth to notice that deep neural network training require thousands of images and case studies in order to ensure a proper and reliable performance. Given the current size of the ongoing dataset and the mentioned requirements, several quite simple tests with different data combinations are considered.

#### 4.1.1 Test 1

In the first test, a small portion of the dataset is used, including just one building, in this case Spedale del Ceppo. For the purpose of training, all the labelled images of the building were randomly divided in 60% for the training (896 images), 20% as validation test (300 images), and 20% as test set (300 images). The training was performed for 30 epochs, yielding a validation accuracy of 94.5%. Once trained, the network was deployed to obtain the prediction on the entire test set. Table 2 reports some of the most significant results. Since just one building was used, the images of training and test set are quite similar, hence this task is very simple and not really significant for the network generalization.

Global Accuracy	Mean Accuracy	Mean IoU	Weighted IoU	Mean BFScore
93,4%	94,1%	82,8%	87,1%	84,2%

Table 2: DeepLabv3+ results on the test set (Test 1).

The performance of the network on the test set is remarkable, yielding to an overall accuracy of 93.4% and an Intersection Over Union (IoU) of 82.8%.

Figure 1 shows a comparison between the input images (A), the ground-truth (B) and the prediction (C) obtained on three of the test set images. Figure 2 shows the confusion matrix.

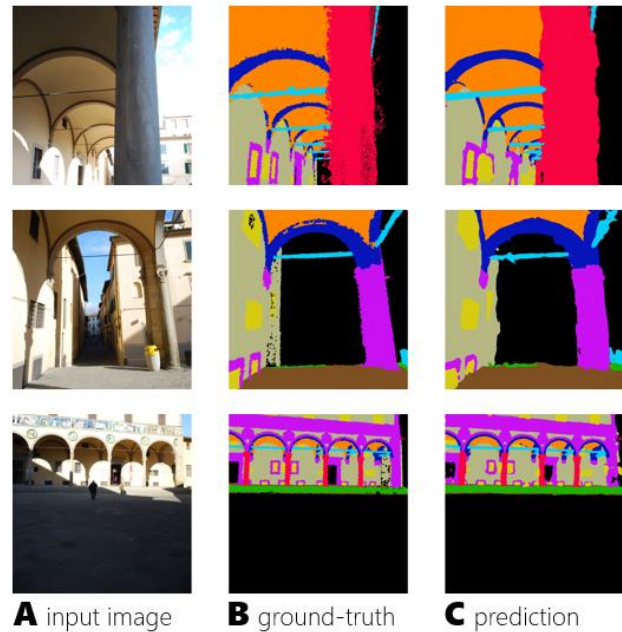


Figure 1: DeepLabv3+ results on the test set (Test 1): A) input image, B) ground-truth, C) prediction.

True Class \ Predicted Class	arch	column	moldings	floor	door_w	indows	wall	stairs	vault	roof	other	none
arch	41.31	0.42*	1.364	0	0.149	0.276	0	3.542	0.006499	1.75	0.118	
column	0.132	96.5*	0.4246	0.5346	0.2048	0.1103	0.2197	0.1368	0.0006537	0.308	10.14	
moldings	1.338	0.5781	90.42	0.3608	2.38	3.4	0.08897	0.1322	0.522	0.09379	0.6857	
floor	0	0.4891	0.2747	95.86	0.7856	0.1671	2.295	0	0	0.03661	0.59	
door_w	0.2288	0.1634	2.055	0.1144	94.2	2.505	0.05631	0.1959	0.003436	0.2421	0.2387	
indows	0.6204	0.4689	3.888	0.09861	1.573	91.53	0.04447	0.5073	0.06285	0.7022	0.509*	
wall	0	0.4369	0.2334	2.159	0.0485	0.04336	96.59	0	0	0.274	0.3655	
stairs	0.194	0.199	0.0991	1.734e-05	0.05698	0.6074	0	90.5	0	2.455	0.04512	
vault	0.007336	0.3599	0.9273	0	0.01866	0.04907	0	0.1003	96.96	2.124	0.2516	
roof	2.706	0.246	0.2607	0.05089	0.2011	0.6057	0.3593	1.852	0.9643	92.04	0.735*	
other	0.3294	1.219	1.132	0.6292	0.2523	1.076	0.3377	0.336	0.2622	0.9465	91.46	
none												

Figure 2: Confusion Matrix for Test 1 - Spedale del Ceppo

This test was repeated for all the other buildings of the dataset, obtaining the results reported in Table 3.

