
PREVIOUS: A METHODOLOGY FOR PREDICTION OF VISUAL INFERENCE PERFORMANCE ON IOT DEVICES

Delia Velasco-Montero
Instituto de Microelectrónica de Sevilla
Universidad de Sevilla-CSIC

Jorge Fernández-Berni
Instituto de Microelectrónica de Sevilla
Universidad de Sevilla-CSIC

Ricardo Carmona-Galán
Instituto de Microelectrónica de Sevilla
Universidad de Sevilla-CSIC

Ángel Rodríguez-Vázquez
Instituto de Microelectrónica de Sevilla
Universidad de Sevilla-CSIC

ABSTRACT

This paper presents PreVIOUS, a methodology to predict the performance of Convolutional Neural Networks (CNNs) in terms of throughput and energy consumption on vision-enabled devices for the Internet of Things. CNNs typically constitute a massive computational load for such devices, which are characterized by scarce hardware resources to be shared among multiple concurrent tasks. Therefore, it is critical to select the optimal CNN architecture for a particular hardware platform according to prescribed application requirements. However, the zoo of CNN models is already vast and rapidly growing. To facilitate a suitable selection, we introduce a prediction framework that allows to evaluate the performance of CNNs prior to their actual implementation. The proposed methodology is based on PreVIOUSNet, a neural network specifically designed to build accurate per-layer performance predictive models. PreVIOUSNet incorporates the most usual parameters found in state-of-the-art network architectures. The resulting predictive models for inference time and energy have been tested against comprehensive characterizations of seven well-known CNN models running on two different software frameworks and two different embedded platforms. To the best of our knowledge, this is the most extensive study in the literature concerning CNN performance prediction on low-power low-cost devices. The average deviation between predictions and real measurements is remarkably low, ranging from 3% to 10%. This means state-of-the-art modeling accuracy. As an additional asset, the fine-grained a priori analysis provided by PreVIOUS could also be exploited by neural architecture search engines.

Keywords

Vision-Enabled IoT, Edge Devices, Deep Learning, Convolutional Neural Networks, Inference Performance, Neural Architecture Search.

1 Introduction

implementation of visual processing at the edge, as opposed to the cloud, presents remarkable advantages such as reduced latency, more efficient use of bandwidth, and lessened privacy issues. These advantages are instrumental for boosting the application scenarios of the Internet-of-Things (IoT) paradigm [1, 2]. Edge vision algorithms must provide enough accuracy for practical deployments while making the most of the limited hardware resources available on embedded devices. Concerning accuracy, Deep Learning (DL) [3] has recently emerged as the reference framework. Deep Neural Networks (DNNs) resulting from training on massive datasets accomplish precise visual inference, greatly improving the performance of classical approaches based on hand-crafted features. However, this accuracy has a cost. The computational and memory requirements of DNNs are much more demanding than those of classical algorithms [4]. This constitutes a challenge when it comes to incorporating DNN-based inference in the processing flow of IoT devices, which is already heavy because of other functions related to networking, power management, additional sensors, etc.

The success of DL in enabling practical vision algorithms and unifying the procedure for a number of tasks such as image recognition, object detection, and pixel segmentation, has prompted research and development at various levels [5]. At software level, various open-source frameworks, both from academia and industry, are accessible on the internet; each of them exploits a particular set of libraries and core system functionalities. At architectural level, new DNN models are ceaselessly reported aiming at enhancing specific aspects, e.g., higher accuracy, faster training, or shorter inference time. Regarding hardware, the pervasiveness of DNNs is forcing the inclusion of ad-hoc strategies that exploit different features of neural layers to speed up their processing. Overall, this extensive DL ecosystem is making the optimal selection of inference components according to prescribed application requirements increasingly difficult in vision-enabled IoT devices.

To assist in the aforementioned selection, we already proposed a methodology based on benchmarking and a companion figure of merit in a previous study [6]. However, benchmarking entails a significant and non-scalable effort because of the complexity and diversity of software libraries, toolchains, DNN models, and hardware platforms. In this paper, we describe PreVIous, a novel methodology that removes the need of comprehensive benchmarking. This methodology is based on the single characterization of PreVIousNet, a Convolutional Neural Network (CNN) specifically designed to encode most of the usual parameters in state-of-the-art DNN architectures for vision. As a result of such characterization on a particular software framework and hardware device, a prediction model is generated. This model provides a precise per-layer estimation of the expected performance for any other CNNs to be eventually run on that software-hardware combination. Seven CNN models on two different software frameworks have been thoroughly characterized to demonstrate the prediction capacity of PreVIous. Regarding hardware, this study is focused on the multi-core Central Processing Units (CPUs) available on two different low-power low-cost platforms, but the methodology could be extended to other types of devices.

