

# Robust Category-Level 6D Pose Estimation with Feature-level Render and Compare (Supplementary Materials)

Wufei Ma<sup>1</sup>, Angtian Wang<sup>1</sup>, Alan Yuille<sup>1</sup>, and Adam Kortylewski<sup>1,2,3</sup>

<sup>1</sup> Johns Hopkins University, Baltimore MD 21218, USA  
{wma27, angtianwang, ayuille1, akortyl1}@jhu.edu

<sup>2</sup> Max Planck Institute for Informatics, Saarbrücken, Germany

<sup>3</sup> University of Freiburg, Germany

## 1 Implementation of Generative 6D Proposals

With the scale-invariant contrastive features, the pose optimization has a clear global minimum near the ground truth location. The loss landscapes are smooth, with decent gradients around the global minimum. This nice property allows us to search for generative 6D proposals from a sparse sampling over six dimensions.

Our generative 6D proposals are predicted in two steps. In the first step we locate the object principal points based on shape-agnostic reconstruction losses. We relax the pose optimization problem by ignoring the interior structure of the 3D model  $\mathcal{M}$  and predict the object principal points by maximizing

$$\max_m \sum_{1 \leq i \leq h, 1 \leq j \leq w} \max_{1 \leq k \leq N} \|\mathbf{F}_{i,j} - C_k\|^2 \quad (1)$$

We solve this optimization problem efficiently by computing and caching the key-point correlation scores  $\|\mathbf{F}_{i,j} - C_k\|^2$  for  $1 \leq i \leq h$ ,  $1 \leq j \leq w$ , and  $1 \leq k \leq N$ . In the second step, we predict an initial 6D pose with sparse pose sampling centered at the coarse 2D locations predicted in the first step. Instead of rendering a feature map for each sampled pose, we compute vertex-level reconstruction losses based on the 2D coordinates of the vertices in the 3D model  $\mathcal{M}$ . Since the structure of the 3D model  $\mathcal{M}$  and the sampled poses are invariant across all testing samples, the 2D coordinates and visibility of the vertices can be pre-computed and cached. This allows us to predict initial 6D poses that are robust to partial occlusion and truncation and easy to optimize with a small complexity overhead.

## 2 Loss Landscapes of Different 6D Proposals

Previous methods for 3D object detection or 6D pose optimization are built on top of a 2D region proposal network or refine predictions from a separate pose estimation network. However, the first-stage networks are unreliable for objects with out-of-distribution textures or shapes or even missing objects if an

