

UNIVERSITY OF ALICANTE

PHD THESIS

**TBD**

*Author*

Alberto GARCIA-GARCIA

*Advisors*

Jose GARCIA-RODRIGUEZ  
Sergio ORTS-ESCOLANO

*A thesis submitted in fulfilment of the requirements  
for the degree of Doctor of Philosophy*

*in the*

3D Perception Lab  
Department of Computer Technology

March 25, 2019



This document was proudly made with L<sup>A</sup>T<sub>E</sub>X and TikZ.

This work is licensed under a Creative Commons  
“Attribution-ShareAlike 4.0 International” license.





*“Will robots inherit the earth? Yes, but they will be our children.”*

Marvin Minsky



# Abstract

Aliquam lectus. Vivamus leo. Quisque ornare tellus ullamcorper nulla. Mauris porttitor pharetra tortor. Sed fringilla justo sed mauris. Mauris tellus. Sed non leo. Nullam elementum, magna in cursus sodales, augue est scelerisque sapien, venenatis congue nulla arcu et pede. Ut suscipit enim vel sapien. Donec congue. Maecenas urna mi, suscipit in, placerat ut, vestibulum ut, massa. Fusce ultrices nulla et nisl.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla ac nisl. Nullam urna nulla, ullamcorper in, interdum sit amet, gravida ut, risus. Aenean ac enim. In luctus. Phasellus eu quam vitae turpis viverra pellentesque. Duis feugiat felis ut enim. Phasellus pharetra, sem id porttitor sodales, magna nunc aliquet nibh, nec blandit nisl mauris at pede. Suspendisse risus risus, lobortis eget, semper at, imperdiet sit amet, quam. Quisque scelerisque dapibus nibh. Nam enim. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Nunc ut metus. Ut metus justo, auctor at, ultrices eu, sagittis ut, purus. Aliquam aliquam.



# Resumen

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

Sed feugiat. Cum sociis natoque penatibus et magnis dis parturient montes, nascentur ridiculus mus. Ut pellentesque augue sed urna. Vestibulum diam eros, fringilla et, consectetur eu, nonummy id, sapien. Nullam at lectus. In sagittis ultrices mauris. Curabitur malesuada erat sit amet massa. Fusce blandit. Aliquam erat volutpat. Aliquam euismod. Aenean vel lectus. Nunc imperdiet justo nec dolor.

Aliquam lectus. Vivamus leo. Quisque ornare tellus ullamcorper nulla. Mauris porttitor pharetra tortor. Sed fringilla justo sed mauris. Mauris tellus. Sed non leo. Nullam elementum, magna in cursus sodales, augue est scelerisque sapien, venenatis congue nulla arcu et pede. Ut suscipit enim vel sapien. Donec congue. Maecenas urna mi, suscipit in, placerat ut, vestibulum ut, massa. Fusce ultrices nulla et nisl.



# Acknowledgements

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.



# Contents

<b>Abstract</b>	vii
<b>Resumen</b>	ix
<b>Acknowledgements</b>	xi
<b>Contents</b>	xiii
<b>List of Figures</b>	xv
<b>List of Tables</b>	xix
<b>List of Acronyms</b>	xxi
<b>1 Introduction</b>	1
1.1 Motivation . . . . .	2
1.2 Approach . . . . .	2
1.3 Contributions . . . . .	2
1.4 Co-Authored Papers . . . . .	2
1.4.1 Chapter 2 . . . . .	3
1.4.2 Chapter 3 . . . . .	3
1.4.3 Chapter 4 . . . . .	3
1.4.4 Other . . . . .	3
1.5 Thesis Structure . . . . .	5
<b>2 Object Classification</b>	7
2.1 Introduction . . . . .	8
2.2 Related Works . . . . .	9
2.2.1 Traditional Approaches . . . . .	9
2.2.2 3D Object Recognition . . . . .	10
2.2.3 Deep Learning . . . . .	11
2.2.4 Volumetric Representations . . . . .	13
2.2.5 Our Proposal in Context . . . . .	16
2.3 Datasets . . . . .	16
2.4 PointNet . . . . .	20
2.4.1 Data Representation . . . . .	20
2.4.2 Network Architecture . . . . .	22
2.4.3 Experiments . . . . .	23

Data Generation . . . . .	23
Implementation and Setup . . . . .	23
Results and Discussion . . . . .	24
2.4.4 Conclusion . . . . .	26
2.5 Noise and Occlusion . . . . .	26
2.5.1 Data Representation . . . . .	26
Tensor Generation . . . . .	26
Occupancy Computation . . . . .	28
2.5.2 Network Architecture . . . . .	29
2.5.3 Experiments . . . . .	31
Data Generation . . . . .	31
Implementation and Setup . . . . .	32
Results and Discussion . . . . .	33
2.5.4 Conclusion . . . . .	38
2.6 LonchaNet . . . . .	38
2.6.1 Data Representation . . . . .	38
2.6.2 Network Architecture . . . . .	39
2.6.3 Experiments . . . . .	40
Data Generation . . . . .	40
Methodology and Setup . . . . .	40
Results and Discussion . . . . .	41
2.6.4 Conclusion . . . . .	42
2.7 Conclusion . . . . .	42
<b>3 Semantic Segmentation</b>	<b>43</b>
3.1 Introduction . . . . .	43
3.2 Related Works . . . . .	43
3.3 The RobotriX . . . . .	43
3.4 UnrealROX . . . . .	43
3.5 2D-3D-SeGCN . . . . .	43
<b>4 Tactile Sensing</b>	<b>45</b>
4.1 Introduction . . . . .	45
4.2 Related Works . . . . .	45
4.3 TactileGCN . . . . .	45
4.4 Conclusion . . . . .	45
<b>5 Conclusion</b>	<b>47</b>
5.1 Findings and Conclusions . . . . .	47
5.2 Limitations . . . . .	47
5.3 Future Work . . . . .	47
<b>Bibliography</b>	<b>49</b>

# List of Figures

2.1	Evolution of the number of academic documents containing the terms 2D, 3D, and Deep Learning together with <i>Computer Vision</i> . Search terms statistics obtained from scopus.com. . . . .	10
2.2	Filters learned by the network proposed by Krizhevsky et al. [35]. Each of the 96 filters shown is of size $11 \times 11 \times 3$ . The filters have clearly learned to detect edges of various orientations and they resemble Gabor filters. Image analysis using that kind of filters is thought to be similar to perception in the human visual system [36]. . . . .	11
2.3	Illustration of the architecture of the aforementioned Convolutional Neural Network (CNN) proposed by Krizhevsky et al.[35] for the ImageNet challenge. Besides the normal components, e.g., convolutional, pooling, and fully connected layers, this network features two different paths. One Graphics Processing Unit (GPU) runs the layers at the top while the other runs the layers at the bottom. . . . .	12
2.4	Common volumetric representations: polygonal mesh (a), point cloud (b), and voxel grid (c) of a chair model (which is color coded by height and depth). . . . .	13
2.5	Effect of the leaf size on binary voxel grids. All grids have the same cubic size: $300 \times 300 \times 300$ units. Leaf sizes vary from 5, 10, and 20 units, resulting in binary grids of $15 \times 15 \times 15$ (a), $30 \times 30 \times 30$ (b), and $60 \times 60 \times 60$ voxels (c) respectively. . . . .	14
2.6	<i>3DShapeNets</i> representation proposed by Wu <i>et al.</i> as shown in their paper [42]. An object (a) is captured from a certain point of view and a depth map is generated (b) which is in turn used to generate a point cloud that will be represented as a voxel grid (c) with empty voxels (in white, not represented), unknown voxels (in blue), and surface or occupied voxels (red). . . . .	14
2.7	Truncated Signed Distance Function (TSDF) representation proposed by Song and Xiao as shown in their paper [43]. An object (a) is captured by a range sensor as a point cloud (b) and then a TSDF grid is generated (red indicates the voxel is in front of surfaces and blue indicates the voxel is behind the surface; the intensity of the color represents the TSDF value). . . . .	15

2.8	Volumetric occupancy grid representation used by <i>VoxNet</i> as shown in their paper [44]. For LIDAR data (a) a voxel size of $0.1\text{m}^3$ is used to create a $32 \times 32 \times 32$ grid (b). For RGB-D data (??), the resolution is chosen so the object occupies a subvolume of $24 \times 24 \times 24$ voxels in a $32 \times 32 \times 32$ grid (d). . . . .	15
2.9	ModelNet10 samples. . . . .	17
2.10	Model distribution per object class or category for both ModelNet-10 and ModelNet-40 training and test splits. . . . .	18
2.11	From CAD models to point clouds. The object is placed in the center of a tessellated sphere, views are rendered placing a virtual camera in each vertex of the icosahedron, the $z$ -buffer data of those views is used to generate point clouds, and the point clouds are transformed and merged at last. . . . .	20
2.12	Various 3D representations for an object. A mesh (a) is transformed into a point cloud (b), and that cloud is processed to obtain a voxelized occupancy grid (c). The occupancy grid shown in this figure is a cube of $30 \times 30 \times 30$ voxels. Each voxel of that cube holds the point density inside its volume. In this case, dark voxels indicate high density whilst bright ones are low density volumes. . . . .	21
2.13	PointNet's 3D CNN architecture. [MISSINGDETAILS] . . . . .	22
2.14	Dataset model processing example to generate the point clouds for PointNet. Some rendered views of a toilet model are shown in (a). The original Object File Format (OFF) mesh is shown in (b). The generated point cloud after merging all points of view is shown in (c), and (d) shows the downsampled cloud using a voxel grid filter with a leaf size of $0.7 \times 0.7 \times 0.7$ . . . . .	23
2.15	Similarity between two objects of different classes: Table and Desk. The point cloud shown in (a) represents an object of the Table class, whilst the point cloud in (b) represents an object whose class is Desk but it is misclassified as a Table due to their resemblance. . . . .	24
2.16	Neuron activations for the output layer of the architecture when classifying all the test samples for both <i>Desk</i> (b) and <i>Table</i> (a) classes. Each row represents an activation vector for a specific sample, so each column is a position of the vector: the activation to that particular class. The first column corresponds to the <i>Desk</i> class, while the second one is the <i>Table</i> . The activation shows the clear confusion between <i>Desk</i> and <i>Table</i> . Although the latter one is much less confused with other classes, many <i>Tables</i> are misclassified as <i>Desks</i> thus lowering the accuracy for this class. . . . .	25
2.17	Comparison of accuracy per class using an unbalanced dataset and a balanced one with a maximum of 400 models per class via random undersampling. Accuracy is harmed in the classes in which models are removed but gained otherwise. . . . .	25
2.18	A fixed occupancy grid ( $8 \times 8 \times 8$ voxels) with 40 units leaf size and 320 units grid size in all dimensions. The grid origin is placed at the minimum $x$ , $y$ , and $z$ values of the point cloud. Front (a), side (b), and perspective (c) views of the grid over a partial view of a segmented table object are shown. . . . .	27

2.19	An adaptive occupancy grid ( $8 \times 8 \times 8$ voxels) with adapted leaf and grid sizes in all dimensions to fit the data. The grid origin is placed at the minimum $x$ , $y$ , and $z$ values of the point cloud. Front (a), side (b), and perspective (c) views of the grid over a partial view of a segmented table object are shown. Notice that the point clouds for the three views are exactly the same for this figure and Figure 2.18, but the grids do change. There is a noticeable difference in the front view. In Figure 2.18, using fixed grids, all voxels are cubic and the point cloud does not fit the grid completely (leftmost column in Figure 2.18a), whilst in this figure, with adaptive grids, the grid is fitted to the cloud. . . . .	27
2.20	Occupied voxels in an adaptive $8 \times 8 \times 8$ grid generated over a partial view point cloud. Those voxels with points inside are shown in a wire-frame representation. Empty voxels are omitted. Occupied voxels must be filled with values which represent the contained shape. . . . .	28
2.21	Binary tensor computed over a point cloud of a partial view of an object (shown in Figure 2.20). Occupied voxels are shown in blue, empty voxels are omitted for the sake of simplicity. . . . .	29
2.22	Normalized density tensor over a point cloud of a partial view of an object (shown in Figure 2.20). Denser voxels are darker and sparse ones are shown in light blue. Empty voxels were removed for visualization purposes. . . . .	30
2.23	CVIU's architecture. [MISSINGDETAILS] . . . . .	30
2.24	Different levels of noise ( $\sigma = 0$ (a), $\sigma = 0.1$ (b), and $\sigma = 1$ (c)) applied to the $z$ -axis of every point of a table partial view. . . . .	31
2.25	Different levels of occlusion ( $\psi = 0\%$ (a), $\psi = 25\%$ (b), and $\psi = 50\%$ (c)) applied randomly to a table partial view. . . . .	32
2.26	Evolution of training and validation accuracy of the model-based CNN using both fixed (a) and adaptive (b) normalized density grids. Different grid sizes (32, 48, and 64) were tested. . . . .	33
2.27	Evolution of validation accuracy of the model-based CNN using both fixed (a) and adaptive (b) normalized density grids as the amount of occlusion in the validation models increases from 0% to 30%. Three grid sizes were tested (32, 48, and 64). . . . .	34
2.28	Evolution of validation accuracy of the model-based CNN using both fixed (a) and adaptive (b) normalized density grids as the standard deviation of the Gaussian noise introduced in the $z$ -axis of the views increases from 0.001 to 10. The common grid sizes were tested (32, 48, and 64). . . . .	35
2.29	Evolution of training and validation accuracy of the model-based CNN using adaptive binary grids (a). Evolution of validation accuracy for the best network weights after training as the amount of occlusion in the validation set increases (b) and different levels of noise are introduced (c). . . . .	36
2.30	Evolution of training and validation accuracy of the 3D CNN using adaptive binary grids with size $32 \times 32 \times 32$ . . . . .	37
2.31	Extracting the slices from a point cloud example [MISSINGDETAILS] . . . . .	39
2.32	A comparison of a slice before and after the dilation process. The dilated image provides a more accurate representation the object. . . . .	40
2.33	Success rate per class for the test split of the ModelNet-10 dataset achieved by LonchaNet after 18300 training iterations using the <i>ModelNet-10</i> dataset (solver type is ADAM, learning rate is 0.00001, $\beta_1$ is 0.9 and $\beta_2$ is 0.999). . . . .	41



# List of Tables

2.1	ModelNet-10 and ModelNet-40 distributions. . . . .	19
2.2	ModelNet leaderboard as of January, 2017. . . . .	42



# List of Acronyms

**2D** two-dimensional

**2.5D** two-and-a-half-dimensional

**3D** three-dimensional

**AMT** Amazon Mechanical Turk

**BRIEF** Binary Robust Independent Elementary Features

**BRISK** Binary Robust Invariant Scalable Keypoints

**BVLC** Berkeley Vision and Learning Center

**CAD** Computer Aided Design

**CDBN** Convolutional Deep Belief Network

**CIFAR** Canadian Institute for Advanced Research

**CNN** Convolutional Neural Network

**FREAK** Fast Retina Keypoint

**GPU** Graphics Processing Unit

**IR** Infrared

**LIDAR** Light Detection and Ranging

**OFF** Object File Format

**ORB** Oriented FAST and Rotated BRIEF

**PCD** Point Cloud Data

**PCL** Point Cloud Library

**POV** Point of View

**ReLU** Rectified Linear Unit

**RGB** Red Green and Blue

**RGB-D** RGB-Depth

**SIFT** Scale Invariant Feature Transform

**SURF** Speeded Up Robust Features

**TSDF** Truncated Signed Distance Function



Chapter **1**

# Introduction

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetur odio sem sed wisi.

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetur odio sem sed wisi.

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetur odio sem sed wisi.

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetur odio sem sed wisi.

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetur odio sem sed wisi.

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetur odio sem sed wisi.

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetur odio sem sed wisi.

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetur odio sem sed wisi.

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetur odio sem sed wisi.

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetur odio sem sed wisi.

## **1.1 Motivation**

## **1.2 Approach**

## **1.3 Contributions**

## **1.4 Co-Authored Papers**

This thesis is the result of continuous effort throughout the last years. Such efforts have sometimes crystallized in form of co-authored publications and conference talks.

### 1.4.1 Chapter 2

- Alberto Garcia-Garcia, Francisco Gomez-Donoso, Jose Garcia-Rodriguez, et al. "PointNet: A 3D Convolutional Neural Network for real-time object class recognition". In: *2016 International Joint Conference on Neural Networks, IJCNN 2016, Vancouver, BC, Canada, July 24-29, 2016*. 2016, pp. 1578–1584. DOI: [10.1109/IJCNN.2016.7727386](https://doi.org/10.1109/IJCNN.2016.7727386). URL: <https://doi.org/10.1109/IJCNN.2016.7727386>
- Alberto Garcia-Garcia, Jose Garcia-Rodriguez, Sergio Orts-Escalano, et al. "A study of the effect of noise and occlusion on the accuracy of convolutional neural networks applied to 3D object recognition". In: *Computer Vision and Image Understanding* 164 (2017), pp. 124–134. DOI: [10.1016/j.cviu.2017.06.006](https://doi.org/10.1016/j.cviu.2017.06.006). URL: <https://doi.org/10.1016/j.cviu.2017.06.006>
- Francisco Gomez-Donoso, Alberto Garcia-Garcia, Jose Garcia-Rodriguez, et al. "LonchaNet: A Sliced-based CNN Architecture for Real-time 3D Object Recognition". In: *2017 International Joint Conference on Neural Networks, IJCNN 2017, Anchorage, Alaska, May 14-19, 2017*. 2017. URL: <https://ieeexplore.ieee.org/document/7965883/>

### 1.4.2 Chapter 3

- Alberto Garcia-Garcia, Jose Garcia-Rodriguez, Sergio Orts-Escalano, et al. "A study of the effect of noise and occlusion on the accuracy of convolutional neural networks applied to 3D object recognition". In: *Computer Vision and Image Understanding* 164 (2017), pp. 124–134. DOI: [10.1016/j.cviu.2017.06.006](https://doi.org/10.1016/j.cviu.2017.06.006). URL: <https://doi.org/10.1016/j.cviu.2017.06.006>
- Alberto Garcia-Garcia, Pablo Martinez-Gonzalez, Sergiu Oprea, et al. "The RobotriX: An eXtremely Photorealistic and Very-Large-Scale Indoor Dataset of Sequences with Robot Trajectories and Interactions". In: *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE. 2018, pp. 6790–6797. URL: <https://ieeexplore.ieee.org/abstract/document/8594495>
- TODO: UnrealROX

### 1.4.3 Chapter 4

- TODO: TactileGCN

### 1.4.4 Other

During the years spent working on the main topics of this thesis, several collaborations and side works were carried out that also were published either as journal papers, conference proceedings, or preprints:

- Sergiu Oprea, Alberto Garcia-Garcia, Jose Garcia-Rodriguez, et al. "A Recurrent Neural Network based Schaeffer Gesture Recognition System". In: *2017 International Joint Conference on Neural Networks, IJCNN 2017, Anchorage, Alaska, May 14-19, 2017*. 2017. URL: <https://ieeexplore.ieee.org/document/7965885/>

