Orientation-boosted Voxel Nets for 3D Object Recognition (Supplementary Material)

BMVC 2017 Submission # 100

11 Auto-Alignment of the Modelnet40 dataset

Modelnet40 [D] consists of more than 12000 *non-aligned* objects in 40 classes. We used the method of Sedaghat & Brox [D] to automatically align the objects class by class.

016

007

008

011

Mesh to Point-Cloud Conversion The auto-alignment method of [2] uses point-cloud
 representations of objects as input. Thus we converted the 3D mesh grids of Modelnet40 to
 point-clouds by assigning uniformly distributed points to object faces.

Hidden faces in the mesh grids needed to be removed, as the so called Hierarchical
Orientation Histogram (HOH) of [I] mainly relies on the exterior surfaces of the objects.
We tackled this issue using the Jacobson's implementation [I] of the "ambient occlusion"
method [I].

We tried to distribute the points roughly with the same density across different faces, regardless of their shape and size, to avoid a bias towards bigger/wider ones. Our basic point-clouds consist of around 50000 points per object, which are then converted to lighter models using the Smooth Signed Distance surface reconstruction method (SSD) [I] as used in [I].

029

Auto-Alignment We first created a "reference-set" in each class, consisting of a random
subset of its objects, with an initial size of 100. This number was then decreased, as the lowquality objects were automatically removed from the reference set, according to [□]. This
reference set was then used to align the remaining objects of the class one by one.

For the HOH descriptor, we used 32 and 8 divisions in ϕ and θ dimensions respectively, for the root component. We also used 8 child components with 16 divisions for ϕ and 4 for θ - see [1].

037

Automatic Assignment of Number of Orientation Classes As pointed out in the main
paper, we do not use the same number of orientation classes for all the object categories.
We implemented the auto-alignment procedure in a way that this parameter is automatically
decided upon for each category: During generation of the reference-set in each class, the
alignment procedure was run with 3 different configurations, for which the search space
spanned over 360, 180 and 90 degrees of rotations respectively. Each run resulted in an error

^{© 2017.} The copyright of this document resides with its authors.

¹⁴⁵ It may be distributed unchanged freely in print or electronic forms.

measure representing the overall quality of the models selected as the reference-set, and we 046 designated respectively 12, 6 and 3 orientation levels to each category, whenever possible. 047 When none of these worked, e.g. for the 'flower_pot' class, we assigned 1 orientation class 048 which is equivalent to discarding the orientation information. 049

2 Analysis

To analyze the behavior of the orientation-boosted network, we compare it to its corresponding baseline network. We would like to know the differences between corresponding filters in the two networks. To find this correspondence, we first train a baseline network, without orientations outputs, for long enough so that it reaches a stable state. Then we use this trained net to initialize the weights of the ORION network, and continue training with a low learning rate. This way we can monitor how the learned features change in the transition from the baseline to the orientation-aware network.

In Figure 1 transition of a single exemplar filter is depicted, and its responses to different 060 rotations of an input object are illustrated. It turns out that the filter tends to become more 061 sensitive to the orientation-specific features of the input object. Additionally some parts of 062 the object, such as the table legs, show stronger response to the filter in the orientation-aware 063 network.

With such an observation, we tried to analyze the overall behavior of the network for 065 specific object classes with different orientations. To this end we introduce the "dominant 066 signal-flow path" of the network. The idea is that, although all the nodes and connections 067 of the network contribute to the formation of the output, in some cases there may exist a 068 set of nodes, which have a significantly higher effect in this process for an specific type of 069 object/orientation. To test this, we take this step-by-step approach: First in a forward pass, 070 the class, c, of the object is found. Then we seek to find the highest contributing node of the 071 last hidden layer:

$$l^{n-1} = \operatorname*{arg\,max}_{k} \{ w^{n-1}_{k,c} a^{n-1}_{k} \} \tag{1} \begin{array}{c} 073\\074 \end{array}$$

where n is the number of layers, a_k^{n-1} are the activations of layer n-1, and $w_{k,c}^{n-1}$ is the 076 weight connecting a_k^{n-1} to the c^{th} node of layer n. This way we naively assume there is a $\frac{1}{0.077}$ significant maximum in the *contributions* and assign its index to l^{n-1} . Later we will see that 078 this assumption proves to be true in many of our observations. We continue "back-tracing" 079 the signal, to the previous layers. Extension of (1) to the convolutional layers is straightforward, as we are just interested in finding the index of the node/filter in each layer. In 081 the end, letting $l^n = c$, gives us the vector l with length equal to the number of network layers, keeping the best contributors' indices in it. Now to depict the "dominant signal-flow path" for a group of objects, we simply obtain l for every member of the group, and plot 084 the histogram of the l^i s as a column. Figure 2(a) shows such an illustration for a specific class-rotation of the objects. It is clearly visible that for many objects of that group, specific nodes have been dominant.

