
DETECTING SOCCER BALLS WITH REDUCED NEURAL NETWORKS: A COMPARISON OF MULTIPLE ARCHITECTURES UNDER CONSTRAINED HARDWARE SCENARIOS

Douglas De Rizzo Meneghetti

Department of Electrical Engineering
FEI University Center
São Bernardo do Campo, SP 09850-901
douglasrizzo@fei.edu.br

Thiago Pedro Donadon Homem

Federal Institute of Education, Science and Technology of São Paulo
São Paulo, SP 05110-000
thiagohomem@ifsp.edu.br

Jonas Henrique Renolfi de Oliveira

Department of Electrical Engineering
FEI University Center
São Bernardo do Campo, SP 09850-901
jonashro@gmail.com

Isaac Jesus da Silva

Department of Electrical Engineering
FEI University Center
São Bernardo do Campo, SP 09850-901
isaacjesus@fei.edu.br

Danilo Hernani Perico

Department of Electrical Engineering
FEI University Center
São Bernardo do Campo, SP 09850-901
dperico@fei.edu.br

Reinaldo Augusto da Costa Bianchi

Department of Electrical Engineering
FEI University Center
São Bernardo do Campo, SP 09850-901
rbianchi@fei.edu.br

Sep. 27, 2020

ABSTRACT

Object detection techniques that achieve state-of-the-art detection accuracy employ convolutional neural networks, implemented to have optimal performance in graphics processing units. Some hardware systems, such as mobile robots, operate under constrained hardware situations, but still benefit from object detection capabilities. Multiple network models have been proposed, achieving comparable accuracy with reduced architectures and leaner operations. Motivated by the need to create an object detection system for a soccer team of mobile robots, this work provides a comparative study of recent proposals of neural networks targeted towards constrained hardware environments, in the specific task of soccer ball detection. We train multiple open implementations of MobileNetV2 and MobileNetV3 models with different underlying architectures, as well as YOLOv3, TinyYOLOv3, YOLOv4 and TinyYOLOv4 in an annotated image data set captured using a mobile robot. We then report their mean average precision on a test data set and their inference times in videos of different resolutions, under constrained and unconstrained hardware configurations. Results show that MobileNetV3 models have a good trade-off between mAP and inference time in constrained scenarios only, while MobileNetV2 with high width multipliers are appropriate for server-side inference. YOLO models in their official implementations are not suitable for inference in CPUs.

Keywords Object detection · Convolutional neural networks · Humanoid robots · Constrained hardware

1 Introduction

The recent successes in the field of object detection are mostly due to the use of deep neural networks, more specifically convolutional neural networks (CNN). The underlying operations that compose CNNs are highly optimized for execution in graphics processing units (GPU). However, in some domains, GPUs may be unavailable and these processes must be executed in CPUs. One such domain is mobile robotics, in which there may be limitations regarding space, weight and energy consumption, constraining the robot’s hardware to contain only CPUs, thus hindering the performance of systems based on deep learning techniques.

With these limitations in mind, this work presents a study of the performance of recent CNN architectures applied to object detection tasks and proposed for constrained hardware settings. We train these models in the task of soccer ball detection and compare their mean average precision (mAP) in a test data set, as well as their inference time in both constrained and unconstrained hardware settings. With this study, we aim to guide readers in their choice of neural network when implementing an object detection system for mobile robots.

We expand preceding work [3] by analyzing the newer MobileNetV2 [22] and MobileNetV3 [6] models, with more combinations of the width multiplier and input resolution parameters, as well as by including the YOLOv3 [20] and YOLOv4 [2] object detection algorithms and their "tiny" counterparts. This work also presents inference time results in both constrained and unconstrained environments and for multiple resolutions of input videos, providing readers with more relevant, up-to-date and comprehensive results to make decisions regarding the choice of neural network model to use in a soccer ball detection system (or any other comparable single-object detection system) under local, mobile, constrained hardware scenarios as well as when remote and unconstrained hardware configurations may be available.

This work is organized as follows: section 2 lists the recent advances in the state-of-the-art in object detection, as well as techniques to create smaller network topologies while still maintaining high detection accuracy. We also describe the network architectures utilized in this work and their detection mechanisms. Section 3 depicts related works. The methodology used throughout the experiments is presented in section 4. Section 5 presents the experiments and the results are discussed in section 6. Lastly, section 7 provides the conclusions and future work.

2 Research Background

In recent years, object detection techniques have advanced at great pace due to the equally fast advances of deep learning and convolutional neural networks applied to computer vision [30]. Two-stage detectors such as Faster R-CNN [21] first generate region proposals and then detect objects only in the selected regions, while one-stage detectors such as YOLO [18] and SSD [7] generate bounding box coordinates and class predictions at the same time, with YOLO using only convolutional layers for this task.

More recently, effort has centered around building strategies for efficiently scaling network models, reaching a trade-off between FLOPS, number of trainable parameters and accuracy. The MobileNetV3 architecture [6] has been partially achieved via hardware-aware neural architecture search techniques [29, 25], while AmoebaNet’s architecture [17], which achieved state-of-the-art classification accuracy on ImageNet [4], was evolved using evolutionary algorithms.