- Francisco Gomez-Donoso, Sergio Orts-Escalano, Alberto Garcia-Garcia, et al. "A robotic platform for customized and interactive rehabilitation of persons with disabilities". In: *Pattern Recognition Letters* 99 (2017), pp. 105–113. DOI: [10.1016/j.patrec.2017.05.027](https://doi.org/10.1016/j.patrec.2017.05.027). URL: <https://doi.org/10.1016/j.patrec.2017.05.027>
- Sergiu Oprea, Alberto GarciaGarcia, Sergio OrtsEscalano, et al. "A long short-term memory based Schaeffer gesture recognition system". In: *Expert Systems* 0.0 (2017), e12247. DOI: [10.1111/exsy.12247](https://doi.org/10.1111/exsy.12247). URL: <https://onlinelibrary.wiley.com/doi/abs/10.1111/exsy.12247>
- Alberto Garcia Garcia, Andreas Beckmann, and Ivo Kabadshow. "Accelerating an FMM-Based Coulomb Solver with GPUs". In: *Software for Exascale Computing-SPPEXA 2013-2015*. Springer, 2016, pp. 485–504. URL: [https://link.springer.com/chapter/10.1007/978-3-319-40528-5\\_22](https://link.springer.com/chapter/10.1007/978-3-319-40528-5_22)
- Alberto Garcia-Garcia, Sergio Orts-Escalano, Sergiu Oprea, et al. "Multi-sensor 3D object dataset for object recognition with full pose estimation". In: *Neural Computing and Applications* 28 (2016), pp. 941–952. ISSN: 1433-3058. DOI: [10.1007/s00521-016-2224-9](https://doi.org/10.1007/s00521-016-2224-9). URL: <http://dx.doi.org/10.1007/s00521-016-2224-9>
- Marcelo Saval-Calvo, Jorge Azorin-Lopez, Andres Fuster-Guillo, et al. "Evaluation of sampling method effects in 3D non-rigid registration". In: *Neural Computing and Applications* 28 (2016), pp. 953–967. ISSN: 1433-3058. DOI: [10.1007/s00521-016-2258-z](https://doi.org/10.1007/s00521-016-2258-z). URL: <http://dx.doi.org/10.1007/s00521-016-2258-z>
- Sergio Orts-Escalano, Jose Garcia-Rodriguez, Miguel Cazorla, et al. "Bioinspired point cloud representation: 3D object tracking". In: *Neural Computing and Applications* 29 (2016), pp. 663–672. ISSN: 1433-3058. DOI: [10.1007/s00521-016-2585-0](https://doi.org/10.1007/s00521-016-2585-0). URL: <https://doi.org/10.1007/s00521-016-2585-0>
- Alberto Garcia-Garcia, Sergio Orts-Escalano, Jose Garcia-Rodriguez, et al. "Interactive 3D object recognition pipeline on mobile GPGPU computing platforms using low-cost RGB-D sensors". In: *Journal of Real-Time Image Processing* 14 (2016), pp. 585–604. ISSN: 1861-8219. DOI: [10.1007/s11554-016-0607-x](https://doi.org/10.1007/s11554-016-0607-x). URL: <https://doi.org/10.1007/s11554-016-0607-x>
- Higinio Mora, Jerónimo M Mora-Pascual, Alberto Garcia-Garcia, et al. "Computational analysis of distance operators for the iterative closest point algorithm". In: *PloS one* 11.10 (2016), e0164694. URL: <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0164694>
- Sergio Orts-Escalano, Jose Garcia-Rodriguez, Vicente Morell, et al. "3D Surface Reconstruction of Noisy Point Clouds Using Growing Neural Gas: 3D Object/Scene Reconstruction". In: *Neural Processing Letters* 43 (2015), pp. 401–423. DOI: [10.1007/s11063-015-9421-x](https://doi.org/10.1007/s11063-015-9421-x). URL: <http://dx.doi.org/10.1007/s11063-015-9421-x>
- Sergio Orts-Escalano, Jose Garcia-Rodriguez, Jose Antonio Serra-Perez, et al. "3D model reconstruction using neural gas accelerated on GPU". in: *Applied Soft Computing* 32 (2014), pp. 87–100. DOI: [10.1016/j.asoc.2015.03.042](https://doi.org/10.1016/j.asoc.2015.03.042). URL: <http://dx.doi.org/10.1016/j.asoc.2015.03.042>

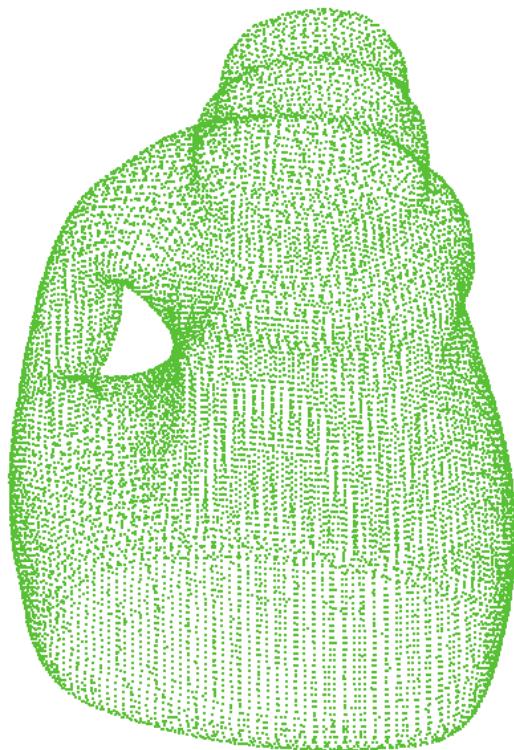
- TODO: ICP

## 1.5 Thesis Structure



# Chapter 2

## Object Classification



### *Abstract*

In this chapter, we address the problem of object classification. To approach this challenge, we rely on the geometric information provided by 3D object representations such as point clouds. Furthermore, we focus on learning-based methods to distinguish objects from different classes while capturing the variability of shape of different objects which belong to the same class. More specifically, we leverage deep learning for such task.

The chapter begins introducing and formulating the object classification task in Section 2.1 followed by a review of the most relevant literature and datasets in Sections 2.2 and 2.3. After that, we present our first proposal for 3D object classification, namely PointNet, in Section 2.4. Later, PointNet is improved and thoroughly tested in adverse conditions with noise and occlusion throughout the study in Section 2.5. Next, LonchaNet is introduced in Section 2.6 as the last iteration of our system that incorporates all the lessons learned by the previous work. Finally, Section 2.7 draws conclusions and sets future lines of research.

## 2.1 Introduction

Object classification is fundamental to computer vision and despite the progress achieved during the last years, it still remains a challenging area of research. Arguably, most of the interest in object classification is due to its usefulness for robotics.

In that regard, recognizing objects is one of the problems that must be solved to achieve total visual scene understanding. Such deeper and better knowledge of the environment eases and enables the execution of a wide variety of more complex tasks. For instance, accurately recognizing objects in a room can be extremely useful for any robotic system that navigates within indoor environments. Due to the unstructured nature of those environments, autonomous robots need to do reasoning grounded in the dynamic real world. In other words, they need to understand the information captured by their sensors to perform tasks such as grasping, navigation, mapping, or even providing humans with information about their surroundings. Identifying the classes to which objects belong is one key step to enhance the aforementioned capabilities.

Despite the easy intuitive interpretation of the problem, its inherent difficulty can be misleading. We humans recognize numerous objects in difficult settings (e.g., different points of view, occlusion, or clutter) with little to no effort. However, approaching that problem is not that easy for a computer and taking into account all the possible settings and combinations of external factors renders this task a difficult one to solve efficiently and with high precision (which is often required in numerous application scenarios).

From a formal point of view, the object classification task can be formulated as follows: given an image  $\mathcal{I}^{H \times W}$  in which an object  $\mathcal{O}$  appears, which can be either a grayscale or RGB array of  $W$  pixels in width and  $H$  pixels in height, the goal is to predict the class of the object  $\mathcal{L}_{\mathcal{O}}$  from a set of  $N$  predefined object classes  $\mathcal{L} = \{\mathcal{L}_0, \mathcal{L}_1, \dots, \mathcal{L}_{N-1}\}$ .

Most of the classic literature of this topic tackled such problem by devising hand-crafted feature descriptors that are extracted on certain keypoints detected over the bidimensional image and later used either to compare them against pre-existing object descriptors in a database to match them to a certain class or either to feed them as input to a shallow machine learning architecture that learns to classify those descriptors to predict the class of the object that appears in the image. That paradigm shifted recently due to the success of deep learning architectures that are able to exploit their feature learning capabilities to avoid the need of hand engineering descriptors while achieving unprecedented accuracy levels. Furthermore, the adoption and spread of depth sensors has also added a literally new dimension to learn from to boost performance. The approaches introduced in this thesis are part of that cutting-edge trend that takes advantage of the additional geometric information facilitated by commodity range scanners to perform learning over them using deep architectures. A more detailed review of the field, from the very beginning to the current trends using 3D data and deep neural networks, is performed in Section 2.2.

Apart from the methods, data also plays a key role in object classification. As methods have evolved, so have datasets. Due to the increasing needs imposed by data-driven approaches, datasets have grown larger, more varied, and richer. That progress has enabled the development of new ways of solving the problem, e.g., using three-dimensional data. Section 2.3 briefly reviews the evolution of such datasets and describes the data that will be used throughout this chapter.

After reviewing the literature and introducing the data we are going to use, we start describing our first approach to perform object classification using 3D data, namely PointNet, capable of learning object classes from point clouds discretized as occupancy grids with uniform voxel grids in the tridimensional space. Section 2.4 describes this

architecture, its data representation, and also benchmarks it on a standard 3D object classification dataset to validate it.

Following that, Section 2.5 analyzes how noise and occlusion impact such 3D deep learning architecture and the importance of the data representation when dealing with such adverse conditions that commonly appear in the real world. In that study, we also propose minor changes to the architecture and the representation themselves that significantly boost accuracy with regard to the originally proposed PointNet.

At last, Section 2.6 takes all the lessons learned from the initial PointNet proposal and the extensive study to introduce a novel slice-based architecture to tackle the 3D object class recognition problem, LonchaNet, which achieved state of the art results in a standard benchmark.

In the end, Section 2.7 concludes this chapter by summarizing the insights gathered while discussing the proposed approaches. Furthermore, we also our main flaws and how to improve upon them besides from proposing new proposing future lines of research for 3D object classification.

## 2.2 Related Works

Since the very beginning of computer vision, a considerable amount of effort has been directed towards achieving robust object recognition systems [16]. This was mainly due to the fact that recognizing objects is a key capability required by robots to operate autonomously in unstructured, real-world environments. That continuous endeavor configured object classification as an ever-evolving area which has followed the general trend of computer vision, i.e., moving from hand-crafted features to trainable feature extractors, and which has also benefited from the improvements on imaging hardware, e.g., depth information from range scanners. In this section, we briefly review that evolution in order to put our proposal in context.

### 2.2.1 Traditional Approaches

Object recognition has been traditionally dominated by feature-based methods. This approach relies on extracting features, i.e., pieces of information which describe simple but significant properties of the objects. Those features are encoded into *descriptors* such as Scale Invariant Feature Transform (**SIFT**)[17], Speeded Up Robust Features (**SURF**)[18], Binary Robust Independent Elementary Features (**BRIEF**)[19], Binary Robust Invariant Scalable Keypoints (**BRISK**)[20], Oriented FAST and Rotated BRIEF (**ORB**)[21], or Fast Retina Keypoint (**FREAK**)[22] to name a few. After extracting those descriptors, machine learning techniques are applied to train a system with them so that it becomes able to classify features extracted from unknown instances. Based on the types of features, these methods can be divided into two categories: global or local feature-based methods. Global ones are characterized by dealing with the object as a whole; they define a set of features which completely encompass the object and describe it effectively. On the other hand, local methods describe local patches of the object, those regions are located around highly distinctive spots named *keypoints*.

Real-world scenes tend to be unstructured environments. This implies that object recognition systems must not be affected by clutter or partial occlusions. In addition, they should be invariant to illumination, transforms, and object variations. Those are the main reasons why local surface feature-based methods have been popular and successful during the last years – since they do not need the whole object to describe it properly, they are able to cope with cluttered environments and occlusions [23].

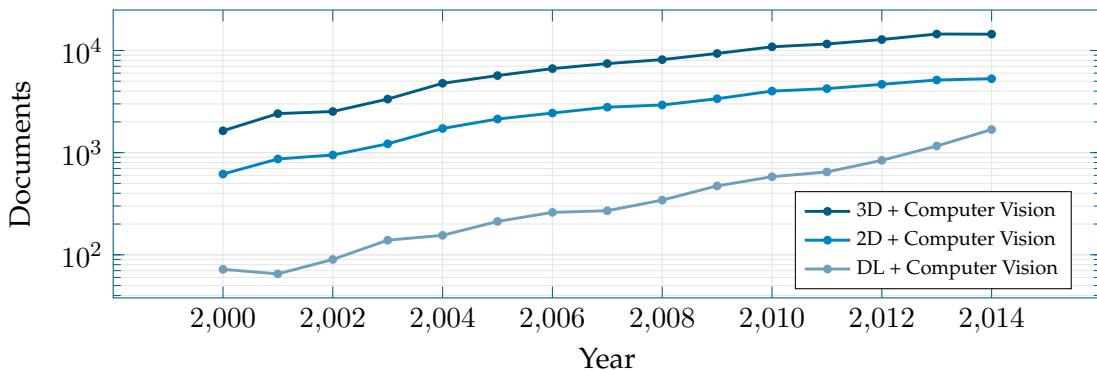
## 2.2.2 3D Object Recognition

Traditionally, those object recognition systems made use of **2D** images with intensity or color information, i.e., Red Green and Blue (**RGB**) images. However, technological advances made during the last years have caused a huge increase in the usage of **3D** information. The field of computer vision in general, and object recognition research in particular, have been slowly but surely moving towards including this richer information into their algorithms.

Nowadays, the use of **3D** information for this task is in a state of continuous evolution, still far behind, in terms of maturity, from the systems that make use of **2D** images. Nevertheless, the use of **2D** information exhibits a handful of problems which hinder the development of robust object recognition systems. Oppositely, the use of range images or point clouds, which provide **2.5D** or **3D** information respectively, presents many significant benefits over traditional **2D**-based systems. Some of the main advantages are the following ones [24]: (1) they provide geometrical information thus removing surface ambiguities, (2) many of the features that can be extracted are not affected by illumination or even scale changes, (3) pose estimation is more accurate due to the increased amount of surface information. Therefore, the use of **3D** data has become a solid choice to overcome the inherent hurdles of traditional **2D** methods.

However, despite all the advantageous assets of **3D** data, researchers had to overcome certain difficulties or drawbacks. On the one hand, sensors capable of providing **3D** were expensive, limited, and performed poorly in many cases. The advent of low-cost **3D** acquisition systems, e.g., Microsoft Kinect, enabled a widespread adoption of these kind of sensors thanks to their accessibility and affordability. On the other hand, **3D** object recognition systems are computationally intensive due to the increased dimensionality. In this regard, advances in computing devices like **GPUs** provided enough computational horsepower to run those algorithms in an efficient manner. In addition, the availability of low-power **GPU** computing devices like NVIDIA's Jetson has supposed a significant step towards deploying robust and powerful object recognition systems in mobile robotic platforms.

The combination of those three factors (the advantages of **3D** data, low-cost sensors, and parallel computing devices) transformed the field of computer vision in general, and object recognition in particular. As we can see in Figure 2.1, there has been a significant dominance of **3D** over **2D** research in computer vision since the year 2000.



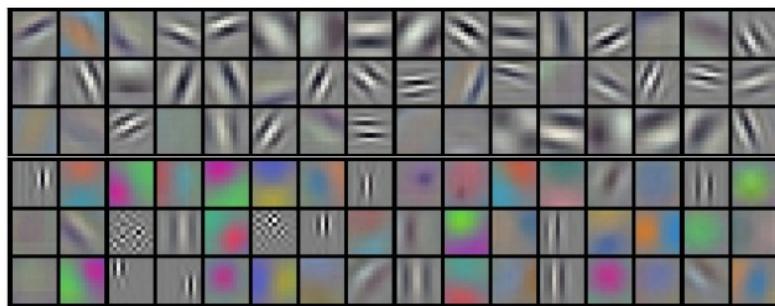
**Figure 2.1:** Evolution of the number of academic documents containing the terms **2D**, **3D**, and Deep Learning together with *Computer Vision*. Search terms statistics obtained from [scopus.com](http://scopus.com).

Therefore, creating a robust 3D object recognition system, which is also able to work in real time, became one of the main goals towards for computer vision researchers [25]. There exist many reviews about 3D object recognition in the literature, including the seminal works of Besl and Rain [26], Brady et al. [27], Arman et al. [28], Campbell and Flynn [29], and Mamic and Bennamoun [30]. All of them perform a general review of the 3D object recognition problem with varying levels of detail and different points of view. The work of Guo et al. [24] is characterized by its comprehensive analysis of different local surface feature-based 3D object recognition methods which were published between the years 1992 and 2013. In that review, they explain the main advantages and drawbacks of each one of them. They also provide an in-depth survey of various techniques used in each phase of a 3D object recognition pipeline, from the keypoint extraction stage to the surface matching one, including the extraction of local surface descriptors. The review is specially remarkable due to its freshness and level of detail. It is important to remark that all the described methods make use of carefully designed feature descriptors by experts in the field.

### 2.2.3 Deep Learning

From the earliest days of computer vision, the aim of researchers has been to replace hand-crafted feature descriptors, which require domain expertise and engineering skills, with multilayer networks able to learn them automatically by using a general-purpose training algorithm [31]. The solution for this problem was discovered during the 1970s and 1980s by different research groups independently [32][33][34]. This gave birth to a whole new branch of machine learning named deep learning.

Deep learning architectures usually consist of a multilayer stack of hierarchical learning modules which compute non-linear input-output mappings. Those modules are just functions of the input with a set of internal weights. The input of each layer in the stack is transformed, using the functions defined by the modules, to increase the selectivity and invariance of the representation. The backpropagation procedure is used to train those multilayer architectures by propagating gradients through all the modules. In the end, deep learning applications use feedforward neural network architectures which learn to map a fixed-size input, e.g., an image, to a fixed-size output, typically a vector containing a probability for each one of the possible categories [31].

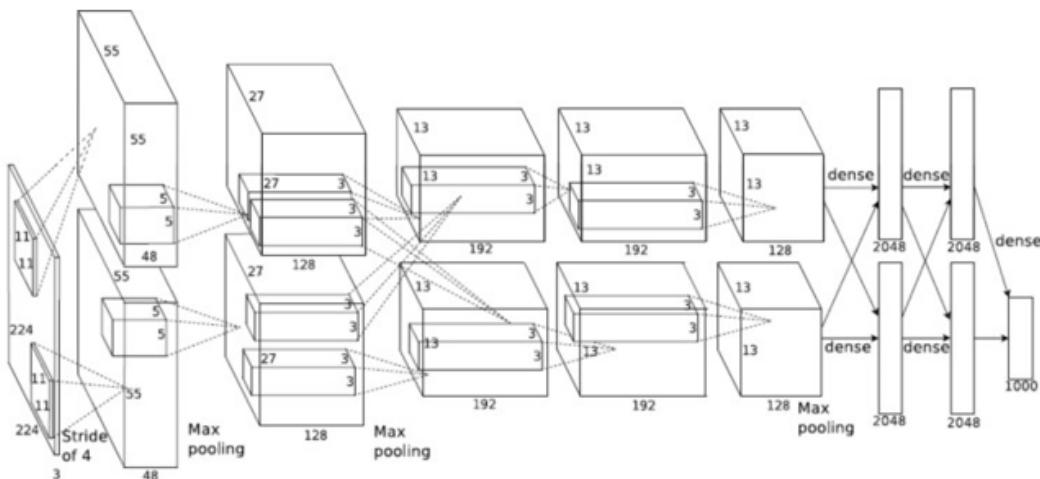


**Figure 2.2:** Filters learned by the network proposed by Krizhevsky et al. [35]. Each of the 96 filters shown is of size  $11 \times 11 \times 3$ . The filters have clearly learned to detect edges of various orientations and they resemble Gabor filters. Image analysis using that kind of filters is thought to be similar to perception in the human visual system [36].

