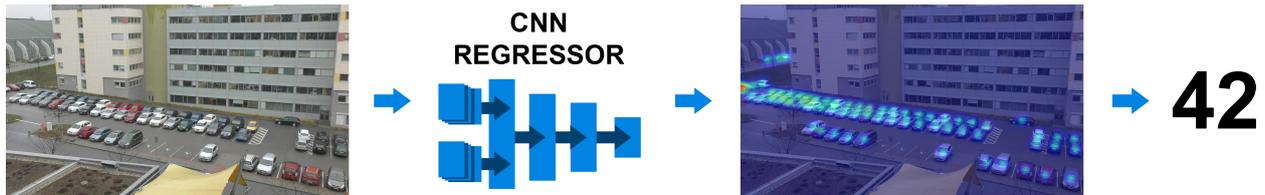


Counting Vehicles in Static Images

Ondřej Zemánek



Abstract

This paper addresses the problem of counting vehicles in static images with no geometric information of the scene. Four convolutional neural network architectures were studied, implemented and trained as a main part of this work. Also, a dataset that consists of 19 310 images in total from 12 views that captures 7 different scenes were taken as part of this work. The trained networks map the appearance of the input sample to its corresponding vehicles density map, which can be easily translated to the vehicle count with keeping the localization of the vehicles in the input image. The main contribution of this work is in an application and a comparison of the state-of-the-art solutions to the problem of object counting. Most of them were mainly designed to count pedestrians in crowded scenes or for medicine images, so the major goal was to adapt these solutions for vehicle counting task. The implemented models were trained on TRANCOS dataset which is a popular benchmark for counting vehicles on annotated low quality highway pictures. Their performance is compared and the results are discussed.

Keywords: visual counting, vehicle counting in static images, car park dataset

Supplementary Material: N/A

*xzeman53@stud.fit.vutbr.cz, Faculty of Information Technology, Brno University of Technology

1. Introduction

Visual counting that aims to accurately estimate the number of vehicles is a hard problem. But its potential is huge in many applications across many industries.

For instance, it can help truck drivers, who need to plan their next break, by monitoring parking capacity near highways. Another application can be long-term analysis of a traffic density on main city roads and highways, so road closures, detours or road expansion can be planned easily and smoothly. Also, solving the vehicle counting problem can bring a cheaper solution for monitoring shopping center parking lots, so instead of using a physical sensor for each parking space, a few cameras can be used to monitor the parking lot.

The most recent state-of-the-art solutions are based mainly on the convolutional neural network model.

Therefore, this work is focused only on these approaches. The best approaches can be divided into objects detection approaches and density map regression based approaches.

The first group uses classification of individual objects in YOLO-like (You Look Only Once) [1] style to detect and count the objects in the input image. Although these approaches can be fast and reliable in trivial cases, in very dense and overcrowded scenes like images with overlapping objects, low-resolution, partly visible objects, images with slightly unseen perspective, the overall performance of these models is limited.

The other group of solutions that reaches much better results in the target scenarios is based on a different idea. Instead of solving this problem by detection of

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32

33 each object in the image, these solutions transform the
 34 visual counting task into object density map estima-
 35 tion from the input image. In other words, they are
 36 using convolutional neural networks to transform an
 37 input image appearance into an object density map in
 38 a certain resolution. From the output of this transfor-
 39 mation, the object count can be easily estimated by the
 40 output density map integration even with keeping the
 41 information of the objects localization.

42 The main contribution of this work is an appli-
 43 cation and a comparison of the existing solutions on
 44 parking lot dataset. To achieve this, I had to analyse
 45 the existing convolutional neural architectures, adapt
 46 these models to vehicle counting problem, create a
 47 large and diverse dataset for training, train them in
 48 with various parameters, and finally, evaluate them.

49 2. Implemented Architectures

50 The state-of-the-art approaches for visual counting are
 51 mainly built upon one of two concepts: density map re-
 52 gression and detection. Thus, the chosen architectures
 53 are designed in this way.