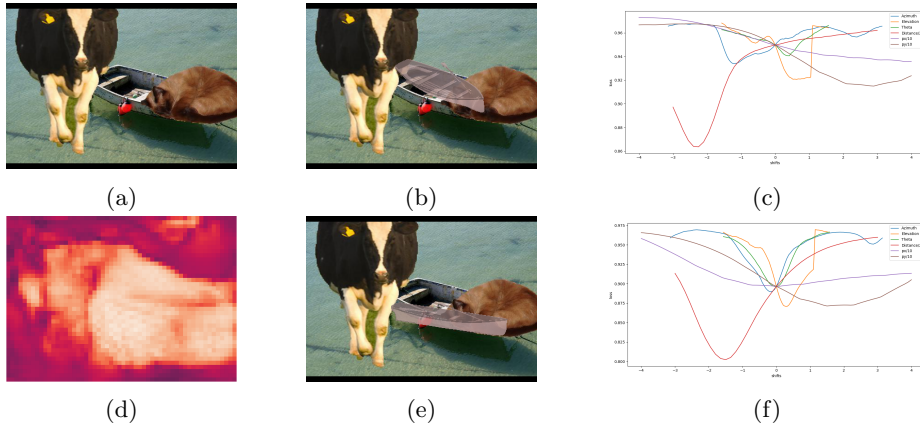


Fig. 1: Our object locating algorithm efficiently locates objects in the image and is robust to partial occlusion and truncation. Moreover, the generative 6D proposals are easier to optimize than using the Faster R-CNN prediction as the initial 6D pose. (b) and (c) shows the Faster R-CNN predicted pose and the loss landscapes over six dimensions centered at the Faster R-CNN prediction. (d) shows the reconstruction loss  $\max_r \|C_r - f_i\|^2$  for each pixel. (e) and (f) shows the generative 6D proposals and the loss landscapes around the proposed 6D poses. In this example with partial occlusion, the initial pose predicted by our object locating algorithm is located at regions with good gradients and hence easy to optimize.

object is partially occluded or truncated. This essentially limits the performance of the hybrid model. Therefore, we propose to search generative 6D proposals from the predicted feature maps using the learned 3D model  $\mathcal{N}$ . The generative 6D proposals are robust to partial occlusion and truncation and are easy to optimize. We compare the loss landscapes of the Faster R-CNN proposals and the generative 6D proposals in Figure 1. The loss landscapes are centered at the proposed 6D poses and show the loss landscapes around the 6D proposals. As we can see, our proposed 6D proposals are located in regions with good gradients pointing to the global minimum, which is easy to optimize.

### 3 Quantitative Results

#### 3.1 Quantitative Results on the PASCAL3D+ dataset and the Occluded PASCAL3D+ dataset

Table 1: Quantitative results of 6D pose estimation on the first 6 categories of PASCAL3D+ dataset.

Subset	Level	Category	Method	Pose Acc ( $\frac{\pi}{6}$ ) $\uparrow$	Pose Acc ( $\frac{\pi}{18}$ ) $\uparrow$	Median Pose Error $\downarrow$	Median ADD $\downarrow$
ImageNet	0	Aeroplane	FRCNN+Cls	69.66	22.60	0.31	0.57
ImageNet	0	Aeroplane	FRCNN+NeMo	58.20	13.83	0.44	1.46
ImageNet	0	Aeroplane	RTM3DExt	49.02	18.68	0.54	1.46
ImageNet	0	Aeroplane	Ours	<b>87.10</b>	<b>46.13</b>	<b>0.19</b>	<b>0.58</b>
ImageNet	0	Bicycle	FRCNN+Cls	66.51	15.35	0.36	0.79
ImageNet	0	Bicycle	FRCNN+NeMo	50.08	10.08	0.52	2.82
ImageNet	0	Bicycle	RTM3DExt	61.86	11.78	0.43	1.61
ImageNet	0	Bicycle	Ours	<b>68.06</b>	<b>22.95</b>	<b>0.33</b>	<b>0.49</b>
ImageNet	0	Boat	FRCNN+Cls	57.79	16.62	0.40	1.22
ImageNet	0	Boat	FRCNN+NeMo	44.85	7.08	0.60	2.34
ImageNet	0	Boat	RTM3DExt	29.75	10.76	1.13	2.91
ImageNet	0	Boat	Ours	<b>60.43</b>	<b>18.22</b>	<b>0.40</b>	<b>1.03</b>
ImageNet	0	Bottle	FRCNN+Cls	<b>84.34</b>	<b>42.70</b>	<b>0.20</b>	1.22
ImageNet	0	Bottle	FRCNN+NeMo	62.38	20.48	0.36	2.32
ImageNet	0	Bottle	RTM3DExt	80.59	43.91	0.21	<b>0.74</b>
ImageNet	0	Bottle	Ours	76.57	8.97	0.34	0.77
ImageNet	0	Bus	FRCNN+Cls	<b>93.61</b>	<b>70.68</b>	0.14	0.49
ImageNet	0	Bus	FRCNN+NeMo	88.91	71.05	0.11	1.32
ImageNet	0	Bus	RTM3DExt	89.66	79.14	<b>0.08</b>	0.43
ImageNet	0	Bus	Ours	82.33	68.05	0.09	<b>0.38</b>
ImageNet	0	Car	FRCNN+Cls	91.74	57.30	0.16	0.48
ImageNet	0	Car	FRCNN+NeMo	77.65	48.19	0.19	1.26
ImageNet	0	Car	RTM3DExt	93.44	65.93	0.13	0.48
ImageNet	0	Car	Ours	<b>96.13</b>	<b>80.68</b>	<b>0.08</b>	<b>0.30</b>

Table 2: Quantitative results of 6D pose estimation on the second 6 categories of PASCAL3D+ dataset.

