

# A Mixed Classification-Regression Framework for 3D Pose Estimation from 2D Images

Siddharth Mahendran<sup>1</sup>  
siddharthm@jhu.edu

Haider Ali<sup>2</sup>  
hali@jhu.edu

René Vidal<sup>1</sup>  
rvidal@cis.jhu.edu

<sup>1</sup> Center for Imaging Science,  
Mathematical Institute for Data Science,  
Johns Hopkins University,  
Baltimore, MD, USA

<sup>2</sup> Department of Computer Science,  
Johns Hopkins University,  
Baltimore, MD, USA

---

We now present more results §1 and implementation details §2 that we were unable to include in the main paper due to space constraints.

## 1 Results

In Tables 2 and 3, we present the results of our eight models and three baselines compared to current state-of-the-art under the two metric of  $MedErr$  and  $Acc_{\frac{\pi}{6}}$  respectively. As we mentioned in the paper, we run all experiments three times and report the mean and standard deviation (in brackets) across these three trials. We also show figures of images where we obtain the least pose estimation error and the most pose estimation error for every object category using one run of model  $\mathcal{M}_{G+}$ . As can be seen from Figs. [1-12], we make the most error under three conditions: (i) when the objects are really blurry (very small in pixel size in the original image), (ii) the shape of the object is uncommon (possibly very few examples seen during training) and (iii) the pose of a test image is very different from common poses observed during training. The first condition is best observed in the bad cases for categories aeroplane and car where almost all the images shown are very blurry. The second condition is best observed in categories boat and chair where the bad cases contain uncommon boats and chairs. The third condition is best observed in categories bottle and tvmonitor where the bad images are in very different poses compared to the best images.

We also present the performance of our models  $\mathcal{M}_G$  and  $\mathcal{M}_{G+}$  across different object categories of the Pascal3D+ dataset during ablation experiments in Tables 4-11. These are detailed tables for the results shown in Tables 3 and 4 of the main paper and an overview of the experiments is shown below.

## 2 Implementation details

We use the ResNet-50 upto layer4 (2048-dim feature output) as our feature network. The pose networks are of the form Input-FC-BN-ReLU-FC-BN-ReLU-FC-Output where FC is a

Expt.	$\mathcal{M}_G$				$\mathcal{M}_{G+}$			
	<i>MedErr</i>	<i>Acc</i> $_{\frac{\pi}{6}}$	<i>MedErr</i>	<i>Acc</i> $_{\frac{\pi}{6}}$	<i>MedErr</i>	<i>Acc</i> $_{\frac{\pi}{6}}$	<i>MedErr</i>	<i>Acc</i> $_{\frac{\pi}{6}}$
K-Means dictionary size $K$	Table 4	Table 5	Table 6	Table 7	Table 6	Table 7	Table 10	Table 11
Weighting parameter $\alpha$	Table 8	Table 9	Table 10	Table 11	Table 10	Table 11	Table 10	Table 11

