

Performance/power assessment of CNN packages on embedded automotive platforms

Paolo Burgio, Gianluca Brilli
University of Modena and Reggio Emilia, Italy
{firstname.lastname}@unimore.it

Abstract—The rise of power-efficient embedded computers based on highly-parallel accelerators opens a number of opportunities and challenges for researchers and engineers, and paved the way to the era of *edge computing*. At the same time, advances in embedded AI for object detection and categorization such as YOLO, GoogleNet and AlexNet reached an unprecedented level of accuracy (mean-Average Precision – mAP) and performance (Frames-Per-Second –FPS). Today, edge computers based on heterogeneous many-core systems are a predominant choice to deploy such systems in industry 4.0, wearable devices, and – our focus– autonomous driving systems. In these latter systems, engineers struggle to make reduced automotive power and size budgets co-exist with the accuracy and performance targets requested by autonomous driving. We aim at validating the effectiveness and efficiency of most recent networks on state-of-the-art platforms with embedded commercial-off-the-shelf System-on-Chips, such as Xavier AGX, Tegra X2 and Nano for NVIDIA and XCZU9EG and XCZU3EG of the Zynq UltraScale+ family, for the Xilinx counterpart. Our work aims at supporting engineers in choosing the most appropriate CNN package and computing system for their designs, and deriving guidelines for adequately sizing their systems.

I. INTRODUCTION

The next generation of automotive systems will be powered by high-performance embedded computers, that exhibit TFLOPs of performance within a reduced Size, Weight and Power consumption (SWaP) [32], [30], [46], [24], [45]. There is a huge interest from car-makers, but also from Tier1s, to use those platforms as a baseline for next-generation Advanced Driving Assistance Systems (ADAS) and Autonomous Vehicles (AV). Notable players are Tesla, BMW, FCA, and Toyota, but also companies that traditionally did not belong to the automotive domain such as NVIDIA [30] and Intel/Mobileye [24]. They are all investing billions of dollars and significant research efforts in this direction. Even Xilinx, which traditionally targeted reconfigurable computing and signal processing systems, recently released its high-performance Ultrascale+ [45] System-on-Module that couples a multi-core host, a graphics co-processor and reconfigurable logics accelerator. Their next platform generation, named Versal [46], is even more promising.

Embedded Deep Neural Networks. Building future AD systems that are safe and reliable rises a number of challenges, because an autonomous vehicle requires to process a huge amount of data coming from cameras, LiDARs, radars and even the internet, to answer to questions such as: *Where am I? Where is everybody else? How do I get from A to B?*, but also: *What is the health status of the driver? Do I*

detect a possible danger for this vehicle, such as a crash with other vehicles or pedestrians?. The current way of answering these questions employs Deep Neural Networks (DNNs) at their state-of-the-art. DNNs have a tremendous effectiveness for tasks as object detection and categorization, which are key modules of the perception and localization systems of tomorrow automated vehicles. For instance, in traffic sign classification, the accuracy of an adequately trained DNNs can reach up to 99% [48], even better than humans!

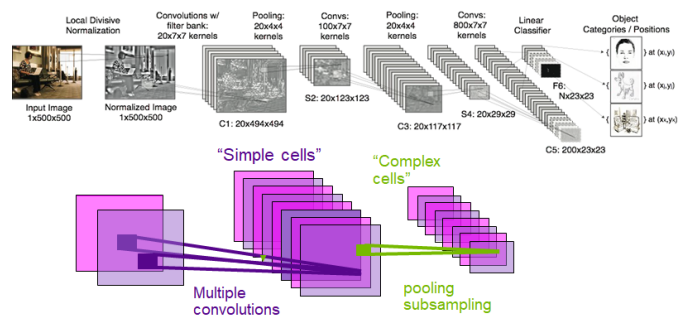


Fig. 1. Scheme of a generic CNN

To sustain the inherent complexity, hence, computational demand, of DNNs on size and power-constrained embedded platforms, typically system designers employ highly-optimized Convolutional Neural Networks (see Figure 1), where convolutional layers extract significant features from images. Nicely, the main operations that CNNs carry are sums and multiplications, more specifically defined as *Multiply-Accumulate Operations (MACs)*, which are highly scalable with the number of neurons, and data-parallelizable, hence perfectly suitable for deployment on many-core accelerators.

What this paper is about. The goal of our research is to explore and assess the performance, accuracy and power profile of state-of-the-art CNN packages (es. YOLO [36] and AlexNet [18]) for image detection and classification, when deployed on embedded computing platforms representative of next-generation automotive domain controllers.

We focus on the task of **inference**, and explore its limits and tradeoff when deployed on next generation of automotive systems, i.e., on tiny, embedded computers, under real-time constraints. We do not focus on the (offline) training problem, and assume the network is already trained. Nonetheless, we show how different datasets might affect performance of the same network on the same machine.

We choose three versions of the NVIDIA Tegra [32], [29], [34], as representative of GPGPU accelerators, and two variants of the Xilinx Ultrascale+ [45], [44], for reconfigurable FPGA-based platforms. For every CNN package, and target platform, we run an extensive set of experiments to assess the performance profile, not only in terms of computing power (where we choose frames-per-second and end-to-end latency on the single frame as representative metrics), but also the power dissipation at voltage and temperature operational ranges, in terms of consumed Watts. Mostly all of these systems are battery-powered, and energy dissipation causes overheating, implying lower SoC lifetime and higher cooling costs. Engineers struggle to find the optimal tradeoff between performance and dissipated power, and this is the reason why they typically measure the performance of these platforms in terms of *TFLOPs-per-watt*. To capture the power efficiency of the system with a single metric for all the experiments, we introduce the FPS/Pow (of the sole SoC/SoM) metric. For measuring network accuracy, we adopt the mean-Average Precision (mAP) ¹, which is a widely adopted and acknowledge way of assessing network precision.

This work follows up, and extends, our preliminary exploratory paper [3], making the best of the results and observations achieved, which ultimately paved the way to this final and complete work. More in details, in our previous work we generically focused on detection and classification CNN, with no in-depth analysis, and with a limited set of (now, obsolete) platforms and networks. Here, we complete our analysis, and go deeper in exploring especially detection networks, building a comprehensive view of the problem. More in details:

- 1). We re-computed our set of experiments on the next generation of computing platforms, including NVIDIA Xavier AGX [32] and the small Jetson NANO [34] and Xilinx Ultra96 [44]. Few previous results (the ones on the “old” NVIDIA TX2) are still reported in Section V for the sake of comparison.
- 2). We extended our analysis to more advanced networks, mainly focusing on the object detection problem. We employ the newest YOLO v3 network ².
- 3). We carried on a detailed, in-depth analysis of how specific features of the network or computing platform affect our metrics, in Section VI. We focus on i) different datasets, ii) network quantization, iii) reduced “tiny” version of the same network, iv) image-level parallelism via batching (on NVIDIA architectures) or multiple threading (on Xilinx’s), v) available power profiles via frequency and power scaling (available only on NVIDIA’s GPGPUs – Section V-A).

Section II some previous works that are related to ours, and gives a brief motivation for the chosen network packages. Section III depicts the architectures considered in this work, and Section IV shows in details the properties and structure of the considered neural networks, and puts previous validation

results in a line. Sections V and VI provide our detailed analysis of such CNN packages on the chosen reference platforms. Finally, Section VII draws some conclusions, and highlights future works.

II. RELATED WORKS, AND NOVELTY

While neural networks are known since decades, doing inference on embedded platform is quite a recent application area, enabled by the IoT era with powerful mobile devices, and advances in automotive and avionics sectors. Every time a new NN package comes to light, authors and interested researchers in the domain carry on an accurate evaluation to assess its advantages and disadvantages, on platforms that are representative of the specific application.

Recently, NVIDIA proposes two articles [31] [33] in which are reported benchmarks of the new Xavier chip compared to the previous version (TX2). The first one contains an exhaustive comparison of the main Convolutional Neural Networks on Jetson AGX Xavier. The tests are made varying the power profile of the Xavier SoC: *10W*, *15W*, *30W* and *EDP* modes (through the *nvpmodel* tool) and varying the size of the batch loaded on the GPU. This comparison is well presented but only compares the Xavier chip and does not consider the whole frequency range allowed by the JetPack. The second article proposed by NVIDIA compares some CNNs on Xavier and TX2, but does not consider others architectures like FPGAs or Many-Cores.

J. Johnson proposes in [16] an analysis that involves some popular CNN models. In particular the results gathered by the author are carried out on several desktop GPUs, like: *Pascal* and *Maxwell Titan X* and *GTX 1080* and *1080 Ti*. This paper differs from what J. Johnson published, because it is focused on desktop GPUs and does not take into account the power consumption.

Eriko Nurvitadhi et al. proposes in [28] a customizable DNN template for FPGA and evaluates some emerging DNN algorithms on Intel Programmable Logics (Arria 10 and Stratix 10) against a desktop powerful GPU: the Pascal Titan X. The results achieved by the authors show that the Stratix 10 is 10%, 50%, and 5.4x better in performance (TOP/sec) than Titan X Pascal GPU on GEMM operations for pruned, Int6, and binarized DNNs, respectively.