	Global Accuracy	Mean Accuracy
2_OSA	75%	66%
3_SS	91%	76%
4_CG	80%	55%
5_CB	74%	49%

Table 3: DeepLabv3+ results on the test set for all the other buildings (Test 1).

### 4.1.2 Test 2

Test 2 considers all the five buildings currently available in the dataset, picking 500 images for each building. The images were randomly shuffled and divided in three sets: 60% as training set (1400), 20% as validation test (466), and 20% as test set (466). As expected, in this case the performance of the network is lower than in Test 1, yielding to an overall accuracy of 88.1% and a IoU of 71.8% (Table 4). Test and training set are still quite similar, but it is remarkable that the model could well generalize the solution even when dealing with several buildings.

Global Accuracy	Mean Accuracy	Mean IoU	Weighted IoU	Mean BFScore
88,1%	88,6%	71,8%	80,9%	72,1%

Table 4: DeepLabv3+ results on the test set (Test 2).

Figure 3 shows a comparison between the input images, the ground-truth and the prediction obtained on three test set images. Figure 4 shows the confusion matrix obtained in Test 2.



Figure 3: DeepLabv3+ results on the test set (Test 2): A) input image, B) ground-truth, C) prediction.

True Class \ Predicted Class	arch	column	moldings	floor	door_w/indows	wall	stairs	vault	roof	other	none
arch	85.31	10.22	1.345	0.000289	0.1117	1.631	0	7.363	0.09019	2.795	0.3314
column	0.7896	89.71	1.772	1.001	0.3385	2.727	0.6842	0.09497	0	0.9443	3.138
moldings	0.9401	1.1	85.99	0.5398	5.497	5.871	0.05365	0.1775	0.4576	0.9011	1.179
floor	0.0002074	0.6445	0.5778	93.94	0.478	0.515	0.814	2.704e-06	0	1.617	1.465
door_w/indows	0.3059	0.2974	3.887	0.4181	90.57	3.641	0.07472	0.1533	0.0466	0.2123	0.3959
wall	1.693	1.259	6.132	0.3901	1.714	85.7	0.07087	0.7205	0.2934	1.242	0.8087
stairs	0	4.869	2.317	3.624	10.56	1.71	80.99	0	0	0.8531	4.652
vault	7.777	0.7906	0.1729	0.000198	0.03275	0.8544	0	88.54	0.03302	2.684	0.1845
roof	0.1885	0.03407	2.01	0.002075	0.7518	0.9987	0	0.7117	95.11	0.7793	0.5374
other	3.848	1.113	1.308	0.9609	0.2034	1.205	0.08648	3.36	0.2295	86.22	1.464
none	0.7825	2.334	2.01	0.6739	0.1805	1.9	0.2129	0.9544	0.265	1.408	89.28

Figure 4: Confusion Matrix for Test 2.

### 4.1.3 Test 3

Finally, in Test 3 the network is tested on the prediction of an unseen scene. Four buildings are used to train the network, picking 500 images for each building, whereas the remaining building is used as test set (1\_SC). Being more challenging, the performance in this test is clearly worse than in Test 1 and 2: the model reached a global accuracy of 54.2% and a IoU of 27.7%.

Global Accuracy	Mean Accuracy	Mean IoU	Weighted IoU	Mean BFScore
54,2%	41,1%	27,7%	37,3%	34,0%

Table 5: DeepLabv3+ results on the test set (Test 3).