Table 1: Overview of Tables of results showing ablation analysis

Method	aero	bike	boat	bottle	bus	car	chair	dtable	mbike	sofa	train	tv	Mean
$\square$	13.8	17.7	21.3	12.9	5.8	9.1	14.8	15.2	14.7	13.7	8.7	15.4	13.59
$\square$	15.4	14.8	25.6	9.3	3.6	6.0	9.7	10.8	16.7	9.5	6.1	12.6	11.68
$\square$	13.6	12.5	22.8	8.3	3.1	5.8	11.9	12.5	12.3	12.8	6.3	11.9	11.15
$\square$	10.0	15.6	19.1	8.6	3.3	5.1	13.7	11.8	12.2	13.5	6.7	11.0	10.88
$\mathcal{R}_E$	14.5 (1.2)	17.7 (0.4)	39.3 (3.4)	7.4 (0.3)	4.0 (0.4)	7.8 (0.5)	15.2 (0.4)	26.6 (2.0)	17.5 (0.2)	10.5 (0.7)	11.5 (2.0)	14.1 (0.6)	15.50 (0.34)
$\mathcal{R}_G$	11.8 (0.6)	15.9 (0.5)	27.2 (2.0)	7.2 (0.6)	2.9 (0.1)	5.2 (0.1)	11.6 (0.4)	15.0 (3.9)	14.3 (0.6)	10.8 (0.1)	5.4 (0.3)	12.4 (0.7)	11.63 (0.51)
$\mathcal{C}$	11.7 (0.7)	15.3 (0.7)	21.5 (2.1)	9.3 (0.3)	4.1 (0.1)	7.4 (0.3)	11.2 (0.2)	17.8 (3.9)	17.0 (0.7)	11.0 (0.4)	7.0 (0.3)	13.1 (0.3)	12.20 (0.44)
$\mathcal{M}_S$	11.0 (0.9)	15.5 (0.8)	21.0 (0.8)	8.8 (0.5)	3.8 (0.2)	7.0 (0.0)	10.8 (0.2)	21.0 (5.8)	16.6 (0.2)	10.7 (1.0)	6.5 (0.3)	13.1 (0.7)	12.14 (0.37)
$\mathcal{M}_G$	12.1 (0.6)	16.0 (0.5)	19.9 (1.1)	8.9 (0.2)	3.4 (0.2)	6.5 (0.3)	10.8 (0.0)	15.2 (6.0)	16.4 (0.5)	9.6 (0.6)	5.9 (0.9)	13.0 (0.7)	11.48 (0.23)
$\mathcal{M}_{LE}$	12.8 (1.0)	15.2 (0.8)	23.4 (3.3)	9.0 (0.2)	4.0 (0.1)	7.4 (0.2)	11.1 (0.1)	16.8 (3.5)	16.1 (0.7)	10.7 (1.4)	6.6 (0.1)	12.3 (0.8)	12.11 (0.32)
$\mathcal{M}_P$	11.4 (0.4)	16.3 (0.6)	25.6 (0.6)	7.0 (0.4)	2.6 (0.1)	5.1 (0.1)	11.3 (0.4)	16.0 (2.1)	13.6 (1.0)	10.2 (0.7)	5.5 (0.3)	12.0 (0.3)	11.38 (0.16)
$\mathcal{M}_{S+}$	12.2 (0.2)	15.7 (1.0)	24.4 (0.6)	9.9 (0.2)	3.6 (0.2)	6.5 (0.2)	12.0 (0.4)	14.8 (0.8)	14.4 (0.8)	11.9 (0.7)	6.4 (0.2)	11.6 (0.3)	11.95 (0.26)
$\mathcal{M}_{G+}$	8.5 (0.0)	14.8 (1.2)	20.5 (1.4)	7.0 (0.4)	3.1 (0.2)	5.1 (0.1)	9.3 (0.1)	11.3 (2.0)	14.2 (0.6)	10.2 (0.3)	5.6 (0.1)	11.7 (0.2)	10.10 (0.38)
$\mathcal{M}_{LE+}$	12.3 (0.4)	16.7 (0.4)	24.7 (1.9)	7.5 (0.2)	3.6 (0.1)	6.5 (0.1)	11.5 (0.5)	15.5 (0.4)	15.1 (0.5)	11.1 (0.5)	7.3 (0.7)	12.1 (0.2)	11.99 (0.15)
$\mathcal{M}_{P+}$	10.6 (0.2)	15.0 (0.2)	23.9 (1.2)	6.7 (0.2)	2.7 (0.1)	4.7 (0.1)	9.8 (0.1)	12.6 (0.5)	13.9 (0.9)	9.7 (0.1)	5.3 (0.1)	11.7 (0.5)	10.54 (0.16)

Table 2: Performance of our models under the *MedErr* metric (lower is better).

fully connected layer, BN is a batch normalization layer and ReLU is the standard rectified linear unit non-linearity. The pose networks for models  $\mathcal{R}_E$  and  $\mathcal{R}_G$  are of size 2048-1000-500-3. The pose network of model  $\mathcal{C}$  is of size 2048-1000-500-100 where 100 is the size of the K-Means dictionary we use to discretize the pose space. The bin and delta networks of models  $\mathcal{M}_S$ ,  $\mathcal{M}_G$ ,  $\mathcal{M}_{LE}$  and  $\mathcal{M}_P$  are of sizes 2048-1000-500-100 and 2048-1000-500-3 respectively. For models  $\mathcal{M}_{S+}$ ,  $\mathcal{M}_{G+}$ ,  $\mathcal{M}_{LE+}$  and  $\mathcal{M}_{P+}$  where we have one delta network per pose-bin per object category, our bin network is of size 2048-1000-500-16 (corresponding to 16 pose-bins) and we use a 2-layer delta network of size 2048-100-3.