Nakahara et. al [25] developed an FPGA version of YOLO where inputs and weights are binarized [13], this network is implemented with SDSoC, an high level synthesis tool provided by Xilinx. They also carry an extensive benchmarking of their network running on Zynq Ultrascale+ and Tegra X2, the comparison shows that the FPGA implementation can reach about 40 frames per second (FPS) with a power dissipation of 4.5 Watt, while the GPU implementation that runs on Tegra X2 achieves only 2 FPS with a consumption of 7 Watt. These latter numbers can certainly be improved.

The Xilinx engineer M. Blott proposed a reduced and quantized version of Tiny-YOLO called Tincy-YOLO [2], which operates at 16 FPS, consuming only 6 Watt for the FPGA accelerator. The benchmark proposed in the article regards a comparison of the accuracy achieved by Tincy-YOLO respect

¹We adopted the following tools to compute the mAP. For COCO: <https://github.com/cocodataset/cocoapi>, for all other tests: <https://github.com/Cartucho/mAP>

²At the time we write, also YOLO v4 and v5 are available, but still not stable nor fully supported on all of our reference platforms (for instance, on Xilinx’s), hence we cannot perform a fair comparison

to other versions of YOLO, but does not include tests with others architectures like GPUs. We aim at completing these results.

Our aim is to assess i) **both** performance and power efficiency, ii) of **multiple** state-of-the-art neural network packages iii) on **multiple** platforms, which are representative of the future embedded systems, not only in the automotive domain (as for positioning of this paper), but also on the embedded domain *tout court*. We believe this work can be a reference for system engineers, to assess and identify the most suitable packages for given hardware system, and the other way around.

Kim et al. [17] perform a comparison of embedded deep learning methods for person detection on embedded platforms such as the Jetson TX2 and Movidius. They investigate the performance of 13 different object detector deep models including variants of YOLO, SSD, RCNN, R-FCN and SqueezeDet. They train and test the NN on their proprietary dataset. In order to deploy the deep models in embedded platforms, they optimize Caffe or Tensorflow implementation using Movidius SDK or TensorRT. This enables the CNN model to utilize the target height/width effectively. The execution times of the methods are computed on the embedded boards. On the other hand, to measure and compare the average precision (AP) and IoU of the deep models, they used a workstation powered by 16 GB of internal memory and Nvidia GTX 1080ti graphics accelerator. Experiments results shows that Tiny YOLO-416 and SSD (VGG-300) are among the fastest models and Faster RCNN (Inception ResNet-v2) and RFCN (ResNet-101) are the most accurate ones. YOLO v3-416 and SSD (VGG-500) are the best tradeoff between Average precision and throughput.

Mittal presents in [23] a survey of works that evaluate and optimize neural network applications on Jetson platform (TK1, TX1, TX2). The survey considers several computer vision applications and offers a good overview of the existing literature. However, it does not take into account the Xavier platform.

Yu et al. [8] compare Real-Time object detection algorithms on embedded platforms. They compare power efficiency, latency and accuracy of Faster RCNN, YOLO and SSD on NVIDIA TK1, XILINX Zynq 7045 and XILINX KU115. However the comparison is not very clear, and different settings are compared, making the comparison unfair.

Chen and Ran [4] present a review of deep learning with edge computing. They report background, measurements and frameworks, consider different applications of deep learning, from computer vision to natural language processing, and take into account several methods for fast inference. Regarding fast inference on device computation, three major efforts have been identified: - Model Design: such methods focus on designing models with a reduced number of parameters in the DNN model, thus reducing memory and execution latency, while aiming to preserve high accuracy. Some examples include MobileNets [12], Single-Shot Detector (SSD)[19], YOLO[36][35], and SqueezeNet [14], with the state of the art that is evolving rapidly. - Model Compression: such methods seek to compress the existing DNN models with minimal accuracy loss compared with the original model. There are several

popular model compression methods: parameter quantization, parameter pruning, and knowledge distillation. - Hardware: To speed up inference of deep learning, hardware manufacturers are leveraging existing hardware such as CPUs and GPUs, as well as producing custom application-specific integrated circuits (ASICs) for deep learning. As an example, consider NVIDIA GPUs, FPGAs or Google TPUs.

Hossain and Lee [11] design a deep learning Real-Time object detection and tracking method for drones with Embedded Device. In this work they compare several SOTA NN for object detections on some GPU based embedded boards. In particular, they consider YOLOv2, YOLOv3, tiny YOLO, SSD and Faster RCNN on TX1, TX2 and AGX Xavier. Moreover, they take into account YOLOv2 and SSDMobileNet on Raspberry Pi, Latte Panda, Odroid, with or without Movidius. The comparison is only performed on the execution times of the neural networks.

Rungsuptaweekoon et al. [37] evaluate the power efficiency of image recognition with YOLO on NVIDIA Jetson TX1, Jetson TX2, and Tesla P40. For this evaluation, they deployed the Low-Power Image Recognition Challenge (LPIRC) system and integrated YOLO, a power meter, and target hardware into the system. The authors compare mAP, accumulated energy consumption, mAP/Energy and frame rate factors. The experimental results show that Jetson TX2 with Max-N mode has the highest throughput; Jetson TX2 with Max-Q mode has the highest power efficiency.

Ding et al. [7] propose REQ-YOLO, a resource aware, systematic weight quantization framework for object detection, considering both algorithm and hardware resource aspects in object detection. They consider Yolo v2 tiny and several boards: NVIDIA Titan X, GTX 1070, Jetson TX2 and Xilinx Virtex-7 485t, Zynq 7020, ADM-7V3 FPGA. They compare the performance on the different platforms in terms of FPS and power (W). They show that, thanks to their framework, they can achieve 314.2 FPS using only 21 W on the ADM-7V3 FPGA.

Lin et al. [6] benchmark several deep learning frameworks (MXNet[1], TensorFlow[9], CNTK[22], Neon[26], and PyTorch[21]) and investigate the FPGA deployment for performing traffic sign classification and detection. They evaluate both training and inference performance. For the latter one, they consider inference speed, accuracy, and power efficiency, by varying different parameters such as floating point precision, batch sizes, etc. They show that TensorFlow is always among the frameworks with the highest inference accuracy. For object detection inference, they compare six SSD models with different base networks, namely, ResNet-18, ResNet-50 [10], MobileNet-v1 [12], SqueezeNet-v1.1 [14], VGG [40], and MobileNet-v2 [38] (with SSDLite), on an NVIDIA GTX 1050 Ti GPU and an Arria 10 FPGA development board. Compared to the reference results on the GPU, they notice that in most of the cases inference speed on the GPU is higher than the FPGA, as well as accuracy. However, FPGA always achieves higher power efficiency than the GPU. Among the FPGA test cases, FP11 is always more power-efficient than FP16, while the data type does not have a clear impact on the power efficiency, but using FP11 bitstreams results in some

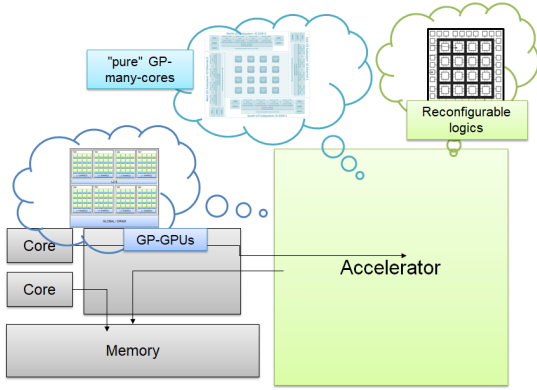


Fig. 2. Modern embedded platform for NN acceleration

decrease in accuracy. Finally, if the bitstream being used is fixed, the data type makes only a very tiny difference on both inference speed and accuracy.

III. AUTOMOTIVE ARCHITECTURES

In this work, we target a quite common architectural template for embedded systems, the *host-accelerator paradigm*, where a multi-core host is coupled with a power-efficient many-core accelerator (See Figure 2). Application is partitioned so that the host subsystem manages data transfer, in our case, the frames from the camera and the accelerator/system memory, and the many-core accelerator runs the highly-parallel workload, that is, in our case, the convolutional neurons, in an energy-efficient manner. We target several different architectural “flavors”, namely three based on NVIDIA GPGPUs accelerators, and two based on reconfigurable FPGA logics. They currently represent the most advanced technologies available on the market in each domain.

GP-GPU: the NVIDIA Tegra family and AGX Xavier

The *Tegra TX2* board embeds an esa-core host with Big.SUPER configuration and a GPU with 2 NVIDIA Streaming Multiprocessors (SM) of the Pascal family, summing up to 256 CUDA cores. The host and accelerators communicate using shared memory banks, employing what NVIDIA calls the *Unified Virtual Memory* space. The platform is claimed to achieve 1 TFLOP of computing power, within approximately 20 Watts of power dissipation. Being explicitly designed for the automotive market, Tegra X2 is qualified according to Functional Safety and Road Vehicles Standard (ISO 26262’s ASIL-B level [15]), and has been marketed in a (family of product) named Drive CX/PX [30], [39].

