

Supplemental Material

In this supplemental material, we provide additional analyses for the classification experiments outlined in the main paper. Furthermore, we provide confusion matrices for the detection experiment on the full image resolution presented in the main paper, as well as a further detection experiment done on a reduced image resolution of 640×480 pixels.

Classification

Figure 6 depicts sample crops used in classification, showing the padding of 20 pixels added to capture relevant context for classification.

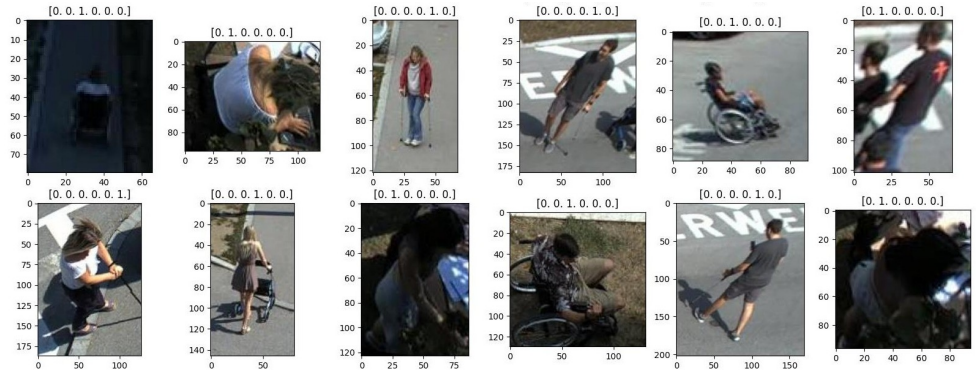


Figure 6: Example crops used for classification, showing the added padding to capture relevant context. The respective title shows the one-hot encoding of the class label (*[pedestrian, wheelchair, rollator, crutch, cane]*).

We provide a comparison of the feature backbones used in the classification experiment in Table 5, providing information on their performance on ImageNet, their complexity as well as final validation accuracy reached in fine-tuning on our dataset. Although the difference in model complexity spans over a decade in terms of number of parameters and almost two decades in floating point operations per second, their final difference in validation accuracy is less than 2 percentage points, allowing to chose smaller, faster implementations.

Figure 7 shows the confusion matrices of the test split for all classifiers. Confusion wise, all classifiers perform the same, with the biggest mix-up existing between classes *crutch* and *cane*. This is not surprising, considering that the visual structures for walking canes and crutches are both very thin and highly similar, differing mostly in the handle.

Detection

We present the confusion matrices for the detection experiments on image size 1280×960 in Figure 8. In each cell, the first line shows the percentage of ground truth labels classified into this cell, the second line denotes the respective number of instances. In general, the biggest confusion exists between the classes *crutch* and *cane*, which both represent small structures mostly distinguished by the grip. A further source for misclassification is these two classes being classified as *pedestrian*, i.e. not using any mobility aid. This is not a surprising result, since both classes represent thin structures, challenging to be captured by a

	Backbone				Fine-tuned
	Acc@1 [↑]	Acc@5 [↑]	Params [↓]	GFLOPs [↓]	ValAcc [↑]
MobileNetV3 L	0.740	0.913	5.5M	0.22	0.954
MobileNetV3 L (V2)	0.740	0.913	5.5M	0.22	0.958
ResNet18	0.698	0.891	<i>11.7M</i>	<i>1.81</i>	0.951
ResNet34	0.733	0.914	21.8M	3.66	0.953
ResNet50	0.761	0.929	25.6M	4.09	0.963
ResNet50 (V2)	<i>0.809</i>	<i>0.954</i>	25.6M	4.09	0.962
DenseNet201	0.770	0.934	20.0M	4.29	<i>0.967</i>
ResNet152	0.783	0.940	60.2M	11.51	0.965
ResNet152 (V2)	0.823	0.960	60.2M	11.51	0.968
ViT-B/16	0.811	0.953	86.6M	17.56	0.955
VGG16	0.716	0.904	138.4M	15.47	0.953

Table 5: Comparison of the classification backbones used with their original performance on ImageNet in terms of top-1 accuracy (Acc@1) and top-5 accuracy (Acc@5), model complexity in terms of millions of model parameters (Params) and billions of floating point operations per second (GFLOPs) as well as final accuracy on the validation split of our dataset (ValAcc). Best values are marked **bold**, second best *italic*.

deep neural network. Furthermore, we can see that the hierarchical training approach vastly reduces the number of false positive detections, especially so for mobility aid use and the thin mobility aids *crutch* and *cane*, irrespective of model size.