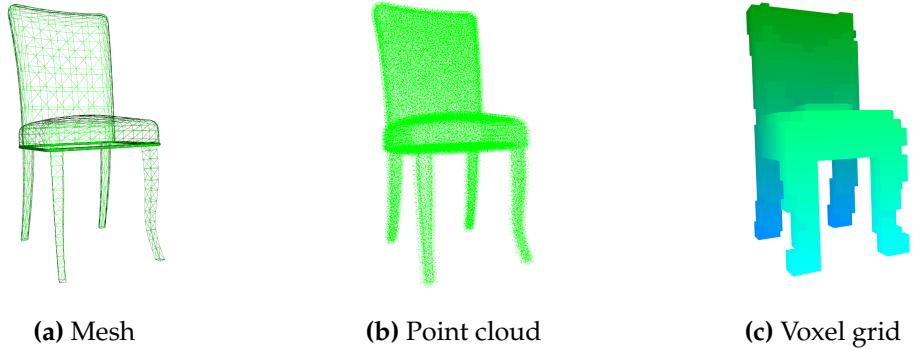
Figure 2.2 shows some sample filter modules automatically learned by training one of the most successful deep learning architectures: the deep convolutional neural network proposed by Krizhevsky et al. [35] to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 [37] contest into 1000 different classes.

In spite of the fact that these kind of architectures showed a huge potential for solving many computer vision problems, they were ignored by the computer vision community. In the latter years, certain breakthrough works revived the interest in deep learning architectures [31]. Recent studies proved that local minima are not an issue with large neural networks. Following a set of seminal works for the field on training deep learning networks [38][39], a group of researchers from the Canadian Institute for Advanced Research (**CIFAR**) introduced unsupervised learning procedures to create layers of feature detectors without labelled data, they also pre-trained several layers and added a final layer of output units; the system was tuned using backpropagation and achieved a remarkable performance when applied to the handwritten digit recognition or pedestrian detection problems [40]. In addition, the advent of **GPUs**, which were easily programmable and extremely efficient for parallel problems, made possible the training of huge networks in acceptable time spans [41].

All those contributions to the field led to the birth of probably the most important milestone regarding deep learning: the Convolutional Neural Network (**CNN**). This special kind of deep network was designed to process data in form of multiple arrays and gained popularity because of its many practical successes. This was due to the fact that they were easier to train and generalized far better than previous models. The architecture of a typical **CNN** is composed by many stages of convolutional layers followed by pooling ones and non-linearity Rectified Linear Unit (**ReLU**) filters; in the end, convolutional and fully connected layers are stacked. The key idea behind using this stack of layers is to exploit the property that many natural signals are compositional hierarchies, in which higher-level features are obtained by composing lower-level ones. Figure 2.3 shows a typical architecture of a **CNN**.



**Figure 2.3:** Illustration of the architecture of the aforementioned **CNN** proposed by Krizhevsky et al.[35] for the ImageNet challenge. Besides the normal components, e.g., convolutional, pooling, and fully connected layers, this network features two different paths. One **GPU** runs the layers at the top while the other runs the layers at the bottom.



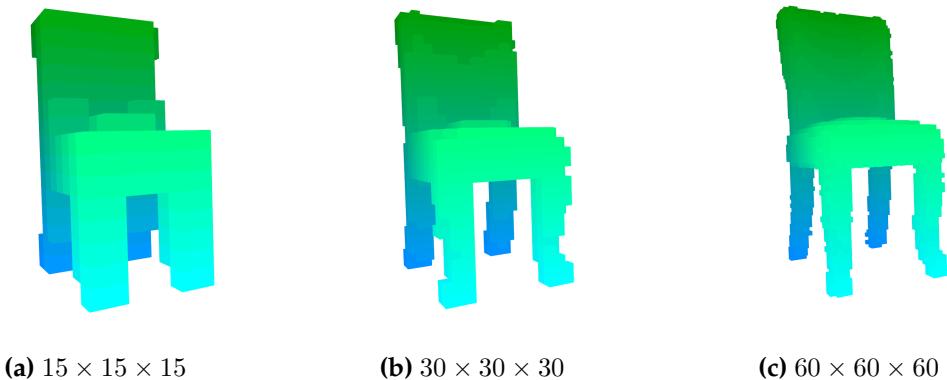
**Figure 2.4:** Common volumetric representations: polygonal mesh (a), point cloud (b), and voxel grid (c) of a chair model (which is color coded by height and depth).

#### 2.2.4 Volumetric Representations

Most of the literature focuses on how **CNNs** can learn filters and recognize objects using **2D** images as input. In this chapter, we will explore the usage of **3D** information to feed the network. For this purpose, we need volumetric representations for such information. Arguably, the most popular representations for volumetric data are **3D** meshes or point clouds, shown in Figure 2.4a and 2.4b respectively. A mesh consists of a collection of vertices (points in a three-dimensional (**3D**) coordinate system), edges (connections between those vertices), and faces (closed sets of edges, usually triangles) that defines the shape of an object. A point cloud is just a set of points defined by  $x$ ,  $y$ , and  $z$  coordinates in a three-dimensional coordinate system that model the surface of an object. However, those representations are unbounded and barely structured since they contain an arbitrary number of components, e.g., vertices or points, and no particular ordering is enforced for those entities.

This fact poses a problem since **CNNs** require a fixed-size representation for the input data. In order to overcome this limitation, alternative volumetric representations must be used to provide samples to the network for both training and testing. The most common volumetric representation which allows a structured and bounded definition of an object shape is the voxel grid. A voxel (word contraction of *volume element* or *volumetric pixel*) is the **3D** equivalent to a **2D** pixel, i.e., it is the minimal unit of a three-dimensional matrix. A volumetric object can be represented as a **3D** matrix of voxels, whose positions are relative to other voxels while points and polygons must be represented explicitly by **3D** coordinates. In this regard, voxels are able to efficiently represent regularly sampled **3D** spaces that are also non-homogeneously filled, while meshes and point clouds are good for representing **3D** shapes with empty or homogeneously filled space. It is important to notice that a voxel is just a data point in a three-dimensional grid, so its value may represent many different properties. The most popular and simple voxel grid type is the binary one (see Figure 2.4c) in which each voxel contains a binary value depending on whether the object's surface intersects or is partially contained in the voxel's volume.

Despite the fact that a binary voxel grid representation allows us to feed a **CNN** with volumetric data coming from different sources (point clouds provided by range sensors or polygonal meshes from **3D** models can be easily converted into voxel grids) a significant amount of information from the original representation is lost. This loss depends on the resolution of the grid, i.e., the voxel size which is usually referred as

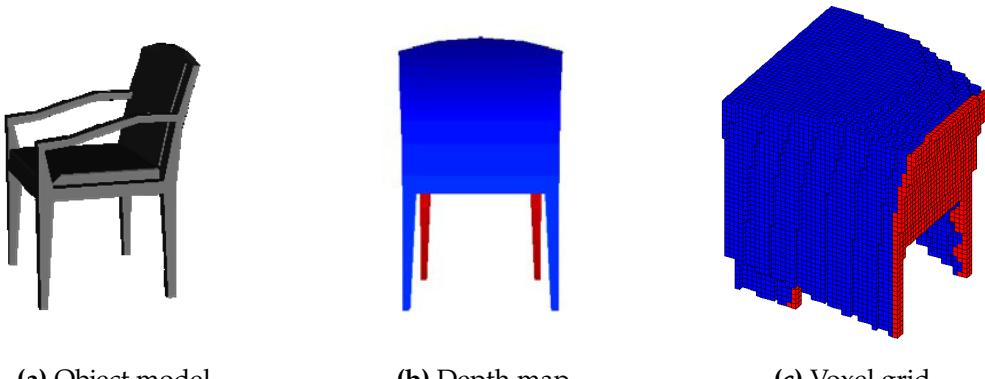


**Figure 2.5:** Effect of the leaf size on binary voxel grids. All grids have the same cubic size:  $300 \times 300 \times 300$  units. Leaf sizes vary from 5, 10, and 20 units, resulting in binary grids of  $15 \times 15 \times 15$  (a),  $30 \times 30 \times 30$  (b), and  $60 \times 60 \times 60$  voxels (c) respectively.

the *leaf size*. Figure 2.5 shows how the resolution of the voxel grid can be tuned to obtain more accurate or more compact volumetric representations. Although an arbitrary precision can be obtained by changing the leaf size, it is hard to determine a specific size which describes with enough detail all the possible inputs for the CNN without sacrificing the compactness of the grid.

Apart from those basic representations, others which maintain the properties of the voxel grid, but include additional information in the values of the cells can be used as well. Here we briefly review the most popular and successful volumetric representations for 3D data that have been used to feed CNNs for object recognition purposes.

The first step was taken by Wu *et al.* [42]; their work 3DShapeNets was the first to apply CNNs to pure 3D representations. Their proposal (shown in Figure 2.6) represents 3D shapes, from captured depth maps that are later transformed into point clouds, as 3D voxel grids of size  $30 \times 30 \times 30$  voxels –  $24 \times 24 \times 24$  data voxels plus 3 extra ones



**Figure 2.6:** 3DShapeNets representation proposed by Wu *et al.* as shown in their paper [42]. An object (a) is captured from a certain point of view and a depth map is generated (b) which is in turn used to generate a point cloud that will be represented as a voxel grid (c) with empty voxels (in white, not represented), unknown voxels (in blue), and surface or occupied voxels (red).

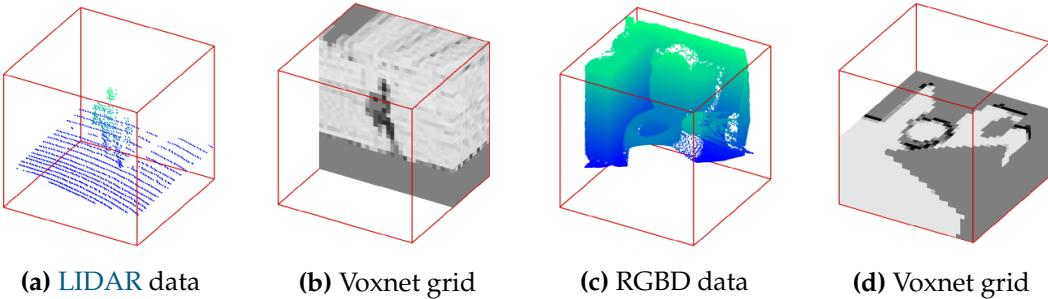


**Figure 2.7:** **TSDF** representation proposed by Song and Xiao as shown in their paper [43]. An object (a) is captured by a range sensor as a point cloud (b) and then a **TSDF** grid is generated (red indicates the voxel is in front of surfaces and blue indicates the voxel is behind the surface; the intensity of the color represents the **TSDF** value).

of padding in both directions to reduce convolution artifacts – which can represent free space, occupied space (the shape itself), and unknown or occluded space depending on the point of view. Neither the grid generation process, nor the leaf size is described but the voxel grid relies on prior object segmentation.

Song and Xiao [43] proposed to adopt a directional **TSDF** encoding which takes a depth map as input and outputs a volumetric representation. They divide a **3D** space using an equally spaced voxel grid in which each cell holds a three-dimensional vector that records the shortest distance between the voxel center and the three-dimensional surface in three directions. In addition, the value is clipped by  $2\delta$ , being  $\delta$  the grid size in each dimension. A  $30 \times 30 \times 30$  voxels grid is fitted to a previously segmented object candidate. Figure 2.7 shows a graphical representation of this approach.

Maturana and Scherer [44] use occupancy grids in *VoxNet* to maintain a probabilistic estimate of the occupancy of each voxel to represent a **3D** shape. This estimate is a function of the sensor data and prior knowledge. They propose three different occupancy models: binary, density, and hit. The binary and density models make use of raytracing to compute the number of hits and pass-throughs for each voxel. The former



**Figure 2.8:** Volumetric occupancy grid representation used by *VoxNet* as shown in their paper [44]. For **LIDAR** data (a) a voxel size of  $0.1m^3$  is used to create a  $32 \times 32 \times 32$  grid (b). For **RGB-D** data (??), the resolution is chosen so the object occupies a subvolume of  $24 \times 24 \times 24$  voxels in a  $32 \times 32 \times 32$  grid (d).

one assumes that each voxel has a binary state, occupied or unoccupied. The latter one assumes that each voxel has a continuous density, based on the probability it will block a sensor beam. The hit grid ignores the difference between unknown and free space, only considering hits; it discards information but does not require the use of raytracing so it is highly efficient in comparison with the other methods. They also propose two different grids for Light Detection and Ranging (**LIDAR**) and **RGB-D** sensor data. For the **RGB-D** case, they use a fixed occupancy grid of  $32 \times 32 \times 32$  voxels, making the object of interest – obtained by a segmentation algorithm or given by a sliding box – occupy a subvolume of  $24 \times 24 \times 24$  voxels. The  $z$  axis of the grid is aligned with the direction of gravity. Figure 2.8 shows the occupancy grids used by VoxNet.

### 2.2.5 Our Proposal in Context

Our proposal builds upon the successes in the literature which suggested that applying deep learning techniques to 3D information to solve the object class recognition problem exhibits potential to raise the bar in terms of performance. Firstly, we introduce a novel way for representing the **3D** input data, which is based on point density occupancy grids, and we integrate it into a **CNN** architecture. The focus of this first iteration is to prove that simple representations and architectures for 3D data can perform reasonably well while keeping computational cost at bay to enable real-time class recognition. Secondly, we carry out an in-depth study of the effect of adverse conditions that characterize real-world scenarios – such as noise caused by the sensor and occlusions due to the positions of the objects in the scene – on the performance of **CNNs** applied to **3D** object class recognition. In this case, our target is the assessment of iterative improvements applied to the previous architecture and representations to show that not only they can perform real-time recognition at reasonable accuracy but also maintain it under adverse circumstances. At last, we take a sideways step to propose a novel approach that uses multiple two-dimensional (**2D**) cross-section views of **3D** models for **3D** object class recognition.

## 2.3 Datasets

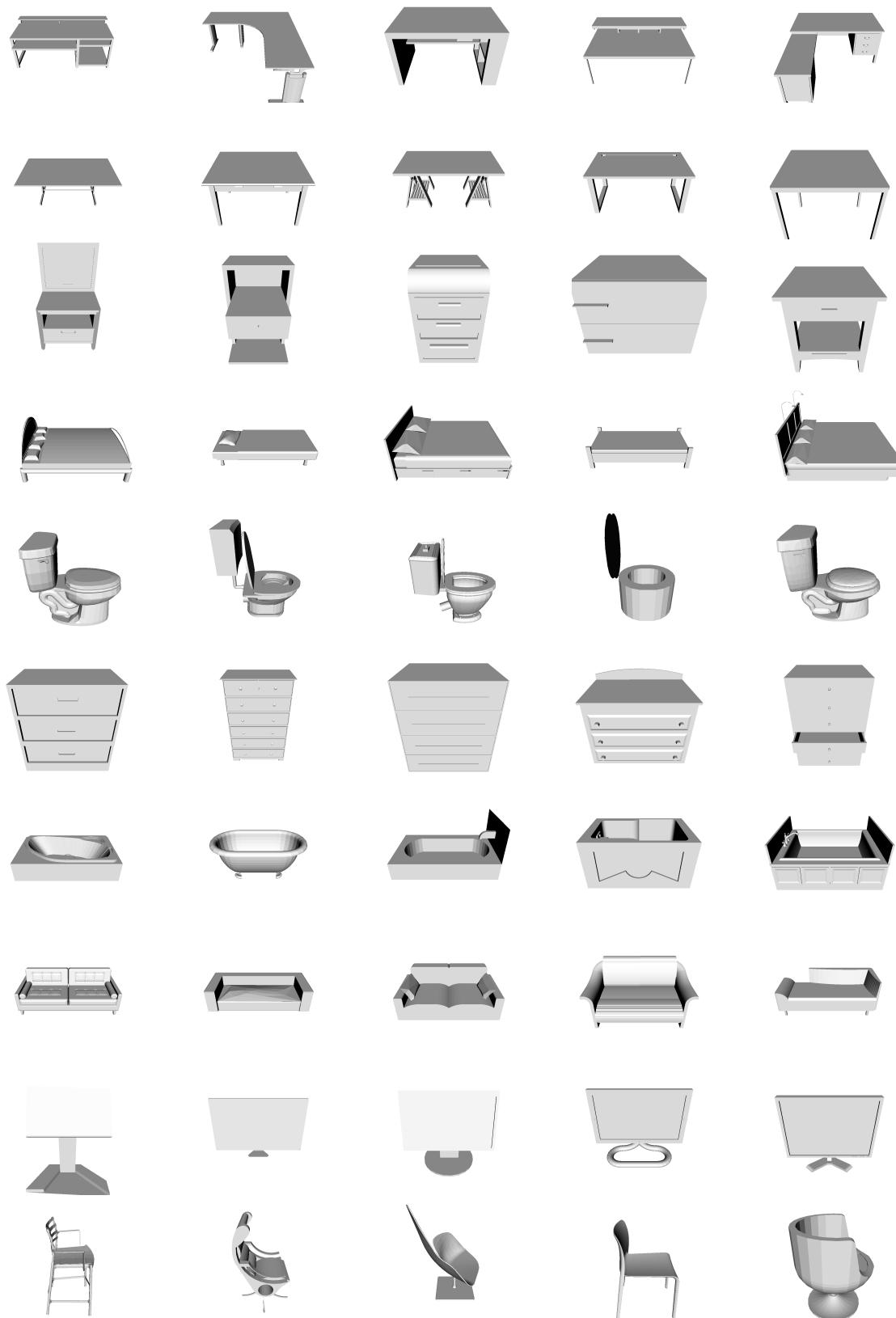
Deep neural network architectures are usually composed by many layers which in turn mean many weights to be learned. Because of that, there is a strong need of large-scale datasets to train those networks in order to avoid overfitting the model to the input data. Nowadays, large-scale databases of real-world **3D** objects are scarce, some of them do not have that high number of objects [45][46][47], or were incomplete by the time this work was performed [48]. A possible workaround to this problem consists of using Computer Aided Design (**CAD**) model databases – which are virtually unlimited – and processing those models to simulate real-world data.

The *Princeton ModelNet* project is one of the most popular large-scale **3D** object dataset. Its goal, as their authors state, is to provide researchers with a comprehensive clean collection of **3D CAD** models for objects, which were obtained via online search engines. Employees from the Amazon Mechanical Turk (**AMT**) service were hired to classify over 150 000 models into 662 different categories.

At the moment, there are two versions of this dataset publicly available for download<sup>1</sup>: *ModelNet-10* and *ModelNet-40*. Those are subsets of the original dataset which

---

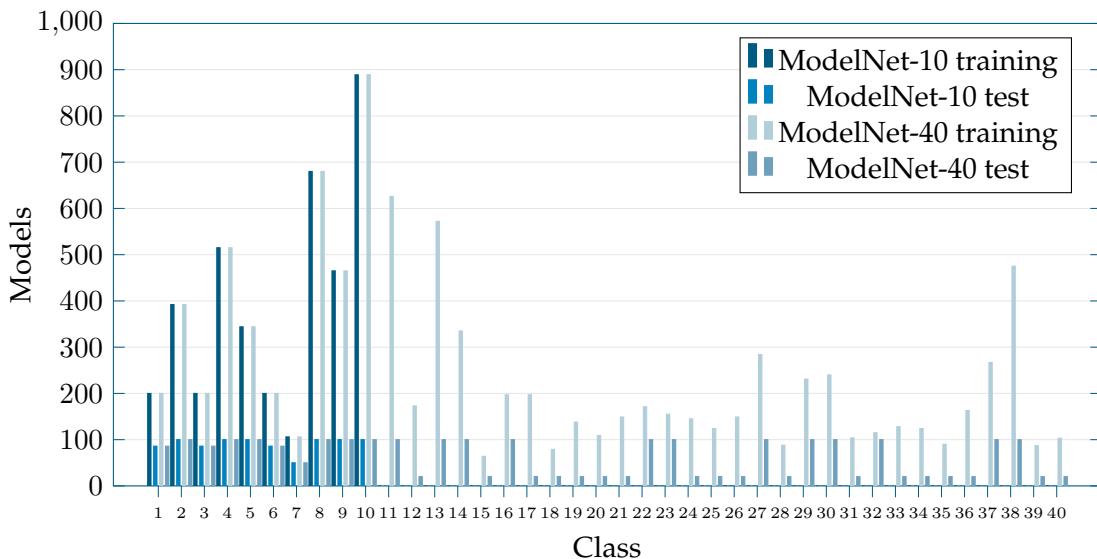
<sup>1</sup><http://modelnet.cs.princeton.edu/>



**Figure 2.9:** ModelNet10 samples.

only provide the 10 and 40 most popular object categories respectively. These subsets are specially clean versions of the complete dataset.