54 First three solutions are based on Density map
 55 regression. So, the input image is regressed into den-
 56 sity map that represents spatial distribution of objects.
 57 Then, this map is integrated into an object count which
 58 corresponds to the input sample.

59 The last studied approach is somewhere in the
 60 middle of these two concepts. It combines density map
 61 regression with classification. More in subsection 2.4.

62 Tensorflow open source platform was used to im-
 63 plement these convolutional neural networks. The
 64 final models of the Counting CNN, the Hydra CNN
 65 and Spatial Division and Conquer Network were in-
 66 spired by the original authors implementations, that
 67 were implemented in the Caffe framework and the Py-
 68 Torch platform respectively. The authors of Stacked
 69 Hourglass model provides only brief description of
 70 the implementation, so the network implementation is
 71 based on this description only.

72 Despite the fact that the authors shares the imple-
 73 mentation details, the training process implementation
 74 is tied its application and the dataset. Thus, the train-
 75 ing process for each used network has to be created
 76 manually based on the deep analysis of the architec-
 77 ture.

78 2.1 From Human Pose Estimation to Visual 79 Counting

80 The Stacked Hourglass architecture proposed by Ne-
 81 well at al. [2] is a composition of multiple modules
 82 called hourglass. Each Hourglass module processes

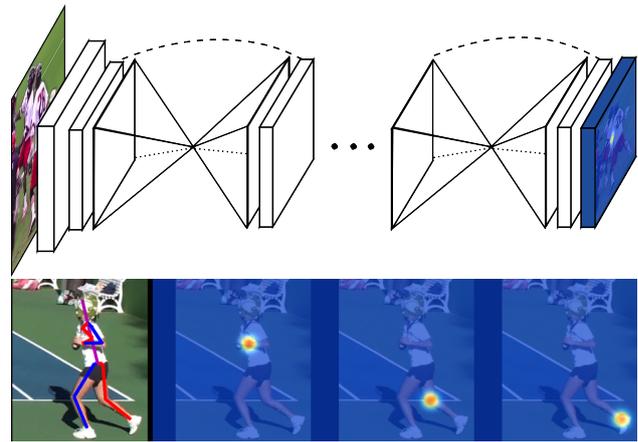


Figure 1. top: Stacked hourglass concept – input image is passed through multiple hourglass modules. Each hourglass model creates an intermediate prediction to improve the final result. **bottom:** Output of the multiple-level hourglass model for Human pose estimation problem [2].

the input features on multiple scales and consolidates
 them to best capture the object landmarks. This is
 done by repeated up-down/bottom-up process. Also,
 it uses an intermediate supervision process, i.e. the
 architecture uses intermediate prediction to improve
 the next prediction. This is done by skipping layers.
 This model (Figure 1) aims on the problem of human
 pose estimation, but as it is shown in the achieved
 results (Sec. 5), it also shows good results in the visual
 counting problem.

2.2 Single-Pipeline CNN with Great Results

The Counting convolutional neural network created by
 Oñoro et al. [3] is a simple sequence of 6 convolutional
 layers and two max-pooling layers, as can be seen in
 Figure 2. The input image patch with size 72×72
 pixels is processed by this sequence and it returns a
 density map with 18×18 pixels size which represents
 the object spatial distribution in the image. Finally,
 the object count is gathered by integrating the density
 map.

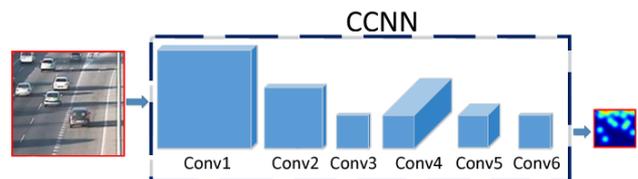


Figure 2. Counting CNN architecture – sequence of 6 convolutional layers in an up-down order [3].