The results for detection on smaller images of size 640×480 are listed in Table 6 in terms of average precision (AP), and in Table 7 in terms of accuracy (ACC), misclassification rate (MIS) and missed detections (MIS). We observe a general performance drop of about 4 – 5 percentage points over all models. Classes *crutch* and *cane* are affected more severe, with *cane* faring even worse than *crutch*. This observation can be explained largely due to the respective mobility aids having only a long and thin footprint within the images, thereby making detection and classification even more challenging with a smaller image resolution. This interpretation is consistent with the fact that classification of class *wheelchair*, which has a comparatively large footprint, does not show any significant performance degradation. Model performance in terms of missed detections (false negatives, FN) and wrong detections (false positives FP in Table 7) stays more or less the same for training with a class hierarchy, models trained with independent classes fare worse for wrong detections, yet a bit better for missed detections. The small model trained with class hierarchy seems to be an outlier in this regard. The confusion matrices presented in Figure 9 show the same qualitative results observed for detectors trained on the large image size, again underlining the benefit of providing the detectors with a safe fall-back by using class hierarchies.

YOLOv5	pedestrian		wheelchair		rollator		crutch		cane		all	
	AP@50	mAP@50-95	AP@50	mAP@50-95	AP@50	mAP@50-95	AP@50	mAP@50-95	AP@50	mAP@50-95	mAP@50	mAP@50-95
small	0.811	0.565	0.979	0.661	0.912	0.639	0.651	0.525	0.383	0.311	0.747	0.540
medium	0.842	0.612	0.980	0.693	0.912	0.648	0.725	0.604	0.471	0.396	0.786	0.591
large	0.857	0.631	0.977	0.693	<i>0.931</i>	0.672	0.742	0.617	0.499	0.413	0.801	<i>0.605</i>
xlarge	<i>0.849</i>	<i>0.626</i>	0.982	0.697	0.938	0.679	0.716	0.599	<i>0.490</i>	0.413	<i>0.795</i>	0.603
h small	0.758	0.508	0.965	0.667	0.868	0.612	0.629	0.529	0.364	0.311	0.717	0.525
h medium	0.812	0.584	0.979	<i>0.716</i>	0.883	0.640	0.661	0.569	0.486	<i>0.420</i>	0.764	0.586
h large	0.812	0.600	<i>0.981</i>	0.734	0.923	0.695	0.671	0.583	0.489	0.428	0.775	0.608
h xlarge	0.799	0.597	0.973	0.689	0.920	<i>0.690</i>	0.632	0.547	0.462	0.401	0.757	0.585

Table 6: Detector performance for models trained on images resolution 640×480 pixels. Models in the top have been trained with independent classes, models prepended with 'h' in the bottom part have been trained with the hierarchical class structure. The **best** model has been marked in **bold face**, the *second best* in *italic*.

YOLOv5	pedestrian			wheelchair			rollator			crutch			cane			all			
	ACC [†]	MCL [‡]	MIS [‡]	ACC [†]	MCL [‡]	MIS [‡]	ACC [†]	MCL [‡]	MIS [‡]	ACC [†]	MCL [‡]	MIS [‡]	ACC [†]	MCL [‡]	MIS [‡]	ACC [†]	MCL [‡]	FP [‡]	FN [‡]
small	0.858	<i>0.089</i>	<i>0.054</i>	0.958	0.016	0.026	0.912	0.075	0.013	0.491	0.500	<i>0.009</i>	0.332	0.652	0.017	0.796	0.171	629	235
medium	0.842	0.097	0.061	0.969	0.008	0.023	0.877	0.103	0.020	<i>0.587</i>	<i>0.400</i>	0.012	0.461	0.522	0.017	0.811	0.152	561	262
large	<i>0.849</i>	0.094	0.057	0.962	<i>0.010</i>	0.028	0.919	0.064	0.018	0.577	0.409	0.014	0.519	0.463	0.019	0.820	<i>0.143</i>	511	259
xlarge	0.842	0.087	0.071	0.964	0.015	<i>0.022</i>	0.935	0.047	0.018	0.590	0.397	0.012	0.454	0.520	0.026	<i>0.816</i>	0.142	521	292
h small	0.802	0.134	0.064	0.947	0.024	0.030	0.856	0.127	0.017	0.486	0.504	0.010	0.332	0.656	0.013	0.762	0.199	353	274
h medium	0.829	0.118	0.053	0.962	0.022	0.016	0.869	0.118	0.013	0.538	0.455	0.007	0.448	0.541	<i>0.011</i>	0.795	0.175	256	212
h large	0.837	0.107	0.056	<i>0.965</i>	0.019	0.016	<i>0.920</i>	0.065	<i>0.015</i>	0.579	0.407	0.014	<i>0.487</i>	<i>0.504</i>	0.009	0.814	0.154	266	229
h xlarge	0.804	0.115	0.080	0.961	0.015	0.024	0.917	<i>0.057</i>	0.025	0.515	0.469	0.016	0.437	0.528	0.015	0.787	0.166	227	330