Other techniques attempt to shrink or expand the dimensions of convolutional layers using hyperparameters. MobileNetV1 [7] introduced the width multiplier and input resolution parameters, discussed later in the text, while EfficientNet [26] and EfficientDet [27] use a compound coefficient to scale all three dimensions of convolutional layers in order to maximize the network’s accuracy. This, allied with neural architecture search, introduced the current state-of-the-art in image classification and object detection using CNNs.

2.1 MobileNets

Introduced in [7], MobileNets are convolutional neural network architectures whose number of trainable parameters can be controlled by two hyperparameters. The first is the *width multiplier* $\alpha \in (0, 1]$, which controls the number of channels in each layer of the network. Smaller values of α reduce the number of parameters in each layer of the network uniformly, also reducing computational cost. The second parameter is the *resolution multiplier* $\rho \in (0, 1]$, which is used to reduce the resolution of the input images and, consequently, the number of operations throughout all layers of the network.

Additional features introduced in MobileNetV1 are batch normalization [9] for learning stabilization, as well as depthwise-separable convolutions [23], a convolution operation that uses fewer parameters to achieve comparable results to regular convolutions.

MobileNetV2 [22] advanced the state-of-the-art by introducing linear bottleneck layers in the network, reducing the size of the inputs in subsequent layers while preventing information from being lost by non-linear activation functions. The ReLU6 non-linearity (whose first appearance was tracked down to [11]) was chosen instead of regular ReLU to prevent loss of information when calculations with low-precision data types are performed.

Finally, MobileNetV3 [6] employs multiple neural architecture search (NAS) algorithms to optimize network architectures for different types of hardware, followed by the manual simplification of the most computationally expensive parts of the generated models. This network employs a variant of the Swish non-linearity [16] called h-swish.

2.2 Single-Shot MultiBox Detector

The Single-shot MultiBox Detector (SSD) [14] is a technique that utilizes a convolutional neural network, called the base network, combined with multiple subsequent convolutional filters of different sizes, to perform detection under different scales and aspect ratios in multiple regions of an input image. The feature maps of the base network may be pretrained in a classification or detection class. When training the network for a detection task, SSD employs techniques such as data augmentation and hard negative mining for faster training, as well as a loss function that is a weighted sum of both localization and classification losses.

2.3 You Only Look Once

You Only Look Once (YOLO) [18] simplified the object detection problem, which was then composed of a region proposal step followed by an image classification step, to a single regression step composed of bounding box coordinates and class probabilities. YOLOv2 [19] introduced the use of anchor boxes to the algorithm, a technique that allowed the detection of multiple objects with different aspect ratios in the same quadrant of an input image, using only convolutional layers.

YOLOv3 [20] introduces the prediction of bounding box coordinates across multiple scales and the use of residual layers [5] to speed up training, while YOLOv4 [2] adapts multiple data augmentation and feature extraction techniques to allow efficient training and inference of a model on a single GPU with 8 to 16 GB of VRAM.

3 Related Work

This section presents recent works that attempt to detect or track soccer balls and other objects relevant to the humanoid soccer scenario.

In [15], a model called JET-Net is proposed for the task of detecting robots and soccer balls in a humanoid soccer environment for NAO soccer games. The network uses building blocks popularized in the MobileNets, such as depthwise separable convolutions, as well as working with images in a single channel in order to reduce data complexity, achieving an inference time of 9 ms per frame.

In [24], ROBO is presented, a network inspired by TinyYOLOv3 which claims to achieve higher accuracy while being 35 times faster. This is achieved by reducing the topology of TinyYOLOv3 further, given that a reduced number of classes will be detected during the robot soccer task; adapting the input resolution of the network to that of the NAO camera; performing downsampling in the input image in the first layers of the network; and using a single anchor box for each class, since the overall shape of object classes is predictable.

A two-stage ball detection algorithm is presented in [28], where region proposals are generated first and then passed to a convolutional network for soccer ball detection, using SSD. The authors report an inference time of 5.13 ms in an Intel NUC with a Core i3 CPU.

DeepBall [10] is a recent fully-convolutional neural network architecture inspired by SSD and YOLO, two techniques covered in this work. DeepBall was created to detect soccer balls in long shot videos of real soccer games.

In [12], neural networks employing temporal convolutions (TCN), ConvLSTM and ConvGRU layers are used to detect and track the movement of soccer balls in a video feed from a humanoid robot. These networks are trained in sequences of images and the authors report the challenge of gathering sequential data to train the network, resorting to synthetic data for pretraining. A final inference time of around 6 ms is reported for TCN networks using an NVIDIA GTX 1050 Ti.

4 Methodology

The final goal of this work is to develop a computer vision system for an autonomous humanoid robot, allowing the robot to detect a soccer ball in a compatible time with the dynamics of the game. The vision system may also run on robots with embedded computers, such as mini-PCs, which are generally CPU-only. To select the object detection technique and underlying neural network that compose the vision system, this paper presents a study that investigates the trade-off between accuracy and inference time of multiple neural network architectures proposed specifically for mobile or embedded hardware settings, with different configurations, when executed on this kind of computer.