Xavier is the new chip developed by NVIDIA with even more computational power. It includes an integrated Volta GPU composed by 512 CUDA-Cores with 64 Tensor Cores, a dual Deep Learning Accelerators (NVDLAs), an octal-core NVIDIA Carmel ARMv8.2 CPU, 16GB 256-bit LPDDR4x with 137GB/s of memory bandwidth, and 650Gbps of high-speed I/O including PCIe Gen 4 and 16 camera lanes of MIPI CSI-2. According to NVIDIA this chip is able to reach a peak performance of 32 TOPS under 30W.

We also consider family’s low-end platform, the Jetson *NANO* [34], a very small compute module that embeds a

quad-cores ARM A57 processor, a GPGPU with only one SM composed of 128 CUDA cores and a 4GB of LPDDR4 main memory. Unfortunately, due to its constrained resources, it does not support bigger networks, such as full YOLO.

Reconfigurable logics: the Xilinx Ultrascale+ The second platform family we consider is Xilinx Zynq UltraScale+ (XU+), a next-generation reconfigurable heterogeneous platform that couples power-efficient host complex (ARM A53 cores with Real-Time cores, the ARM Cortex R5), and reconfigurable logics. The SoC also features a Mali GPU, that anyway can be used only for graphics. In our setup we used two different versions of the XU+ family: the XCZU9EG and the XCZU3EG. For the XCZU9EG SoC we used the Xilinx Zynq *ZCU102* [45] a very big carrier board with a really rich I/O connectivity, that can become quite energy-consuming, for this reason we expect its power consumption to be higher than the other ones. For the XCZU3EG, we adopted the Avnet *Ultra96* development platform [44] a carrier board with a very small form factor, for the latter we expect a power consumption close to Jetson Nano. The XCZU3EG have about the 25% of the total reconfigurable resources of the XCZU9EG.

Expected Results. We expect to see better performance on GPGPUs in general, mainly due to the fact that the frameworks used are widely tested and therefore highly optimized. More in details, we expect an increase of about $3\times$ factor across different platform generations, since NVIDIA declares a computational power of about 1 TFLOP for TX2 and about 2.8 TFLOPS for Xavier. With regard to the performance on FPGA-based architectures, we expect lower throughput in terms of FPS due to the lower diffusion of the aforementioned chips and the lower life cycle of the frameworks used. Finally, regarding the power consumption, we expect lower consumption on FPGA, and an excellent energy efficiency on XU+ and NVIDIA Xavier.

IV. CONVOLUTIONAL NEURAL NETWORKS

Convolutional Neural Networks or *CNNs* are typologies of neural networks particularly suitable for computer vision, image recognition and in general for processing 2D-input data. The basic structure of a CNN consists of a stack of layers in which an input image flows from the input layer toward the output layer. The input image is called *input feature map* (*fmap*) and it is mapped directly into the input neurons; this is beneficial in image recognition, because the spatial information inside images is preserved. The input feature map is represented as an object in three dimensions called *tensor*, whose dimensions are height H_{in} , width W_{in} and depth D_{in} . Each layer of which the CNN is composed, produces a stack of D_{out} —images of $W_{out} \times H_{out}$ —pixels, this output tensor is called *output feature map*. The most important and one of the most compute-intensive operation performed by CNNs is the discrete 2D-convolution, it is performed by convolutional layers and it is defined as:

$$o_{i,j}^k = \sum_{c=0}^{D_{in}} \sum_{h=0}^{K_H} \sum_{w=0}^{K_W} (w_{h,w,c}^k x_{i+h,j+w,c}) + b_k$$

Due to its mathematical structure, the convolutional layers, are suitable to be implemented on parallel hardware like Many-Cores platforms.

Image Classification is probably the most well-known problem in computer vision and it consists of classifying an image into one or more different categories. The problem of Object Detection is more complex, as it consists of classifying several objects inside an input image with some spatial information called *bounding boxes*, and for this reason, this operation requires more computing power and energy to be carried out.

Detection: You-Only-Look-Once (YOLO). The model of the CNNs for objects detection analyzed in this paper is the YOLO [36], the state of the art for objects detection via Convolutional Neural Networks. Almost the whole detection systems apply its model to an image at multiple locations and scales and then returns the computed bounding boxes that achieves the maximum scores. YOLO uses a totally different approach, it applies a single neural network to the full input image, the network divides the image into smaller regions and predicts bounding boxes and classes for each region, the network also compute weights for all regions and finally it returns only the bounding boxes which exceeds a fixed threshold. The YOLO full model requires a lot of memory and computations, then the authors proposed some different and reduced versions of the original model called YOLO-Small and Tiny-YOLO.

Classification: AlexNet. AlexNet is a large scale deep convolutional neural network proposed by A. Krizhevsky et al. in [1], AlexNet is designed for image classification and trained on ImageNet ILSVRC. This CNN is composed by eight trained layers: five convolutional layers and three fully connected layers, the network starts performing an image reduction to match the size of the input neurons that is $227 \times 227 \times 3$. AlexNet introduces some new layers, for example a new normalization layer called *Local Response Normalization (LRN)*.

V. ARCHITECTURAL EXPLORATION

We recall that we are in the automotive domain, where systems not only deliver high peek performance, but also do it in a power efficient manner. For this reason, we measure performance of considered CNNs under two major metrics, that are, the *number of frames* processed in one second and the *total power* dissipated by the SoC in the inference phase. We also introduce a last metric that is composed by FPS and power, hence captures the *power-efficiency* of the system. Unless specified, all throughput measures are in Frames-per-Second, latency measures are in milliseconds, and power consumption is in Watt.

Since NVIDIA boards provide frequency scaling mechanism both for the host/CPU and accelerator/GPGPU, for Tegra X2 and Xavier we show 3D graphs where the x-axis reports the GPU frequency ranging from about 100MHz up to 1,3GHz. In the y-axis is reported the whole host frequencies, that ranging from about 300MHz up to about 2GHz. Finally, the metric of interest is reported in z-axis.

For the Host part, here we have a single core for feeding video frames to the network and one core for reading the

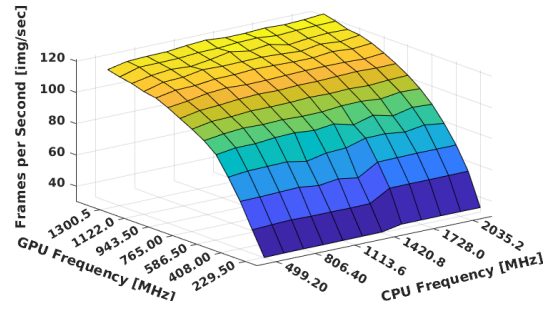


Fig. 3. AlexNet on TX2, FPS @ 19V

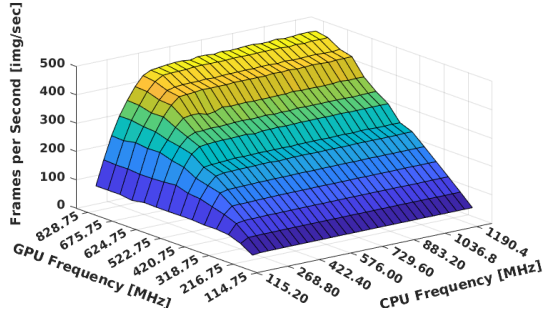


Fig. 4. AlexNet on Xavier, FPS @ 19V

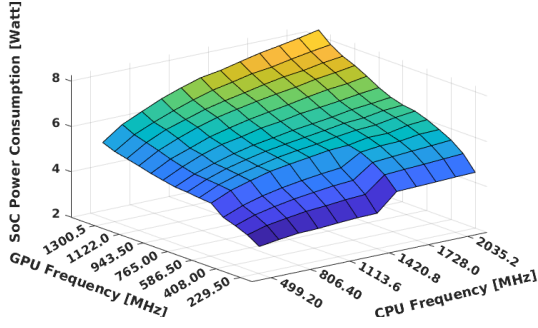
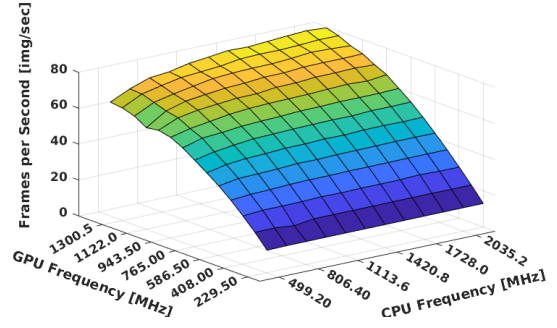
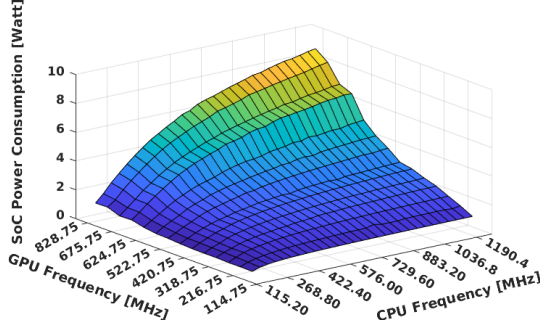
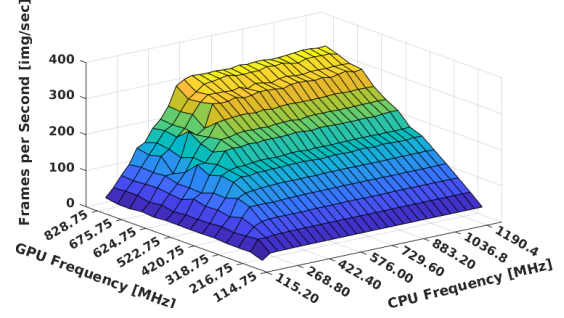
PS Freq. [MHz]	PL Freq. [MHz]	T.Put [FPS]	Pboard [Watt]	Pinf [Watt]	Esys [FPS/Watt]
2400	200	61.554	23.1165	0.5097	120.7652
1110	99	0.2462	20.2231	0.4603	0.534900