Table 7: Detector performance in terms of accuracy (ACC), misclassification (MCL) and missed (MIS) detections as well as overall false negatives (FN) and false positive (FP) de- tections for detectors trained on image size 640×480 .

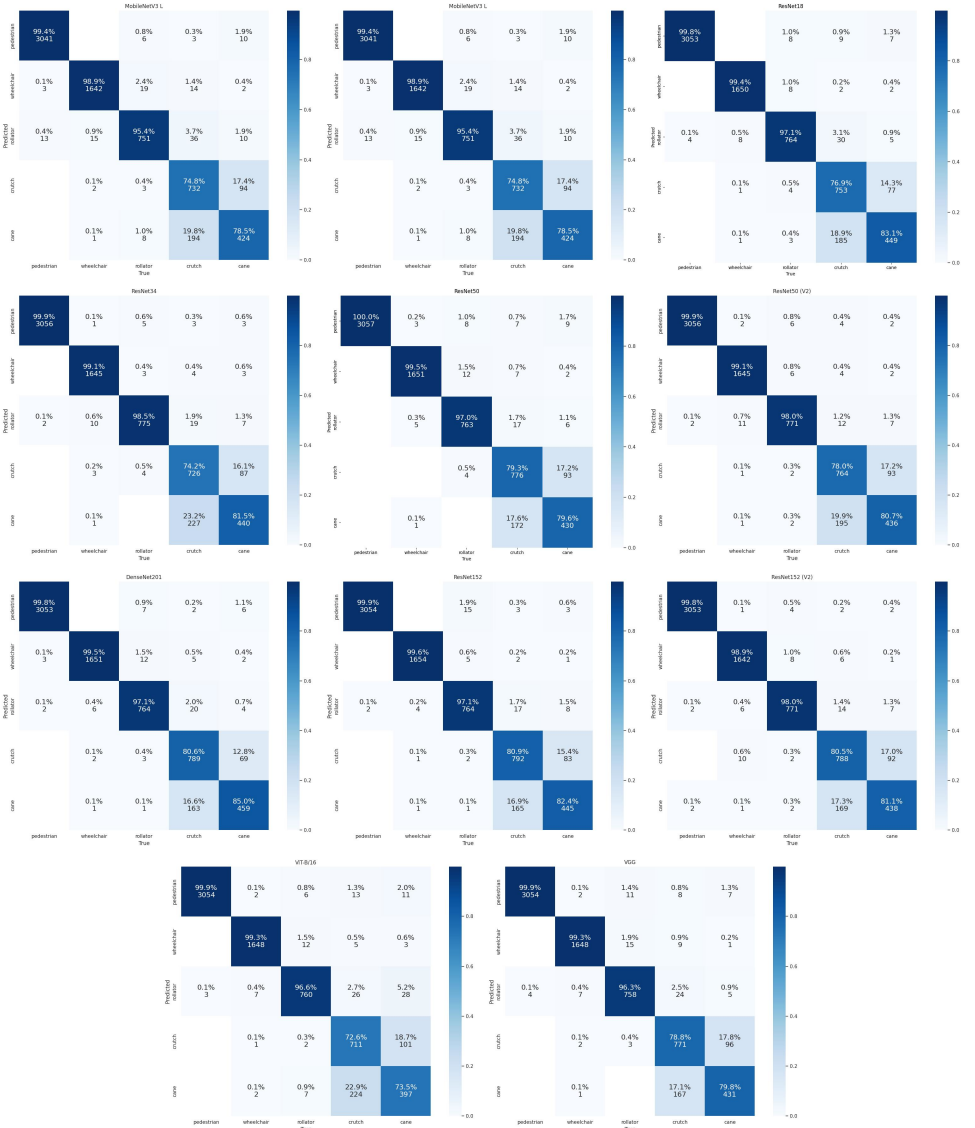


Figure 7: Confusion matrices of the classifiers compared. The biggest mix-up exists for the classes *Crutch* and *Cane* across all classifiers.



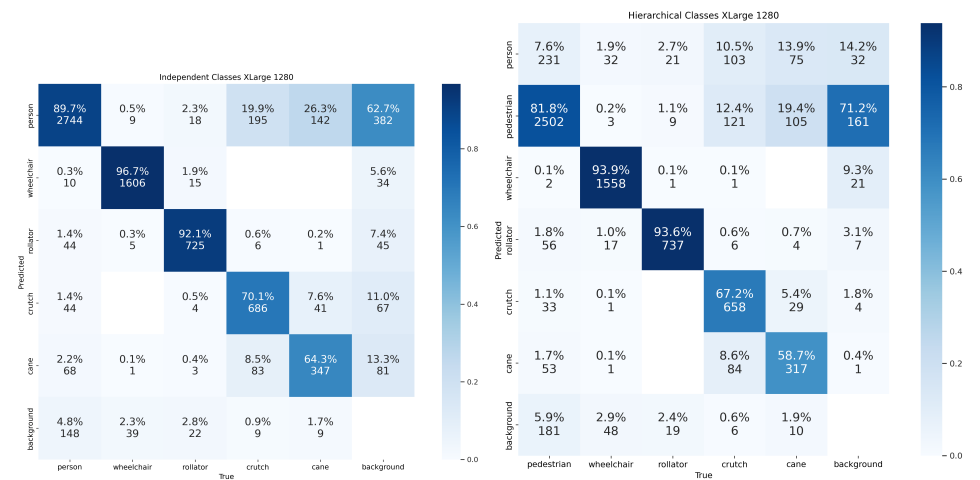