Subset	Level	Category	Method	Pose Acc ( $\frac{\pi}{6}$ ) $\uparrow$	Pose Acc ( $\frac{\pi}{18}$ ) $\uparrow$	Median Pose Error $\downarrow$	Median ADD $\downarrow$
ImageNet	0	Chair	FRCNN+Cls	80.67	27.42	0.25	1.33
ImageNet	0	Chair	FRCNN+NeMo	48.92	5.92	0.53	2.80
ImageNet	0	Chair	RTM3DExt	83.23	30.97	0.26	0.93
ImageNet	0	Chair	Ours	<b>89.55</b>	<b>51.08</b>	<b>0.17</b>	<b>0.59</b>
ImageNet	0	Diningtable	FRCNN+Cls	65.39	27.84	0.29	1.07
ImageNet	0	Diningtable	FRCNN+NeMo	59.15	17.87	0.40	3.02
ImageNet	0	Diningtable	RTM3DExt	72.68	28.97	0.28	1.85
ImageNet	0	Diningtable	Ours	<b>82.83</b>	<b>59.32</b>	<b>0.13</b>	<b>0.76</b>
ImageNet	0	Motorbike	FRCNN+Cls	62.08	10.74	0.39	1.01
ImageNet	0	Motorbike	FRCNN+NeMo	51.51	11.91	0.51	2.43
ImageNet	0	Motorbike	RTM3DExt	62.75	13.42	0.41	2.10
ImageNet	0	Motorbike	Ours	<b>71.98</b>	<b>21.14</b>	<b>0.31</b>	<b>0.58</b>
ImageNet	0	Sofa	FRCNN+Cls	<b>91.51</b>	37.82	0.20	0.74
ImageNet	0	Sofa	FRCNN+NeMo	82.53	37.82	0.23	1.57
ImageNet	0	Sofa	RTM3DExt	91.99	<b>48.08</b>	0.18	0.56
ImageNet	0	Sofa	Ours	90.06	47.60	<b>0.18</b>	<b>0.45</b>
ImageNet	0	Train	FRCNN+Cls	<b>94.27</b>	<b>61.61</b>	<b>0.15</b>	<b>0.70</b>
ImageNet	0	Train	FRCNN+NeMo	87.77	48.45	0.18	1.03
ImageNet	0	Train	RTM3DExt	90.09	54.95	0.16	1.25
ImageNet	0	Train	Ours	58.36	30.50	0.36	1.24
ImageNet	0	Tvmonitor	FRCNN+Cls	78.14	22.03	0.30	1.00
ImageNet	0	Tvmonitor	FRCNN+NeMo	66.88	17.20	0.37	2.60
ImageNet	0	Tvmonitor	RTM3DExt	<b>82.15</b>	22.99	0.28	0.66
ImageNet	0	Tvmonitor	Ours	78.62	<b>30.06</b>	<b>0.27</b>	<b>0.62</b>
ImageNet	0	Mean	FRCNN+Cls	78.90	37.35	0.22	0.74
ImageNet	0	Mean	FRCNN+NeMo	66.06	28.44	0.33	1.84
ImageNet	0	Mean	RTM3DExt	74.94	39.56	0.23	0.92
ImageNet	0	Mean	Ours	<b>81.45</b>	<b>47.68</b>	<b>0.19</b>	<b>0.53</b>

Table 3: Quantitative results of 6D pose estimation on the first 6 categories of Occluded PASCAL3D+ dataset (occlusion level 1).