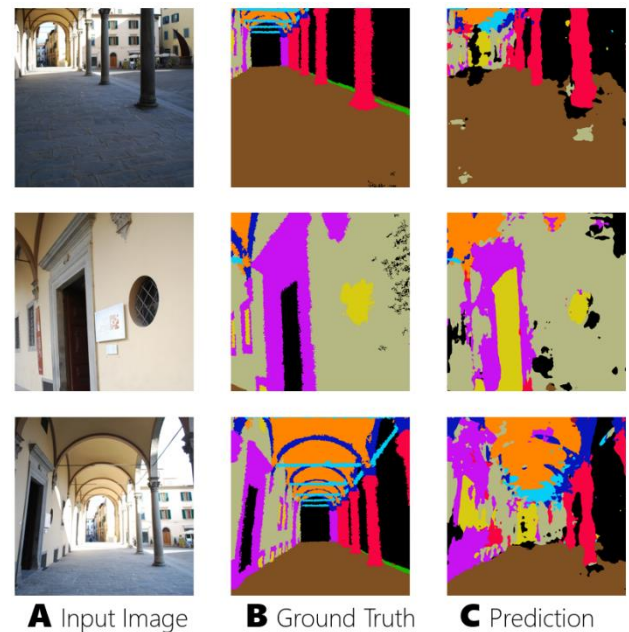


Figure 5: DeepLabv3+ results on the test set (Test 3): A) input image, B) ground-truth, C) prediction.

True Class \ Predicted Class	arch	column	moldings	floor	door_w/indows	wall	stairs	vault	roof	other	none
arch	37.67	5.585	3.445	0.07957	0.4583	11.01	4.644e-05	26.76	0.179	3.641	9.336
column	2.341	73.22	2.85	2.947	6.917	7.782	0.06343	0.5362	0.1142	0.7267	8.731
moldings	3.347	5.984	26.7	0.8164	5.234	70.45	0.04944	1.089	0.9771	1.295	74
floor	0.009957	0.3165	0.6337	85.51	0.7731	7.84	0.2065	0.002156	0.0005669	0.08146	10.16
door_w/indows	1.342	1.673	7.722	0.8306	26.4	18.38	0.07406	0.7151	0.3593	0.9394	41.58
wall	1.895	2.877	4.845	0.5135	10.31	52.53	0.05157	3.975	0.1676	1.156	39.99
stairs	0.009097	3.271	3.145	54.63	0.2827	4.54	0.6669	0.005298	0.002426	0.92	32.52
vault	9.136	0.734	1.697	0.03375	0.07574	11.73	1.088e-05	85.03	0.1289	4.224	7.219
roof	1.863	0.8921	5.116	0.03039	10.93	1.108	0.0004729	2.59	7.873	2.082	67.52
other	1.41	2.579	1.875	0.3603	1.24	10.9	0.01542	28.58	0.7682	17.96	24.31
none	0.9249	2.237	2.386	16.2	3.179	7.054	0.08035	0.6603	0.0711	0.8458	66.16

Figure 6: Confusion Matrix for Test 3.

The confusion matrix, in Figure 6, shows that the most common classes like "floor", "vault", "wall" and "column" are well segmented. This proves that the strong class imbalance, which characterize the ongoing dataset, is a critical factor for the network performance.

#### 4.2 Back-projection procedure

Finally, this subsection shows the results of the back-projection procedure on two dataset buildings, (*1\_SC*) Spedale del Ceppo and (*2\_OSA*) Ospedale Sant'Antonio. Two back-projection cases are considered for comparison.

First, a point cloud is built using a set of masks created directly from the ground truth images used for the training phase of the network. Secondly, a point cloud is built using the set of masks created from the prediction output by the network. In both the cases the image prediction is obtained by means of the neural network trained in Test 1.

The comparison between such two reconstructions allows to evaluate the correct functioning of the procedure and the influence of the prediction accuracy on the final result.

Figure 8 shows A) the ground-truth point cloud, B) the point cloud reconstruction from the image ground-truth, and C) the point cloud reconstruction from the image prediction output by the neural network for the Spedale del Ceppo building.