2.3 Even Better Results with Multi-Scale In-puts

Next implemented architecture is called Hydra CNN,
 like the mythological creature Hydra with nine heads

83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102

103
104
105
106

107 and a big body. It is mainly based on multiple Counting
 108 CNN's that are used as the input heads of the creature
 109 as illustrated in Figure 3. Each head processes the
 110 input sample on a different scale, so the final architec-
 111 ture is a scale-aware solution for visual counting. The
 112 patch scale (crop ratio of the original sample) for head
 113 H_i is defined as follows,

$$H_s = 1 - \frac{1}{C} \cdot H_i, \quad (1)$$

114 where C is the number of heads. In case of the first
 115 level head, the input patch corresponds to the original
 116 sample. The heads' intermediate outputs are fused by
 117 three fully connected layers, so even the input array
 118 contains the image in multiple scales, the output is a
 119 single density map like in the case of Counting CNN.
 120 The concept can be easily extended by adding a new
 121 head or simplified by removing one. The results are
 122 much better than in the case of the simple counting
 123 CNN.

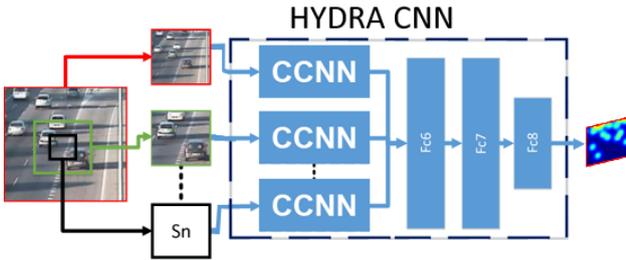


Figure 3. Hydra CNN scheme – input sample is up-sampled multiple times and processed by the Counting CNNs. The sub-results are merged into a single output map by fully connected layers.

124 2.4 Open-Set to Closed-Set Transformation as 125 the State of the Art

126 Previously described methods are modeled in a regres-
 127 sion manner. Xiong et al. [4] proposed a new different
 128 approach with their Spatial Divide-and-Conquer Net-
 129 work (S-DCNet) that is more complex and it uses mul-
 130 tiple modules with different purposes. The main idea
 131 is spatial division of the input image into small regions,
 132 each with a closed set of a defined range, so they can
 133 transform quantity to intervals which the network can
 134 easily classify.

135 It is based on the main idea that the visual counting
 136 problem is an open-set problem by nature. But only
 137 limited and closed set labeled counts can be observed
 138 in reality. So the goal is to transform the original
 139 problem into a close-set one.

140 This is done by spatial division of the image until
 141 the count of every part is in a specified range. For
 142 example, the experimented range for vehicle counting I

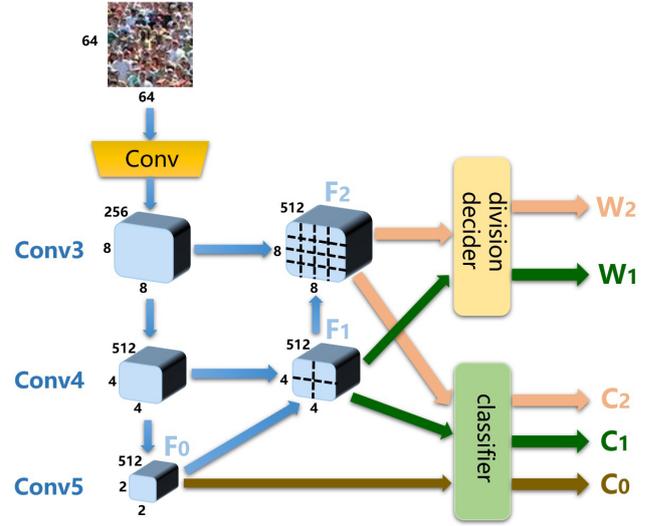


Figure 4. The architecture of Spatial division-and-conquer network. The first 5 layers are convolutional and follow the VGG16 concept [5]. These layers are then connected to the classifier and division decoder straight forward or preceded by fusion with another layer. Finally, the output weights and classification are processed to get the object density map.

143 am using is $[0, 5]$. That means the architecture spatially
 144 divides the input image, until every part reaches 5
 145 object at max.