TABLE I
ALEXNET ON ZYNQ ULTRASCALE+

values of the sensor; the others cores are off. Multi-threading and batching are explored later, in Section VI-E. We don't want our experiments to be disturbed by software interference from other application, and for this reason we set the *realtime* priority of the task that runs the network to the maximum, using (the native) Ubuntu Linux API `chrt`, so we avoid unwanted context switches and through `taskset -c` we set a cpu affinity for the task that handles the network and the one that read the sensor.

Classification: AlexNet. In the first comparison we show the AlexNet model, that as mentioned above is a CNN for image classification. Figures 3, 4, and V show the frames-per-second achieved, respectively on TX2, Xavier, and UltraScale+. For inference on this network, GPUs outperforms the other two accelerators by an order of magnitude, in particular the TensorRT implementation, surprisingly was able to reach up to about 120 FPS on TX2 and about 300 FPS on Xavier, while the other are more or less comparable: about 60 FPS on UltraScale+.

It is interesting to note how, in case of TX2 platform (Figure 3), performance does not seem to be affected much by host frequency scaling, but only by GPGPU scaling, and this means that the CPU can keep the pace of the accelerator also at lower frequencies. This is a positive effect on power

Fig. 5. AlexNet on TX2, P_{SoC} @ 19VFig. 7. Tiny-YOLO on TX2, FPS @ 19VFig. 6. AlexNet on Xavier, P_{SoC} @ 19VFig. 8. Tiny-YOLO on Xavier, FPS @ 19V

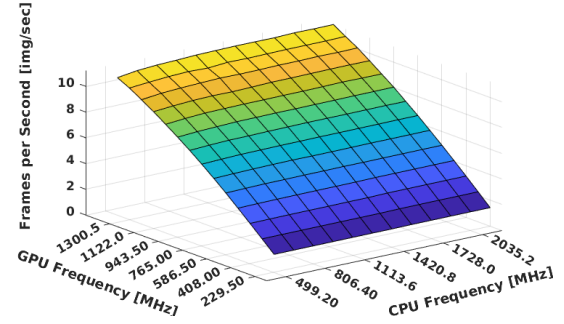
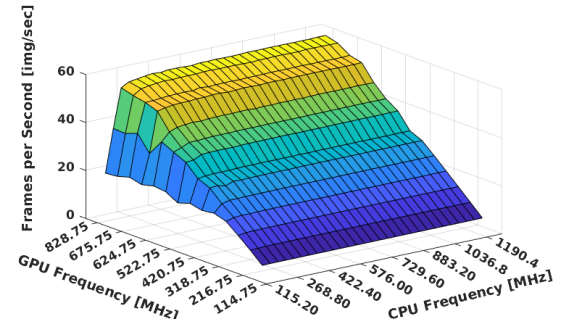
Model	PS Freq. [MHz]	PL Freq. [MHz]	T.Put [FPS]	Pboard [Watt]	Pinf [Watt]	Esys [FPS/Watt]
Tiny-YOLO	2400	200	22.6807	24.1798	1.1465	19.7826
Small-YOLO	2400	200	7.93110	23.7137	0.6804	11.6565
YOLO	2400	200	6.67288	23.5545	0.5217	12.7906

TABLE II
YOLOS ON ZYNQ ULTRASCALE+

efficiency, as shown in [3]³. Also, the GPGPU keeps good performance scaling with frequency, meaning that we can potentially run additional workload (in case, exploiting the missing ARM cores as controller) without performance loss. Similarly, the Xavier board (Figure 4), on the other hand, shows a tremendous peak performance of 300 FPS, which can potentially be improved increasing GPGPU frequency. Also, here there is poor we can do increasing host core frequency.

Regarding the power consumption, the P_{SoC} parameter, Figures 5 and 6 show comparable performance on GPGPUs (about 1.5 and 4.5 Watts on average for TX2 and Xavier). In figure V we report the power consumption P_{inf} of the FPGA-based accelerator, that was able to run AlexNet with about 500 milliWatts.

Detection: The YOLO-family. In this section, we show results for the detection network YOLO in its reduced version Tiny-YOLO, that compared to the full model, has less memory footprint and less complexity. Figures 7, 8, and V, show a big difference in terms of frame-per-second on GPUs compared to the others: in particular, this model on TX2 and Xavier can reach about 70 and 300 FPS respectively, while on UltraScale+ reach about 22 FPS. Regarding the

Fig. 9. YOLO on TX2, FPS @ 19VFig. 10. YOLO on Xavier, FPS @ 19V

power consumption, we have experimented the same trend as AlexNet: a power-hungry behavior on GPUs (about 7 and 10 Watts on average), to sustain the tremendous peak performance. UltraScale+ still consumes about an order of magnitude less power respect to the others architectures. In this case we have good performance on Xavier, about 55 FPS

³For reasons of space, we do not report here energy charts, as they are already published

(Figure 10) and comparable performance on TX2 and XU+: about 11 and 7 FPS respectively (Figures 9 and V).

Reasoning on results. In general, in terms of FPS, we can see that beyond particular frequencies of the CPU, the number of frames processed per second does not increase much, this makes us understand that in terms of power-efficiency it is not convenient to maximize the frequency of the processor, but there are fixed speeds in which the power-efficiency is maximized. For example, with reference to Figures 9 and 10 (where the frames processed by the YOLO network on TX2 and Xavier are shown) we can see that, on TX2, scaling the host core from about 350MHz to 2035MHz we have an increase of really few FPS, about 2 – 3. On Xavier (shown in Figure 10) we have the same behavior, supposing to not consider the 115, 192 and 268MHz frequencies that are excessively low to be able to feed the accelerator constantly with video frames. This behavior can be seen mainly on computationally-intensive neural networks, such as *Tiny-YOLO* and even more *YOLO*.

The charts in Figure 11 to 14 show the SoC power consumption of inference phase of YOLO variants on the considered platforms.

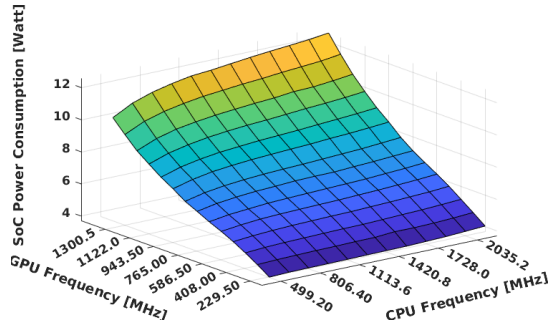


Fig. 11. Tiny-YOLO on TX2, P_{SoC} @ 19V

Observation: *For the reasons mentioned above, it is not always convenient to maximize the frequencies and in this way the power-efficiency surface takes an inverted 3D parabolic shape.*

Discussing about the accelerator frequency, we will discuss only the GPGPU case, because only Tegras expose a reasonably high number of frequency slots to let results be representative. In this case, with reference to the power graphs, for example 13 and 14, (where we can see the power dissipated by the *YOLO* network) we notice a steep growth in power consumption at high accelerator frequencies. In terms of FPS, on the other hand, we have noticed an opposite trend, and the framerate achieved increases slowly at high frequencies. This behavior leads to a maximum point in the power-efficiency function, in which the accelerator frequency is not at its higher value.