Figure 8: Confusion matrices for YOLOv5 detectors trained on image size 1280×960 . The left column shows the results for models trained with independent classes, the right column shows the results for those trained with hierarchical classes.

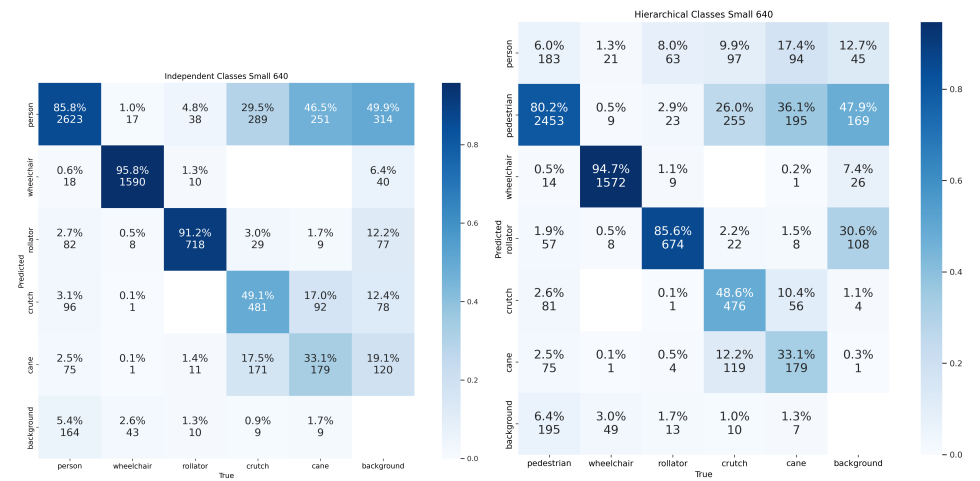




Figure 9: Confusion matrices for YOLOv5 detectors trained on image size 640×480 . The left column shows the results for models trained with independent classes, the right column shows the results for those trained with hierarchical classes.