146 The division decoder in Figure 4 is trained to decide
 147 whether it is necessary to divide the input or not. It
 148 returns W_i weights in range $[0, 1]$, where greater value
 149 means higher need for division. The weight is then
 150 used to compute the division result DIV_i . So, higher
 151 resolution count maps C_i are applied if the division
 152 weight W_i is higher.

153 Next, the classifier module predicts the object counts
 154 C_i for each output feature map. The classification is
 155 done on specified intervals in the closed-set range, i.e.
 156 for each feature map, a single object count is obtained.

157 After the model prediction, following post-process
 158 is applied:

$$DIV_i = (1 - W_i) \circ avg(C_{i-1}) + W_i \circ C_i, \quad (2)$$

159 where “ \circ ” denotes the Hadamard product and avg
 160 is an averaging re-distribution operator. Finally, pre-
 161 dicted object count is represented by last DIV_i map.
 162 The i level corresponds to max division time which is
 163 the value that limits the input image division. So, if
 164 the i is 2, then the model takes into account up-to 2
 165 times divided input image.

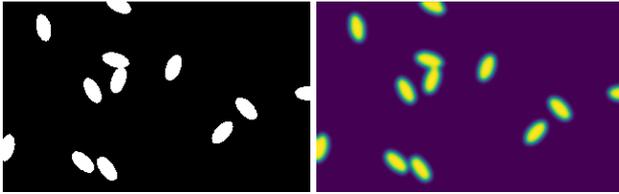


Figure 5. Toy dataset example of input sample (left) and output density map (right).



Figure 6. Example of the labeling style. Scribbles are used for high-resolution vehicles and dots for partially visible cars or vehicles in distant.

166 3. Training Details

167 The authors of the proposed solutions are proving
168 their network performance with pretrained model that
169 has been successfully trained on available benchmark
170 dataset. Unfortunately, they do not always share train-
171 ing details.

172 For instance, the paper about proposed Spatial
173 division-and-conquer architecture contains very vague
174 information about the training process. Therefore, to
175 train this model, it was necessary to understand the
176 architecture in depth and experiment with model pa-
177 rameters.

178 3.1 Iterative Training Process With a Toy Da- 179 taset

180 To train the models on large and diverse dataset, it is
181 necessary to know how the models performs with dif-
182 ferent training settings and parameters. As the fastest
183 approach to make the training process work in the way
184 how it was designed is training on a simplified, “toy”
185 dataset.

186 Toy dataset, as can be seen in Figure 5 is a custom
187 generated dataset which can be easily created in a short
188 time. Its parameters, like Gaussian blur variance σ ,
189 resolution and object shape, can be modified for each
190 convolutional neural network architecture. The main
191 benefit of this simple dataset is that the training process
192 is much faster than training on real images.

193 3.2 Output Activation and Loss Functions

194 Even with a fast training process, there is another im-
195 portant factor influencing training success. That is
196 the combination of the last activation function and the
197 loss function. The last activation function gives the
198 transformation of the linear output value and the loss
199 function is used to compute the trained model error.

200 Unfortunately, there is no general-purpose com-
201 bination of these functions. So, as the implemented
202 architectures do not have the same output format, it
203 was necessary to understand the model pipeline and
204 decide which combination is the best for each output.

205 For classification problem is common to use sig-
206 moid or soft-max function as the last activation func-
207 tion and some type of cross-entropy function. Also,

208 There are multiple cross-entropy loss functions for dif-
209 ferent purpose, like multi-class or binary classification.

210 For regression problem, it is typical to use linear
211 function as the activation function and use L1-norm
212 (mean absolute error) or L2-norm (mean square error)
213 as loss functions (density map). Otherwise, the
214 activation function can be set to sigmoid with use of
215 quadratic loss function so the regressed value is in
216 range $[0, 1]$ (weight, normalized output)¹.

217 The Counting CNN, Hydra-CNN model and the
218 Stacked Hourglass model are using a combination of
219 linear activation function and mean square error loss
220 function.