4.1 Network architectures and training

Twenty MobileNetV2 configurations were tested by modifying the values for the width and resolution multipliers. For the width multiplier, the values 1, 0.75, 0.5 and 0.35 were used, with the resulting networks having 3.47, 2.61, 1.95 and 1.66 million trainable parameters, respectively.

The values used for the resolution multiplier were chosen so that the input resolution of the network is equal to 224, 192, 160, 128 and 96. The combined values of both hyperparameters resulted in a total of twenty models that were trained using the soccer ball data set, described on section 4.4.

MobileNetV3 models are composed of the “Large” and “Small” variants presented in [6], both with width multipliers of 1 and 0.75, possessing 5.4 (Large, $\alpha = 1$), 4 (Large, $\alpha = 0.75$), 2.9 (Small, $\alpha = 1$) and 2.4 million (Small, $\alpha = 0.75$) trainable parameters, as well as minimalistic versions of both variants with $\alpha = 1$, possessing 3.9 and 2 million parameters. Minimalistic models do not contain the more advanced squeeze-and-excite units, hard-swish, and 5×5 convolutions operations from the non-minimalistic counterparts. YOLO models are composed of the v3 and v4 versions of the neural networks, presented in [20] and [2] respectively, as well as their “tiny” counterparts.

All models used pretrained weights learned in the COCO dataset [13].¹

4.2 Training procedure

The models were trained in a server with Intel Xeon Gold 5118@2.3 GHz processors totaling 48 CPUs, 192 GB of RAM and an NVIDIA Tesla V100-PCIE with 16 GB of memory, running CentOS 7.6.1810. The MobileNet models were trained for 50000 training steps. The YOLO and TinyYOLO models were trained for a total of 6000 training steps, following recommendations from the original developers of the model, given the number of classes to be detected.

All MobileNetV2 models were trained using batches of 32 images and the RMSProp optimizer with initial learning rate of $4 \cdot 10^{-3}$, an exponential decay schedule with a decay factor of 0.95 and a momentum coefficient of 0.9. All MobileNetV3 models were trained using batches of 32 images, stochastic gradient descent with initial learning rate of 0.4, a cosine decay schedule and a momentum coefficient of 0.9.

YOLO and TinyYOLO models were trained using batches of 64 images, stochastic gradient descent with a momentum coefficient of 0.9. YOLOv3 and TinyYOLOv3 models used a learning rate of 10^{-3} , while YOLOv4 and TinyYOLOv4 used a learning rate of $2 \cdot 10^{-3}$. Two step decays at 80% and 90% of the training were applied to these learning rates.

4.3 Humanoid Robot

The humanoid robot to which the object detection system is geared towards weighs about 5.9 kg and measures 81 cm in height. It is composed of 19 Dynamixel servomotors (a combination of MX-64, MX-106 and XM430 models), totaling 19 degrees of freedom. The humanoid robot uses a Genius WideCam F100 (Full HD) camera for image capture and a CH Robotics UM7 orientation sensor. The center of mass has a height of 36.1 cm and the robot has a foot area of 174 cm^2 . Other measurements include 39.5 cm of shoulder length, 38.5 cm of leg height, 18.8 cm of neck height and 38.5 cm of arm length. The robot is equipped with an Intel NUC Core i7 mini-PC. A picture of the robot is presented in figure 1a.

¹In order to facilitate the replication of these results and encourage the development of similar approaches by other researchers, the software used in this paper is available for download at: <https://github.com/douglasrizzo/JINT2020-ball-detection>.