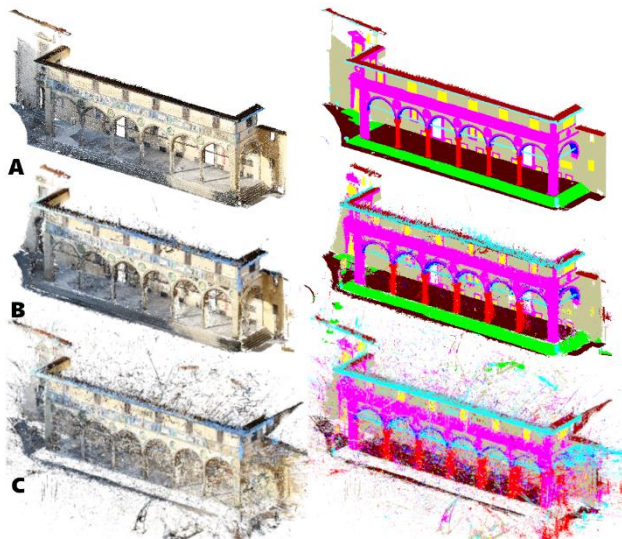


Figure 8: Back-Projection results for *1\_SC*: A) ground-truth point cloud, B) image ground-truth reconstruction, C) prediction reconstruction.

Table 6 reports the Global Accuracy, the Mean Intersection Over Union and the DICE coefficient for the two reconstructions.

TYPE OF RECONSTRUCTION	Global Accuracy	Mean IoU	Mean DICE
Ground Truth Images	82,6%	56,7%	68,9%
Prediction Images	33,4%	25,6%	29,7%

Table 6: Evaluation metrics for the back-projection of *1\_SC*.

Figure 8 shows that the projection procedure works quite well with the ground-truth images, but its performance degrades when working with the prediction images.

Figure 9 shows the results for the Ospedale Sant'Antonio building.



Figure 9: Back-Projection results for *2\_OSA*: A) ground-truth point cloud, B) image ground-truth reconstruction, C) prediction reconstruction.

Table 7 reports some of the most remarkable evaluation metrics.

TYPE OF RECONSTRUCTION	Global Accuracy	Mean IoU	Mean DICE
Ground Truth Images	69,7%	49,7%	61,2%
Prediction Images	57,3%	37,7%	52,7%

Table 7: Evaluation metrics for the back-projection of *2\_OSA*.

In this second case, the projection procedure reached a lower overall accuracy, but the performance on the cloud built with the prediction images has less degradation compared with the previous case. As expected, the results strictly depend on the quality of image segmentation prediction and hence on the accuracy of the neural network.

Table 8 compares the image segmentation with the 3D point cloud segmentation results.

	Prediction	Global Acc	Mean IoU
<i>1_SC</i>	Images	93,4%	82,8%
	Point Cloud	33,4%	25,6%
<i>2_OSA</i>	Images	75,0%	66,0%
	Point Cloud	57,3%	37,7%

Table 8: Comparison between the neural network-based image segmentation results and the corresponding results on the 3D point cloud.

Despite the high value of the accuracy in the image prediction for the first case (*1\_SC*) the Global Accuracy on the back-projection is quite low, with a percentage loss of 60%. On the second case the accuracy on the images is lower, but the percentage loss on the back-projection is less than 20%. Considering only these two case studies is not possible to correctly evaluate and generalize the performance of the procedure: future tests will be performed to this aim.

However, these issues are related at least to two critical aspects: (i) the lack of accuracy of the prediction, in particular for what concerns object edges, and (ii) the problem of label overlapping. Issue (i) is related also to the dataset characteristics: a larger number of buildings and images should be considered, and the class imbalance should be reduced. In addition to increasing the dataset size increasing the number of the buildings, using data augmentation and class weighting during training should be useful as well.

Instead, the label overlapping issue (ii) should be at least partially compensated by introducing a regularization step in the back-projection procedure.

Since the results obtained using the ground truth segmented images are promising (low percentage loss: 10% for Spedale del Ceppo and 6% for Spedale Sant'Antonio), an improvement in the automatic image segmentation shall lead to a significant increase of the overall point cloud classification performance.

## 5. CONCLUSIONS

This paper presented a procedure for the semantic segmentation of 3D point clouds of heritage buildings. The procedure is composed by two main steps: (i) the labelling of the intermediate images, using a pre-trained (deep learning-based) neural network, and then (ii) the projection of the image features on the 3D point cloud, thanks to a masking method. The results are quite promising, but more in depth tests shall be performed to assess the general performance of the entire workflow. Furthermore, the foreseen generalization and size increase of the training dataset shall improve the performance of the automatic image segmentation step. Finally, an improvement of back-projection procedure, in particular to better deal with the label overlapping issue, will be investigated as well. Such changes are expected to improve the point cloud segmentation results.