221 The Spatial Division-and-Conquer model is much
222 more complex. The network output consist of the divi-
223 sion weights and quantity interval classifications. The
224 division weight is a regression to values between 1 and
225 0, so the combination of sigmoid activation function
226 and mean square error loss function are implemented.
227 The count interval problem is a multi-class classifica-
228 tion, so the soft-max activation function with mean
229 absolute error loss function is applied.

230 3.3 Ground Truth Labeling

231 For ground truth labels is used dotted annotation blur-
232 red by 2D Gaussian function. So, the vehicles in
233 images are labeled by only single dot in first step.
234 Then the dotted map is then blurred with Gaussian
235 blur. Even after blurring the dots out, an integration of
236 the blurred map still corresponds to the vehicle count.
237 During the training process this map is used as the
238 ground truth density map.

239 4. Diverse and Robust Training Dataset

240 To train the implemented network to count vehicles
241 in images, it is crucial to have a big enough dataset

¹<https://medium.com/@phuctrt/189815343d3f>

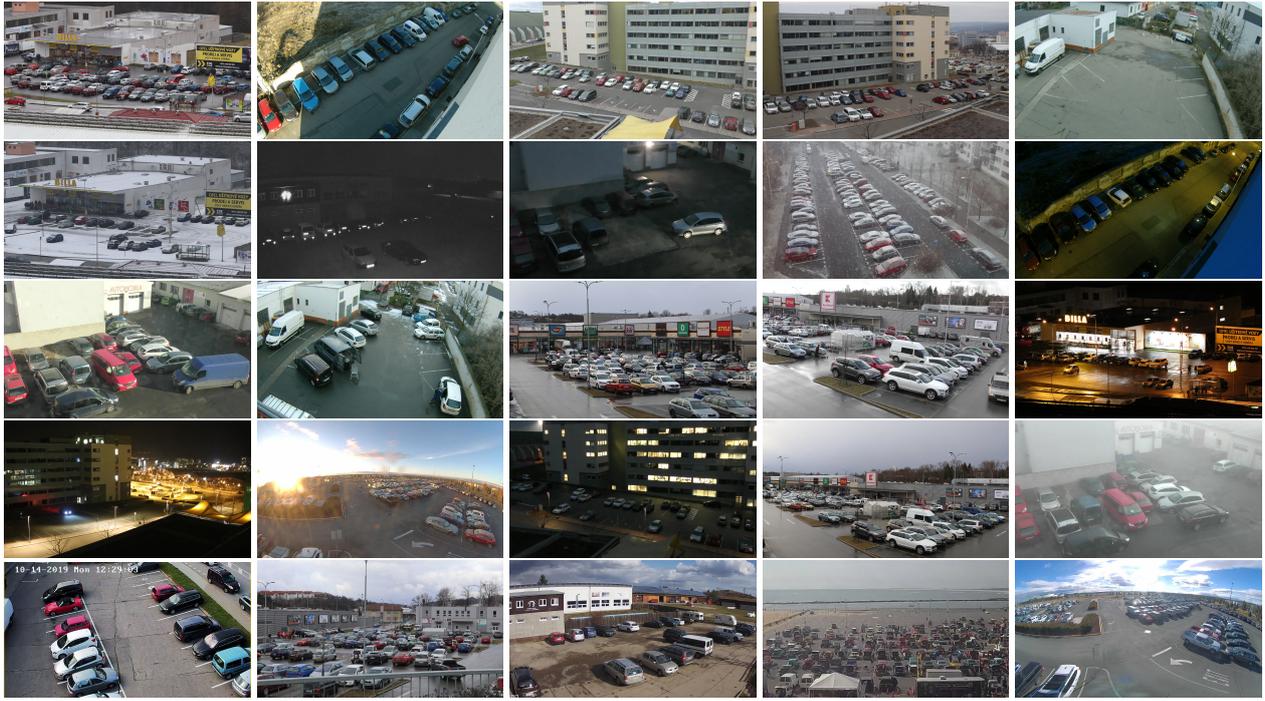


Figure 7. Custom dataset samples.