The manuscript is organized as follows. Section 2 summarizes related work and sets the context to point out the contribution of PreVIous to the state of the art. An overview of CNNs is provided in Section 3, where their usual layers and fundamental characteristics are briefly described. Section 4 elaborates on the main elements defining PreVIous and how it has been applied in practical terms in this study. The core of PreVIous, i.e., PreVIousNet, is further described in Section 5. The vast set of experimental results that confirm the modeling capacity of PreVIous is reported in Section 6. Finally, we draw the most relevant conclusions arising from these results in Section 7.

2 Related Work

As previously mentioned, the implementation of CNNs on resource-constrained devices is a remarkable challenge that has been addressed through various approaches. The common objective of all of them is to maximize throughput and inference accuracy while minimizing energy consumption.

Architecture design. Several strategies have been investigated to boost inference performance of CNNs. For instance, the well-known SqueezeNet model [7] features a massive reduction of parameters with respect to previous models while still achieving a notable accuracy. Other architectures tailored for embedded devices have also been proposed [8–14]. These models were designed to alleviate their computational burden as a whole. In other words, they were not specifically adapted for a particular platform. Therefore, their performance significantly varies depending on the host system [15]. A preliminary evaluation could be conducted by simply comparing the number of operations required for each network. However, this direct assessment does not usually translate into accurate values of measured performance metrics, in particular power consumption [16].

Performance benchmarking. To select the optimal CNN according to performance requirements on a particular platform, recent studies have carried out systematic benchmarking on several hardware systems [6, 17–19]. Notwith-

standing, benchmarking has a limited scope, considering the rapid evolution and massive number of CNNs, platforms, and software tools presently available.

Architecture optimization. Manual design of computationally efficient CNNs is time-consuming and requires experience. Several approaches for network compression have been investigated, ranging from network quantization [20–23] and channel pruning [24–27] to special network implementations [28–33]. Alternatively, CNN design is currently moving from manual tuning towards automatic algorithms. Neural Architecture Search (NAS) or sequential model-based optimization algorithms have been proposed to optimize metrics such as latency, energy, and accuracy on prescribed platforms [34–44]. The incorporation of specific platform constraints to such automatic approaches involves modeling how the network architecture relates with the optimization target. Inference performance modeling is highly valuable in this context.

Performance modeling. Various reported methods aimed at an objective similar to that of PreVIous in terms of performance prediction. We can classify them according to their particular research focus:

1. *Hardware accelerators:* to leverage the energy efficiency of CNN accelerators such as Eyeriss [45], an energy estimation methodology was proposed [46], [47]. It relies on the energy costs of memory accesses and multiply-accumulate operations (MACs) at each level in the memory hierarchy. This methodology also includes an energy-aware pruning process for network optimization.
2. *GPU-based systems:* these platforms usually include power monitor tools, which facilitate modeling performance. On the basis of energy measurements extracted with such tools and taking network metrics and layer configurations as inputs, machine learning models targeting energy consumption and latency have been investigated [48–50].
3. *Cloud-mobile computing scenarios:* characterization of layerwise network performance enables finding the optimal DNN partition between the cloud and the mobile device while optimizing resources such as energy consumption, latency, and bandwidth [51–53].
4. *Automatic optimization algorithms:* modeling is especially valuable to drive the design of highly efficient convolutional and fully connected layers, thus accelerating optimization algorithms based on latency and energy metrics [35, 37, 42–44]. Performance models leveraged by these algorithms comprise from simple look-up tables to more complex machine learning models.
5. *Training optimization:* to reduce training costs, some works modeled the performance of high-end GPUs during training [54, 55].

Most of the performance modeling studies mentioned above revolve around high-end systems or energy-demanding platforms. However, IoT application scenarios clearly benefit from low-cost low-power devices. In addition, the reported models present limitations. They either focus on particular types of layers or involve an extensive benchmark. With PreVIous, we have addressed these key points. First, we have worked with inexpensive devices featuring enough computational power to perform CNN-based inference. Second, the performance models provide fine-grained per-layer information, covering many different types of layers with the characterization of a single CNN.

All in all, the main contributions of this study are:

- A methodology that allows to evaluate the performance of CNN models accurately layer by layer in terms of throughput and energy consumption. The CNNs do not need to be actually run to obtain this evaluation, which is automatically generated from a previously built predictive model.
- A neural network whose characterization enables the construction of the aforementioned prediction model. This network incorporates a large variety of CNN layers and interconnections between them in order to achieve fine-grained CNN profiling.
- Rapid identification of layers whose execution time or energy consumption is distinctively higher than others on a particular software-hardware combination. Such layers would be the first ones to be modified by an optimization procedure or a NAS engine.
- A broad analysis of CNNs running on different software frameworks and hardware platforms. This analysis, which has served the purpose of gauging the goodness of the proposed methodology, is intrinsically valuable as an extensive set of measured performance metrics.

3 Overview of CNNs

Next, as a basis for subsequent sections, we summarize fundamental aspects of CNNs, including architectural details, key network metrics, and typical implementation strategies.