## REFERENCES

- Badrinarayanan, V., Kendall, A., Cipolla, R., 2015. Segnet: a deep convolutional encoder-decoder architecture for scene segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39 (12), 2481–2495. [10.1109/TPAMI.2016.2644615](https://doi.org/10.1109/TPAMI.2016.2644615).
- Boulch, A.; Guerry, J.; Le Saux, B.; Audebert, N., 2018. SnapNet: 3D point cloud semantic labeling with 2D deep segmentation networks. *Computers and Graphics (Pergamon)*, 71, 189–198.
- Chen, L. C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A. L., 2017. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4), pp. 834–848.
- Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B., 2016. The cityscapes dataset for semantic urban scene understanding. *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 3213–3223.
- Grilli, E., Remondino, F., 2019. Classification of 3D digital heritage. *Remote Sensing*, 11(7).
- Grilli, E., Dininno, D., Petrucci, G., Remondino, F., 2018. From 2D to 3D supervised segmentation and classification for cultural heritage applications. *International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences - ISPRS Archives*, 42(2), 399–406.
- López, F.J., Lerones, P.M., Llamas, J., Gómez-García-Bermejo, J., Zalama, E., 2018. A Review of Heritage Building Information Modeling (H-BIM). *Multimodal Technologies and Interaction*, 2, 21. <https://doi.org/10.3390/mti2020021>.
- Malinverni, E. S., Pierdicca, R., Paolanti, M., Martini, M., Morbidoni, C., Matrone, F., Lingua, A., 2019. Deep Learning for Semantic Segmentation of 3D Point Cloud. *International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences - ISPRS Archives*, 42(2/W15), 735–742.
- Matrone, F., Grilli, E., Martini, M., Paolanti, M., Pierdicca, R., Remondino, F., 2020. Comparing machine and deep learning methods for large 3D heritage semantic segmentation. *ISPRS International Journal of Geo-Information*, 9(9). <https://doi.org/10.3390/ijgi9090535>.
- Murtiyoso, A., Lhenry, C., Landes, T., Grussenmeyer, P., Alby, E., 2021. Semantic segmentation for building façade 3D point cloud from 2D orthophoto images using transfer learning. *International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences - ISPRS Archives*, 43(B2-2021), 201–206.
- Murtiyoso, A., Pellis, E., Grussenmeyer, P., Landes, T., Masiero, A., 2022. Towards semantic photogrammetry: generating semantically rich point clouds from close range photogrammetry. *Sensors*, 22.
- Pellis, E., Masiero, A., Tucci, G., Betti, M., Grussenmeyer, P., 2021. Assembling an Image and Point Cloud Dataset for Heritage Buildings Semantic Segmentation. *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences, XLVI-M-1–2021*, 539–546.
- Qi, C. R., Su, H., Mo, K., Guibas, L. J., 2017a. Pointnet: Deep learning on point sets for 3D classification and segmentation. *In Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 652–660).
- Qi, C. R., Yi, L., Su, H., Guibas, L. J., 2017b. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413*.
- Shan, J., Toth, C. K. (Eds.), 2018. *Topographic laser ranging and scanning: principles and processing*. CRC press.
- Su, H.; Maji, S.; Kalogerakis, E.; Learned-Miller, E., 2015. Multi-view convolutional neural networks for 3D shape recognition. *Proceedings of the IEEE International Conference on Computer Vision*, Vol. 2015 Inter, pp. 945–953.
- Yastikli, N., 2007. Documentation of cultural heritage using digital photogrammetry and laser scanning. *Journal of Cultural heritage*, 8(4), 423–427.
- Zhang, J., Zhao, X., Chen, Z., Lu, Z., 2019. A Review of Deep Learning-Based Semantic Segmentation for Point Cloud. *IEEE Access*, vol. 7, pp. 179118–179133. doi: [10.1109/ACCESS.2019.2958671](https://doi.org/10.1109/ACCESS.2019.2958671).