In Figure 2(b), the dominant paths of the baseline and ORION networks for some sample ⁰⁸⁷ object categories of the Modelnet10 dataset are illustrated. It can be seen that in the baseline ⁰⁸⁸ network, the dominant paths among various rotations of a class mostly share a specific set ⁰⁸⁹ of nodes. This is mostly visible in the convolutional layers – e.g. see the red boxes. On the ⁰⁹⁰ contrary, the dominant paths in the ORION network rarely follow this rule and have more ⁰⁹¹



- higher values \longrightarrow

127 Figure 1: The picture illustrates the activations of one of the nodes of the first layer, while the 128 network transitions from a baseline network to ORION. The input is always the same object, 129 which has been fed to the network in its possible discretized rotations (columns) at each step (row). We simulated this transition by first training the baseline network and then fine-tuning 131 our orientation-aware architecture on top of the learned weights. To be able to depict the 3D 132 feature maps, we had to cut out values below a specific threshold. It can be seen that the encoded filter detects more orientation-specific aspects of the object, as it moves forward in 134 learning the orientations. In addition, it seems that the filter is becoming more sensitive to a *table* rather than only a horizontal surface – notice the table legs appearing in the below 136 rows.

4 AUTHOR(S): ORIENTATION-BOOSTED VOXEL NETS FOR 3D OBJECT RECOGNITION



Figure 2: (a) shows the "dominant signal-flow path" of the network, for an exemplar object 175 category-orientation. Each column contains the activations of one layer's nodes. Obviously 176 the columns are of different sizes. Higher intensities show dominant nodes for the specific 177 group of objects. Details of the steps taken to form such an illustration are explained in the 178 text. In (b), rows represent object classes, while in different columns we show rotations of 179 the objects. So each cell is a specific rotation of a specific object category. It can be seen that 180 in the baseline network, many of the rotations of a class, share nodes in their dominant path 181 (e.g. see the red boxes), whereas, in the ORION network the paths are more distributed over 182 all the nodes.

distributed path nodes. We interpret this as one of the results of orientation-boosting, and ahelping factor in better classification abilities of the network.

Extended Architecture

	Conv1	Conv2	Conv3	Conv4	Pool4
# of filters	32	64	128	256	-
kernel size	3x3x3	3x3x3	3x3x3	3x3x3	2x2x2
stride	2	1	1	1	2
padding	0	0	0	0	-
dropout ratio	0.2	0.3	0.4	0.6	-
batch normalization	\checkmark	\checkmark	\checkmark	\checkmark	-
	fc1	fc2:class	fc2:orientation		
# of outputs	128	10/40	variable [†]		
dropout ratio	0.4	-	-		
batch normalization	×	×	×		

Table 1: Details of the extended architecture introduced in Tables 1 & 2 of the main article. [†] The number of nodes dedicated to the orientation output varies in different experiments.

4 Orientation Estimation Results

Although the orientation estimation was used merely as an auxiliary task, here in Table 2 we report the accuracies of the estimated orientation classes. Note that getting better results on orientation estimation would be possible by emphasizing on this task – e.g. see the detection experiment in the main article.

	Accuracy %				
	Sydney	NYUv2	ModelNet10	ModelNet40	
ORION	71.5	51.9	89.0	86.5	
ORION - Extended	70.1	54.5	89.3	87.6	

Table 2: Orientation estimation accuracies on different datasets. The extended architecture of the second row, is the one introduced in the main article and detailed in Table 1 of this document.

- 22:

6 AUTHOR(S): ORIENTATION-BOOSTED VOXEL NETS FOR 3D OBJECT RECOGNITION

References	230
[1] Fatih Calakli and Gabriel Taubin. Ssd: Smooth signed distance surface reconstruction <i>Computer Graphics Forum</i> , volume 30, pages 1993–2002. Wiley Online Library, 20	231 1. In 232 11. 233
[2] Alec Jacobson et al. gptoolbox: Geometry processing toolbox, 20 http://github.com/alecjacobson/gptoolbox.)16. 234 235 236
 [3] Gavin Miller. Efficient Algorithms for Local and Global Accessibility Shading. In <i>I ceedings of the 21st Annual Conference on Computer Graphics and Interactive T niques</i>, SIGGRAPH '94, pages 319–326, New York, NY, USA, 1994. ACM. IS 978-0-89791-667-7. doi: 10.1145/192161.192244. URL http://doi.acm.or/10.1145/192161.192244. 	Pro- 237 ech- 238 BN 239 eg/ 240 241
[4] Nima Sedaghat and Thomas Brox. Unsupervised Generation of a Viewpoint Annot Car Dataset from Videos. In <i>Proceedings of the IEEE International Conference on C</i> <i>puter Vision (ICCV)</i> , 2015.	ated 242 om- 243 244 245
[5] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Ta and Jianxiong Xiao. 3D ShapeNets: A Deep Representation for Volumetric Sha In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognit</i>	ang, 246 pes. 247 <i>ion</i> , 248
pages 1912–1920, 2015.	249 250 251
	252 253 254
	255 256 257
	258 259 260
	261 262 263
	264 265 266
	267 268
	269 270 271
	272 273 274