Subset	Level	Category	Method	Pose Acc ( $\frac{\pi}{6}$ ) $\uparrow$	Pose Acc ( $\frac{\pi}{18}$ ) $\uparrow$	Median Pose Error $\downarrow$	Median ADD $\downarrow$
ImageNet	1	Aeroplane	FRCNN+Cls	50.89	12.83	0.51	0.90
ImageNet	1	Aeroplane	FRCNN+NeMo	37.33	7.89	0.74	1.61
ImageNet	1	Aeroplane	RTM3DExt	28.71	6.31	2.64	4.01
ImageNet	1	Aeroplane	Ours	<b>61.62</b>	<b>18.09</b>	<b>0.39</b>	<b>0.90</b>
ImageNet	1	Bicycle	FRCNN+Cls	39.75	6.47	0.76	1.10
ImageNet	1	Bicycle	FRCNN+NeMo	24.29	5.21	1.37	2.53
ImageNet	1	Bicycle	RTM3DExt	20.50	3.63	3.14	8.03
ImageNet	1	Bicycle	Ours	<b>53.00</b>	<b>13.88</b>	<b>0.49</b>	<b>0.72</b>
ImageNet	1	Boat	FRCNN+Cls	34.20	9.01	1.08	1.66
ImageNet	1	Boat	FRCNN+NeMo	25.95	4.56	1.39	2.43
ImageNet	1	Boat	RTM3DExt	10.21	2.93	3.11	4.87
ImageNet	1	Boat	Ours	<b>39.52</b>	<b>10.10</b>	<b>0.72</b>	<b>1.51</b>
ImageNet	1	Bottle	FRCNN+Cls	<b>67.49</b>	30.60	<b>0.28</b>	1.74
ImageNet	1	Bottle	FRCNN+NeMo	45.77	14.62	0.74	2.72
ImageNet	1	Bottle	RTM3DExt	65.98	<b>32.10</b>	0.29	1.60
ImageNet	1	Bottle	Ours	62.98	7.92	0.42	<b>0.95</b>
ImageNet	1	Bus	FRCNN+Cls	<b>86.34</b>	<b>54.46</b>	<b>0.17</b>	0.82
ImageNet	1	Bus	FRCNN+NeMo	78.37	41.18	0.22	1.36
ImageNet	1	Bus	RTM3DExt	55.03	34.16	0.38	3.42
ImageNet	1	Bus	Ours	74.76	49.53	0.18	<b>0.58</b>
ImageNet	1	Car	FRCNN+Cls	78.65	40.59	0.20	0.69
ImageNet	1	Car	FRCNN+NeMo	66.93	33.45	0.28	1.33
ImageNet	1	Car	RTM3DExt	56.05	28.04	0.38	2.20
ImageNet	1	Car	Ours	<b>84.69</b>	<b>60.20</b>	<b>0.14</b>	<b>0.41</b>

Table 4: Quantitative results of 6D pose estimation on the second 6 categories of Occluded PASCAL3D+ dataset (occlusion level 1).