On the one hand, ModelNet-10 is composed of a collection of over 5000 models classified into 10 categories and divided into training and test splits. In addition, the orientation of all **CAD** models of the dataset was manually aligned. On the other hand, ModelNet-40 features over 9800 models classified into 40 categories, also including training and test sets. However, the orientations of its models are not aligned as they are in ModelNet-10. Figure 2.9 shows some model examples from ModelNet-10. Figure 2.10 and Table 2.1 show the model distribution per each class of both subsets taking into account the training and test splits.



**Figure 2.10:** Model distribution per object class or category for both ModelNet-10 and ModelNet-40 training and test splits.

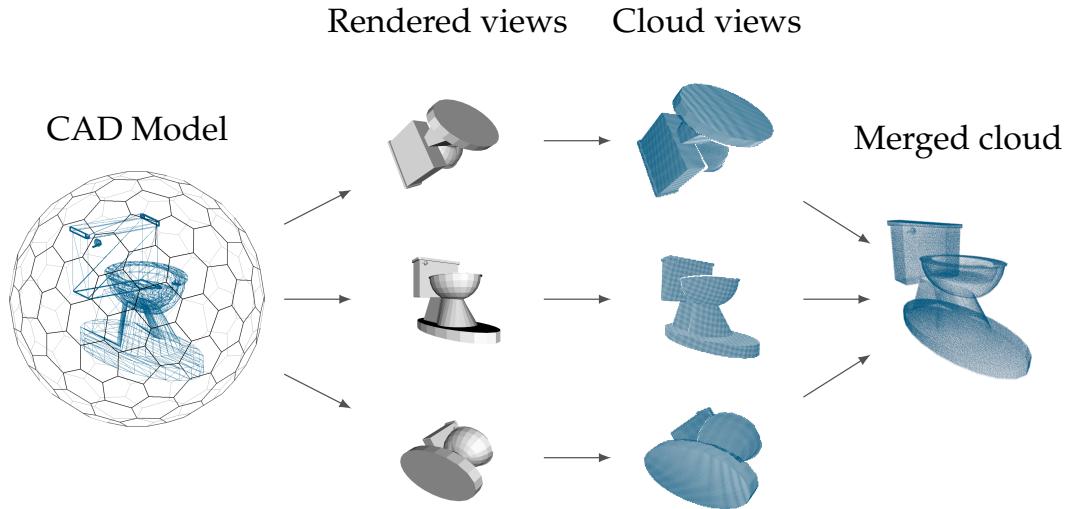
Since the final goal of the **CNNs** is to provide means to recognize objects onboard a mobile robotic platform which features RGB-Depth (**RGB-D**) sensors, it is logical to transform the full mesh representation provided by the dataset into the representation that our networks will deal with. **RGB-D** cameras output depth maps which can be used to generate **3D** point clouds of the scene viewed by the camera. In this regard, we will transform each **CAD** object of the dataset into partial point clouds from different Point of Views (**POVs**).

For this purpose, we converted the **OFF** models into Point Cloud Data (**PCD**) clouds using a raytracing-based process. The object is placed in the center of a **3D** sphere, which is tessellated to a certain level, and a virtual camera pointing to the center of the sphere is placed in each vertex of that truncated icosahedron. Then those partial views are rendered and their *z*-buffer data, which contains the depth information, is used to generate point clouds from each **POV**. In the end, those views are translated and rotated, depending on their **POV**, and merged into a cloud for the full object.

Figure 2.11 shows a diagram of the aforementioned process. For the conversion, we used the first tessellation level of the sphere, which generates 42 vertices or **POVs**. A resolution of  $256 \times 256$  pixels was used for rendering the views. A voxel grid filter with a leaf size of  $0.7 \times 0.7 \times 0.7$  units is applied to the merged cloud to equalize the point density, which is higher in certain zones due to view overlapping.

Category	ModelNet-10		ModelNet-40	
	Training Set	Test set	Training Set	Test set
1	Desk	200	86	200
2	Table	392	100	392
3	Nighstand	200	86	200
4	Bed	515	100	515
5	Toilet	344	100	344
6	Dresser	200	86	200
7	Bathtub	106	50	106
8	Sofa	680	100	680
9	Monitor	465	100	465
10	Chair	889	100	889
11	Airplane	-	-	626
12	Bench	-	-	173
13	Bookshelf	-	-	572
14	Bottle	-	-	335
15	Bowl	-	-	64
16	Car	-	-	197
17	Cone	-	-	167
18	Cup	-	-	79
19	Curtain	-	-	138
20	Door	-	-	109
21	Flower Pot	-	-	149
22	Glass Box	-	-	171
23	Guitar	-	-	155
24	Keyboard	-	-	145
25	Lamp	-	-	124
26	Laptop	-	-	149
27	Mantel	-	-	284
28	Person	-	-	88
29	Piano	-	-	231
30	Plant	-	-	240
31	Radio	-	-	104
32	Range Hood	-	-	115
33	Sink	-	-	128
34	Stairs	-	-	124
35	Stool	-	-	90
36	Tent	-	-	163
37	TV Stand	-	-	267
38	Vase	-	-	475
39	Wardrobe	-	-	87
40	X-Box	-	-	103

**Table 2.1:** ModelNet-10 and ModelNet-40 distributions.



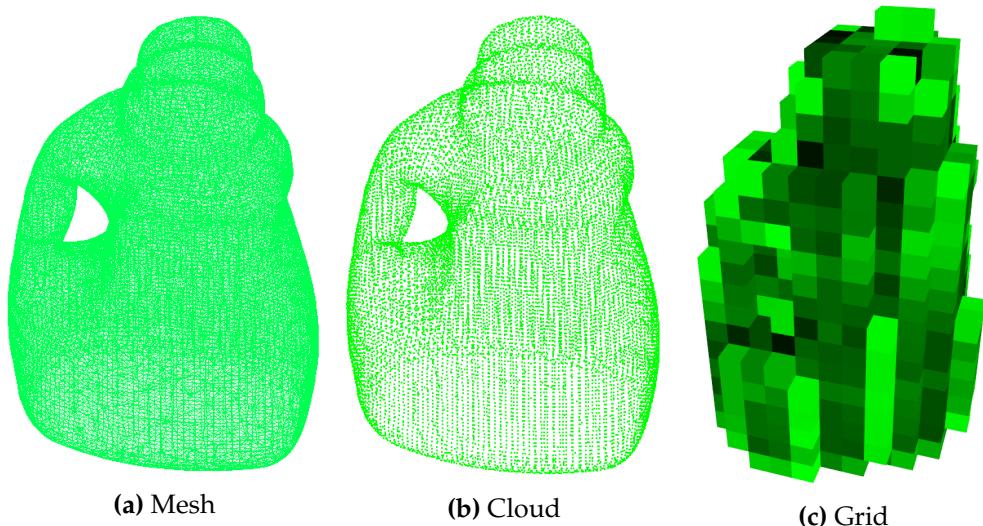
**Figure 2.11:** From **CAD** models to point clouds. The object is placed in the center of a tessellated sphere, views are rendered placing a virtual camera in each vertex of the icosahedron, the *z*-buffer data of those views is used to generate point clouds, and the point clouds are transformed and merged at last.

## 2.4 PointNet

In this first iteration of our efforts towards an object class recognition solution, we propose a system that takes a point cloud of an object as input and predicts its class label by leveraging two novel components: (1) a point density volumetric grid to estimate spatial occupancy inside each voxel, and (2) a **3D-CNN** which is trained to predict object classes from those grids. The occupancy grid – inspired by VoxNet’s [44] occupancy models that rely on probabilistic estimates – provides a compact representation of the object’s **3D** information originally present in the point cloud that can be fed to the **CNN** architecture, which in turn computes a label for that sample, i.e., predicts the class of the object. The components of this approach, namely *PointNet*, will be described throughout this section. Firstly, the proposed data representation is introduced in Section 2.4.1. Secondly, the network itself is presented in Section 2.4.2. Next, we carry out a set of experiments to validate the system in Section 2.4.3. At last, in Section 2.4.4, we draw conclusions about this work and set the stage for the next iteration.

### 2.4.1 Data Representation

As we mentioned before, our proposed architecture takes a point cloud of an object as input to recognize it. However, point clouds are unstructured representations that cannot be easily handled by common **CNN** architectures due to the lack of a matrix-like organization. The most straightforward way to apply formal convolutions to that unstructured space is to impose a certain organization into it. Occupancy grids expose a compact representation of the volumetric space; they stand between meshes or clouds, which offer rich but unstructured information, and voxelized representations,



**Figure 2.12:** Various 3D representations for an object. A mesh (a) is transformed into a point cloud (b), and that cloud is processed to obtain a voxelized occupancy grid (c). The occupancy grid shown in this figure is a cube of  $30 \times 30 \times 30$  voxels. Each voxel of that cube holds the point density inside its volume. In this case, dark voxels indicate high density whilst bright ones are low density volumes.

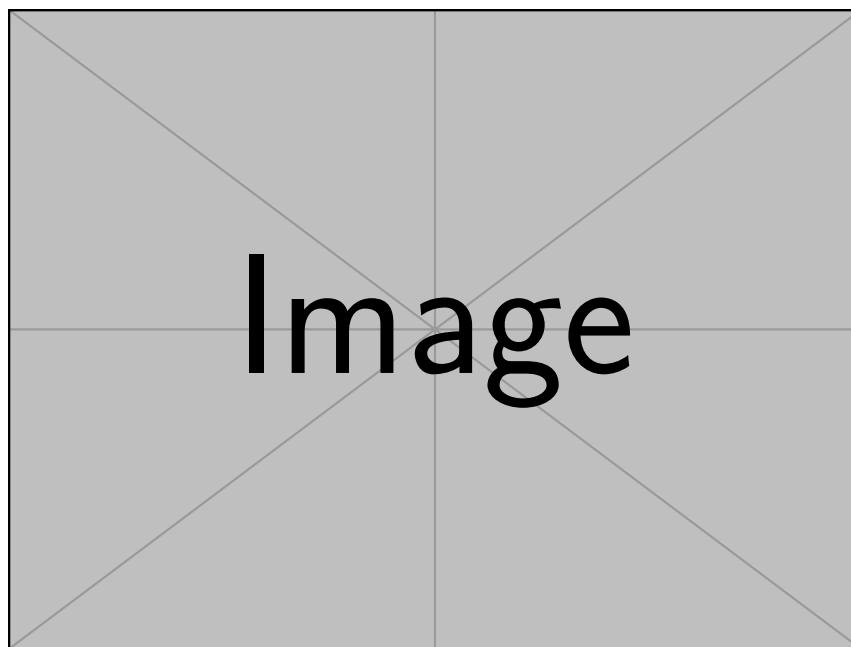
with packed but poor information. At that midpoint, occupancy grids provide important shape cues to perform learning while enabling an efficient processing of that information thanks to their matrix-like organization.

As we previously reviewed in Section 2.2, certain 3D deep learning architectures make use of occupancy grids as a representation for the input data. For instance, 3D ShapeNets [42] is a Convolutional Deep Belief Network (**CDBN**) which represents a 3D shape as a  $30 \times 30 \times 30$  binary tensor in which a one indicates that a voxel intersects the mesh surface, and a zero represents empty space. VoxNet [44] introduces three different occupancy grids that employ 3D ray tracing to compute the number of beams hitting or passing each voxel and then use that information to compute the value of each voxel depending on the chosen model: a binary occupancy grid using probabilistic estimates, a density grid in which each voxel holds a value corresponding to the probability that it will block a sensor beam, and a hit grid that only considers hits thus ignoring empty or unknown space. The binary and density grids proposed by Maturana *et al.* [44] differentiate unknown and empty space, whilst the hit grid and the binary tensor do not. VoxNet’s occupancy grid outperforms 3D ShapeNets in terms of accuracy in the ModelNet challenge for the 3D-centric approaches described above. However, ray tracing grids considerably harmed performance in terms of execution time so that other approaches must be considered for a real-time implementation.

With PointNet, we propose an occupancy grid inspired by the aforementioned successes but aiming to maintain a reasonable accuracy while enabling a real-time implementation. In our volumetric representation, each point of a cloud is mapped to a voxel of a fixed-size occupancy grid. Before performing that mapping, the object cloud is scaled to fit the grid. Each voxel will hold a value representing the number of points mapped to itself. At last, the values held by each cell are normalized. Figure 2.12 shows the derivation of the proposed occupancy grid representation from other typical tridimensional representations of a sample object.

### 2.4.2 Network Architecture

As we have previously stated, **CNNs** have proven to be very useful for recognizing and classifying objects in **2D** images. A convolutional layer can recognize basic patterns such as corners or planes and if we stack several of them they can learn a topology of hierarchical filters. Furthermore, the composition of several of these regions can define a feature of a more complex object. By doing that, a combination of various filters is able to recognize a full object. We apply this approach used in **2D** images to **3D** recognition. The deep architecture featured by *PointNet* is represented in Figure 2.13.



**Figure 2.13:** PointNet’s 3D **CNN** architecture. [MISSINGDETAILS]

### 2.4.3 Experiments

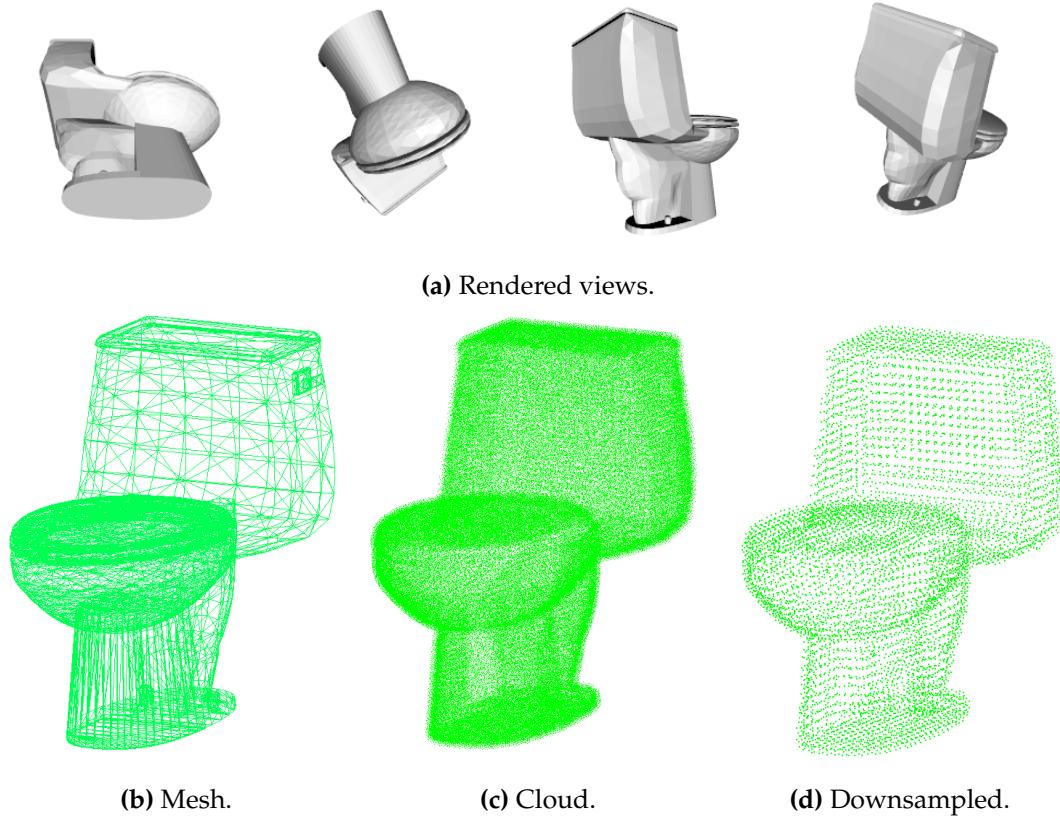
#### Data Generation

We generated point clouds to feed PointNet following the procedure described in Section 2.3. Figure 2.14 illustrates the result of the aforementioned process. After that, the resulting point clouds are used to train, randomizing the order of the models, and test the system taking into account the corresponding splits provided by ModelNet-10.

#### Implementation and Setup

This architecture was implemented using the Point Cloud Library ([PCL](#)) [49][50] which provides state-of-the-art algorithm implementations for 3D point cloud processing and Caffe [51], a deep learning framework developed and maintained by the Berkeley Vision and Learning Center ([BVLC](#)) and an active community of contributors on GitHub <sup>2</sup>. This BSD-licensed C++ library enables researchers to design, train, and deploy [CNN](#) architectures efficiently, mainly thanks to its drop-in integration of NVIDIA cuDNN [52] to take advantage of [GPU](#) acceleration.

<sup>2</sup><https://github.com/BVLC/caffe>



**Figure 2.14:** Dataset model processing example to generate the point clouds for PointNet. Some rendered views of a toilet model are shown in (a). The original [OFF](#) mesh is shown in (b). The generated point cloud after merging all points of view is shown in (c), and (d) shows the downsampled cloud using a voxel grid filter with a leaf size of  $0.7 \times 0.7 \times 0.7$ .

All the timings and results were obtained by performing the experiments in the following test setup: Intel Core i5-3570 with 8 GB of 1600 MHz DD3 RAM on an ASUS P8H77-M PRO motherboard (Intel H77 chipset). Additionally, the system includes an NVIDIA Tesla K20 GPU, and a Seagate Barracuda 7200.14 secondary storage. Caffe RC2 was run over ElementaryOS Freya 0.3.1, an Ubuntu-based Linux distribution. It was compiled using CMake 2.8.7, g++ 4.8.2, CUDA 7.0, and cuDNN v3.

## Results and Discussion

As a result of training PointNet with a learning rate of 0.0001 and a momentum of 0.9 during 200 iterations using the ModelNet-10 dataset, it obtained a success rate of 77.6%. As shown in Figure [MISSINGREF], the confusion matrix reveals the stability of the system, mainly confusing items that look alike such as desk and table. Because of the nature of **CNNs**, which heavily rely on detecting combinations of features, these kind of errors are common. As we can observe in Figure 2.15, the visual features that define a desk and a table are almost the same, making it hard to distinguish between both classes. Figure 2.16 shows the neuron activations for the output layer of the architecture, proving that *Desk* and *Table* are consistently confused during the tests. In light of these experiments, and taking into account the knowledge of the **CNNs** principles, it is conceivable to think that a deeper network would provide better results so more experiments were carried out.

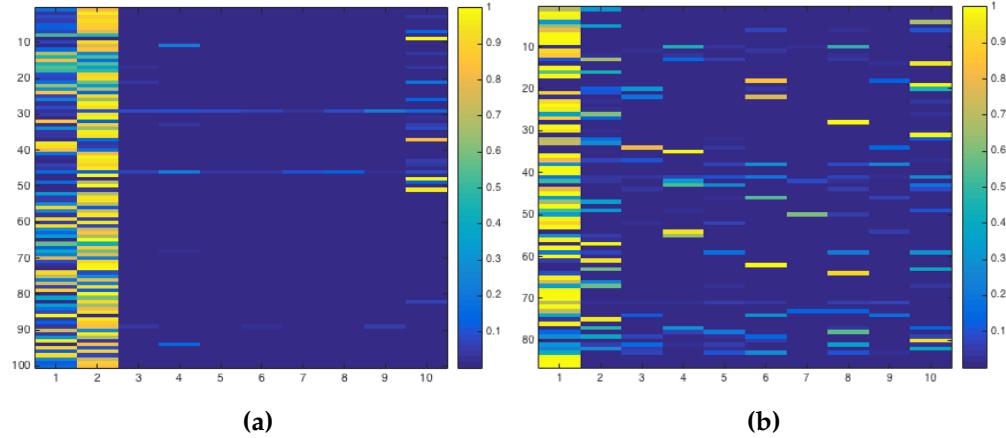
In the deeper network experiment we added several layers to the PointNet architecture. One more convolutional layer was added since these layers are coupled to the detection of the features of the objects, so the more layers there are, a better or more expressive model is produced. An Inner Product layer was also added. Since these layers make the classification possible, adding more of them would theoretically provide better classification results.