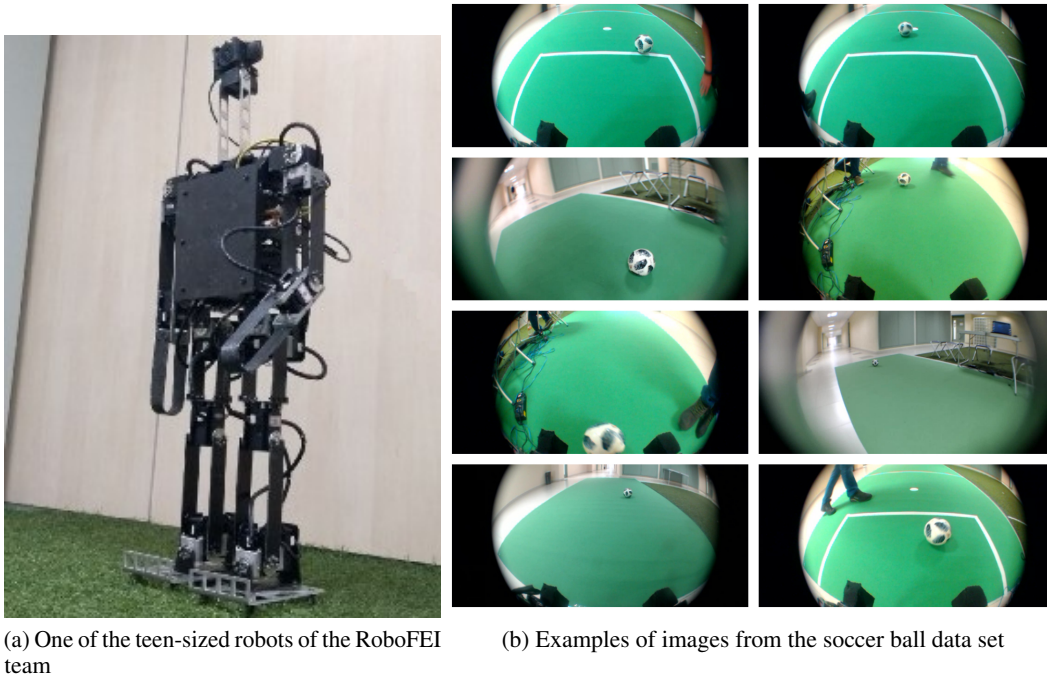


Figure 1: Image data set and humanoid robot used to collect it.

4.4 Image data set

The data set used in this work [1]² consists of 4364 images in 1920×1080 resolution, collected from the point-of-view of the humanoid robot described in section 4.3. A fish-eye lens is used to maximize the field of view of the robot and so all images in the data set also inherit this feature.

Of the 4364 images, 4014 compose the annotated training set and 250, the annotated test set. In these sets, the soccer balls visualized by the robot have been marked with bounding boxes. The training and test sets were collected from different sets of videos, with the purpose of minimizing data correlation.

Each image contains a single soccer ball, captured under multiples lighting conditions, as well as at different angles and distances from the camera. There are pictures of both stationary and moving soccer balls. Figure 1b presents examples of the data set.

5 Experiments

In order to recommend the best neural network for the soccer ball detection system, we trained multiple MobileNet and YOLO models in the training set of the data described in section 4.4 and measured the mean average precision (mAP) in the test set during training, taking the final mAP at the end of training as well as the *de facto* measure of precision for all models. Another factor measured by the experiments was inference time, which was measured using a 30-second video taken from the robot’s point-of-view in 1920×1080 resolution, which was also scaled down to 1280×720 , 640×480 and 480×360 pixels. All networks then processed these four versions of the same video and the mean inference time over all frames of each video for each network was recorded.

The MobileNet implementations selected for this work are provided in the TensorFlow Object Detection API [8], while the YOLO implementations are provided by [2]. The constrained hardware configuration in which the inference time of the models was captured is equipped with an i5-4210U CPU @ 1.70GHz and 8 GB of RAM and no GPU, in line with the hardware typically used by an autonomous mobile robot. For comparison, the same experiments were performed in the training computer, under GPU and CPU-only settings. In all cases, when CPU-only experiments were executed, all CPU cores were allowed to be utilized.

²The dataset was also made available at <http://ieee-dataport.org/open-access/open-soccer-ball-dataset>

Table 1: Average precision (higher is better) and inference time in milliseconds (lower is better) of the trained models. The five best in each category are marked in bold.

Network	Width mult.	Input res.	mAP	Inference time (ms)		
				Core i5	V100	Xeon
MobileNetV2	0.35	96	0.4065	99.964	65.37	64.959
		128	0.7095	101.166	51.062	49.642
		160	0.6304	88.651	52.642	50.97
		192	0.6756	87.006	53.613	52.232
		224	0.4984	78.553	42.852	41.703
	0.5	96	0.4065	91.403	58.919	57.626
		128	0.6986	81.08	44.529	43.388
		160	0.3361	91.775	58.288	58.616
		192	0.0944	116.076	65.629	64.828
		224	0.3253	78.624	42.759	43.18
	0.75	96	0.7284	86.569	51.065	51.528
		128	0.6954	84.159	42.097	41.866
		160	0.6679	81.351	41.883	42.309
		192	0.6952	78.699	42.347	41.776
		224	0.7874	85.186	48.343	47.854
	1	96	0.8133	122.853	56.992	57.83
		128	0.7672	82.277	46.921	47.799
		160	0.8597	88.886	52.278	52.569
		192	0.3632	110.75	61.263	60.052
		224	0.8177	79.547	42.438	42.183
MobileNetV3 (large min.)	1	224	0.6007	85.808	58.706	59.581
MobileNetV3 (large)	0.75	224	0.8847	89.362	63.515	63.703
MobileNetV3 (large)	1	224	0.6875	120.045	88.017	91.369
MobileNetV3 (small min.)	1	224	0.6024	79.142	48.68	49.236
MobileNetV3 (small)	0.75	224	0.7067	60.654	49.328	47.741
MobileNetV3 (small)	1	224	0.8651	96.975	70.689	70.017
TinyYOLOv3			0.3381	588.235	33.557	85.47
TinyYOLOv4			0.3504	714.286	29.851	119.048
YOLOv3			0.1355	5000	44.248	588.235
YOLOv4			0.1419	5000	50	833.333

6 Results

Table 1 displays the mAP of all trained models in the test set, as well as the inference time in milliseconds for different hardware configurations. In this table, the inference time was calculated with the videos in their native 1920×1080 resolution. These results allow us to conclude that the canonical YOLOv3 and YOLOv4 implementations, as well as their tiny counterparts, are not optimized for inference on CPUs, achieving the highest inference times of all models in the Intel Core i5-4210U. In fact, it is stated in [2] that the model implementations are optimized for inference on single GPUs. This can be seen by the comparatively low times achieved in the Tesla V100 GPU, especially by the TinyYOLO models, which achieved the lowest inference times of all models.