For the models,  $\mathcal{M}_G$  and  $\mathcal{M}_{G+}$ , we initialize the network weights with 1 epoch of training over the models  $\mathcal{M}_S$  and  $\mathcal{M}_{S+}$ . All other models are initialized using pre-trained networks on the ImageNet image classification problem. The models  $\mathcal{M}_S$ ,  $\mathcal{M}_G$ ,  $\mathcal{M}_P$ ,  $\mathcal{M}_{S+}$  and  $\mathcal{M}_{P+}$  were trained with  $\alpha = 1$ . For the model  $\mathcal{M}_{G+}$ , we use a value of  $\alpha = 10$  and for the models  $\mathcal{M}_{LE}$  and  $\mathcal{M}_{LE+}$ , we use  $\alpha = 0.1$ .

Method	aero	bike	boat	bottle	bus	car	chair	dtable	mbike	sofa	train	tv	Mean
$\mathcal{I}$	0.81	0.77	0.59	0.93	0.98	0.89	0.80	0.62	0.88	0.82	0.80	0.80	0.8075
$\mathcal{I}$	0.74	0.83	0.52	0.91	0.91	0.88	0.86	0.73	0.78	0.90	0.86	0.92	0.8200
$\mathcal{I}$	0.78	0.83	0.57	0.93	0.94	0.90	0.80	0.68	0.86	0.82	0.82	0.85	0.8103
$\mathcal{I}$	0.83	0.82	0.64	0.95	0.97	0.94	0.80	0.71	0.88	0.87	0.80	0.86	0.8392
$\mathcal{R}_E$	0.77 (0.01)	0.75 (0.01)	0.41 (0.03)	0.96 (0.00)	0.91 (0.01)	0.83 (0.01)	0.72 (0.02)	0.56 (0.04)	0.75 (0.01)	0.90 (0.02)	0.75 (0.03)	0.87 (0.00)	0.7656 (0.0015)
$\mathcal{R}_G$	0.80 (0.01)	0.78 (0.00)	0.54 (0.03)	0.97 (0.00)	0.95 (0.01)	0.93 (0.01)	0.83 (0.01)	0.59 (0.02)	0.82 (0.01)	0.91 (0.01)	0.81 (0.01)	0.86 (0.01)	0.8166 (0.0041)
$\mathcal{C}$	0.84 (0.01)	0.77 (0.01)	0.60 (0.01)	0.95 (0.01)	0.97 (0.01)	0.95 (0.01)	0.90 (0.01)	0.63 (0.06)	0.78 (0.02)	0.94 (0.01)	0.81 (0.01)	0.87 (0.01)	0.8350 (0.0045)
$\mathcal{M}_S$	0.83 (0.01)	0.78 (0.01)	0.61 (0.00)	0.96 (0.00)	0.96 (0.00)	0.94 (0.01)	0.90 (0.01)	0.56 (0.04)	0.79 (0.01)	0.95 (0.02)	0.82 (0.00)	0.87 (0.02)	0.8303 (0.0014)
$\mathcal{M}_G$	0.83 (0.02)	0.76 (0.03)	0.63 (0.01)	0.96 (0.01)	0.97 (0.02)	0.93 (0.00)	0.91 (0.00)	0.57 (0.04)	0.78 (0.02)	0.95 (0.02)	0.82 (0.01)	0.88 (0.01)	0.8335 (0.0067)
$\mathcal{M}_{LE}$	0.83 (0.02)	0.77 (0.02)	0.58 (0.03)	0.96 (0.00)	0.96 (0.01)	0.94 (0.01)	0.91 (0.00)	0.71 (0.10)	0.81 (0.00)	0.93 (0.01)	0.81 (0.01)	0.87 (0.01)	0.8410 (0.0025)
$\mathcal{M}_P$	0.80 (0.01)	0.77 (0.01)	0.56 (0.00)	0.97 (0.01)	0.97 (0.00)	0.93 (0.01)	0.82 (0.01)	0.57 (0.04)	0.81 (0.03)	0.92 (0.00)	0.82 (0.01)	0.88 (0.01)	0.8185 (0.0035)
$\mathcal{M}_{S+}$	0.82 (0.01)	0.80 (0.01)	0.59 (0.01)	0.94 (0.01)	0.97 (0.01)	0.94 (0.01)	0.91 (0.02)	0.63 (0.04)	0.81 (0.00)	0.97 (0.00)	0.83 (0.01)	0.87 (0.00)	0.8387 (0.0044)
$\mathcal{M}_{G+}$	0.87 (0.01)	0.81 (0.01)	0.64 (0.01)	0.96 (0.00)	0.97 (0.01)	0.95 (0.01)	0.92 (0.01)	0.67 (0.10)	0.85 (0.01)	0.97 (0.01)	0.82 (0.01)	0.88 (0.00)	0.8588 (0.0111)
$\mathcal{M}_{LE+}$	0.81 (0.02)	0.77 (0.02)	0.56 (0.03)	0.96 (0.01)	0.97 (0.01)	0.92 (0.00)	0.86 (0.01)	0.73 (0.06)	0.79 (0.02)	0.93 (0.02)	0.80 (0.01)	0.89 (0.01)	0.8329 (0.0032)
$\mathcal{M}_{P+}$	0.84 (0.00)	0.82 (0.00)	0.59 (0.01)	0.97 (0.00)	0.97 (0.01)	0.95 (0.00)	0.88 (0.01)	0.68 (0.04)	0.84 (0.01)	0.93 (0.01)	0.81 (0.01)	0.89 (0.01)	0.8470 (0.0044)