Subset	Level	Category	Method	Pose Acc ( $\frac{\pi}{6}$ ) $\uparrow$	Pose Acc ( $\frac{\pi}{18}$ ) $\uparrow$	Median Pose Error $\downarrow$	Median ADD $\downarrow$
ImageNet	1	Chair	FRCNN+Cls	54.03	16.73	0.47	1.97
ImageNet	1	Chair	FRCNN+NeMo	24.19	2.62	1.01	3.30
ImageNet	1	Chair	RTM3DExt	38.10	8.67	0.76	3.29
ImageNet	1	Chair	Ours	<b>66.33</b>	<b>21.77</b>	<b>0.36</b>	<b>0.93</b>
ImageNet	1	Diningtable	FRCNN+Cls	47.60	21.05	0.60	1.44
ImageNet	1	Diningtable	FRCNN+NeMo	34.24	6.64	1.02	2.80
ImageNet	1	Diningtable	RTM3DExt	35.37	10.39	1.64	4.54
ImageNet	1	Diningtable	Ours	<b>65.68</b>	<b>32.58</b>	<b>0.28</b>	<b>1.07</b>
ImageNet	1	Motorbike	FRCNN+Cls	36.51	5.36	0.84	1.20
ImageNet	1	Motorbike	FRCNN+NeMo	26.64	2.42	1.34	2.47
ImageNet	1	Motorbike	RTM3DExt	15.05	1.38	3.14	9.17
ImageNet	1	Motorbike	Ours	<b>53.29</b>	<b>12.11</b>	<b>0.48</b>	<b>0.90</b>
ImageNet	1	Sofa	FRCNN+Cls	81.18	<b>31.73</b>	<b>0.25</b>	0.99
ImageNet	1	Sofa	FRCNN+NeMo	69.19	19.37	0.33	1.61
ImageNet	1	Sofa	RTM3DExt	62.18	17.16	0.37	1.92
ImageNet	1	Sofa	Ours	<b>82.10</b>	27.49	0.28	<b>0.79</b>
ImageNet	1	Train	FRCNN+Cls	<b>80.09</b>	<b>42.34</b>	<b>0.20</b>	<b>1.11</b>
ImageNet	1	Train	FRCNN+NeMo	73.78	32.39	0.25	1.37
ImageNet	1	Train	RTM3DExt	69.98	37.60	0.25	1.88
ImageNet	1	Train	Ours	50.39	23.38	0.51	1.61
ImageNet	1	Tvmonitor	FRCNN+Cls	57.63	15.42	0.45	1.55
ImageNet	1	Tvmonitor	FRCNN+NeMo	42.69	7.63	0.60	2.69
ImageNet	1	Tvmonitor	RTM3DExt	52.60	12.01	0.50	1.75
ImageNet	1	Tvmonitor	Ours	<b>66.40</b>	<b>16.40</b>	<b>0.40</b>	<b>0.93</b>
ImageNet	1	Mean	FRCNN+Cls	61.48	26.11	0.33	1.07
ImageNet	1	Mean	FRCNN+NeMo	48.34	17.46	0.55	1.90
ImageNet	1	Mean	RTM3DExt	43.55	17.68	0.82	3.29
ImageNet	1	Mean	Ours	<b>66.63</b>	<b>30.84</b>	<b>0.31</b>	<b>0.77</b>

Table 5: Quantitative results of 6D pose estimation on the first 6 categories of Occluded PASCAL3D+ dataset (occlusion level 2).