242 of images with similar parameters. As the parameters
 243 depend mainly on the final application, it is necessary
 244 to define it.

245 Parking lot occupancy monitoring was chosen as
 246 the main application. Therefore, the dataset images
 247 should capture parking areas or similar places like
 248 highways or streets with cars.

249 4.1 Existing Datasets

250 Currently, there are three suitable datasets for this ap-
 251 plication. The first one is the TRAffic ANd COnges-
 252 tionS (TRANCOS) [6] dataset with more than thou-
 253 sand labeled low-quality images from highway. The
 254 next dataset is called CNRPark+EXT² and it captures
 255 occupancy of parking lots with roughly 4300 labeled
 256 images. Lastly, the CARPK³ is a collection of drone
 257 images of huge parking lots with 1500 annotated im-
 258 ages where only a part of it is suitable for our applica-
 259 tion.

260 4.2 Newly Created Dataset

261 As the main goal of this work is an robust real-world
 262 application for vehicle counting problem, we need
 263 much more diverse and robust dataset to train the net-
 264 works. Therefore, a new and more diverse parking lot
 265 dataset was collected as part of this work. Figure 7
 266 shows examples of this custom dataset. It consist of
 267 19310 images in total from 12 views that capture 7
 268 different scenes.

²<http://www.cnrpark.it/>

³<https://lafi.github.io/LPN/>

Each location was captured from a similar angle 269
 to the ground to simulate the common monitoring 270
 cam position. The recording process took place from 271
 September to March, so diverse weather and lighting 272
 conditions were captured. Also, three online webcams 273
 were recorded as part of the custom dataset. This adds 274
 another few thousands of images to this dataset. 275

4.3 Annotation 276

The training cannot be done without the ground truth 277
 labels for the dataset pictures. Several annotation 278
 styles were tested to label the images, like bounding 279
 box, silhouette, scribble. Although the bounding box 280
 annotation is common approach for object detection 281
 and silhouette annotation is even more precise, these 282
 two label styles are too time-consuming and unneces- 283
 sary for our application. Thus, the faster and sufficient 284
 scribble and dots labeling styles were chosen as can 285
 be seen in Figure 6. 286

So far, more than 3500 images were annotated as 287
 part of this work and the labeling process still contin- 288
 ues. 289

5. Achieved Results 290

The presented networks were trained on the TRAN- 291
 COS dataset to demonstrate their performance so far. 292
 The dataset contains training, validation and test sets, 293
 so the results can be accurately compared on samples 294
 that were not used for training. Visual comparison of 295
 the trained models on TRANCOS dataset can be seen 296
 in Figure 8. 297