Table 3: Performance of our models under the  $Acc_{\frac{\pi}{6}}$  metric (higher is better).

Method	aero	bike	boat	bottle	bus	car	chair	dtable	mbike	sofa	train	tv	Mean
K=24	12.4 (0.3)	16.3 (0.1)	23.5 (0.3)	8.7 (0.1)	2.7 (0.0)	5.4 (0.3)	11.6 (0.1)	17.4 (2.6)	16.3 (0.4)	13.7 (1.4)	6.2 (0.5)	15.6 (0.5)	12.48 (0.35)
K=50	12.9 (0.9)	16.0 (0.9)	21.1 (1.3)	8.3 (0.2)	3.3 (0.2)	6.1 (0.1)	11.2 (0.4)	22.2 (5.1)	17.7 (1.1)	11.8 (0.9)	5.8 (0.2)	13.9 (0.2)	12.53 (0.45)
K=100	12.1 (0.6)	16.0 (0.5)	19.9 (1.1)	8.9 (0.2)	3.4 (0.2)	6.5 (0.3)	10.8 (0.0)	15.2 (6.0)	16.4 (0.5)	9.6 (0.6)	5.9 (0.9)	13.0 (0.7)	11.48 (0.23)
K=200	10.6 (0.2)	16.4 (1.4)	21.6 (0.2)	8.1 (0.4)	3.2 (0.1)	6.0 (0.2)	9.9 (0.2)	14.6 (1.2)	16.0 (0.8)	11.1 (0.6)	6.3 (0.9)	13.4 (0.4)	11.44 (0.20)

Table 4: Ablation analysis of the size of K-Means dictionary in model  $\mathcal{M}_G$  under the  $MedErr$  metric (lower is better).

Method	aero	bike	boat	bottle	bus	car	chair	dtable	mbike	sofa	train	tv	Mean
K=24	0.84 (0.01)	0.82 (0.01)	0.58 (0.01)	0.95 (0.01)	0.97 (0.01)	0.94 (0.00)	0.89 (0.01)	0.57 (0.04)	0.77 (0.01)	0.91 (0.02)	0.81 (0.01)	0.86 (0.01)	0.8266 (0.0047)
K=50	0.82 (0.01)	0.80 (0.02)	0.59 (0.03)	0.95 (0.00)	0.97 (0.01)	0.94 (0.01)	0.92 (0.01)	0.54 (0.02)	0.78 (0.02)	0.91 (0.02)	0.83 (0.01)	0.89 (0.01)	0.8281 (0.0063)
K=100	0.83 (0.02)	0.76 (0.03)	0.63 (0.01)	0.96 (0.01)	0.97 (0.02)	0.93 (0.00)	0.91 (0.00)	0.57 (0.04)	0.78 (0.02)	0.95 (0.02)	0.82 (0.01)	0.88 (0.01)	0.8335 (0.0067)
K=200	0.84 (0.01)	0.76 (0.02)	0.62 (0.00)	0.96 (0.00)	0.98 (0.00)	0.94 (0.01)	0.92 (0.02)	0.65 (0.04)	0.80 (0.02)	0.96 (0.02)	0.82 (0.00)	0.87 (0.01)	0.8439 (0.0054)

Table 5: Ablation analysis of the size of K-Means dictionary in model  $\mathcal{M}_G$  under the  $Acc_{\frac{\pi}{6}}$  metric (higher is better).