A. Power profiles (GPGPU only)

Besides core frequency scaling, NVIDIA GPGPUs offer a very powerful mechanism, that are, **power models**. This platform-specific feature enables different CPU-GPU frequency scaling in a programmer-friendly manner, without having to manually tune their cores' frequencies. They are

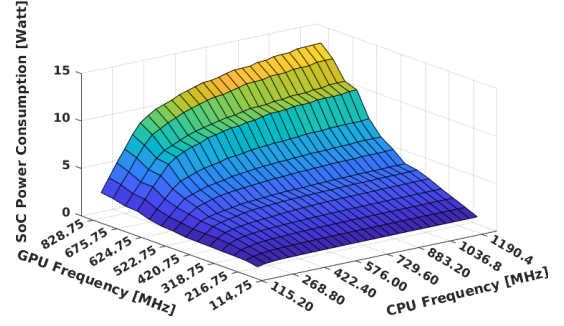


Fig. 12. Tiny-YOLO on Xavier, P_{SoC} @ 19V

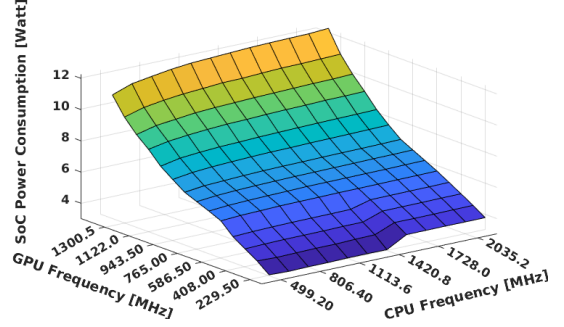


Fig. 13. YOLO on TX2, P_{SoC} @ 19V

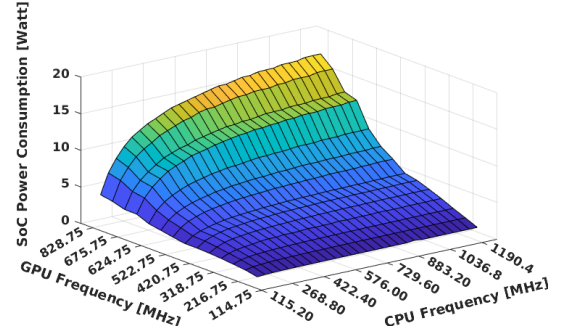


Fig. 14. YOLO on Xavier, P_{SoC} @ 19V

shown in Tables V-A and V-A. For all platforms, '0' is the most performing, hence power expensive, mode, while '1' is the other way around. Other models numbers, where available, are in between this spectrum, in increasing performance. We carried on an extensive set of experiments on two networks (namely, Yolo v2 and Yolo v3) on the three reference boards. Due to reasons of space, we report in Tables only results for the COCO dataset (and not YoloV2-tiny), but the trend is similar for all networks and datasets, and this is probably the most interesting result.

Observation: *power modes don't seem to significantly affect network precision (see mAP column), and do not significantly impact on any network or dataset more than others.*

Comparison matrix: all results together. This section is full of charts, and it's possible that even expert readers might get lost on them. For the sake of readability, tables in Figures V-A and V-A groups results that assess the throughput and power efficiency of the most interesting CNNs for each

Profile	TX2						Xavier						Nano					
	E2E latency	T.Put	Pow (board)	Pow (SoC)	mAP	P.eff	E2E latency	T.Put	Pow (board)	Pow (SoC)	mAP	P.eff	E2E latency	T.Put	Pow (board)	Pow (SoC)	mAP	P.eff
0 (perf)	145.10	6.88	12.54	8.873	23.90%	0.775	62.07	16.11	37.354	26.98	23.90%	0.597	405.39	2.466	8.9303	6.7476	0.239	0.365
1 (power)	228.70	4.37	6.346	3.952	23.90%	1.106	264.75	3.777	11.438	3.6746	23.90%	1.028	586.48	1.705	5.2426	3.1873	0.239	0.535
2	176.70	5.65	6.194	5.605	23.90%	1.008	121.06	8.26	15.8422	7.8147	23.90%	1.057	-	-	-	-	-	-
3	164.10	6.09	7.125	6.574	23.90%	0.926	97.17	10.29	20.2578	11.6641	23.90%	0.882	-	-	-	-	-	-
4	159.80	6.25	7.201	6.631	23.90%	0.943	94.39	10.59	20.9475	12.3842	23.90%	0.855	-	-	-	-	-	-
5	-	-	-	-	-	-	91.21	10.96	21.1356	12.5894	23.90%	0.871	-	-	-	-	-	-
6	-	-	-	-	-	-	89.40	11.19	21.2648	12.7357	23.90%	0.879	-	-	-	-	-	-

TABLE III
POWERMODES ON YOLO v2

Profile	TX2						Xavier						Nano					
	E2E latency	T.Put	Pow (board)	Pow (SoC)	mAP	P.eff	E2E latency	T.Put	Pow (board)	Pow (SoC)	mAP	P.eff	E2E latency	T.Put	Pow (board)	Pow (SoC)	mAP	P.eff
0 (perf)	291.48	3.43	15.105	11.438	31.80%	0.3	106.59	9.38	38.399	28.025	31.80%	0.335	-	-	-	-	-	-
1 (power)	431.77	2.32	7.904	5.51	31.80%	0.4	514.25	1.944	11.362	3.5986	31.80%	0.540	-	-	-	-	-	-
2	332.70	3.00	11.438	10.849	31.80%	0.3	211.50	4.72	16.283	8.2555	31.80%	0.572	-	-	-	-	-	-
3	333.63	2.99	11.362	10.811	31.80%	0.3	163.60	6.11	21.47	12.8763	31.80%	0.475	-	-	-	-	-	-
4	334.56	2.98	11.248	10.678	31.80%	0.3	163.35	6.12	21.337	12.7737	31.80%	0.479	-	-	-	-	-	-
5	-	-	-	-	-	-	161.45	6.19	21.66	13.1138	31.80%	0.472	-	-	-	-	-	-
6	-	-	-	-	-	-	162.24	6.16	21.717	13.1879	31.80%	0.467	-	-	-	-	-	-

TABLE IV
POWERMODES ON YOLO v3

SoC	Accel. Freq. [MHz]	Host Freq. [MHz]	T.Put [FPS]	Esys [FPS/Watt]
Xavier	624.7	499.2	164.2	54.093
XU+	200.0	2400	22.68	19.782
TX2	765.0	499.2	58.64	8.1613

TABLE V
TINY YOLO POWER EFFICIENCY COMPARISON

SoC	Accel. Freq. [MHz]	Host Freq. [MHz]	T.Put [FPS]	Esys [FPS/Watt]
Xavier	675.7	192.0	162.5	194.59
XU+	200.0	2400	61.55	120.77
TX2	943.5	343.6	107.6	22.906

TABLE VI
ALEXNET POWER EFFICIENCY COMPARISON

	Type	#classes	First layer (Yolo)	First layer (Tiny-Yolo)	Validation test size
COCO VOC	Generic	80	608 × 608	416 × 416	1000
	Generic	20	608 × 608	416 × 416	1000
Berkeley A (reduced)	Automotive	3	512 × 512	-	1000
Berkeley B (original)	Automotive	10	512 × 288	-	1000

TABLE VII
DATA SETS THAT WE USED IN OUR EXPLORATION

platform, namely Tiny-YOLO and AlexNet, at the best of their performance.

VI. EXPERIMENTAL SETTINGS

A. Varying datasets

Table VI-A shows the main features of the four datasets we employ in our analysis. Datasets highly affect network precision, and (how we will now show) also performance. Often, network engineers undergo a complete retraining of their network with small variation of the same dataset, e.g., removing unnecessary classes, to obtain highly-tuned optimal

performance. In our work we decided not to explore this dimension for all the networks and dataset, as to do so would make the design space explode. We nonetheless analyzed a reduced version of “Berkeley, with only three classes, hence, potentially faster and more accurate than the reference full version of the dataset. Table(s) VI-A and VI-A show how target datasets perform on selected networks and platforms. Looking at mAP, we see it is in line with what expected and declared in the reference papers and guides for each network and dataset. Also, there is a slight difference in terms of precision between GPGPUs and FPGAs, due to low-level optimizations such as data types (FPGAs use INT8, while GPGPUs use different optimization for different networks, most of which are undisclosed in the production environment).

One interesting thing to note is that different networks and datasets affect system performance (throughput and E2E latency), which in turn affects power consumption.

Observation: *there is an (in)direct correlation between the dataset and the power consumption.*

Also here, we can note how FPGA-based platforms are at least twice power efficient as GPGPUs. Among GPGPUs, strangely, TX2 outperforms Xavier in power efficiency, and also Nano offer similar (≈ 4 vs. ≈ 6) numbers, while consuming $\approx 1/4$ of power.