This architecture was trained during 1,000 iterations and tested every 200 iterations. The best result was provided by the 800 iterations test with an accuracy of 76.7%, while the 1,000 iterations test dropped the performance to a 75.9% due to overfitting.

It is well known that training using an unbalanced dataset tends to harm those classes with the least number of examples and to benefit those with the most, as stated by [MISSINGREF]. Having this in mind, and knowing that ModelNet-10 is highly unbalanced as shown in Table [MISSINGTABLE], the dataset was balanced by limiting

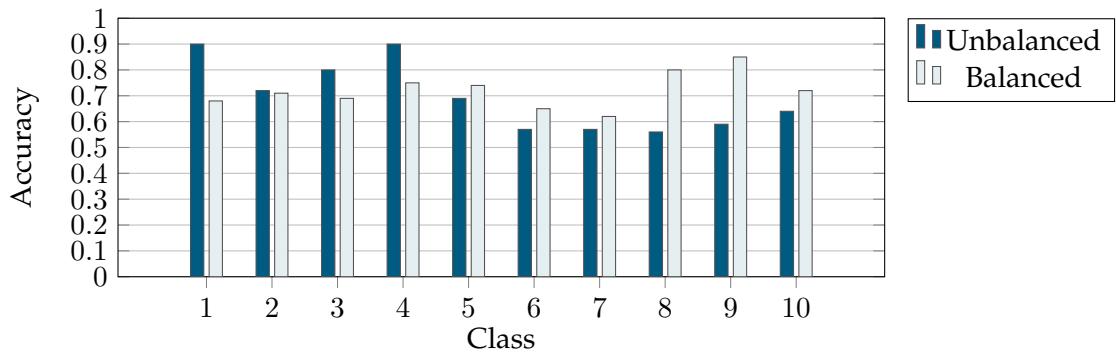


**Figure 2.15:** Similarity between two objects of different classes: Table and Desk. The point cloud shown in (a) represents an object of the Table class, whilst the point cloud in (b) represents an object whose class is Desk but it is misclassified as a Table due to their resemblance.



**Figure 2.16:** Neuron activations for the output layer of the architecture when classifying all the test samples for both *Desk* (b) and *Table* (a) classes. Each row represents an activation vector for a specific sample, so each column is a position of the vector: the activation to that particular class. The first column corresponds to the *Desk* class, while the second one is the *Table*. The activation shows the clear confusion between *Desk* and *Table*. Although the latter one is much less confused with other classes, many *Tables* are misclassified as *Desks* thus lowering the accuracy for this class.

the number of examples of each class to 400 using random undersampling. This does not fully solve the problem but improves the difference between the classes with the least number of examples and those with the most. The network was trained and tested with this more balanced dataset and it achieved an accuracy of 72.9%. The fact is that balancing the training set makes the accuracy of the classes with less examples higher, but it harms the success rate on classes with more instances as seen in Figure 2.17.



**Figure 2.17:** Comparison of accuracy per class using an unbalanced dataset and a balanced one with a maximum of 400 models per class via random undersampling. Accuracy is harmed in the classes in which models are removed but gained otherwise.

After analyzing the results, it can be stated that neither a deeper network nor balancing the dataset increase accuracy. In fact, the experiments of the original architecture with the unbalanced ModelNet-10 exhibited better performance with a 77.6% success rate. In addition, PointNet takes an average time of 24.6 miliseconds to classify

an example (in comparison with VoxNet, which can take up to half a second for its raytracing-based implementation). These results prove the system as a fast and accurate **3D** object class recognition tool.

#### 2.4.4 Conclusion

PointNet is a brand new kind of CNN for object class recognition that handles tridimensional data, inspired by VoxNet and 3D ShapeNets but using density occupancy grids as inner representation for input data. It was implemented in Caffe and provides a faster method than the state of art ones yet obtaining a high success rate as the experiments over the ModelNet10 dataset. This fact enlightens a promising future in real-time 3D recognition tasks.

Following on this work, we plan to improve the inner representation by using adaptable occupancy grids instead of fixed-size ones. In addition, we will integrate the system in an object recognition pipeline for 3D scenes. Our network will receive a point cloud segment of the scene where the object lies, produced by a preprocessing method, and that segments will be used to generate the occupancy grids that will be learned by the system. This implies adapting the system for learning partial views of the objects and dealing with occlusions and scale changes. As an additional feature, we will include pose estimation in that pipeline, all of this with goal of developing an end-to-end 3D object recognition system.

### 2.5 Noise and Occlusion

#### 2.5.1 Data Representation

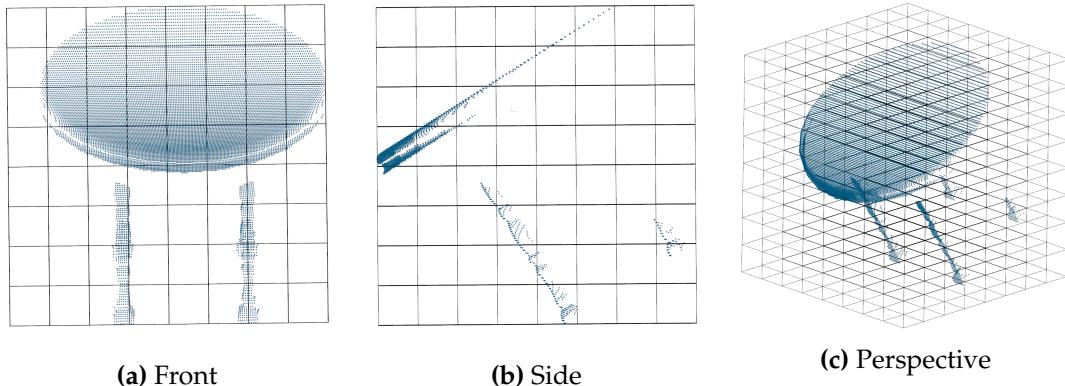
As is clear from the previous sections, a volumetric representation to be fed to a **2.5D** or **3DCNN** must encode the **3D** shape of an object as a **3D** tensor of binary or real values. This is due to the fact that raw **3D** data is sparse, i.e., a **3D** shape is only defined on its surface, and **CNNs** are not engineered for this kind of data.

In this regard, our proposal for the study is twofold. First, we implemented two different ways of generating the structure of the tensor – position, grid size, and leaf size – using a fixed grid and an adaptive one. Second, we developed two possible occupancy measures for the volumetric elements of the tensor.

#### Tensor Generation

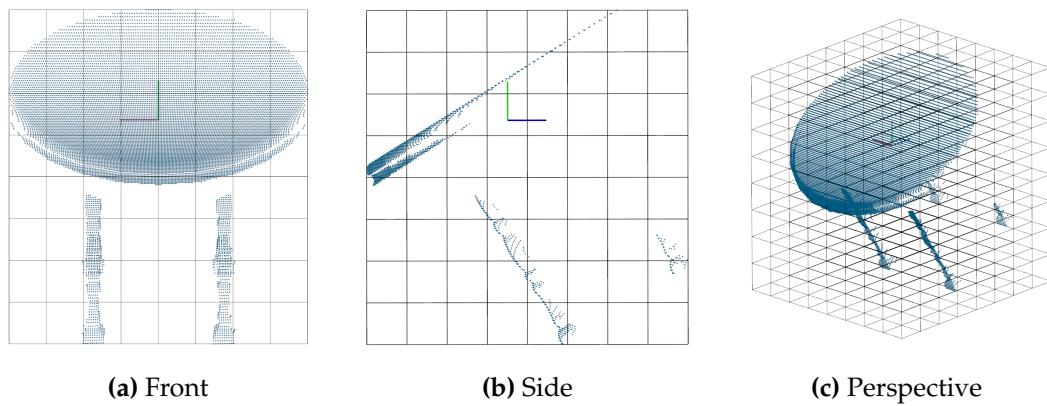
Providing that the input to our network consists of point clouds generated from the information provided by **RGB-D** sensors, we need to generate a discretized representation of the unbounded **3D** data to feed the network. Each cloud will be represented as a **3D** tensor. For that purpose, we need to spawn a grid to subdivide the space occupied by the point clouds. Two types are proposed: one with fixed leaf and grid sizes, and another one which will adapt those sizes to fit the data.

**Fixed** This kind of grid sets its origin at the minimum  $x$ ,  $y$ , and  $z$  values of the point cloud. Then the grid is spawned, with fixed and predefined sizes for both grid and voxels. After that, the cloud is scaled up or down to fit the grid. The scale factor is computed with respect to the dimension of maximum difference between the cloud and the grid. The cloud is scaled with that factor in all axes to maintain the original ratios. As a result, a cubic grid is generated as shown in Figure 2.18.



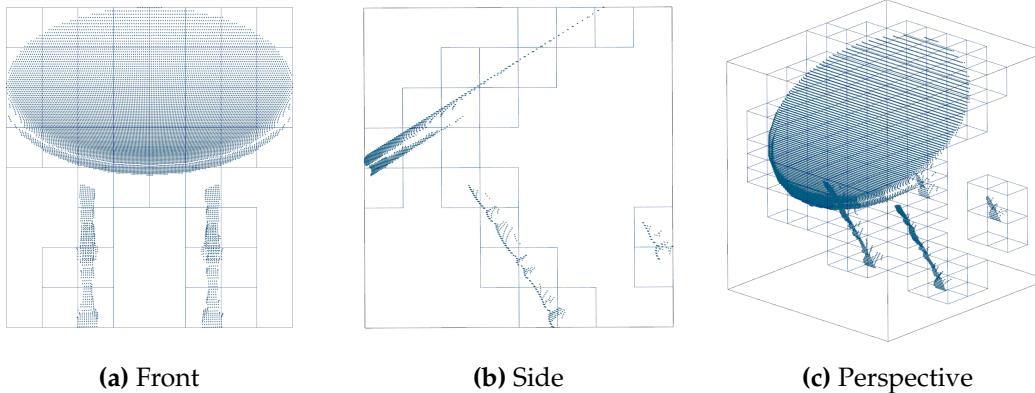
**Figure 2.18:** A fixed occupancy grid ( $8 \times 8 \times 8$  voxels) with 40 units leaf size and 320 units grid size in all dimensions. The grid origin is placed at the minimum  $x$ ,  $y$ , and  $z$  values of the point cloud. Front (a), side (b), and perspective (c) views of the grid over a partial view of a segmented table object are shown.

**Adaptive** The adaptive grid also sets its origin at the minimum  $x$ ,  $y$ , and  $z$  values of the point cloud. Next, the grid size is adapted to the cloud dimensions. The leaf size is also computed in function of the grid size. Knowing both parameters, the grid is spawned, fitting the point cloud data. As a result, a non-cubic grid is generated. As shown in Figure 2.19, all voxels have the same size, but they are not necessarily cubic.



**Figure 2.19:** An adaptive occupancy grid ( $8 \times 8 \times 8$  voxels) with adapted leaf and grid sizes in all dimensions to fit the data. The grid origin is placed at the minimum  $x$ ,  $y$ , and  $z$  values of the point cloud. Front (a), side (b), and perspective (c) views of the grid over a partial view of a segmented table object are shown. Notice that the point clouds for the three views are exactly the same for this figure and Figure 2.18, but the grids do change. There is a noticeable difference in the front view. In Figure 2.18, using fixed grids, all voxels are cubic and the point cloud does not fit the grid completely (leftmost column in Figure 2.18a), whilst in this figure, with adaptive grids, the grid is fitted to the cloud.

It is important to remark that, in both cases (fixed and adaptive), the number of voxels in the grid is fixed. Figures 2.18 and 2.19 show examples for both types using



**Figure 2.20:** Occupied voxels in an adaptive  $8 \times 8 \times 8$  grid generated over a partial view point cloud. Those voxels with points inside are shown in a wireframe representation. Empty voxels are omitted. Occupied voxels must be filled with values which represent the contained shape.

$8 \times 8 \times 8$  voxels for the sake of a better visualization.

It is also important to notice that each representation serves a purpose. The fixed grid will not always fit the data perfectly so it might end up having sparse zones with no information at all (as seen in Figure 2.18a on the first column). However, it can be used right away for sliding box detection. On the contrary, the adaptive grid fits the data to achieve a better representation. Nonetheless, it relies on a proper segmentation of the object to spawn the grid.

### Occupancy Computation

After spawning the grid to generate a discrete space, we need to determine the values for each cell or voxel of the 3D tensor. In order to do that, we must encode the geometric information of the point cloud into each occupied cell (see Figure 2.20). In other words, we have to summarize as a single value, the information of all points which lie inside a certain voxel. One way to do that is using occupancy measures. For that purpose, we propose two different alternatives: binary occupancy, normalized density.

**Binary** The binary tensor is the simplest representation that can be conceived to encode the shape. Voxels will hold binary values, they will be considered occupied if at least a point lies inside, and empty otherwise. Figure 2.21 shows an example of this tensor.

**Normalized Density** Binary representations are simple and require low computational power. However, complex shapes may get oversimplified so useful shape information gets lost. This representation can be improved by taking into account more shape information. A possible alternative consists of computing the point density inside each voxel, i.e., counting the number of points that fall within each cell.

It is important to notice that point density directly depends on the cloud resolution which in turn depends on many factors involving the camera and the scene, e.g., it is common for **RGB-D** to generate denser shapes in closer surfaces. To alleviate this problem, we can normalize the density inside each voxel dividing each value by the

maximum density over the whole tensor. An example of normalized density tensor is shown in Figure 2.22.

### 2.5.2 Network Architecture

In this section, we will describe the main layers that compose the CNN that will be used for the study. Figure 2.23 shows a diagram of the chosen architecture. It is highly inspired by *Voxnet* [44] and *PointNet* [12]. The network was implemented using Caffe. It features 2D convolutions and takes full 3D object model point clouds as input.

The input layer is a custom data layer implemented in Caffe which takes object point clouds as inputs and generates the corresponding discrete volumetric representation as discussed in the previous section.

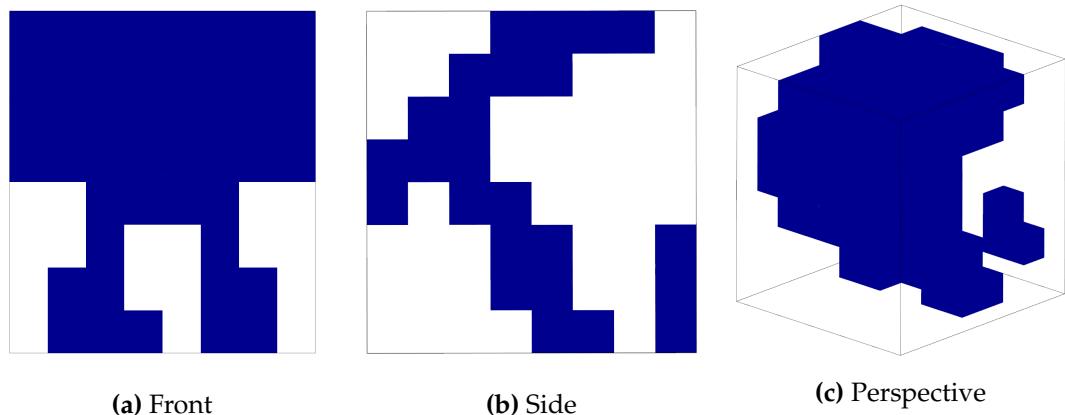
Next, we can find a *convolution layer* or  $C(m, n, d)$ . This layer applies  $m$  filters of size  $n \times n$  and a stride of  $d \times d$  voxels. In our case, this first convolution layer learns 48  $3 \times 3$  filters using a stride of  $1 \times 1$  voxels. This convolution layer is followed by a *ReLU* activation to introduce non-linearities to the model.

After that, another convolution layer is found. In this case, it will learn 128  $5 \times 5$  filters with a stride of  $1 \times 1$  voxels again. This layer is also followed by a *ReLU* activation one.

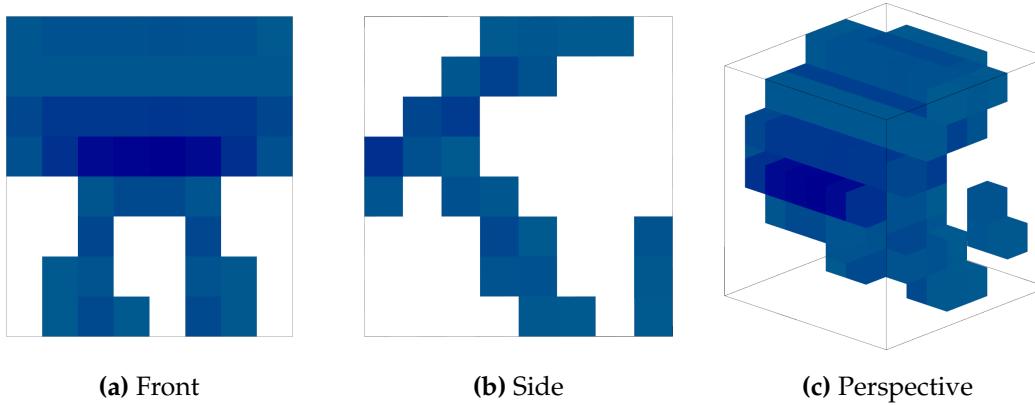
A *pooling layer* or  $P(n, d)$  takes place after those blocks. It performs a max-pooling process to summarize the input data, taking the maximum value of a fixed local spatial region of  $n \times n$  which is滑ed across the input volume using a stride of  $d \times d$  voxels. In this case, a pooling region of  $2 \times 2$  voxels with the same stride was chosen.

At last, we can find an *inner product layer* or  $IP(n)$ . It is just a fully connected layer, a traditional neural network architecture which consists of  $n$  neurons (1024 in this case). It is followed by a *ReLU* activation and a *dropout layer* Srivastava2014 or  $DP(r)$ . The function of the dropout layer is to avoid overfitting, randomly dropping connections with a probability  $r$  (0.5 in our case). In the end, another fully connected layer represents the output of the network, with as many output neurons as classes has our classification problem. Since our dataset has 10 classes (see Section ??) this layer has 10 neurons.

We use the term 2.5D to refer to this network due to the fact that it processes 3D data using 2D convolutions. This means that, in the end, its convolutions do not fully

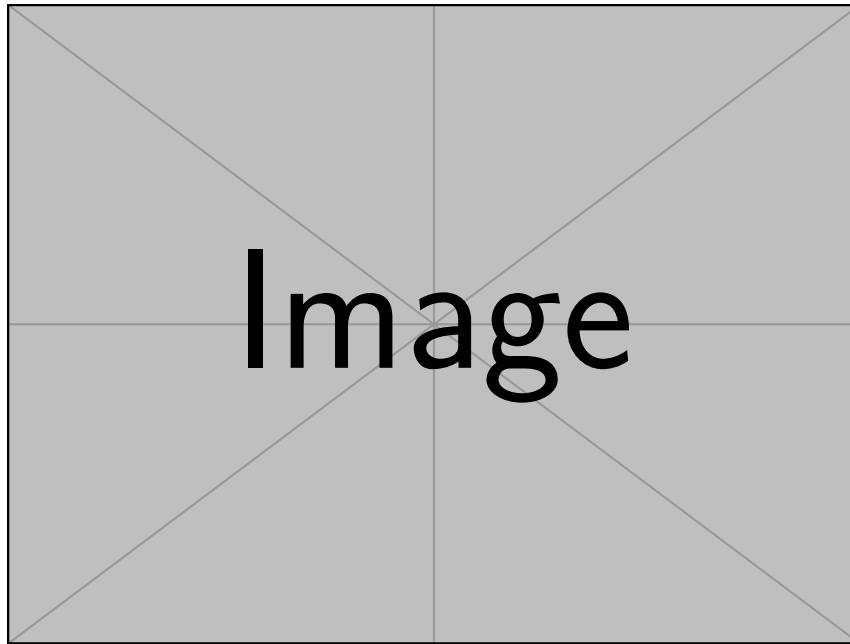


**Figure 2.21:** Binary tensor computed over a point cloud of a partial view of an object (shown in Figure 2.20). Occupied voxels are shown in blue, empty voxels are omitted for the sake of simplicity.