Subset	Level	Category	Method	Pose Acc ( $\frac{\pi}{6}$ ) $\uparrow$	Pose Acc ( $\frac{\pi}{18}$ ) $\uparrow$	Median Pose Error $\downarrow$	Median ADD $\downarrow$
ImageNet	2	Aeroplane	FRCNN+Cls	26.68	3.74	1.07	1.26
ImageNet	2	Aeroplane	FRCNN+NeMo	22.52	3.09	1.21	1.68
ImageNet	2	Aeroplane	RTM3DExt	10.67	1.49	3.14	7.47
ImageNet	2	Aeroplane	Ours	<b>35.86</b>	<b>5.98</b>	<b>0.83</b>	<b>1.35</b>
ImageNet	2	Bicycle	FRCNN+Cls	18.58	2.26	1.69	1.67
ImageNet	2	Bicycle	FRCNN+NeMo	15.35	1.29	2.28	2.52
ImageNet	2	Bicycle	RTM3DExt	6.14	0.65	3.14	18.65
ImageNet	2	Bicycle	Ours	<b>36.03</b>	<b>8.40</b>	<b>0.73</b>	<b>1.04</b>
ImageNet	2	Boat	FRCNN+Cls	20.00	<b>3.87</b>	1.98	2.16
ImageNet	2	Boat	FRCNN+NeMo	14.81	1.44	2.17	2.64
ImageNet	2	Boat	RTM3DExt	4.75	1.66	3.14	12.19
ImageNet	2	Boat	Ours	<b>27.18</b>	3.76	<b>1.40</b>	<b>2.16</b>
ImageNet	2	Bottle	FRCNN+Cls	<b>39.31</b>	<b>13.38</b>	2.80	3.09
ImageNet	2	Bottle	FRCNN+NeMo	24.41	5.66	3.03	3.87
ImageNet	2	Bottle	RTM3DExt	33.24	15.31	3.14	7.90
ImageNet	2	Bottle	Ours	38.21	4.14	<b>0.65</b>	<b>1.88</b>
ImageNet	2	Bus	FRCNN+Cls	<b>70.48</b>	<b>33.52</b>	<b>0.27</b>	0.88
ImageNet	2	Bus	FRCNN+NeMo	60.95	20.76	0.38	1.27
ImageNet	2	Bus	RTM3DExt	20.38	10.29	3.14	32.98
ImageNet	2	Bus	Ours	55.43	25.14	0.41	<b>0.86</b>
ImageNet	2	Car	FRCNN+Cls	58.31	23.47	0.36	0.88
ImageNet	2	Car	FRCNN+NeMo	51.30	19.33	0.50	1.31
ImageNet	2	Car	RTM3DExt	25.46	9.49	3.14	8.43
ImageNet	2	Car	Ours	<b>70.02</b>	<b>35.76</b>	<b>0.26</b>	<b>0.68</b>

Table 6: Quantitative results of 6D pose estimation on the second 6 categories of Occluded PASCAL3D+ dataset (occlusion level 2).

Subset	Level	Category	Method	Pose Acc ( $\frac{\pi}{6}$ ) $\uparrow$	Pose Acc ( $\frac{\pi}{18}$ ) $\uparrow$	Median Pose Error $\downarrow$	Median ADD $\downarrow$
ImageNet	2	Chair	FRCNN+Cls	31.06	8.28	1.08	2.60
ImageNet	2	Chair	FRCNN+NeMo	13.66	1.66	2.45	3.79
ImageNet	2	Chair	RTM3DExt	22.36	3.93	3.14	7.38
ImageNet	2	Chair	Ours	<b>42.03</b>	<b>9.52</b>	<b>0.64</b>	<b>1.50</b>
ImageNet	2	Diningtable	FRCNN+Cls	31.98	<b>14.63</b>	0.99	1.92
ImageNet	2	Diningtable	FRCNN+NeMo	25.38	5.06	1.32	2.76
ImageNet	2	Diningtable	RTM3DExt	15.90	4.34	3.14	7.85
ImageNet	2	Diningtable	Ours	<b>43.45</b>	13.91	<b>0.64</b>	<b>1.57</b>
ImageNet	2	Motorbike	FRCNN+Cls	18.60	1.58	1.72	1.60
ImageNet	2	Motorbike	FRCNN+NeMo	13.16	1.40	2.06	2.48
ImageNet	2	Motorbike	RTM3DExt	3.51	0.18	3.14	30.67
ImageNet	2	Motorbike	Ours	<b>34.56</b>	<b>5.44</b>	<b>0.77</b>	<b>1.35</b>
ImageNet	2	Sofa	FRCNN+Cls	<b>64.68</b>	<b>23.22</b>	<b>0.34</b>	1.35
ImageNet	2	Sofa	FRCNN+NeMo	55.85	11.52	0.47	1.63
ImageNet	2	Sofa	RTM3DExt	38.39	8.25	0.98	4.12
ImageNet	2	Sofa	Ours	56.05	11.90	0.45	<b>1.20</b>
ImageNet	2	Train	FRCNN+Cls	<b>60.03</b>	<b>23.98</b>	<b>0.36</b>	<b>1.86</b>
ImageNet	2	Train	FRCNN+NeMo	56.12	19.09	0.43	1.93
ImageNet	2	Train	RTM3DExt	46.82	23.00	0.68	4.59
ImageNet	2	Train	Ours	39.15	12.56	1.75	2.39
ImageNet	2	Tvmonitor	FRCNN+Cls	41.03	<b>9.47</b>	0.71	2.12
ImageNet	2	Tvmonitor	FRCNN+NeMo	28.74	5.15	1.19	2.84
ImageNet	2	Tvmonitor	RTM3DExt	31.40	6.98	3.14	5.03
ImageNet	2	Tvmonitor	Ours	<b>47.34</b>	8.80	<b>0.55</b>	<b>1.47</b>
ImageNet	2	Mean	FRCNN+Cls	41.94	14.74	0.75	1.47
ImageNet	2	Mean	FRCNN+NeMo	34.32	9.64	1.05	2.03
ImageNet	2	Mean	RTM3DExt	21.27	7.24	3.14	9.00
ImageNet	2	Mean	Ours	<b>47.95</b>	<b>16.25</b>	<b>0.56</b>	<b>1.22</b>