B. Quantization (FPGA only)

Quantization is one of the most significant optimization that can be employed on artificial neural networks, as they affect both performance and precision. The aim of quantization phase is to employ fixed point integers for network’s weights and activation to better exploit the characteristics of the underlying hardware, without heavily impact on average precision [20] Table(s) VI-A shows how a selected network and dataset performs under quantization. In this work we used INT8 quantization only on FPGA-based systems, because GPGPUs are highly efficient on performing floating-point computing, so we do not expect big improvement in terms of mAP/efficiency

Model	Training set	TX2						Xavier						Nano					
		E2E latency	T.Put	Pow (Board)	Pow (SoC)	mAP	P.eff	E2E latency	T.Put	Pow (Board)	Pow (SoC)	mAP	P.eff	E2E latency	T.Put	Pow (Board)	Pow (SoC)	mAP	P.eff
Yolo v2	COCO	145.10	6.88	12.54	8.873	23.90%	0.775	62.07	16.11	37.354	26.98	23.90%	0.597	405.39	2.466	8.9303	6.7476	23.90%	0.365
Yolo v2	VOC	70.10	14.3	10.678	7.011	58.82%	2.040	32.46	30.8	32.87	22.496	58.77%	1.369	175.69	5.691	9.5826	7.3999	58.77%	0.769
YOlo v2-Tiny	COCO	31.40	31.77	7.999	4.332	5.30%	7.334	11.81	84.67	22.8162	12.442	5.30%	6.805	60.84	16.43	5.94	3.7573	5.30%	4.373
Yolo v3-Tiny	COCO	13.29	75.26	10.309	4.116	6.48%	18.287	4.97	201.34	17.2258	6.4750	6.48%	31.095	42.01	23.80	7.97	2.8877	6.54%	8.243
Yolo v3	VOC	147.98	6.76	14.712	4.116	78.35%	1.642	73.22	13.658	32.255	6.475	78.35%	2.109	-	-	-	-	-	-
Yolo v3	COCO	291.48	3.43	15.105	11.438	31.80%	0.300	106.59	9.38	38.399	28.025	31.80%	0.335	-	-	-	-	-	-
Yolo v3	Berkeley	129.90	7.69	14.174	10.507	38.50%	0.732	61.52	16.25	41.515	31.141	38.57%	0.522	-	-	-	-	-	-

TABLE VIII
FLOATING POINT VERSIONS OF YOLO V3, VARYING DATASETS ON GPGPU

Model	Training set	Ultra96						ZCU102					
		E2E latency	T.Put	Pow (Board)	Pow (SoC)	mAP	P.eff	E2E latency	T.Put	Pow (Board)	Pow (SoC)	mAP	P.eff
Yolo v2	COCO	-	-	-	-	-	-	-	-	-	-	-	-
Yolo v2	VOC	85.956	11.5	9.18	2.604	63.78%	4.416	40.064	24.96	37.824	5.46	64.56%	4.571
Yolo v2-Tiny	COCO	-	-	-	-	-	-	-	-	-	-	-	-
Yolo v3-Tiny	COCO	83.542	11.97	7.836	1.26	4.60%	9.500	71.994	13.89	34.44	2.076	4.60%	6.691
Yolo v3	VOC	150.82	6.63	9.912	3.336	73.17%	1.987	69.444	14.4	38.364	6.00	75.07%	2.400
Yolo v3	COCO	950.57	1.052	8.088	1.512	31.30%	0.696	319.3	3.1318	34.656	2.292	31.30%	1.366
Yolo v3	Berkeley (A)	22.492	44.46	9.576	3.00	19.44%	14.820	11.76	85.03	36.552	4.188	19.44%	20.303
Yolo v3	Berkeley (B)	157.9	6.333	9.912	3.336	55.20%	1.898	72.463	13.80	37.848	5.484	55.20%	2.516

TABLE IX
QUANTIZED VERSIONS OF YOLO V3, VARYING DATASETS ON FPGA

Quantization	Ultra96						ZCU102					
	E2E latency	T.Put	Pow (Board)	Pow (SoC)	mAP	P.eff	E2E latency	T.Put	Pow (Board)	Pow (SoC)	mAP	P.eff
INT8	86.956	11.5	9.180	2.604	63.78%	4.41628	40.064	24.96	37.824	5.46	64.56%	4.57140
INT8 Pruned 22g	31.847	31.4	9.744	3.168	62.94%	9.91161	19.109	52.33	37.092	4.728	63.10%	11.0681
INT8 Pruned 24g	28.116	35.56	9.684	3.108	61.41%	11.4436	17.143	58.33	36.756	4.392	61.96%	13.2809
INT8 Pruned 26g	24.473	40.86	9.702	3.126	61.59%	13.0710	15.213	65.73	36.984	4.62	61.25%	14.2272

TABLE X
QUANTIZED AND PRUNED VERSION OF THE YOLO V2 PROVIDED WITH XILINX SDK

on these latter architectures. We show results for a single network and dataset (namely, Yolo v2 trained with VOC), as we want to show the variation of performance compared to baseline (in this case, INT8), however, similar results hold also for other configuration.

C. Pruning (FPGA only)

Pruning a neural network means reducing its size and complexity by removing parameters and activations that "do not play a big role" inside the whole network [5]. Also in this case we did not exploit this technique on NVIDIA GPGPUs, because their inference engine framework, called *TensorRT*, does not allow network pruning. We show results in Table VI-B. The 'g' of names in the first column ('22g', '24g', '26g'), is used by Xilinx to label the declared performance in GFLOPs of the network. Intuitively, the higher the number, the higher the throughput, and this is reflected in our tables. Results are available only for FPGA platforms, as such a fine-tuning of the network is not possible in the NVIDIA boards. We see how pruning and quantization affects performance, but this does not reflect in higher power consumption, nor harnesses network precision.

Observation: *This apparently counter-intuitive result is due to low-level network optimization (remember that FPGA-based system enable engineers coding neurons nearly in hardware),*

		E2E latency	T.Put	Pow (Board)	Pow (SoC)	mAP	P.eff
Full Yolo V3	TX2	291.48	3.43	15.11	11.44	31.80%	0.30
	Xavier	106.59	9.38	38.40	28.03	31.80%	0.33
	NANO	-	-	-	-	-	-
	Ultra96	950.57	1.052	8.09	1.51	31.30%	0.70
	ZCU102	319.00	3.1318	34.66	2.29	31.30%	1.37
Full Yolo V2	TX2	145.10	6.88	12.54	8.87	23.90%	0.78
	Xavier	62.07	16.11	37.35	26.98	23.90%	0.60
	NANO	405.39	2.466	8.93	6.75	23.90%	0.37
	Ultra96	86.96	11.5	9.18	2.60	63.78%	4.42
	ZCU102	400.64	24.96	37.82	5.46	64.56%	4.57
Tiny Yolo V2	TX2	31.40	31.77	8.00	4.33	5.30%	7.33
	Xavier	11.81	84.67	22.82	12.44	5.30%	6.81
	NANO	60.84	16.43	5.94	3.76	5.30%	4.37
	Ultra96	-	-	-	-	-	-
	ZCU102	-	-	-	-	-	-
Tiny Yolo V3	TX2	13.29	75.262	10.31	4.12	6.48%	18.29
	Xavier	4.97	201.34	17.23	6.48	6.48%	31.10
	NANO	42.01	23.804	7.97	2.89	6.54%	8.24
	Ultra96	835.42	11.97	7.84	1.26	4.60%	9.50
	ZCU102	719.94	13.89	34.44	2.08	4.60%	6.69

TABLE XI
TINY VERSION OF THE YOLO V2 AND V3 ON REFERENCE PLATFORMS

and enables higher power efficiency almost "for free" at the same precision.

D. Tiny networks

Small network variants are extremely suitable for low-end, resource-constrained platforms ⁴. For instance, the NVIDIA NANO cannot support a full YOLO v3 network. On the other

⁴Numbers shown refer to experiments with the COCO dataset

hand, it might not be convenient to deploy a reduced network version on some platforms. This explains the empty cells in Table VI-D, where we show how different network variants perform on different platforms.

Observation: *Intuitively, while smaller platforms achieve highest power efficiency in running smaller networks, it might not always be worth, such as for instance in the case of Tiny-Yolo on Ultra96, that exhibit very poor mAP of $\approx 4.6\%$.*

E. Batching/Threading

The most effective/important optimization we employ consists of loading multiple images into the accelerator, this technique allows to improve the utilization of the accelerator. For example, each layer of the network will have some amount of overhead and synchronization required to compute forward inference. By computing more results in parallel, this overhead is paid off more efficiently. In NVIDIA platforms, this feature is called *batching*, and it is typically performed by a single application thread that “packs” images in a *batch* transfer to the GPGPU accelerator. On the other hand, the FPGA-based platforms that we target provided *multi-threaded* version of the host-side application, which is capable of delivering multiple images to the DNN IP exploiting the multicore host and multiple DNN cores. In our experiments, batching/threading does not affect network precision (mAP), hence, we don’t report it.

Detection networks. We carried on experiments on the most heavyweight and powerful network (YOLO v3 trained with COCO), and results are reported in Table VI-E Surprisingly, they show how, for the same network, performance do not vary much, and this is because of the size of network itself and of input frames (80 classes, frames are 608x608, see Table VI-A). On the other hand, on FPGA-based platforms we see interesting results. High parallelism in Ultra96 only slightly increases performance, either for Full or Tiny Yolo V3. The reason is that, the board is already overloaded by computation on this heavyweight network. On the other hand, because of the higher number of DNN cores instantiated, ZCU102 can still keep the scaling pace up to 4-8 thread (resp. for the two datasets), but then it also “surrenders” to the heavy workload.