As for the MobileNet models, we can see that MobileNetV3 and MobileNetV2 with $\alpha = 1$ achieved the highest mAP in the test data set. However, with the high variation in inference times between all combinations of hyperparameters, the results from Table 1 alone do not provide enough information to compare model performances. To remedy that, we calculate a performance score for each neural network in each hardware setting $p_{m,h} = \frac{mAP_m}{t_{m,h}}$, where mAP_m represents the mAP of network m (hardware independent) and h , the hardware setting the inference time t of model m was gathered from. Then, the performance scores of all models in the same hardware are normalized by the highest

Table 2: Normalized scores of the trained models under different hardware. Higher is better. A normalized score of 1 indicates that the model performed the best under that hardware, compared with the others.

Network	Width mult.	Input res.	Normalized score		
			Core i5	V100	Xeon
MobileNetV2	0.35	96	0.349	0.323	0.323
		128	0.602	0.721	0.737
		160	0.61	0.622	0.638
		192	0.666	0.654	0.667
		224	0.545	0.604	0.617
	0.5	96	0.382	0.358	0.364
		128	0.74	0.814	0.831
		160	0.314	0.299	0.296
		192	0.07	0.075	0.075
		224	0.355	0.395	0.389
	0.75	96	0.722	0.74	0.729
		128	0.709	0.857	0.857
		160	0.705	0.828	0.814
		192	0.758	0.852	0.858
		224	0.793	0.845	0.849
	1	96	0.568	0.741	0.726
		128	0.8	0.849	0.828
		160	0.83	0.853	0.844
		192	0.281	0.308	0.312
		224	0.882	1	1
MobileNetV3 (large min.)	1	224	0.601	0.531	0.52
MobileNetV3 (large)	0.75	224	0.85	0.723	0.716
MobileNetV3 (large)	1	224	0.492	0.405	0.388
MobileNetV3 (small min.)	1	224	0.653	0.642	0.631
MobileNetV3 (small)	0.75	224	1	0.744	0.764
MobileNetV3 (small)	1	224	0.766	0.635	0.637
TinyYOLOv3			0.049	0.523	0.204
TinyYOLOv4			0.042	0.609	0.152
YOLOv3			0.002	0.159	0.012
YOLOv4			0.002	0.147	0.009

performance score in that hardware, leading to the normalized score $s_{m,h} = \frac{p_{m,h}}{\max_{\eta} p_{\eta,h}}$. Achieving a normalized score $s_{m,h} = 1$ means that model m had the highest mAP/inference time ratio off all models in that hardware.³

Table 2 presents the normalized scores of all models. Overall, MobileNetV2 models with width multipliers $\alpha \in \{0.75, 1\}$ had the best scores of all MobileNetV2 models in all hardware settings. We can also see that the five best models in unconstrained hardware settings are the same for both CPU and GPU. However, when operating under the Intel Core i5 4210U processor, both MobileNetV3 models (large and small) with $\alpha = 0.75$ achieved the highest scores, making MobileNetV3 models a viable option for an object detection system that operates under constrained hardware settings, whereas they did not exhibit the same performance in GPUs.

6.1 Performance under different input resolutions

This section presents the inference times in milliseconds of all model implementations when processing input videos of multiple resolutions (1920×1080 , 1280×720 , 640×480 , 480×360). Table 3 presents the mean and standard deviation of the inference time of all YOLO models for the tested hardware configurations. Overall, all YOLO and

³This score may easily break or be less informative if one neural network in the sample has disproportionately low inference time or high mAP. However, given the well-behaved values presented in Table 1, we consider the use of the proposed score appropriate for the purposes of our analysis.

Table 3: Inference time of YOLO and TinyYOLO models when processing videos of multiple input resolutions.

Network	Hardware	Inference time	
		mean	std. dev.
TinyYOLOv3	Tesla V100	41.957	10.022
	Xeon Gold 5118	88.983	3.369
	i5-4210U	588.235	0.000
TinyYOLOv4	Tesla V100	38.808	9.136
	Xeon Gold 5118	114.529	6.094
	i5-4210U	714.286	0.000
YOLOv3	Tesla V100	46.726	4.081
	Xeon Gold 5118	588.235	0.000
	i5-4210U	5000.000	0.000
YOLOv4	Tesla V100	50.385	0.676
	Xeon Gold 5118	817.308	32.051
	i5-4210U	5000.000	0.000

TinyYOLO models achieved low standard deviation in this test, implying that their implementation [2] is indifferent to the input resolution of images.