Table 7: Quantitative results of 6D pose estimation on the first 6 categories of Occluded PASCAL3D+ dataset (occlusion level 3).

Subset	Level	Category	Method	Pose Acc ( $\frac{\pi}{6}$ ) $\uparrow$	Pose Acc ( $\frac{\pi}{18}$ ) $\uparrow$	Median Pose Error $\downarrow$	Median ADD $\downarrow$
ImageNet	3	Aeroplane	FRCNN+Cls	10.41	<b>2.33</b>	1.68	<b>1.41</b>
ImageNet	3	Aeroplane	FRCNN+NeMo	8.64	0.78	1.82	1.72
ImageNet	3	Aeroplane	RTM3DExt	3.32	0.11	3.14	12.39
ImageNet	3	Aeroplane	Ours	<b>14.40</b>	1.66	<b>1.46</b>	1.61
ImageNet	3	Bicycle	FRCNN+Cls	7.49	1.33	2.06	2.20
ImageNet	3	Bicycle	FRCNN+NeMo	6.16	0.67	2.64	2.68
ImageNet	3	Bicycle	RTM3DExt	0.83	0.00	3.14	83.21
ImageNet	3	Bicycle	Ours	<b>19.80</b>	<b>2.83</b>	<b>1.28</b>	<b>1.63</b>
ImageNet	3	Boat	FRCNN+Cls	11.27	<b>3.04</b>	2.53	<b>2.70</b>
ImageNet	3	Boat	FRCNN+NeMo	7.78	0.79	2.83	2.96
ImageNet	3	Boat	RTM3DExt	2.03	0.34	3.14	49.79
ImageNet	3	Boat	Ours	<b>12.40</b>	1.24	<b>2.10</b>	3.13
ImageNet	3	Bottle	FRCNN+Cls	14.31	4.03	3.10	3.86
ImageNet	3	Bottle	FRCNN+NeMo	9.03	1.81	3.10	4.33
ImageNet	3	Bottle	RTM3DExt	20.14	<b>8.89</b>	3.14	15.41
ImageNet	3	Bottle	Ours	<b>20.97</b>	0.83	<b>1.06</b>	<b>3.33</b>
ImageNet	3	Bus	FRCNN+Cls	46.45	12.67	0.57	1.34
ImageNet	3	Bus	FRCNN+NeMo	37.62	8.64	0.79	1.54
ImageNet	3	Bus	RTM3DExt	1.92	0.58	3.14	38.54
ImageNet	3	Bus	Ours	33.21	7.68	1.04	1.82
ImageNet	3	Car	FRCNN+Cls	30.86	6.84	1.07	<b>1.14</b>
ImageNet	3	Car	FRCNN+NeMo	26.89	5.52	1.28	1.52
ImageNet	3	Car	RTM3DExt	7.42	2.18	3.14	56.89
ImageNet	3	Car	Ours	<b>42.21</b>	<b>11.08</b>	<b>0.74</b>	1.20

Table 8: Quantitative results of 6D pose estimation on the second 6 categories of Occluded PASCAL3D+ dataset (occlusion level 3).