One consideration we must draw on the different approaches to parallelism offered by the two platform families. While, apparently, the effect is the same (we increase accelerator utilization by offloading multiple data frames), as shown, the two approaches are deeply different.

Classification networks. In Figure 15 and 16 we can see the effect of batching multiple input images on GPGPU system, in terms of FPS and power consumption. From the graphs shown, we can see that through batching, TX2 and Xavier were able to reach about 500 and 2500 FPS respectively in 13 and 25 Watts, and performance in Xavier could potentially scale more with the batch size.

Observation: *GPGPU batching might be ineffective with big detection networks, because the host system might not be able to process and pack all camera frames in real-time. Threading, instead, might play a role, where available, because it enables multiple concurrent stream of data.*

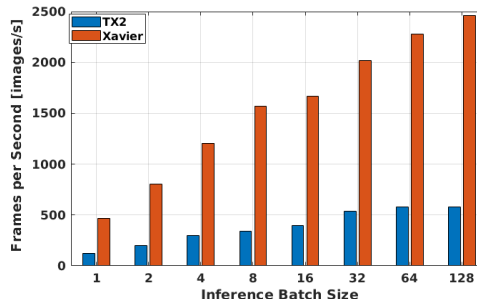


Fig. 15. Framerate achieved by AlexNet on TX2 and Xavier, quantized with 16-bit floating point numbers, varying the *batch size* parameter

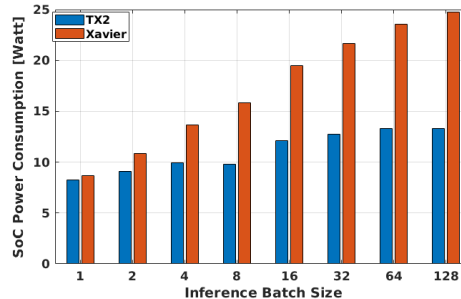


Fig. 16. Power consumption of AlexNet on TX2 and Xavier, quantized with 16-bit floating point numbers, varying the *batch size* parameter

Finally, few very important final consideration follow.

Observation: *in none of our tests we caused a bottleneck neither in system bus nor in system memory, when transferring frames back and forth the accelerator. This seems to indicate that the complex memory/bus hierarchy of modern platforms is such powerful, that, at least for the considered workloads, its limits are still far away to be reached.*

Reasoning on results. In GPGPU systems, the application code must collect a number of data frame (i.e., images from multiple cameras), before being able to pack and offload it to the accelerator. Especially, in NVIDIA systems, the NN package and inference engine drivers give poor freedom to the system designer to access camera in parallel from within the same application. This aspect could certainly can be improved, if the system were not proprietary and closed. A possible enhancement for reconfigurable platforms, and especially Xilinx’s SoMs, is that it is possible for engineers to direct the video stream across the FPGA to get through the DRAM banks. This point is really interesting, as it potentially enables on-the-flight detection on the video stream without involving the host cores, and deserves more investigation from system engineers and experts of hardware design.

VII. CONCLUSIONS

In this work, we evaluated the power-efficiency and performance of image classification and objects detection Convolutional Neural Networks, on embedded platforms representative of next generation automotive domain controllers. Our results assess higher peek performance, in terms of frame-per-second on GPGPU-based accelerators compared to the reconfigurable

Network	Batch size /Nthreads	TX2		Xavier		Nano		Ultra96		ZCU102	
		E2E latency	T.Put	E2E latency	T.Put	E2E latency	T.Put	E2E latency	T.Put	E2E latency	T.Put
Full Yolo	1	303.41	3.296	153.38	6.520	-	-	950.57	1.052	319.30	3.132
Full Yolo	2	604.76	1.654	293.42	3.408	-	-	950.57	1.959	319.30	6.132
Full Yolo	4	1215.53	0.823	583.09	1.715	-	-	950.57	3.007	319.30	11.459
Full Yolo	8	2397.03	0.417	1156.41	0.865	-	-	950.57	3.051	319.30	15.662
Full Yolo	16	4826.71	0.207	2362.29	0.423	-	-	950.57	3.068	319.30	15.586
Yolo Tiny	1	13.29	75.262	4.97	201.341	42.01	23.804	83.54	11.97	71.99	13.890
Yolo Tiny	2	26.38	37.911	9.57	104.520	82.74	12.086	83.54	21.189	71.99	27.530
Yolo Tiny	4	52.54	19.032	18.79	53.230	166.37	6.011	83.54	36.548	71.99	52.930
Yolo Tiny	8	106.28	9.410	37.44	26.708	338.02	2.958	83.54	39.219	71.99	59.410
Yolo Tiny	16	211.17	4.736	74.50	13.422	678.90	1.473	83.54	35.96	71.99	59.860

TABLE XII
EFFECT OF BATCHING/MULTI-THREADING ON FULL AND TINY VERSION OF THE YOLO V3

logics. Regarding the AlexNet network, with the batch size parameter set at 128, the new NVIDIA Xavier is resulted faster than TX2 of about $5\times$, with an increase factor of $1.92\times$ in power consumption. We experienced comparable performance for the object detection networks, due to their complexity and size, while for classification network, GPGPUs still outperform (non-optimized) FPGAs by at least one order of magnitude. On the other hand, FPGAs are better from the power consumption and power-efficiency viewpoint. Figures V-A and V-A show what are probably the most interesting (and useful) outcomes of our effort, a comparison of the two main CNN models, namely AlexNet and Tiny-YOLO, in terms of power-efficiency (Formula 3). They show the highest power-efficiency of XU+, because the power consumption due to the FPGA fabric is very low, which was somehow expected.

We also explored both platform optimizations, namely power scaling/profiles, batching and threading, and network-specific features such as quantization and multiple datasets, assessing not only the performance/power dimension, but also how they do affect network precision. Interesting results act as guidelines for engineers, and as starting point for future works, when new platforms and networks come to light.

REFERENCES

- [1] Apache Software Foundation. Apache MXNet: A Flexible and Efficient Library for Deep Learning.
- [2] M. Blott (Xilinx, Inc.). Tiny YOLO: a real-time, low-latency, low-power object detection system running on a Zynq UltraScale+ MPSoC, 2017.
- [3] G. Brilli, P. Burgio, and M. Bertogna. Convolutional neural networks on embedded automotive platforms: A qualitative comparison. In *2018 International Conference on High Performance Computing & Simulation (HPCS)*, pages 496–499, 2018.
- [4] J. Chen and X. Ran. Deep learning with edge computing: A review. *Proceedings of the IEEE*, 107(8):1655–1674, Aug 2019.
- [5] B. D., G. O. JJ., F. J., and G. J. What is the state of neural network pruning. *Conference on Machine Learning and Systems*, 2020.
- [6] W. Dai and D. Berleant. Benchmarking contemporary deep learning hardware and frameworks: a survey of qualitative metrics. *IEEE Transactions of intelligent vehicles*, 07 2019.
- [7] C. Ding, S. Wang, N. Liu, K. Xu, Y. Wang, and Y. Liang. REQ-YOLO: A Resource-Aware, Efficient Quantization Framework for Object Detection on FPGAs, 2019.
- [8] J. Y. et al. Real-time object detection towards high power efficiency. In *2018 Design, Automation Test in Europe Conference Exhibition (DATE)*, pages 704–708, 2018.
- [9] Google Brain. TensorFlow.
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [11] S. Hossain and D.-J. Lee. Deep Learning-Based Real-Time Multiple-Object Detection and Tracking from Aerial Imagery via a Flying Robot with GPU-Based Embedded Devices. *Sensors (Basel, Switzerland)*, 19(15), July 2019.
- [12] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [13] I. Hubara, et al. Binarized neural networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4107–4115. Curran Associates, Inc., 2016.
- [14] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <0.5mb model size. *arXiv:1602.07360*, 2016.
- [15] International Organization for Standardization / Technical Committee 22 (ISO/TC 22). ISO/DIS 26262-1 - Road vehicles ? Functional safety. Technical report, International Organization for Standardization / Technical Committee 22 (ISO/TC 22), Geneva, Switzerland, July 2009.
- [16] J. Johnson. Benchmarks for popular convolutional neural network models on CPU and different GPUs, with and without cuDNN, 2017.
- [17] C. E. Kim, M. M. D. Oghaz, J. Fajtl, V. Argyriou, and P. Remagnino. A comparison of embedded deep learning methods for person detection. *CoRR*, abs/1812.03451, 2018.
- [18] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *Proceedings of NIPS'12*, pages 1097–1105, USA, 2012. Curran Associates Inc.
- [19] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*, pages 21–37. Springer, 2016.
- [20] C. M., B. Y., and D. J.P. Binaryconnect: Training deep neural networks with binary weights during propagations. *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- [21] Meta AI. PyTorch.
- [22] Microsoft. Microsoft Cognitive Toolkit.
- [23] S. Mittal. A Survey on Optimized Implementation of Deep Learning Models on the NVIDIA Jetson Platform. *Journal of Systems Architecture*, 12 2018.
- [24] Mobileye. Mobileye an intel company.
- [25] H. Nakahara, H. Yonekawa, T. Fujii, and S. Sato. A Lightweight YOLOv2: A Binarized CNN with A Parallel Support Vector Regression for an FPGA. In *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays, FPGA '18*, pages 31–40, New York, NY, USA, 2018. ACM.
- [26] Nervana Systems. Neon.
- [27] N. Ni and V. Kathail (Xilinx, Inc.). Caffe to Zynq: State-of-the-Art Machine Learning Inference Performance in Less Than 5 Watts, 2017.
- [28] E. Nurvitadhi, G. Venkatesh, J. Sim, D. Marr, R. Huang, J. Ong Gee Hock, Y. T. Liew, K. Srivatsan, D. Moss, S. Subhaschandra, and G. Boudoukh. Can fpgas beat gpus in accelerating next-generation deep neural networks? In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays, FPGA '17*, page 5–14, New York, NY, USA, 2017. Association for Computing Machinery.
- [29] NVIDIA. Jetson TX2 Module, 2017.
- [30] NVIDIA. NVIDIA DRIVE PX: scalable AI Supercomputer For Autonomous Driving, 2017.