As for the MobileNet results, we first discuss the distributions of results in the three hardware settings. In total, 104 values were collected in each setting (twenty V2 and six V3 models applied to videos in four resolutions), with the following statistics: $\mu_{i5} = 68.292$, $\sigma_{i5} = 17.374$, $\mu_{V100} = 47.24$, $\sigma_{V100} = 8.756$, $\mu_{Xeon} = 47.009$ and $\sigma_{Xeon} = 8.981$.

Due to the similarity in the inference times collected from the NVIDIA Tesla V100 GPU and the Intel Xeon Gold 5118, we executed a two-sample Kolmogorov-Smirnov test between the two samples, with a result of $p = 0.97371$, indicating that the measurements collected in the 48 CPUs and the single GPU are similar with high statistical relevance. Because of that, in this section we only report results for the MobileNets in the NVIDIA Tesla V100 GPU and the Intel i5-4210U processor.

The distance between μ_{i5} and μ_{V100} is indicative of the performance lost by executing deep learning models in constrained CPUs, while a larger standard deviation on the CPU (σ_{i5}) indicate that there is a larger variability in network performance, given the resolution of the input video.

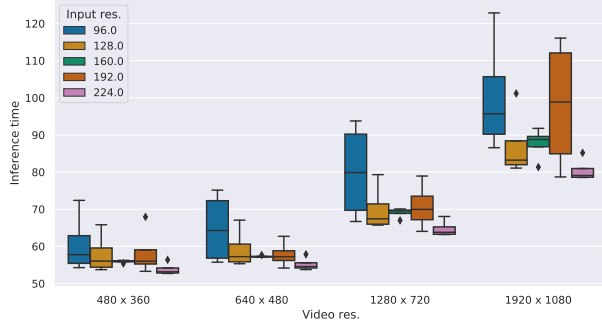
Figure 2 presents the inference time by frame in milliseconds for the MobileNet models when processing the same video under multiple resolutions in the Intel i5-4210U CPU and an NVIDIA Tesla V100 GPU. Unlike the results for the YOLO models, all MobileNets show a significant speedup when applied to input frames of lower resolutions. This information may be relevant, as the implementations by [8] already operate in downscaled images, with a resolution of 300×300 , indicating that downscaling the input feed prior to executing object detection is a valid strategy to achieve lower inference times. This speedup is visualized in both constrained CPU and GPU settings.

7 Conclusions

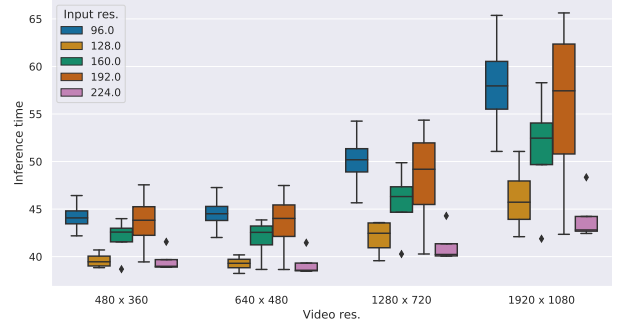
This work presented a comparative study of the performance of multiple neural network architectures, designed for fast inference with a reduced number of trainable parameters while maintaining high precision, when applied in the task of soccer ball detection under constrained and unconstrained hardware scenarios, as well as when processing input images of different resolutions.

Results have shown that MobileNetV2 models with high width multipliers have the best trade-off between mAP and inference time in unconstrained hardware settings, being suitable when executing inference in remote servers is an option. However, under a local, constrained, CPU-only scenario, MobileNetV3 models have shown the best scores, while not having remarkable performance when operating in a GPU. Lastly, the official implementations of YOLO and TinyYOLO, being optimized for inference in GPUs, displayed poor results in our low-end Intel Core i5-4210U processor.

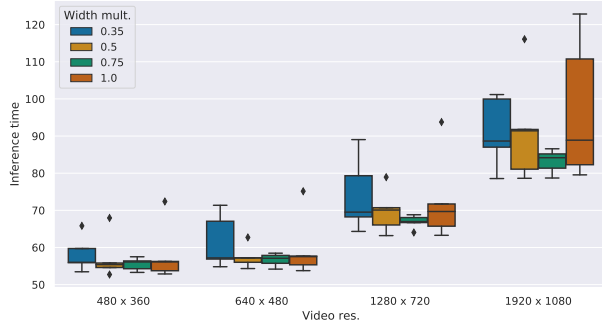
In future work, we aim to evaluate the performance of state-of-the-art reduced object detection models in embedded systems with GPUs, seem them as the next step in hardware for mobile robotics.