Method	aero	bike	boat	bottle	bus	car	chair	dtable	mbike	sofa	train	tv	Mean
K=4	10.4 (0.5)	13.3 (0.4)	21.9 (1.4)	7.2 (0.3)	2.9 (0.2)	5.3 (0.1)	9.9 (0.1)	16.3 (1.9)	14.1 (0.0)	10.4 (0.7)	5.0 (0.3)	12.5 (0.5)	10.78 (0.35)
K=8	10.5 (0.5)	14.8 (0.5)	21.5 (1.5)	6.8 (0.4)	2.7 (0.1)	4.9 (0.1)	9.7 (0.4)	16.1 (2.9)	14.9 (0.1)	10.2 (0.3)	5.6 (0.3)	12.3 (0.5)	10.85 (0.38)
K=16	9.9 (0.4)	14.3 (0.5)	21.3 (0.7)	7.3 (0.1)	2.7 (0.1)	4.9 (0.1)	9.6 (0.3)	13.0 (3.8)	14.7 (0.5)	10.8 (1.2)	5.2 (0.3)	11.7 (0.3)	10.46 (0.28)
K=24	9.7 (0.2)	15.3 (0.8)	23.5 (1.0)	7.1 (0.4)	2.9 (0.1)	5.0 (0.2)	10.0 (0.5)	13.3 (2.6)	14.4 (0.3)	11.3 (0.5)	5.3 (0.5)	13.1 (0.4)	10.91 (0.22)

Table 6: Ablation analysis of the size of K-Means dictionary in model  $\mathcal{M}_G+$  under the *MedErr* metric (lower is better).

Method	aero	bike	boat	bottle	bus	car	chair	dtable	mbike	sofa	train	tv	Mean
K=4	0.85 (0.00)	0.80 (0.00)	0.61 (0.02)	0.97 (0.00)	0.97 (0.01)	0.95 (0.00)	0.87 (0.02)	0.67 (0.08)	0.84 (0.01)	0.93 (0.02)	0.83 (0.00)	0.86 (0.01)	0.8453 (0.0085)
K=8	0.83 (0.01)	0.79 (0.02)	0.60 (0.01)	0.97 (0.00)	0.96 (0.01)	0.95 (0.00)	0.91 (0.01)	0.62 (0.04)	0.81 (0.02)	0.95 (0.00)	0.83 (0.00)	0.89 (0.00)	0.8427 (0.0017)
K=16	0.84 (0.01)	0.82 (0.02)	0.61 (0.02)	0.96 (0.00)	0.98 (0.00)	0.96 (0.00)	0.92 (0.01)	0.67 (0.07)	0.82 (0.01)	0.97 (0.01)	0.82 (0.01)	0.90 (0.02)	0.8553 (0.0035)
K=24	0.87 (0.01)	0.80 (0.01)	0.60 (0.01)	0.96 (0.00)	0.97 (0.01)	0.95 (0.00)	0.90 (0.01)	0.65 (0.04)	0.83 (0.01)	0.94 (0.01)	0.82 (0.01)	0.87 (0.01)	0.8467 (0.0054)

Table 7: Ablation analysis of the size of K-Means dictionary in model  $\mathcal{M}_G+$  under the *Acc $\frac{\pi}{6}$*  metric (higher is better).

Method	aero	bike	boat	bottle	bus	car	chair	dtable	mbike	sofa	train	tv	Mean
$\alpha = 0.1$	11.8 (0.1)	16.1 (0.8)	20.8 (0.5)	8.4 (0.2)	3.3 (0.1)	6.4 (0.1)	10.6 (0.2)	28.5 (11.7)	15.0 (1.0)	11.2 (0.7)	6.0 (0.1)	12.3 (0.7)	12.53 (1.00)
$\alpha = 1$	12.1 (0.6)	16.0 (0.5)	19.9 (1.1)	8.9 (0.2)	3.4 (0.2)	6.5 (0.3)	10.8 (0.0)	15.2 (6.0)	16.4 (0.5)	9.6 (0.6)	5.9 (0.9)	13.0 (0.7)	11.48 (0.23)
$\alpha = 10$	12.1 (0.8)	14.5 (0.1)	22.8 (0.5)	8.7 (0.3)	3.1 (0.0)	6.5 (0.3)	10.9 (0.4)	15.1 (0.5)	16.3 (0.8)	10.6 (0.5)	6.0 (0.5)	13.0 (0.3)	11.63 (0.24)

Table 8: Ablation analysis of the weighting parameter  $\alpha$  in model  $\mathcal{M}_G$  under the *MedErr* metric (lower is better).