Subset	Level	Category	Method	Pose Acc ( $\frac{\pi}{6}$ ) $\uparrow$	Pose Acc ( $\frac{\pi}{18}$ ) $\uparrow$	Median Pose Error $\downarrow$	Median ADD $\downarrow$
ImageNet	3	Chair	FRCNN+Cls	15.04	3.60	2.61	3.29
ImageNet	3	Chair	FRCNN+NeMo	6.36	0.64	2.91	4.32
ImageNet	3	Chair	RTM3DExt	13.98	2.12	3.14	9.18
ImageNet	3	Chair	Ours	<b>24.58</b>	<b>3.60</b>	<b>1.02</b>	<b>2.36</b>
ImageNet	3	Diningtable	FRCNN+Cls	19.91	<b>7.44</b>	1.69	<b>2.15</b>
ImageNet	3	Diningtable	FRCNN+NeMo	14.70	2.14	2.60	3.02
ImageNet	3	Diningtable	RTM3DExt	11.53	3.26	3.14	8.83
ImageNet	3	Diningtable	Ours	<b>24.37</b>	4.37	<b>1.06</b>	2.17
ImageNet	3	Motorbike	FRCNN+Cls	6.13	0.36	2.39	2.47
ImageNet	3	Motorbike	FRCNN+NeMo	6.31	0.72	2.77	3.06
ImageNet	3	Motorbike	RTM3DExt	0.90	0.18	3.14	67.59
ImageNet	3	Motorbike	Ours	<b>14.95</b>	<b>1.26</b>	<b>1.46</b>	<b>2.41</b>
ImageNet	3	Sofa	FRCNN+Cls	<b>38.07</b>	<b>8.48</b>	0.95	1.75
ImageNet	3	Sofa	FRCNN+NeMo	30.37	5.13	1.11	1.79
ImageNet	3	Sofa	RTM3DExt	27.42	6.31	3.14	6.09
ImageNet	3	Sofa	Ours	34.91	4.14	<b>0.62</b>	<b>1.48</b>
ImageNet	3	Train	FRCNN+Cls	<b>37.74</b>	9.43	2.55	<b>2.35</b>
ImageNet	3	Train	FRCNN+NeMo	34.31	7.72	2.56	2.40
ImageNet	3	Train	RTM3DExt	28.13	<b>13.72</b>	3.14	9.78
ImageNet	3	Train	Ours	23.50	6.69	<b>2.46</b>	3.56
ImageNet	3	Tvmonitor	FRCNN+Cls	21.84	<b>5.63</b>	3.00	3.00
ImageNet	3	Tvmonitor	FRCNN+NeMo	17.06	1.71	2.98	3.53
ImageNet	3	Tvmonitor	RTM3DExt	20.14	4.27	3.14	8.43
ImageNet	3	Tvmonitor	Ours	<b>33.11</b>	4.10	<b>0.77</b>	<b>2.35</b>
ImageNet	3	Mean	FRCNN+Cls	22.42	<b>5.58</b>	2.01	1.95
ImageNet	3	Mean	FRCNN+NeMo	18.17	5.30	2.40	2.35
ImageNet	3	Mean	RTM3DExt	10.17	3.11	3.14	19.92
ImageNet	3	Mean	Ours	<b>27.43</b>	5.30	<b>1.07</b>	<b>1.94</b>