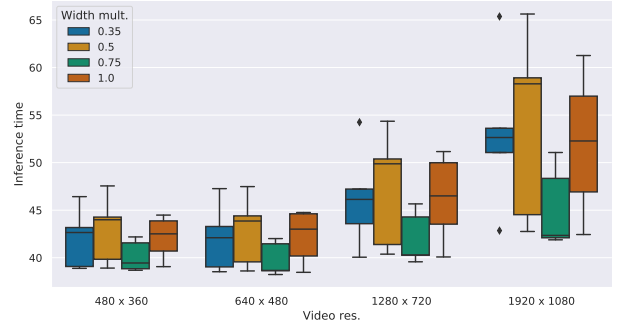
(a) MobileNetV2 by input resolution on i5-4210U



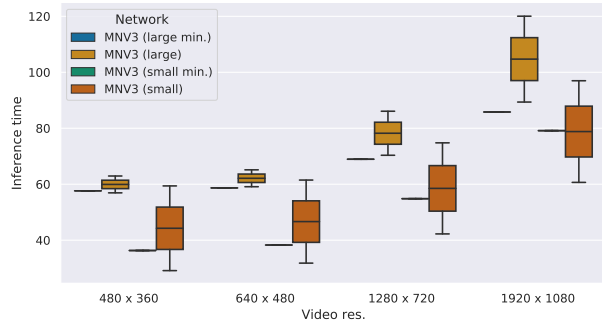
(b) MobileNetV2 by input resolution on V100



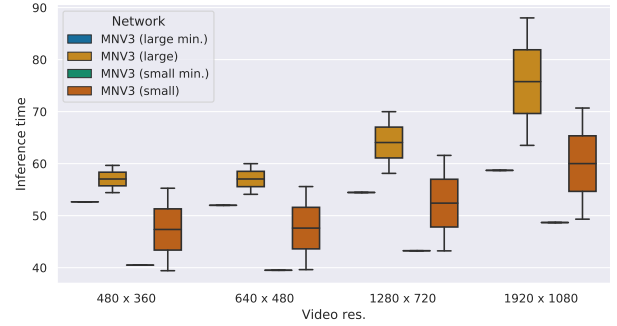
(c) MobileNetV2 by width multiplier on i5-4210U



(d) MobileNetV2 by width multiplier on V100



(e) MobileNetV3 on i5-4210U



(f) MobileNetV3 on V100

Figure 2: Inference time of MobileNetV2 and V3 models in Intel i5-4210U (left) and NVIDIA Tesla V100 (right).

Acknowledgements

The authors acknowledge the São Paulo Research Foundation (FAPESP Grant 2019/07665-4) for supporting this project. This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior – Brasil (CAPES) – Finance Code 001. This is a preprint of an article published in the Journal of Intelligent & Robotic Systems. The final authenticated version is available online at: <https://doi.org/10.1007/s10846-021-01336-y>.

Conflict of interest

The authors declare that they have no conflict of interest.

References

- [1] Bianchi, R.A.d.C., Perico, D.H., Homem, T.P.D., da Silva, I.J., Meneghetti, D.D.R.: Open Soccer Ball Dataset (2020). DOI 10/ghcfxn
- [2] Bochkovskiy, A., Wang, C.Y., Liao, H.Y.M.: YOLOv4: Optimal Speed and Accuracy of Object Detection. arXiv:2004.10934 [cs, eess] (2020). URL <http://arxiv.org/abs/2004.10934>
- [3] de Oliveira, J.H.R., da Silva, I.J., Homem, T.P.D., Meneghetti, D.D.R., Perico, D.H., Bianchi, R.A.d.C.: Object detection under constrained hardware scenarios: A comparative study of reduced convolutional network architectures. In: 2019 XVI Latin American Robotics Symposium and VII Brazilian Robotics Symposium (LARS/SBR). IEEE (2019)
- [4] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Li Fei-Fei: ImageNet: A large-scale hierarchical image database. In: 2009 IEEE Conference on Computer Vision and Pattern Recognition, pp. 248–255 (2009). DOI 10/cvc7xp
- [5] He, K., Zhang, X., Ren, S., Sun, J.: Deep Residual Learning for Image Recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778 (2016). DOI 10.1109/cvpr.2016.90
- [6] Howard, A., Sandler, M., Chu, G., Chen, L.C., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., Le, Q.V., Adam, H.: Searching for MobileNetV3. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 1314–1324 (2019). URL https://openaccess.thecvf.com/content_ICCV_2019/html/Howard_Searching_for_MobileNetV3_ICCV_2019_paper.html
- [7] Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., Adam, H.: MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. CoRR **abs/1704.04861** (2017). URL <http://arxiv.org/abs/1704.04861>
- [8] Huang, J., Rathod, V., Sun, C., Zhu, M., Korattikara, A., Fathi, A., Fischer, I., Wojna, Z., Song, Y., Guadarrama, S., Murphy, K.: Speed/Accuracy Trade-Offs for Modern Convolutional Object Detectors. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE (2017). DOI 10.1109/cvpr.2017.351
- [9] Ioffe, S., Szegedy, C.: Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In: International Conference on Machine Learning, pp. 448–456 (2015). URL <http://proceedings.mlr.press/v37/ioffe15.html>
- [10] Komorowski, J., Kurzejamski, G., Sarwas, G.: DeepBall: Deep Neural-Network Ball Detector. In: International Conference on Computer Vision Theory and Applications. Prague, Czech Republic (2019). DOI 10/gg94mp
- [11] Krizhevsky, A.: Convolutional deep belief networks on CIFAR-10. Tech. rep. (2010)
- [12] Kukleva, A., Khan, M.A., Farazi, H., Behnke, S.: Utilizing Temporal Information in Deep Convolutional Network for Efficient Soccer Ball Detection and Tracking. In: S. Chalup, T. Niemueller, J. Suthakorn, M.A. Williams (eds.) RoboCup 2019: Robot World Cup XXIII, Lecture Notes in Computer Science, pp. 112–125. Springer International Publishing, Cham (2019). DOI 10/gg94mm
- [13] Lin, T.Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., Perona, P., Ramanan, D., Zitnick, C.L., Dollár, P.: Microsoft COCO: Common Objects in Context (2014)
- [14] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: SSD: Single Shot MultiBox Detector. In: Computer Vision – ECCV 2016, vol. 9905, pp. 21–37. Springer International Publishing (2016). DOI 10.1007/978-3-319-46448-0_2
- [15] Poppinga, B., Laue, T.: JET-Net: Real-Time Object Detection for Mobile Robots. In: S. Chalup, T. Niemueller, J. Suthakorn, M.A. Williams (eds.) RoboCup 2019: Robot World Cup XXIII, vol. 11531, pp. 227–240. Springer International Publishing, Cham (2019). DOI 10.1007/978-3-030-35699-6_18