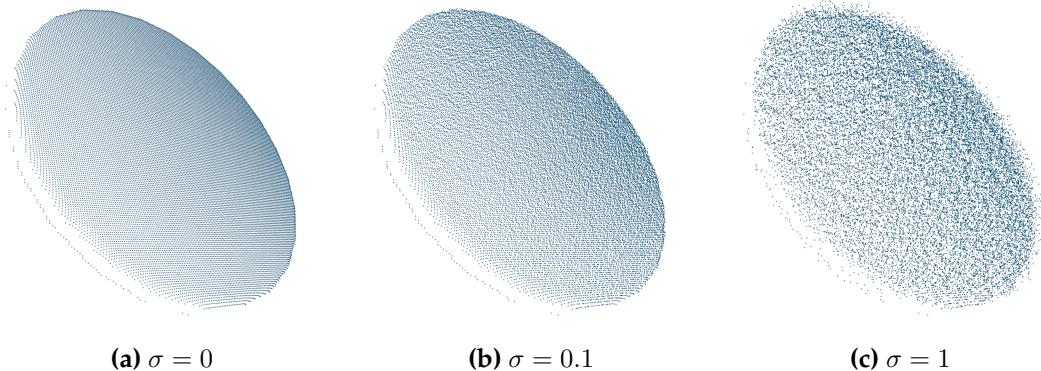


**Figure 2.22:** Normalized density tensor over a point cloud of a partial view of an object (shown in Figure 2.20). Denser voxels are darker and sparse ones are shown in light blue. Empty voxels were removed for visualization purposes.

take into account the depth spatial dimension of the input as if we were using pure 3D convolution filters. It is intuitive to think that a 3D CNN would yield better results due to that extra spatial dimension. However, a 3D CNN has some disadvantages that made us consider using a 2.5D CNN instead for the experimentation: (1) higher computational cost, (2) memory footprint is also much higher, (3) more parameters thus harder training. For those reasons, the main body of the experiments were carried out using the 2.5D approach.



**Figure 2.23:** CVIU's architecture. [MISSINGDETAILS]



**Figure 2.24:** Different levels of noise ( $\sigma = 0$  (a),  $\sigma = 0.1$  (b), and  $\sigma = 1$  (c)) applied to the  $z$ -axis of every point of a table partial view.

### 2.5.3 Experiments

In order to assess the performance of the proposed model-based [CNN](#) we carried out an extensive experimentation to determine the accuracy of the model and its robustness against occlusions and noise – situations that often occur in real-world scenes. For that purpose we started using the normalized density grids since they offer a good balance between efficiency and representation. We also investigated the effect of both fixed and adaptive grids using different sizes. Further experimentation was performed to compare the normalized density grids with the binary ones. We also carried out a brief experiment using a [3D CNN](#) to compare its performance with the [2.5D](#) counterpart.

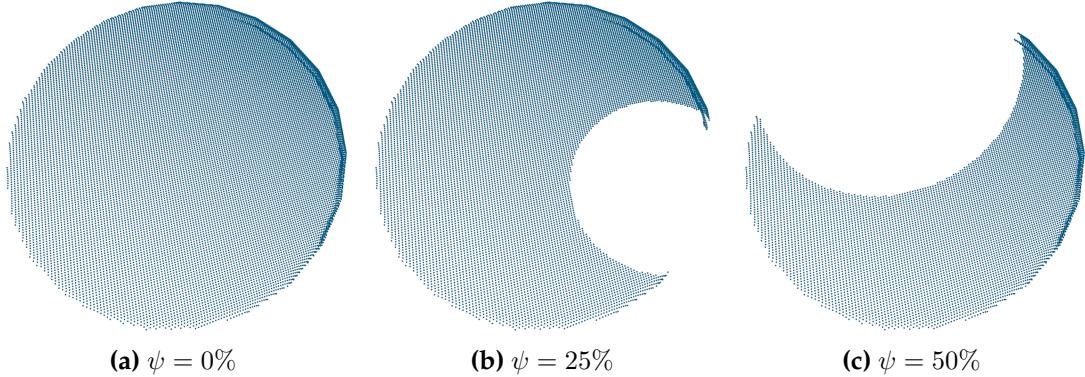
#### Data Generation

Data was generated in a similar fashion as we did for PointNet. The meshes from the ModelNet-10 dataset were converted to point clouds following the procedure described in Section 2.3. However, this time we applied a set of transformations to simulate noise and occlusions in our data.

**Noise Simulation** The partial views generated using the previously described process are not a good simulation of the result that we would obtain by using a low-cost [RGB-D](#) sensor. Those systems are noisy, so the point clouds produced by them are not a perfect representation of the real-world objects.

In order to properly simulate the behavior of a sensor, a model is needed. In our case, we are dealing with low-cost [RGB-D](#) sensors such as Microsoft Kinect and Primesense Carmine. A complete noise model for those sensors, specifically for the Kinect device, must take into account occlusion boundaries due to distance between the Infrared ([IR](#)) projector and the [IR](#) camera, 8-bit quantization,  $9 \times 9$  pixel correlation window smoothing, and  $z$ -axis or depth Gaussian noise [Gschwandtner2011](#).

We will make use of a simplification of this model, only taking into account the Gaussian noise since it is the most significant one for the generated partial views. In this regard, the synthetic views are augmented by adding Gaussian noise to the  $z$  dimension of the point clouds with mean  $\mu = 0$  and different values for the standard deviation  $\sigma$  to quantify the noise magnitude. Figure 2.24 shows the effect of this noise over a synthetic partial view of one object of the dataset.



**Figure 2.25:** Different levels of occlusion ( $\psi = 0\%$  (a),  $\psi = 25\%$  (b), and  $\psi = 50\%$  (c)) applied randomly to a table partial view.

**Occlusion Simulation** In addition to modelling the sensor to improve our synthetic data, it is important to also take the environment into account. In a real-world scenario, objects are not usually perfectly isolated and easily segmented; in fact, it is common for them to be occluded by other elements of the scene.

The occlusion simulation process consists of picking a random point of the cloud with a uniform probability distribution. Then, a number of closest neighbors to that point are picked. At last, both the neighbors and the point are considered occluded surface and removed from the point cloud. The number of neighbors to pick depends on the amount of occlusion  $\psi$  we want to simulate. For instance, for an occlusion  $\psi = 25\%$  we will remove neighbors until the rest of the cloud contains a 75% of the original amount of points, i.e., we will remove a 25% of the original cloud. Figure ?? shows the effect of the random occlusion process with different occlusion factors  $\psi$  over a synthetic partial view of a table object of the dataset.

It is important to notice the randomness of the occlusion process. This means that even with a high  $\psi$  it is possible not to remove any important surface information and vice versa. In other words, it is possible for some objects to remove a 50% of their points and still be recognizable because the removed region was not significant at all, e.g., a completely flat surface. However it is possible to render an object unrecognizable by removing a small portion of its points if the randomly picked surface is significant for its geometry. This remark is specially important when testing the robustness of the system. In order to guarantee that an appropriate measure of the robustness against missing information is obtained, a significant amount of testing sets must be generated and their results averaged so that it is highly probable to test against objects which have been occluded all over their surface across the whole testing set.

### Implementation and Setup

Results were obtained using the following test setup: Intel Core i7-5820K with 32 GiB of Kingston HyperX 2666MHz and CL13 DDR4 RAM on an Asus X99-A motherboard (Intel X99 chipset). Additionally, the system included an NVIDIA Tesla K40c GPU used for training and inference. The framework of choice was Caffe RC2 running on Ubuntu 14.04.02. It was compiled using CMake 2.8.7, g++ 4.8.2, CUDA 7.5, and cuDNN v3.

## Results and Discussion

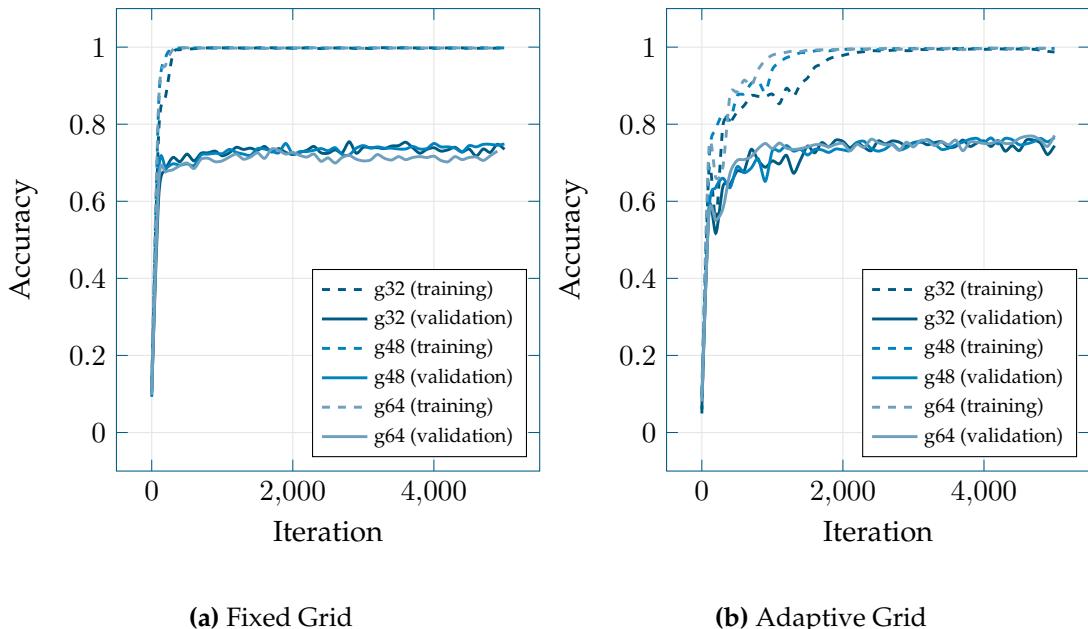
The networks were trained for a maximum of 5000 iterations – weights were snapshoted every 100 iterations so in the end we selected the best sets of them as if we were early stopping – using Adadelta as optimizer with  $\delta = 1 \cdot 10^{-8}$ . The regularization term or weight decay in Caffe was set to  $5 \cdot 10^{-3}$ . A batch size of 32 training samples was chosen.

After describing the experimentation setup, the dataset that was used to train and test the networks, and the ways of simulating noise and occlusion for the test sets, we will present and discuss the results of the experiments. Firstly, the normalized density tensor results – using the [2.5D CNN](#) – will be presented. After that, we will proceed with the binary tensor ones. Furthermore, we will report the experiments which produced the best results with a pure [3D CNN](#) with fully [3D](#) convolutions. At last, we will perform a comparison with the state of the art.

**Density Tensor** Figure 2.26 shows the accuracy results of the network for both grid types and increasing sizes. The peak accuracies for the fixed grids are  $\approx 0.75$ ,  $\approx 0.76$ , and  $\approx 0.73$  for sizes 32, 48, and 64 respectively. In the case of the adaptive one, the peak accuracies are  $\approx 0.77$ ,  $\approx 0.78$ , and  $\approx 0.79$  for the sizes 32, 48, and 64 respectively.

Taking those facts into account, we can extract two conclusions. First, the adaptive grid is able to achieve a slightly better peak accuracy in all cases; however, the fixed grid takes less iterations to reach accuracy values close to the peak in all cases. Second, there is no significant difference in using a bigger grid size of 64 voxels instead of a smaller one of 32.

The most important fact that can be observed in the aforementioned figures is that there is a considerable gap between training and validation accuracy in all situations. As we can observe, all networks reach maximum accuracy for the training set whilst



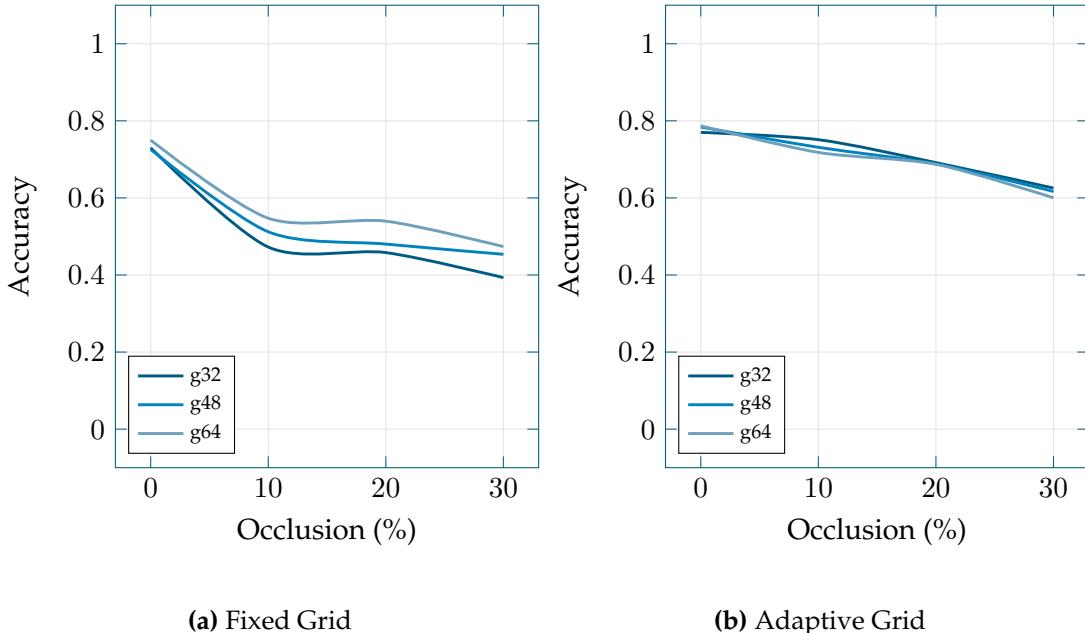
**Figure 2.26:** Evolution of training and validation accuracy of the model-based [CNN](#) using both fixed (a) and adaptive (b) normalized density grids. Different grid sizes (32, 48, and 64) were tested.

the validation one hits a glass ceiling at approximately 0.80. We hypothesize that the network suffers overfitting even when we thoroughly applied measures to avoid that. The most probable cause for that problem is the reduced number of training examples. In the case of ModelNet10 the training set consists of only 3991 models. Considering the complexity of the [CNN](#), it is reasonable to think that the lack of a richer training set is causing overfitting.

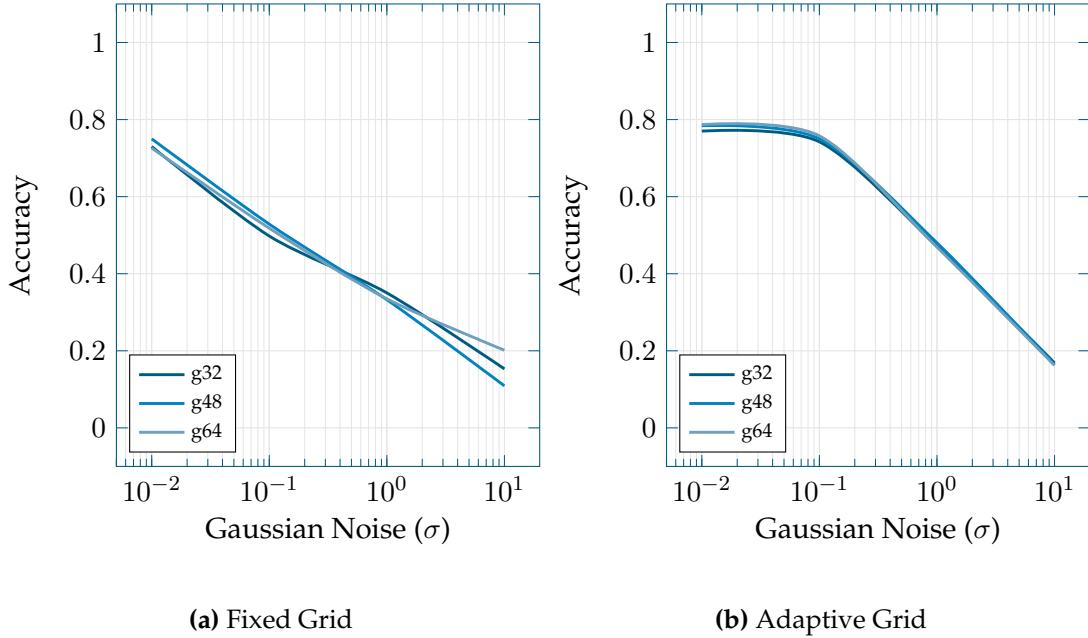
Concerning the robustness against occlusion, we took the best networks after training and tested them using the same validation sets as before but introducing occlusions in them (up to a 30%). Figure 2.27 shows the accuracy of both grid types with different sizes as the amount of occlusion in the validation model increases. As we can observe, occlusion has a significant and negative impact on the fixed grid – bigger grid sizes are less affected – going down from  $\approx 0.75$  accuracy to 0.40 – 0.50 approximately in the worst and best case respectively when a 30% of the model is occluded. On the contrary, the adaptive grid does not suffer that much – it goes down from  $\approx 0.78$  to  $\approx 0.60$  in the worst case – and there is no significant difference between grid sizes. In conclusion, the adaptive grid is considerably more robust to occlusion than the fixed one.

Regarding the resilience to noise, we also tested the best networks obtained from the aforementioned training process using validation sets with different levels of noise (ranging from  $\sigma = 1 \cdot 10^{-2}$  to  $\sigma = 1 \cdot 10^1$ ). Figure 2.28b shows the results of those experiments. It can be observed that adding noise has a significant impact on the fixed grid, even small quantities, reducing the accuracy from  $\approx 0.75$  to  $\approx 0.60$ ,  $\approx 0.4$ , and  $\approx 0.2$  for  $\sigma = 1 \cdot 10^{-1}$ ,  $\sigma = 1 \cdot 10^0$ , and  $\sigma = 1 \cdot 10^1$  respectively. On the other hand, the adaptive one shows remarkable robustness against low levels of noise (up to  $\sigma = 1 \cdot 10^{-1}$ ), barely diminishing its accuracy.

In the end, both grids suffer huge penalties in accuracy when noise levels higher than  $\sigma = 1 \cdot 10^{-1}$  are introduced, being the adaptive one less affected. The grid size has



**Figure 2.27:** Evolution of validation accuracy of the model-based [CNN](#) using both fixed (a) and adaptive (b) normalized density grids as the amount of occlusion in the validation models increases from 0% to 30%. Three grid sizes were tested (32, 48, and 64).



**Figure 2.28:** Evolution of validation accuracy of the model-based CNN using both fixed (a) and adaptive (b) normalized density grids as the standard deviation of the Gaussian noise introduced in the  $z$ -axis of the views increases from 0.001 to 10. The common grid sizes were tested (32, 48, and 64).