- [31] NVIDIA. Jetson AGX Xavier: Deep Learning Inference Benchmarks, 2018.
- [32] NVIDIA. Jetson AGX Xavier Developer Kit, 2018.
- [33] NVIDIA. NVIDIA Jetson AGX Xavier Delivers 32 TeraOps for New Era of AI in Robotics, 2018.
- [34] NVIDIA. Jetson Nano Developer Kit, 2019.
- [35] J. Redmon. YOLO: Real-Time Object Detection.
- [36] J. Redmon and A. Farhadi. YOLO9000: better, faster, stronger. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 6517–6525, 2017.
- [37] K. Rungsuptaweekoon, V. Visoottiviseth, and R. Takano. Evaluating the power efficiency of deep learning inference on embedded GPU systems. In *2017 2nd International Conference on Information Technology (INCIT)*, pages 1–5, 2017.
- [38] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018.
- [39] Shri Sundaram. Building autonomous vehicles using Drive PX 2, 2017.
- [40] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [41] D. Wang, K. Xu, and D. Jiang. PipeCNN: An OpenCL-based open-source FPGA accelerator for convolution neural networks. In *2017 International Conference on Field Programmable Technology (ICFPT)*, pages 279–282, Dec 2017.
- [42] Xilinx. PetaLinux Tools.
- [43] I. Xilinx. HLS based Deep Neural Network Accelerator Library for Xilinx Ultrascale+ MPSoCs, 2018.
- [44] Xilinx, Inc., . Xilinx Zynq Ultra96 Evaluation Kit.
- [45] Xilinx, Inc., . Xilinx Zynq UltraScale+ MPSoC ZCU102 Evaluation Kit.
- [46] Xilinx inc. Versal AI Core Series.
- [47] Y. Jia et.al. Caffe: Convolutional Architecture for Fast Feature Embedding. In *Proceedings of the 22Nd ACM International Conference on Multimedia, MM '14*, pages 675–678, New York, NY, USA, 2014. ACM.
- [48] Z. Zhu, D. Liang, S. Zhang, X. Huang, B. Li, and S. Hu. Traffic-sign detection and classification in the wild. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2110–2118, 2016.

APPENDIX

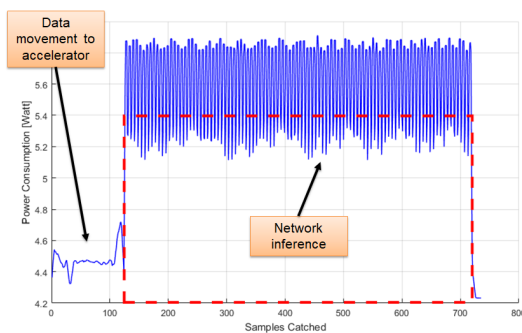


Fig. 17. approximation of the power consumption during the inference phase

Accurately measuring the power consumption is a challenging task, because the different hardware platforms considered in this work have different set of sensors and primitives.

The methodology we used for the **Zynq UltraScale+** is based on measuring the current that flows from the power supply into the Zynq board, because there are power sensors/counters that are not available on that board. This current flow was measured through a multimeter (configured as Amperometer) connected in series as the *VDD* and *GND* cables of the board power supply.

The total current absorption of the development board was measured with the precision of 40 samples-per-second, that is the highest for the multimeter that we used. Based on the samples caught by the amperometer, the dynamic power consumption has been approximated by multiplying the mean current absorbed and the voltage supply with the formula:

$$Pow_{(board)} = V_{DD} \times I_{mean} \quad (1)$$

Then we derived P_{inf} a power measure obtained by algebraic subtraction of the power computed with the previous formula and the power measured when the *System-on-Chip (SoC)* is in idle:

$$Pow_{(SoC)} = Pow_{(board)} - Pow_{(idle)} \quad (2)$$

This approximate metric roughly captures the power actually consumed in the inference process. Subsequently we approximate the power-efficiency of the system, combining FPS and power dissipation as in the following formula:

$$P.eff = \frac{FPS}{P_{mean}} \left[\frac{images/s}{Watt} = images/J \right] \quad (3)$$

For the power consumption of the FPGA, we compute the power consumption with the formulas 1 and 2. On the other hand, **NVIDIA Tegras** boards have on-board sensors (and related APIs directly accessible from software), hence for these platforms we directly access the power consumption using sensors.

Finally for NVIDIA TX2 and Jetson AGX Xavier, like the previous case, we used a sensor-based approach, through the *INA3221* sensor⁵. This sensor has three channels that can

⁵The `sysfs` nodes to read power on Jetson are located in `/sys/bus/i2c/drivers/ina3221x`.

monitor respectively the *VDD_IN*, *VDD_CPU* and *VDD_DDR* power rails, through an *I2C* connection.

For these experiments we need to monitor the *VDD_IN* rail, because is the rail that provides energy to the SoC itself. In Figure 17 we can see the typical trend in power consumption, that for example we can compute trough the Formula 1 in which we can recognize the *image loading phase* (on the left in the graph) and the *inference phase* (on the right), that consists of a sequence of power peaks. The dashed rectangle shows an approximation of the total energy needed during the inference phase.

It is important to note that we also run a reduced set of experiments by applying the methodology we used for the XU+ board i.e., physically cutting rails) also to the NVIDIA boards, and measured power was comparable to the ones from sensors. This proves that the approach of physically acting on the board is accurate enough for our purposes (besides the fact that it damages the board itself).

A. Evaluated Frameworks

Inference engines play a big role in maximising the performance of a network, on every platform, and for this reason they usually come with the SDK by platform provider. In our exploration, we employed the reference framework for every specific network and platform. This is the typical choice of system engineers, because it represents the best tradeoff among the metrics, and it's supported by the community and platform provider.

TensorRT is an optimized Deep Neural Network library for accelerating the inference phase, to exploit TensorRT we used tkDNN⁶, a framework written in C++ and CUDA. tkDNN is optimized to execute mainly on embedded SoC like for example the NVIDIA Tegra family, rather than GPU-based servers. We also used *TensorRT* directly, (through the *trtexec* tool), for testing AlexNet on TX2 and Xavier.

DPU ecosystem the Xilinx DPU is an optimized engine to accelerate DNNs on top of Xilinx's SoCs. DPU is released as IP-core with Xilinx SDK, and it can be easily inserted in a custom hardware design. The DPU operation requires the application processing unit (APU) to service interrupts to coordinate data transfer. DNNDK is a framework provided by Xilinx in the Ultrascale SDK, that allows to quantize and compile a deep learning model for a target DPU architecture. A machine learning engineer can define a custom neural network model with some of the most popular engines, like Caffe or Tensorflow and after the training phase it is possible to convert the model to fixed point and deploy it to a DPU engine. At the moment the DNNDK, framework can quantize a NN with an arbitrary fixed point precision, but the DPUv2 in this current release supports only the INT8 quantization, for these reason all the networks quantized with this framework have this type of quantization.

Other Frameworks tested in this work are Caffe, a deep learning framework developed by the Berkeley Vision and Learning Center (BVLC) and proposed in [47]. Caffe is a standard for implementing CNNs and it is well integrated in

⁶<https://github.com/ceccocats/tkDNN>

several others frameworks and libraries for embedded accelerators. xfDNN is a framework based on Caffe for accelerating deep neural networks on Xilinx UltraScale+ MPSoCs. This framework supports 16-bit integer data type. It became an opensource project called CHaiDNN [43].

PipeCNN is another CNN framework, written in OpenCL and proposed by D. Wang et al. in [41] that we used for implementing AlexNet on Xilinx UltraScale+. This framework is based on *pipes*, a new feature introduced in OpenCL 2.x.