Method	aero	bike	boat	bottle	bus	car	chair	dtable	mbike	sofa	train	tv	Mean
$\alpha = 0.1$	0.84 (0.01)	0.77 (0.02)	0.62 (0.01)	0.96 (0.01)	0.96 (0.01)	0.94 (0.01)	0.91 (0.02)	0.51 (0.04)	0.82 (0.01)	0.96 (0.02)	0.81 (0.01)	0.88 (0.01)	0.8306 (0.0049)
$\alpha = 1$	0.83 (0.02)	0.76 (0.03)	0.63 (0.01)	0.96 (0.01)	0.97 (0.02)	0.93 (0.00)	0.91 (0.00)	0.57 (0.04)	0.78 (0.02)	0.95 (0.02)	0.82 (0.01)	0.88 (0.01)	0.8335 (0.0067)
$\alpha = 10$	0.82 (0.02)	0.79 (0.01)	0.59 (0.01)	0.96 (0.00)	0.97 (0.01)	0.94 (0.01)	0.91 (0.01)	0.67 (0.00)	0.81 (0.02)	0.95 (0.02)	0.82 (0.00)	0.88 (0.01)	0.8424 (0.0060)

Table 9: Ablation analysis of the weighting parameter  $\alpha$  in model  $\mathcal{M}_G$  under the *Acc $\frac{\pi}{6}$*  metric (higher is better).

Method	aero	bike	boat	bottle	bus	car	chair	dtable	mbike	sofa	train	tv	Mean
$\alpha = 0.1$	10.3 (0.6)	16.0 (2.0)	24.0 (2.1)	7.1 (0.1)	3.2 (0.5)	5.5 (0.8)	10.3 (1.3)	11.8 (0.9)	15.2 (0.9)	10.6 (1.1)	6.0 (0.9)	12.3 (0.8)	11.01 (0.92)
$\alpha = 1$	9.9 (0.4)	14.3 (0.5)	21.3 (0.7)	7.3 (0.1)	2.7 (0.1)	4.9 (0.1)	9.6 (0.3)	13.0 (3.8)	14.7 (0.5)	10.8 (1.2)	5.2 (0.3)	11.7 (0.3)	10.46 (0.28)
$\alpha = 10$	8.5 (0.0)	14.8 (1.2)	20.5 (1.4)	7.0 (0.4)	3.1 (0.2)	5.1 (0.1)	9.3 (0.1)	11.3 (2.0)	14.2 (0.6)	10.2 (0.3)	5.6 (0.1)	11.7 (0.2)	10.10 (0.38)

Table 10: Ablation analysis of the weighting parameter  $\alpha$  in model  $\mathcal{M}_G+$  under the *MedErr* metric (lower is better).

Method	aero	bike	boat	bottle	bus	car	chair	dtable	mbike	sofa	train	tv	Mean
$\alpha = 0.1$	0.86 (0.02)	0.81 (0.02)	0.59 (0.04)	0.96 (0.01)	0.98 (0.01)	0.95 (0.00)	0.90 (0.03)	0.65 (0.04)	0.80 (0.01)	0.96 (0.01)	0.82 (0.01)	0.89 (0.00)	0.8473 (0.0102)
$\alpha = 1$	0.84 (0.01)	0.82 (0.02)	0.61 (0.02)	0.96 (0.00)	0.98 (0.00)	0.96 (0.00)	0.92 (0.01)	0.67 (0.07)	0.82 (0.01)	0.97 (0.01)	0.82 (0.01)	0.90 (0.02)	0.8553 (0.0035)
$\alpha = 10$	0.87 (0.01)	0.81 (0.01)	0.64 (0.01)	0.96 (0.00)	0.97 (0.01)	0.95 (0.01)	0.92 (0.01)	0.67 (0.10)	0.85 (0.01)	0.97 (0.01)	0.82 (0.01)	0.88 (0.00)	0.8588 (0.0111)

Table 11: Ablation analysis of the weighting parameter  $\alpha$  in model  $\mathcal{M}_G+$  under the  $Acc_{\frac{\pi}{6}}$  metric (higher is better).