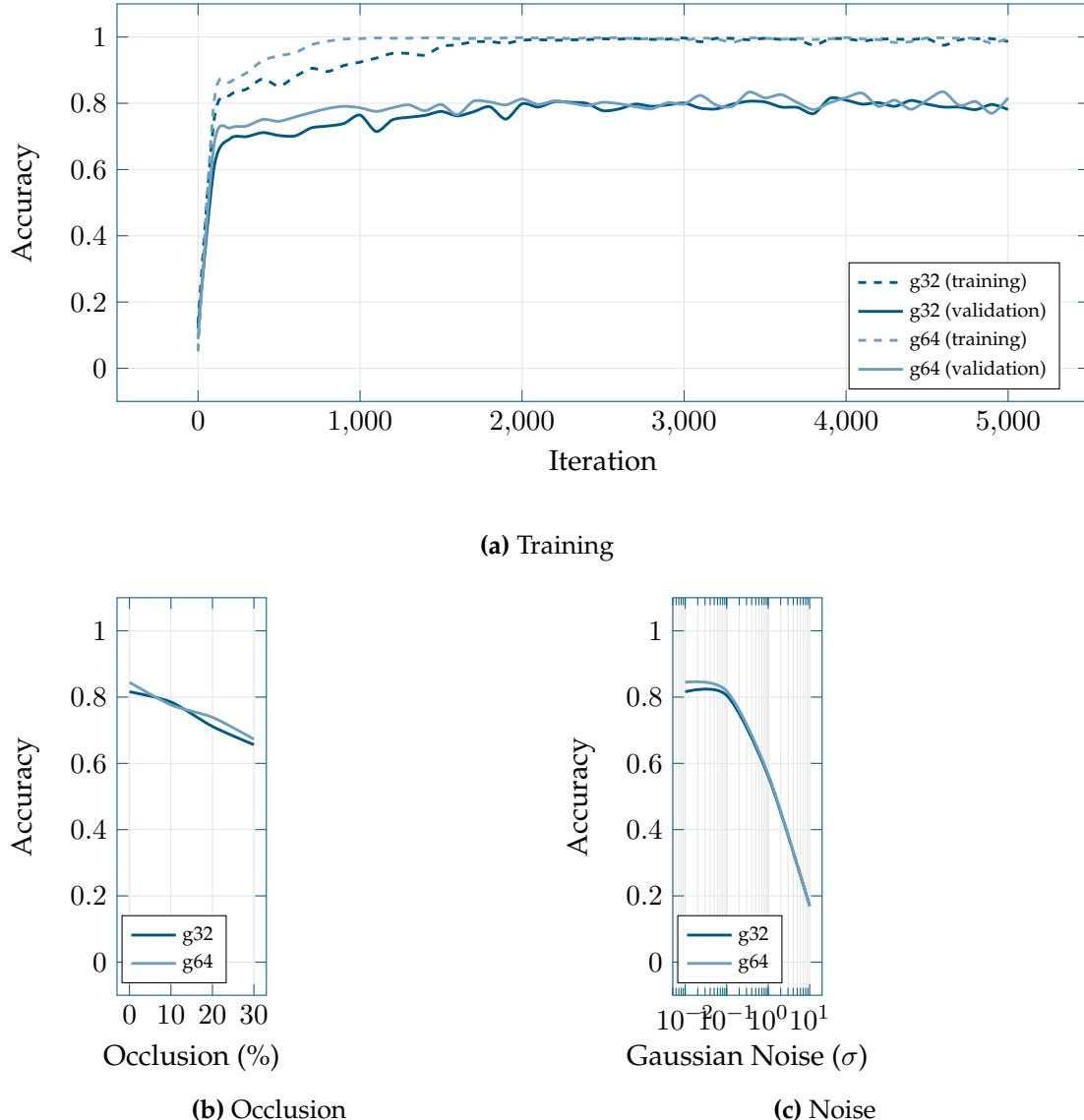
little to no effect in both cases, only in the fixed grid bigger sizes are slightly more robust when intermediate to high levels of noise are introduced. In conclusion, the adaptive grid is significantly more resilient to low levels of noise, and slightly outperforms the fixed one when dealing with intermediate to high ones.

**Binary Tensor** After testing the performance of the normalized density grid, we also trained and assessed the accuracy of the binary one in the same scenarios. This test intended to show whether there is any gain in using representations which include more information about the shape – at a small penalty to execution time.

For this experimentation we picked the best performer in the previous sections: the adaptive grid. We also discarded the intermediate size (48 voxels) since there was no significant difference between it and the others. Figure ?? shows the accuracy results of the network trained using binary grids. As we can observe, there is no significant difference between grid sizes neither. However, using this representation we achieved a peak accuracy of approximately 0.85, using 64 voxels grids, which is better to some extent than the normalized density one shown in Figure 2.26.

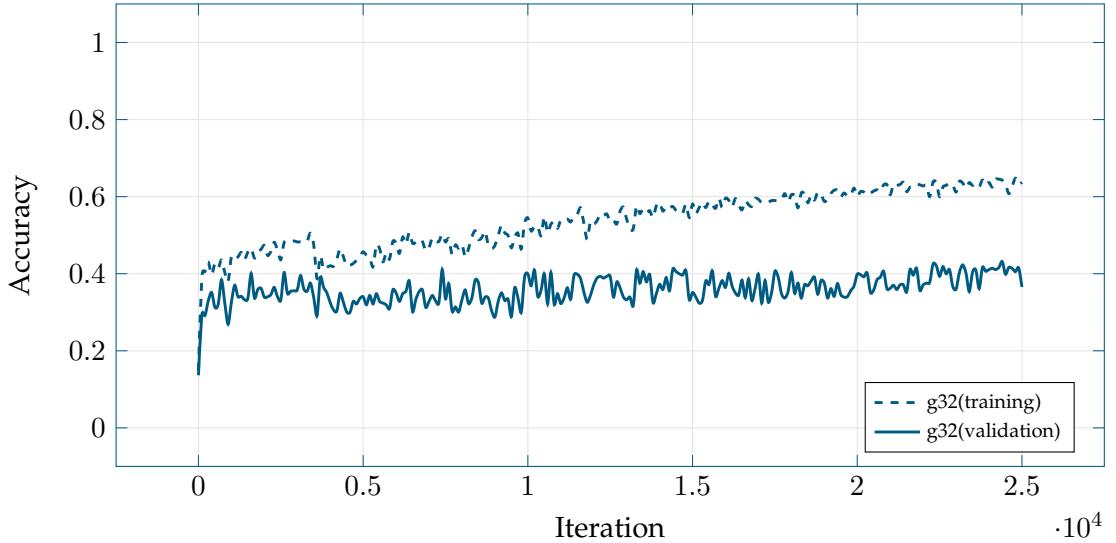
Occlusion and noise tolerance (shown in Figures 2.29b and 2.29c respectively) is mostly similar to the robustness shown by the normalized density adaptive grid (see Figures 2.27b and 2.28b) except from a small offset caused by the higher accuracy of the binary grid network.

In conclusion, the less-is-better effect applies in this situation and turns out that the simplification introduced by the binary representation helps the network during the learning process. It is pending to check if this statement is still valid if the validation accuracy is not bounded by network overfitting.



**Figure 2.29:** Evolution of training and validation accuracy of the model-based CNN using adaptive binary grids (a). Evolution of validation accuracy for the best network weights after training as the amount of occlusion in the validation set increases (b) and different levels of noise are introduced (c).

**3D CNN** At last, we tested the best configuration – binary adaptive grids – with a 3D CNN architecture with pure 3D convolutions. We kept the same architecture we introduced in Section ??, but extended its convolution and pooling layers to three dimensions. We then trained the network using adaptive binary grids as the volumetric representation of choice and monitored validation and training errors. Due to memory limitations on the GPU we could only experiment with grids of  $32 \times 32 \times 32$  voxels.



**Figure 2.30:** Evolution of training and validation accuracy of the 3D CNN using adaptive binary grids with size  $32 \times 32 \times 32$ .

Figure 2.30 shows the results of this experiment. As we can observe, we trained the network for five times more iterations than before and even then we couldn't achieve a proper convergence. The training set accuracy kept increasing slowly up to approximately 0.65 whilst the validation one got stuck around 0.40 for the whole experiment.

In conclusion, porting the 2.5D network directly to 3D just by extending its convolution and pooling layers to slide along the depth axis did not produce good results using the same dataset and setup that produced a significantly good outcome with the 2.5D architecture. We hypothesize various causes for this problem.

On the one hand, the data representation might not be adequate for such fine-grained convolutions. It is presumable that bigger grids, e.g.,  $64 \times 64 \times 64$ , would yield better results. However, given the size of the model, they could not be tested in the available GPU.

On the other hand, the complexity of the network increased considerably after including that extra dimension in convolution and pooling layers. This means that the number of parameters of the network gets increased significantly, making it harder to train with so few samples due to overfitting. This hypothesis is backed up by the fact that training accuracy kept increasing slowly while validation one got stuck. This would eventually lead to a perfect fit on the training set but low accuracy on the validation split.

**Discussion** To sum up, we determined that the adaptive grid slightly outperforms the fixed one in normal conditions. The same happens with the grid size, obtaining marginally better results with bigger sizes. However, when it comes down to noise and occlusion robustness, the adaptive grid exceeds the accuracy of the fixed grid by a large margin for low levels of occlusion and noise, whilst for intermediate and high levels the impact on both grids is somewhat similar. In other words, the adaptive grid is better than the fixed one and it is preferable to use a bigger grid size if the performance impact can be afforded.

It is important to remark that the binary occupancy measure performed better than the normalized density one, both using adaptive grids, while maintaining similar resilience against noise and occlusions. The best network trained with normalized density grids reached a peak accuracy of approximately 0.79 while the best binary one achieved approximately a 0.85 accuracy on the validation set.

Another remarkable fact was that all networks exhibited a considerable amount of overfitting, i.e., training accuracy was almost perfect whilst validation was far away from it by a considerable margin. We hypothesize that this was due to the fact that the dataset has few training examples considering the complexity of the network. Besides, we also inspected the confusion matrix shown in Table ?? to gain insight about the behavior of our network. As we can observe, there are many misclassified samples of classes that are similar. If we take a closer look at some of the misclassified samples (see Figures ??, ??, and ??) it is reasonable to think that the network is not able to classify them properly because they are extremely similar. In this regard, the dataset must be augmented introducing noise, translations, rotations, and variations of the models to avoid overfitting and learn better those models that can be easily misclassified.

## 2.5.4 Conclusion

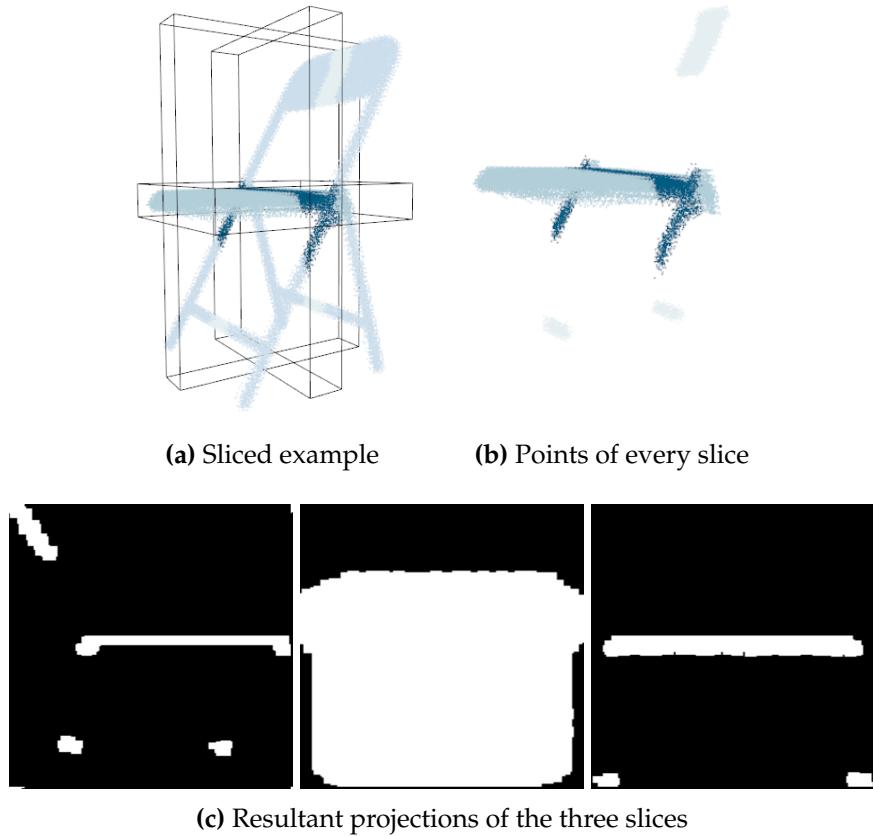
## 2.6 LonchaNet

In this work, we present a novel approach that uses multiple 2D views from 3D models for 3D object recognition. The proposed 2D renderings are based on cross sections of 3D models. The proposed method outperforms most existing approaches that participated in the ModelNet challenge. Currently, we are in the second position of the leaderboard, obtaining 94.37% classification accuracy. Our proposal is focused on learning a discriminative representation that is able to distinguish among most categories of the ModelNet dataset. We also contribute with a new architecture that uses the existing GoogLeNet network **Szegedy2015**. We use three independent GoogLeNet networks for learning specific features of each cross-section or slice from the 3D model.

As stated before, we propose a 3D object recognition method using deep learning. First, for every example in the dataset, we took three sections from an object, one for each 3D axis, and project the 3D points to a plane, so we obtain three images from every example. Each of these three images that shape a single example are then fed to a deep Convolutional Neural Network. Our novel deep architecture features three GoogLeNets, one for each image, joined in a layer prior to the classification layer. This provides us with great expressiveness and a high success rate.

### 2.6.1 Data Representation

LonchaNet takes point clouds as input, but the neural architecture itself uses three images corresponding to three slices, so first we need to extract the slices from the 3D point cloud. To that end, we load the point clouds and calculate the central point for each axis. Then we make a slice across the XY, XZ, and YZ planes with a thickness of 5% of the model size. Those points that fall inside these sections are isolated and projected in their planes to generate an image of  $500 \times 500$  pixels. These images are binary maps in which the background is black and the projected points are white. This process is shown in Figure 2.31.



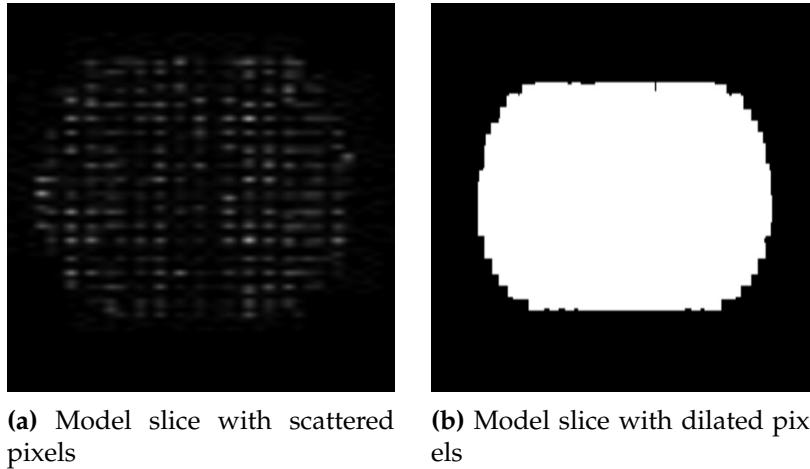
**Figure 2.31:** Extracting the slices from a point cloud example [MISSINGDETAILS]

Due to the inconsistent point density produced by the sampling process, there are some slices in which the points are very scattered, so the projection does not represent the object faithfully. To deal with this problem, a post-process is carried out for each projection in which we apply dilations of 10 pixels using a square as structuring element. This post-processing step fills the gap between the scattered points and produces a more suitable representation of the object, as shown in Figure 2.32. This process is performed for each example present in the dataset, so for each point cloud there are three corresponding images, one per slice.

This sliced representation of the object allows us to train and test a 3D recognition system in a 2D fashion, that provides the high success rate and the fast execution time that characterizes deep image recognition networks yet preserving and taking advantage of the 3D information implicitly materialized in the slicing method.

## 2.6.2 Network Architecture

The main architecture of LonchaNet is composed of three isolated GoogLeNet that are joined at the bottom in a Concatenation Layer prior a Fully Connected Layer which is the final classifier, as shown in Figure ??.



**Figure 2.32:** A comparison of a slice before and after the dilation process. The dilated image provides a more accurate representation the object.

As stated before, GoogLeNet is the state-of-the-art deep Convolutional Neural Network for image recognition tasks, providing the highest accuracy in several challenges, so we chose it over the other topologies.

In this architecture, all convolutions, including those inside the inception modules, use **ReLU** activation. The size of the receptive field in this network is  $224 \times 224$  taking RGB channels with mean subtraction, although in LonchaNet ensemble we use binary maps with no mean normalization. GoogLeNet network consists on 22 layers deep when considering only layers with parameters (or 27 layers if we also consider pooling ones). The overall number of layers (independent building blocks) used for the construction of the network is about 100. However, this number depends on the machine learning infrastructure system used. The use of average pooling step prior the classifier is based on [Lin2013](#), although this implementation differs in the use of an extra linear layer. This enables adapting and fine-tuning the network for other datasets.

LonchaNet features three independent GoogLeNet (we will reference them as "branches") learning features that define an object for each slice, so we force each branch to specialize the filters in particular features of every slice. Finally, the responses of each branch are concatenated in a single output that is fed to a fully connected layer that acts as a classifier.

## 2.6.3 Experiments

### Data Generation

### Methodology and Setup

All timings and results were obtained by conducting the experiments in the following test setup: Intel Core i7-5820K with 32 GiB of Kingston HyperX 2666 MHz and CL13 DDR4 RAM on an Asus X99-A motherboard (Intel X99 chipset). Secondary storage was provided by a Samsung 850 EVO SSD. Additionally, the system included two NVIDIA Tesla K40c **GPUs** used for training and inference.

The framework of choice was Caffe RC2 running on Ubuntu 14.04.02. It was compiled using CMake 2.8.7, g++ 4.8.2, CUDA 7.5, and cuDNN v3.

## Results and Discussion

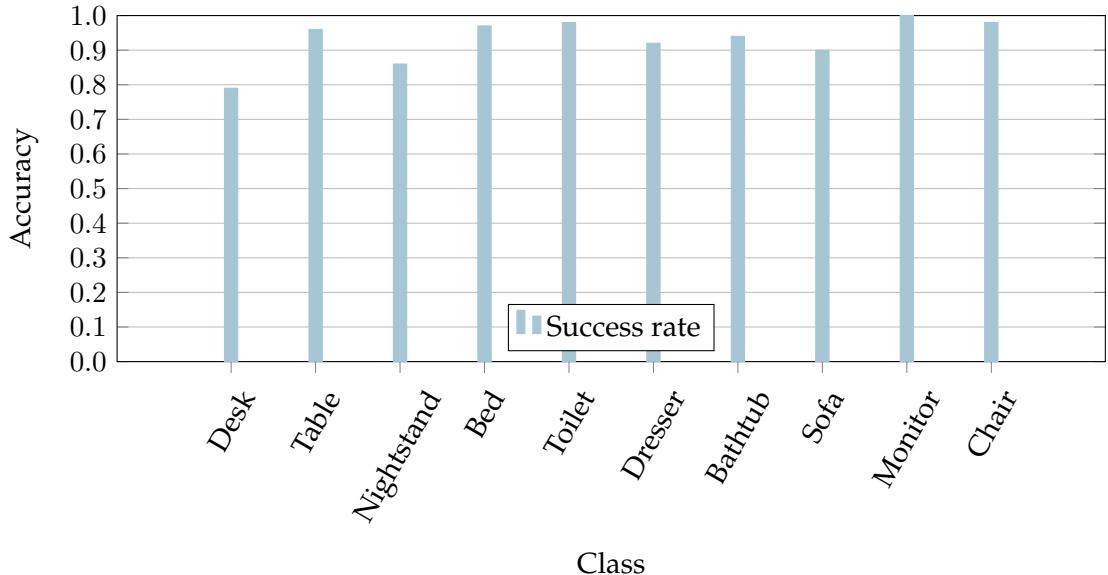
We trained from scratch and tested LonchaNet with the ModelNet-10 dataset, as described earlier in the subsection ???. It is remarkable that no transfer learning or fine tuning over a pre-trained set of parameters was used. It is remarkable that the training and test splits are defined by the dataset itself and also it is worth saying that the number of examples per class are not balanced, as seen in Table ???. This fact could harm the accuracy of the system, biasing the learning and the classification to the classes with more number of examples, as stated by **He2009**.