- [16] Ramachandran, P., Zoph, B., Le, Q.V.: Searching for Activation Functions. arXiv:1710.05941 [cs] (2017). URL <http://arxiv.org/abs/1710.05941>
- [17] Real, E., Aggarwal, A., Huang, Y., Le, Q.V.: Regularized Evolution for Image Classifier Architecture Search. In: 33rd AAAI Conference on Artificial Intelligence, AAAI 2019 (2019). URL <http://arxiv.org/abs/1802.01548>
- [18] Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You Only Look Once: Unified, Real-Time Object Detection. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE (2015). DOI 10.1109/cvpr.2016.91
- [19] Redmon, J., Farhadi, A.: YOLO9000: Better, Faster, Stronger (2016). URL <http://arxiv.org/abs/1612.08242v1>
- [20] Redmon, J., Farhadi, A.: YOLOv3: An Incremental Improvement (2018). URL <http://arxiv.org/abs/1804.02767v1>
- [21] Ren, S., He, K., Girshick, R.B., Sun, J.: Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. CoRR **abs/1506.01497** (2015). DOI 10.1109/tpami.2016.2577031
- [22] Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C.: MobileNetV2: Inverted Residuals and Linear Bottlenecks. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4510–4520 (2018). DOI 10/gfxgjz
- [23] Sifre, L.: Rigid-motion scattering for image classification. Ph. D. Thesis, Ecole Polytechnique, CMAP, Palaiseau, France (2014)
- [24] Szemenyei, M., Estivill-Castro, V.: ROBO: Robust, Fully Neural Object Detection for Robot Soccer. In: S. Chalup, T. Niemueller, J. Suthakorn, M.A. Williams (eds.) RoboCup 2019: Robot World Cup XXIII, vol. 11531, pp. 309–322. Springer International Publishing, Cham (2019). DOI 10.1007/978-3-030-35699-6_24
- [25] Tan, M., Chen, B., Pang, R., Vasudevan, V., Sandler, M., Howard, A., Le, Q.V.: MnasNet: Platform-Aware Neural Architecture Search for Mobile. arXiv:1807.11626 [cs] (2019). URL <http://arxiv.org/abs/1807.11626>
- [26] Tan, M., Le, Q.V.: EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. arXiv:1905.11946 [cs, stat] (2020). URL <http://arxiv.org/abs/1905.11946>
- [27] Tan, M., Pang, R., Le, Q.V.: EfficientDet: Scalable and Efficient Object Detection. arXiv:1911.09070 [cs, eess] (2020). URL <http://arxiv.org/abs/1911.09070>
- [28] Teimouri, M., Delavaran, M.H., Rezaei, M.: A Real-Time Ball Detection Approach Using Convolutional Neural Networks. In: S. Chalup, T. Niemueller, J. Suthakorn, M.A. Williams (eds.) RoboCup 2019: Robot World Cup XXIII, vol. 11531, pp. 323–336. Springer International Publishing (2019). DOI 10.1007/978-3-030-35699-6_25
- [29] Yang, T.J., Howard, A., Chen, B., Zhang, X., Go, A., Sandler, M., Sze, V., Adam, H.: NetAdapt: Platform-aware neural network adaptation for mobile applications. In: European Conference on Computer Vision (ECCV) (2018). URL https://openaccess.thecvf.com/content_ECCV_2018/papers/Tien-Ju_Yang_NetAdapt_Platform-Aware_Neural_ECCV_2018_paper.pdf
- [30] Zhao, Z.Q., Zheng, P., Xu, S.T., Wu, X.: Object Detection With Deep Learning: A Review. IEEE Transactions on Neural Networks and Learning Systems pp. 1–21 (2019). DOI 10.1109/tnnls.2018.2876865