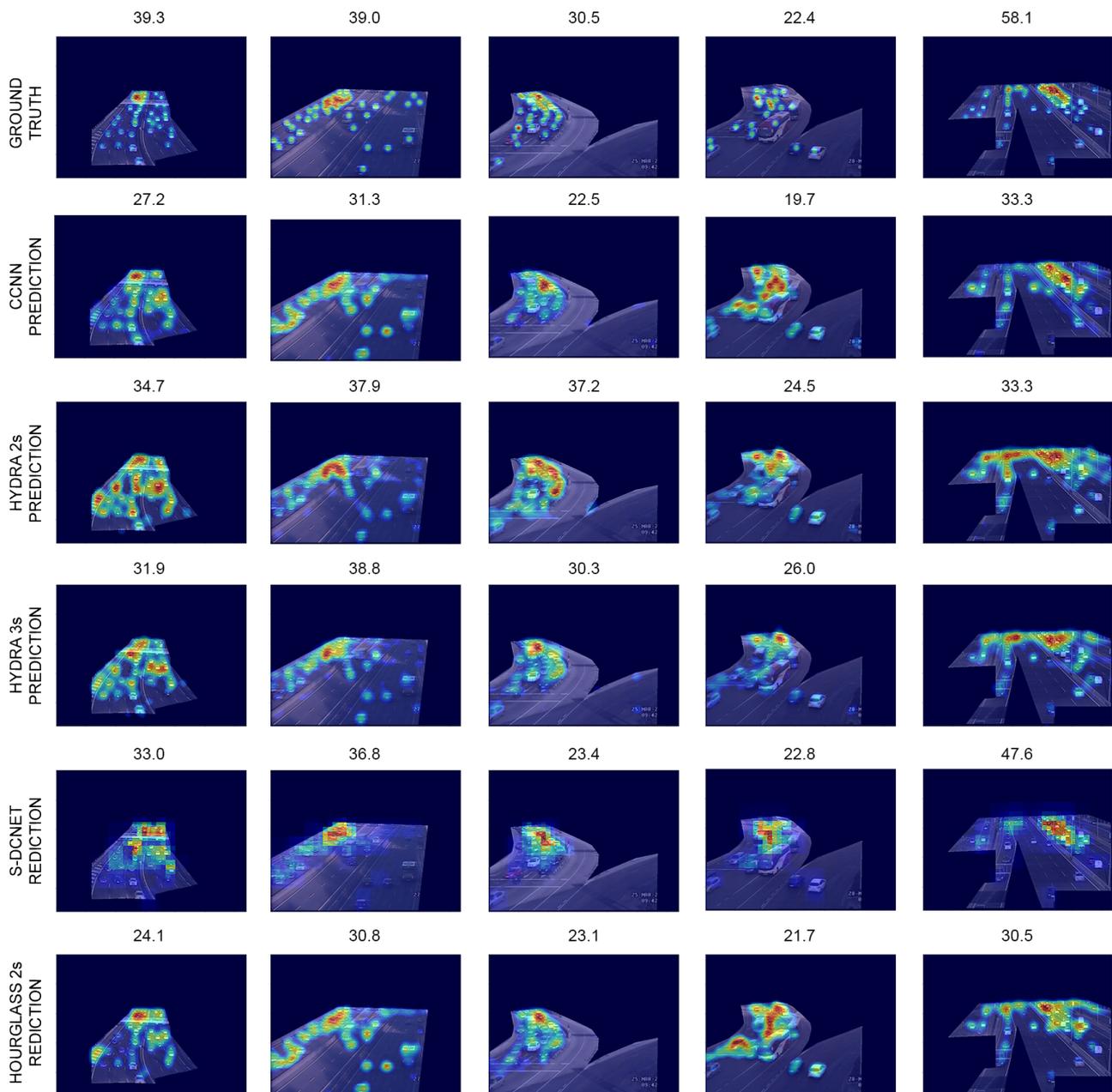


Figure 8. Five test samples from TRANCOS dataset with trained models predictions. Top row corresponds to the target image with ground-truth. “Hydra 2s”, “Hydra 3s”, “Hourglass 2s” stands for the Hydra CNN with 2 heads, 3 heads and the Stacked Hourglass model with 2 stacks respectively. The ground-truth counts are slightly different because Gaussian blur was applied to the label maps.

298 The comparison of the used architectures on full
 299 TRANCOS dataset can be seen in table 1. The evaluation
 300 was done with Grid Average Mean Absolute Error
 301 (GAME) [7] metric defined by

$$GAME(L) = \frac{1}{N} \sum_{n=1}^N \left(\sum_{l=1}^{4^L} |C_{pre}^l - C_{gr}^l| \right), \quad (3)$$

302 where N denotes the number of images, C_{pre}^l and
 303 C_{gr}^l are the predicted and ground-truth count of the
 304 L-th subregion, respectively.