Regarding the parameters that affects to the learning process, we trained the architecture with a base learning rate of 0.00001, multiplying the current learning rate by 0.75 every 10000 iterations. To compute the weights update we use the ADAM **KingmaB14** solver with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . The training process executed for 20000 iterations, but the best weights set was produced on the iteration #18300 which threw a test accuracy of 94.3709%, the second best score in the leaderboard of the ModelNet10 challenge.

It is also worth noticing the fast computations of our architecture. One training iteration with a batch size of 30 examples takes an average of 2.25 seconds. Also, classify a new examples only takes 0.0896 seconds.

Figure 2.33 shows the success rate per class. It can be seen that the classes with the less success rates are the very same classes with the less number of examples, namely the classes desk and nightstand. In light of this fact we can expect that if we could have a balanced dataset, the error rate of these classes would decrease significantly.

In Figure ???, we can observe the confusion matrix of the classification accuracy for the test split of ModelNet-10. This table, alongside the precision-recall curves presented in Figure ???, confirms the stability and reliability of the system, which confuses only objects of classes that heavily look alike from a visual perspective.



**Figure 2.33:** Success rate per class for the test split of the ModelNet-10 dataset achieved by LonchaNet after 18300 training iterations using the *ModelNet-10* dataset (solver type is ADAM, learning rate is 0.00001,  $\beta_1$  is 0.9 and  $\beta_2$  is 0.999).

That is the case of the classes desk and table which the system confuses, yet not the other way around. This could be possibly caused by the fact that a desk is a type of table, meaning the desks have minor visual features that go unnoticed in some examples, making those examples hard to distinguish from the table ones, as shown in Figures ?? and ?? . In addition, in the Table ?? , it could be seen that the class desk have a reduced number of examples compared with other classes such as sofa or chair.

Also, there is a confusion problem among the classes nightstand and dresser. This kind of confusion is very common in CNN architectures because, again, they rely their classification capabilities on visual features of an object and, as seen in Figures ?? and ?? , these two classes are visually similar.

Figure ?? remarks the ambiguity of the visual features between the classes desk and table, and nightstand and dresser and point up the difficulty of the problem.

Using our current architecture, we reached the second place in the leaderboard of the challenge, with an accuracy of 94.37% in the ModelNet-10 Accuracy task. The state of the leaderboard (as of January 2017) is shown in Table ?? .

#### 2.6.4 Conclusion

In this last iteration, we introduced a novel architecture for 3D object recognition, LonchaNet. Our system takes three slices of the input point cloud (one per 3D axis), projects the points to a plane, generating three images, and sends them to the neural network. The architecture consist in three independent GoogLeNet branches which activations are concatenated and fed to a fully connected layer. Each of this branches learns particular features of a slice.

This method allows us to take advantage of the fast 2D computation whilst preserving the 3D information. LonchaNet posits as the second place in the ModelNet challenge with a succes rate of 94.37% in the ModelNet-10 accuracy test, yet providing a extremely fast computation: once the model is trained, classifying a 3D object only takes 0.0896 seconds.

### 2.7 Conclusion

Method	ModelNet10 Accuracy	ModelNet40 Accuracy
VRN Ensemble [53]	97.14%	95.54%
<b>LonchaNet</b>	<b>94.37%</b>	N/A
ORION [54]	93.80%	N/A
FusionNet [55]	93.11%	90.80%
Pairwise [56]	92.80%	90.70%
GIFT [57]	92.35%	83.10%
VoxNet [44]	92.00%	83.00%
3D-GAN [58]	91.00%	83.30%
DeepPano [59]	85.45%	77.63%
3DShapeNets [42]	83.50%	77.00%
MVCNN [60]	N/A	90.10%

**Table 2.2:** ModelNet leaderboard as of January, 2017.

Chapter **3**

# Semantic Segmentation

## **3.1 Introduction**

## **3.2 Related Works**

## **3.3 The RobotriX**

## **3.4 UnrealROX**

## **3.5 2D-3D-SeGCN**



Chapter **4**

# Tactile Sensing

## **4.1 Introduction**

## **4.2 Related Works**

## **4.3 TactileGCN**

## **4.4 Conclusion**



Chapter **5**

## Conclusion

### **5.1 Findings and Conclusions**

### **5.2 Limitations**

### **5.3 Future Work**



# Bibliography

- [1] Alberto Garcia-Garcia, Francisco Gomez-Donoso, Jose Garcia-Rodriguez, et al. "PointNet: A 3D Convolutional Neural Network for real-time object class recognition". In: *2016 International Joint Conference on Neural Networks, IJCNN 2016, Vancouver, BC, Canada, July 24-29, 2016*. 2016, pp. 1578–1584. DOI: [10.1109/IJCNN.2016.7727386](https://doi.org/10.1109/IJCNN.2016.7727386). URL: <https://doi.org/10.1109/IJCNN.2016.7727386>.
- [2] Alberto Garcia-Garcia, Jose Garcia-Rodriguez, Sergio Orts-Escalano, et al. "A study of the effect of noise and occlusion on the accuracy of convolutional neural networks applied to 3D object recognition". In: *Computer Vision and Image Understanding* 164 (2017), pp. 124–134. DOI: [10.1016/j.cviu.2017.06.006](https://doi.org/10.1016/j.cviu.2017.06.006). URL: <https://doi.org/10.1016/j.cviu.2017.06.006>.
- [3] Francisco Gomez-Donoso, Alberto Garcia-Garcia, Jose Garcia-Rodriguez, et al. "LonchaNet: A Sliced-based CNN Architecture for Real-time 3D Object Recognition". In: *2017 International Joint Conference on Neural Networks, IJCNN 2017, Anchorage, Alaska, May 14-19, 2017*. 2017. URL: <https://ieeexplore.ieee.org/document/7965883/>.
- [4] Alberto Garcia-Garcia, Pablo Martinez-Gonzalez, Sergiu Oprea, et al. "The RobotriX: An eXtremely Photorealistic and Very-Large-Scale Indoor Dataset of Sequences with Robot Trajectories and Interactions". In: *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE. 2018, pp. 6790–6797. URL: <https://ieeexplore.ieee.org/abstract/document/8594495>.
- [5] Sergiu Oprea, Alberto Garcia-Garcia, Jose Garcia-Rodriguez, et al. "A Recurrent Neural Network based Schaeffer Gesture Recognition System". In: *2017 International Joint Conference on Neural Networks, IJCNN 2017, Anchorage, Alaska, May 14-19, 2017*. 2017. URL: <https://ieeexplore.ieee.org/document/7965885/>.
- [6] Francisco Gomez-Donoso, Sergio Orts-Escalano, Alberto Garcia-Garcia, et al. "A robotic platform for customized and interactive rehabilitation of persons with disabilities". In: *Pattern Recognition Letters* 99 (2017), pp. 105–113. DOI: [10.1016/j.patrec.2017.05.027](https://doi.org/10.1016/j.patrec.2017.05.027). URL: <https://doi.org/10.1016/j.patrec.2017.05.027>.
- [7] Sergiu Oprea, Alberto GarciaGarcia, Sergio OrtsEscolano, et al. "A long short-term memory based Schaeffer gesture recognition system". In: *Expert Systems* 0.0 (2017), e12247. DOI: [10.1111/exsy.12247](https://doi.org/10.1111/exsy.12247). URL: <https://onlinelibrary.wiley.com/doi/abs/10.1111/exsy.12247>.

- [8] Alberto Garcia Garcia, Andreas Beckmann, and Ivo Kabadshow. "Accelerating an FMM-Based Coulomb Solver with GPUs". In: *Software for Exascale Computing-SPPEXA 2013-2015*. Springer, 2016, pp. 485–504. URL: [https://link.springer.com/chapter/10.1007/978-3-319-40528-5\\_22](https://link.springer.com/chapter/10.1007/978-3-319-40528-5_22).
- [9] Alberto Garcia-Garcia, Sergio Orts-Escalano, Sergiu Oprea, et al. "Multi-sensor 3D object dataset for object recognition with full pose estimation". In: *Neural Computing and Applications* 28 (2016), pp. 941–952. ISSN: 1433-3058. DOI: [10.1007/s00521-016-2224-9](https://doi.org/10.1007/s00521-016-2224-9). URL: <http://dx.doi.org/10.1007/s00521-016-2224-9>.
- [10] Marcelo Saval-Calvo, Jorge Azorin-Lopez, Andres Fuster-Guillo, et al. "Evaluation of sampling method effects in 3D non-rigid registration". In: *Neural Computing and Applications* 28 (2016), pp. 953–967. ISSN: 1433-3058. DOI: [10.1007/s00521-016-2258-z](https://doi.org/10.1007/s00521-016-2258-z). URL: <http://dx.doi.org/10.1007/s00521-016-2258-z>.
- [11] Sergio Orts-Escalano, Jose Garcia-Rodriguez, Miguel Cazorla, et al. "Bioinspired point cloud representation: 3D object tracking". In: *Neural Computing and Applications* 29 (2016), pp. 663–672. ISSN: 1433-3058. DOI: [10.1007/s00521-016-2585-0](https://doi.org/10.1007/s00521-016-2585-0). URL: <https://doi.org/10.1007/s00521-016-2585-0>.
- [12] Alberto Garcia-Garcia, Sergio Orts-Escalano, Jose Garcia-Rodriguez, et al. "Interactive 3D object recognition pipeline on mobile GPGPU computing platforms using low-cost RGB-D sensors". In: *Journal of Real-Time Image Processing* 14 (2016), pp. 585–604. ISSN: 1861-8219. DOI: [10.1007/s11554-016-0607-x](https://doi.org/10.1007/s11554-016-0607-x). URL: <https://doi.org/10.1007/s11554-016-0607-x>.
- [13] Higinio Mora, Jerónimo M Mora-Pascual, Alberto Garcia-Garcia, et al. "Computational analysis of distance operators for the iterative closest point algorithm". In: *PloS one* 11.10 (2016), e0164694. URL: <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0164694>.
- [14] Sergio Orts-Escalano, Jose Garcia-Rodriguez, Vicente Morell, et al. "3D Surface Reconstruction of Noisy Point Clouds Using Growing Neural Gas: 3D Object/Scene Reconstruction". In: *Neural Processing Letters* 43 (2015), pp. 401–423. DOI: [10.1007/s11063-015-9421-x](https://doi.org/10.1007/s11063-015-9421-x). URL: <http://dx.doi.org/10.1007/s11063-015-9421-x>.
- [15] Sergio Orts-Escalano, Jose Garcia-Rodriguez, Jose Antonio Serra-Perez, et al. "3D model reconstruction using neural gas accelerated on GPU". In: *Applied Soft Computing* 32 (2014), pp. 87–100. DOI: [10.1016/j.asoc.2015.03.042](https://doi.org/10.1016/j.asoc.2015.03.042). URL: <http://dx.doi.org/10.1016/j.asoc.2015.03.042>.
- [16] Alexander Andreopoulos and John K. Tsotsos. "50 Years of object recognition: Directions forward". In: *Computer Vision and Image Understanding* 117.8 (2013), pp. 827–891.
- [17] David G. Lowe. "Distinctive image features from scale-invariant keypoints". In: *International journal of computer vision* 60.2 (2004), pp. 91–110.
- [18] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. "Surf: Speeded up robust features". In: *Computer vision-ECCV 2006*. Springer, 2006, pp. 404–417.
- [19] Michael Calonder, Vincent Lepetit, Christoph Strecha, et al. "Brief: Binary robust independent elementary features". In: *Computer Vision-ECCV 2010* (2010), pp. 778–792.

- [20] Stefan Leutenegger, Margarita Chli, and Roland Y. Siegwart. "BRISK: Binary robust invariant scalable keypoints". In: *Computer Vision (ICCV), 2011 IEEE International Conference on*. IEEE. 2011, pp. 2548–2555.
- [21] Ethan Rublee, Vincent Rabaud, Kurt Konolige, et al. "ORB: an efficient alternative to SIFT or SURF". In: *Computer Vision (ICCV), 2011 IEEE International Conference on*. IEEE. 2011, pp. 2564–2571.
- [22] Alexandre Alahi, Raphael Ortiz, and Pierre Vandergheynst. "FREAK: Fast retina keypoint". In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. Ieee. 2012, pp. 510–517.
- [23] David G. Lowe. "Object recognition from local scale-invariant features". In: *Computer vision, 1999. The proceedings of the seventh IEEE international conference on*. Vol. 2. Ieee. 1999, pp. 1150–1157.
- [24] Yulan Guo, Mohammed Bennamoun, Ferdous Sohel, et al. "3D object recognition in cluttered scenes with local surface features: A survey". In: *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 36.11 (2014), pp. 2270–2287.
- [25] Jean Ponce, Svetlana Lazebnik, Fredrick Rothganger, et al. "Toward true 3D object recognition". In: *Reconnaissance de Formes et Intelligence Artificielle*. 2004.
- [26] Paul J. Besl and Ramesh C. Jain. "Three-dimensional object recognition". In: *ACM Computing Surveys (CSUR)* 17.1 (1985), pp. 75–145.
- [27] Jim P. Brady, Nagaraj Nandhakumar, and Jake K. Aggarwal. "Recent progress in object recognition from range data". In: *image and vision computing* 7.4 (1989), pp. 295–307.
- [28] Farshid Arman and Jake K. Aggarwal. "Model-based object recognition in dense-range imagesa review". In: *ACM Computing Surveys (CSUR)* 25.1 (1993), pp. 5–43.
- [29] Richard J. Campbell and Patrick J. Flynn. "A survey of free-form object representation and recognition techniques". In: *Computer Vision and Image Understanding* 81.2 (2001), pp. 166–210.
- [30] George Mamic and Mohammed Bennamoun. "Representation and recognition of 3D free-form objects". In: *Digital Signal Processing* 12.1 (2002), pp. 47–76.
- [31] Yann Le Cun, Yoshua Bengio, and Geoffrey E. Hinton. "Deep learning". In: *Nature* 521.7553 (2015), pp. 436–444.
- [32] Paul J. Werbos. "Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences". PhD thesis. Harvard University, 1974.
- [33] Yann Le Cun. "A learning scheme for asymmetric threshold networks". In: *Proceedings of Cognitiva 85* (1985), pp. 599–604.
- [34] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. "Learning representations by back-propagating errors". In: *Cognitive modeling* 5 (1988), p. 3.
- [35] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks". In: *Advances in neural information processing systems*. 2012, pp. 1097–1105.
- [36] Stjepan Marčelja. "Mathematical description of the responses of simple cortical cells". In: *JOSA* 70.11 (1980), pp. 1297–1300.
- [37] Alex Berg, Jia Deng, and Fei-Fei Li. *ImageNet large scale visual recognition challenge 2010*. 2010.

- [38] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. "A fast learning algorithm for deep belief nets". In: *Neural computation* 18.7 (2006), pp. 1527–1554.
- [39] Yoshua Bengio, Pascal Lamblin, Dan Popovici, et al. "Greedy layer-wise training of deep networks". In: *Advances in neural information processing systems* 19 (2007), p. 153.
- [40] Pierre Sermanet, Koray Kavukcuoglu, Sandhya Chintala, et al. "Pedestrian detection with unsupervised multi-stage feature learning". In: *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*. IEEE. 2013, pp. 3626–3633.
- [41] Rajat Raina, Anand Madhavan, and Andrew Y. Ng. "Large-scale deep unsupervised learning using graphics processors". In: *Proceedings of the 26th annual international conference on machine learning*. ACM. 2009, pp. 873–880.
- [42] Zhirong Wu, Shuran Song, Aditya Khosla, et al. "3D ShapeNets: A Deep Representation for Volumetric Shapes". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2015. URL: <http://arxiv.org/abs/1406.5670>.
- [43] Shuran Song and Jianxiong Xiao. "Deep Sliding Shapes for Amodal 3D Object Detection in RGB-D Images". In: *CoRR* abs/1511.02300 (2015). URL: <http://arxiv.org/abs/1511.02300>.
- [44] Daniel Maturana and Sebastian Scherer. "Voxnet: A 3d convolutional neural network for real-time object recognition". In: *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*. IEEE. 2015, pp. 922–928.
- [45] Kevin Lai, Liefeng Bo, Xiaofeng Ren, et al. "A large-scale hierarchical multi-view rgb-d object dataset". In: *Robotics and Automation (ICRA), 2011 IEEE International Conference on*. IEEE. 2011, pp. 1817–1824.
- [46] Ashutosh Singh, Jin Sha, Karthik S. Narayan, et al. "Bigbird: A large-scale 3d database of object instances". In: *Robotics and Automation (ICRA), 2014 IEEE International Conference on*. IEEE. 2014, pp. 509–516.
- [47] Bo Li, Yijuan Lu, Chunyuan Li, et al. "Shrec14 track: extended large scale sketch-based 3D shape retrieval". In: *Eurographics Workshop on 3D Object Retrieval*. 2014, pp. 121–130.
- [48] Sungjoon Choi, Qian-Yi Zhou, Stephen Miller, et al. "A Large Dataset of Object Scans". In: *CoRR* abs/1602.02481 (2016).
- [49] Radu B. Rusu. "Point Cloud Library". In: *IEEE International Conference on Robotics and Automation (ICRA)*. 2011, pp. 1–4.
- [50] Aitor Aldoma, Zoltan-Csaba Marton, Federico Tombari, et al. "Tutorial: Point cloud library: Three-dimensional object recognition and 6 dof pose estimation". In: *IEEE Robotics & Automation Magazine* 19.3 (2012), pp. 80–91.
- [51] Yangqing Jia, Evan Shelhamer, Jeff Donahue, et al. "Caffe: Convolutional Architecture for Fast Feature Embedding". In: *Proceedings of the ACM International Conference on Multimedia*. 2014, pp. 657–678.
- [52] Sharan Chetlur, Cliff Woolley, Philippe Vandermersch, et al. "cuDNN: Efficient Primitives for Deep Learning". In: (2014), pp. 1–9. ISSN: 08876266. DOI: [10.1002/polb.23894](https://doi.org/10.1002/polb.23894). eprint: [1410.0759](https://arxiv.org/abs/1410.0759).
- [53] Andrew Brock, Theodore Lim, JM Ritchie, et al. "Generative and Discriminative Voxel Modeling with Convolutional Neural Networks". In: *arXiv preprint arXiv:1608.04236* (2016).

- [54] Nima Sedaghat, Mohammadreza Zolfaghari, and Thomas Brox. "Orientation-boosted Voxel Nets for 3D Object Recognition". In: *arXiv preprint arXiv:1604.03351* (2016).
- [55] Vishakh Hegde and Reza Zadeh. "FusionNet: 3D Object Classification Using Multiple Data Representations". In: *CoRR* abs/1607.05695 (2016). URL: <http://arxiv.org/abs/1607.05695>.
- [56] Edward Johns, Stefan Leutenegger, and Andrew J Davison. "Pairwise Decomposition of Image Sequences for Active Multi-View Recognition". In: *arXiv preprint arXiv:1605.08359* (2016).
- [57] Song Bai, Xiang Bai, Zhichao Zhou, et al. "Gift: A real-time and scalable 3d shape search engine". In: *arXiv preprint arXiv:1604.01879* (2016).
- [58] Jiajun Wu, Chengkai Zhang, Tianfan Xue, et al. "Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling". In: *arXiv preprint arXiv:1610.07584* (2016).
- [59] Baoguang Shi, Song Bai, Zhichao Zhou, et al. "Deeppano: Deep panoramic representation for 3-d shape recognition". In: *IEEE Signal Processing Letters* 22.12 (2015), pp. 2339–2343.
- [60] Hang Su, Subhransu Maji, Evangelos Kalogerakis, et al. "Multi-view convolutional neural networks for 3d shape recognition". In: *Proceedings of the IEEE International Conference on Computer Vision*. 2015, pp. 945–953.