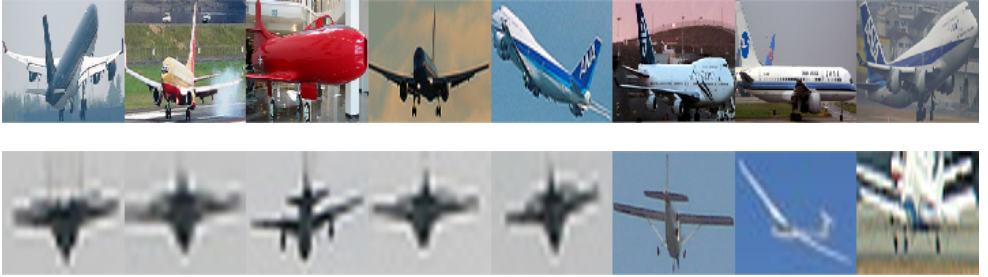


Figure 1: Best (top row) and Worst (bottom row) images for Category: Aeroplane

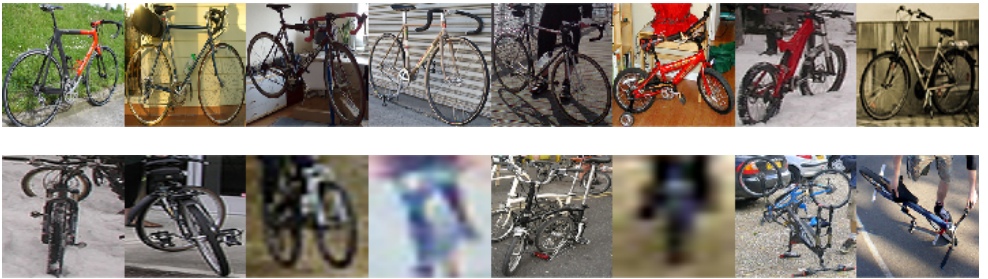


Figure 2: Best (top row) and Worst (bottom row) images for Category: Bicycle



Figure 3: Best (top row) and Worst (bottom row) images for Category: Boat



Figure 4: Best (top row) and Worst (bottom row) images for Category: Bottle



Figure 5: Best (top row) and Worst (bottom row) images for Category: Bus



Figure 6: Best (top row) and Worst (bottom row) images for Category: Car



Figure 7: Best (top row) and Worst (bottom row) images for Category: Chair

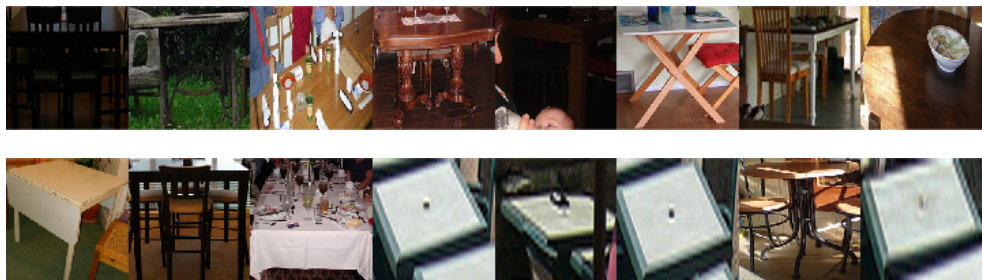


Figure 8: Best (top row) and Worst (bottom row) images for Category: Diningtable



Figure 9: Best (top row) and Worst (bottom row) images for Category: Motorbike



Figure 10: Best (top row) and Worst (bottom row) images for Category: Sofa



Figure 11: Best (top row) and Worst (bottom row) images for Category: Train



Figure 12: Best (top row) and Worst (bottom row) images for Category: Tvmonitor

## References

- [1] Alexander Grabner, Peter M. Roth, and Vincent Lepetit. 3d pose estimation and 3d model retrieval for objects in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [2] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Kosecka. 3d bounding box estimation using deep learning and geometry. In *CVPR*, 2017.
- [3] Hao Su, Charles R. Qi, Yangyan Li, and Leonidas J. Guibas. Render for cnn: Viewpoint estimation in images using cnns trained with rendered 3d model views. In *IEEE International Conference on Computer Vision*, 2015.
- [4] Shubham Tulsiani and Jitendra Malik. Viewpoints and keypoints. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.