305 The Counting CNN achieves good results in to-

Method	GAME 0	GAME 1	GAME 2	GAME 3
CCNN	12.18	16.44	20.35	22.97
Hydra 2s	10.77	14.28	17.69	21.13
Hydra 3s	11.02	14.37	17.01	20.64
SHG 2s	14.30	15.84	18.23	22.81
S-DCNet	8.56	9.357	10.40	11.83

Table 1. Trained model evaluation with GAME metric on the TRANCOS dataset. The best performance is in boldface. “SHG 2s” stands for the Stacked Hourglass model with 2 stacks

total count prediction, but in higher levels on GAME 306

307 metric shows some false prediction. In case of the
308 Hydra CNN the results are better but the problem with
309 noisy prediction is remains. The best predictions gives
310 the Spatial Divide-and-Conquer Network, which has
311 accurate prediction across all GAME levels, so the
312 spatial prediction is very precise. Lastly, the Stacked
313 hourglass model with 2 stacks shows great results in
314 spatial prediction, but the total density of the predicted
315 map is lower than ground-truth.

316 Next step in comparison of these architectures is
317 the custom dataset evaluation with the GAME metric.
318 However, the training process is still in progress and I
319 don't want to present temporally results.

320 6. Conclusions

321 In this paper, I have shown four different architectures
322 for visual counting that has been implemented, de-
323 scribed the training details of these models. Also, The
324 custom car park dataset with 19 300 images and 3 500
325 already labeled pictures have been presented.

326 The architectures were evaluated on the popular
327 TRANCOS dataset and the results were shown in last
328 chapter.

329 The newly created dataset is being continually up-
330 dated with new locations and more importantly it is
331 being labeled.

332 Also extended version of the Spatial division and
333 conquer network has been recently released by the
334 authors and I am finishing the implementation of this
335 network.

336 Finally, all the presented networks are being trained
337 on the custom dataset and will be evaluated.

338 Acknowledgements

339 I would like to thank my supervisor prof. Ing. Adam
340 Herout, Ph.D. for great support and constructive criti-
341 cism during my work on this paper. Also, I would like
342 to thank to Ing. Jakub Špaňhel for his advice and help
343 with a dataset acquisition.

344 Access to computing and storage facilities owned
345 by parties and projects contributing to the National
346 Grid Infrastructure MetaCentrum provided under the
347 programme "Projects of Large Research, Development,
348 and Innovations Infrastructures" (CESNET LM2015-
349 042), is greatly appreciated.

350 References

351 [1] Joseph Redmon, Santosh Divvala, Ross Girshick,
352 and Ali Farhadi. You only look once: Unified,
353 real-time object detection. In *Proceedings of the*
354 *IEEE conference on computer vision and pattern*
355 *recognition*, pages 779–788, 2016.

[2] Alejandro Newell, Kaiyu Yang, and Jia Deng. 356
Stacked hourglass networks for human pose es- 357
timation. In *European conference on computer* 358
vision, pages 483–499. Springer, 2016. 359

[3] Daniel Onoro-Rubio and Roberto J López-Sastre. 360
Towards perspective-free object counting with 361
deep learning. In *European Conference on Com-* 362
puter Vision, pages 615–629. Springer, 2016. 363

[4] Haipeng Xiong, Hao Lu, Chengxin Liu, Liang Liu, 364
Zhiguo Cao, and Chunhua Shen. From open set 365
to closed set: Counting objects by spatial divide- 366
and-conquer. In *Proceedings of the IEEE Inter-* 367
national Conference on Computer Vision, pages 368
8362–8371, 2019. 369

[5] Karen Simonyan and Andrew Zisserman. Very 370
deep convolutional networks for large-scale im- 371
age recognition. *arXiv preprint arXiv:1409.1556*, 372
2014. 373

[6] Roberto López-Sastre Saturnino Maldon- 374
ado Bascón Ricardo Guerrero-Gómez-Olmedo, 375
Beatriz Torre-Jiménez and Daniel Oñoro-Rubio. 376
Extremely overlapping vehicle counting. In 377
Iberian Conference on Pattern Recognition and 378
Image Analysis (IbPRIA), 2015. 379

[7] Ricardo Guerrero-Gómez-Olmedo, Beatriz 380
Torre-Jiménez, Roberto López-Sastre, Saturnino 381
Maldonado-Bascón, and Daniel Onoro-Rubio. 382
Extremely overlapping vehicle counting. In 383
Iberian Conference on Pattern Recognition and 384
Image Analysis, pages 423–431. Springer, 2015. 385