3.1 Common layers

CNN architectures usually comprise a heterogeneity of layers. In most architectures, the first layer is fed with a 3-channel input image, and the network progressively reduces the spatial dimensions ($H \times W$) of feature maps (*fmaps*), while increasing the number of channels C . Thus, each layer takes a 3D input tensor \mathbf{I} , performs operations that involve a set of learnt weights \mathbf{W} , and generates output data \mathbf{O} for the next layer. Typical layers included in the majority of state-of-the-art CNNs, and covered by PreVIOUSNet, are briefly described below.

- **Convolutional (CONV).** Input data are convolved with a 4D tensor \mathbf{W} composed of N kernels of dimensions $k_h \times k_w \times C_{in}$. The n -th kernel $\mathbf{W}^{(n)}$, $n = 1, \dots, N$ yields the n -th 2D output feature map in which the activations are obtained from:

$$\mathbf{O}_{x,y,n} = \sum_{c=1}^{C_{in}} \sum_{i=1}^{k_w} \sum_{j=1}^{k_h} \mathbf{W}^{(n)}_{i,j,c} \mathbf{I}_{x+i,y+j,c}. \quad (1)$$

This convolutional layer requires $k_h k_w C_{in} H_{out} W_{out} N$ MACs. In general, learnt biases, denoted by $b^{(n)}$, are also added to each output, adding $H_{out} W_{out} N$ MACs to the computation. The operation of this layer is illustrated in Fig. 1.

- **Fully Connected (FC).** These layers are usually located at the end of the network to perform classification on the extracted feature maps. Similar to classical neural networks, the operating data are arranged in 1D vectors. A weight factor is applied to each connection between input and output activations. Additional biases can be added. Generally assuming N_{in} inputs that yields N_{out} outputs, a FC layer involves a computational cost of $N_{in} N_{out}$ MACs.
- **Pooling.** This type of layer lowers the spatial dimensions of *fmaps* by applying a simple operation to each $k_h \times k_w$ patch with a stride s . A total of $k_h k_w H_{out} W_{out} C_{out}$ operations – not necessarily MACs – are performed, with $H_{out} = \lfloor \frac{H_{in} - k_h}{s} + 1 \rfloor$. *Maximum* and *average* are the most usual functions employed to reduce dimensionality.
- **Rectified Linear Unit (ReLU).** To introduce non-linearities between layers, various functions are applied, among which ReLU is the most popular one. It performs the simple operation of selecting the maximum between each input activation $I_{i,j,c}$ and 0. This simplicity speeds up the calculation of non-linearities with respect to activation functions such as *sigmoid* and *tanh*. In addition, ReLU is more suitable for rapid training convergence [56].
- **Batch-Normalization (BN).** Currently, this is the most popular normalization layer implemented in state-of-art CNNs for training acceleration. It normalizes activations on the i -th channel in terms of zero-mean and one-variance across the training batch [57]. Two weights per channel are learnt (scale and shift), and two operations per activation are performed.
- **Concatenation (Concat)** of data from multiple layers, usually along the channel dimension, is convenient for merging branches in the network. For instance, this is the last layer within the *Inception* module included in a number of CNNs [58–60]. No mathematical operation is performed, only data reorganization.
- **Element-wise Operation (Eltwise).** This layer performs element-wise operations such as addition, product, or maximum on multiple input activations.

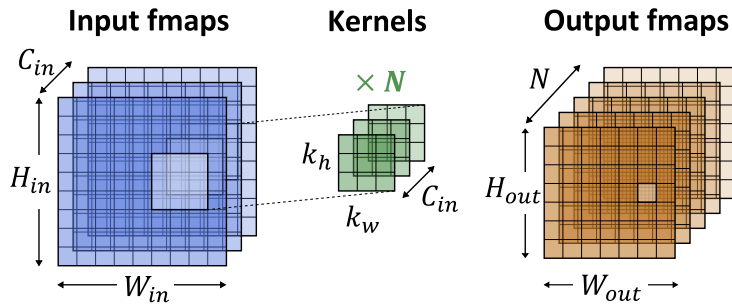


Figure 1: Convolutional layers constitute the core operation of CNNs. Each kernel filter – depicted in green – operates on sliding local regions of the input fmaps – *receptive fields* in light blue – to produce the corresponding output fmaps.

- **Scale.** This layer multiplies each input activation by a factor, thus requiring *HWC* MACs. Optionally, biases can be added.
- **Softmax** is the most notable loss function for classification tasks. It outputs a normalized probability distribution from a vector of N class predictions by applying the function $softmax(I_c) = \frac{e^{I_c}}{\sum_{p=1}^N I_p}$.

3.2 Intrinsic CNN Metrics

For practical deployment of a network on resource-constrained devices, it is worth considering the following parameters:

1. **Accuracy.** Once trained on a dataset for a specific application, a CNN provides a particular inference precision. In this regard, *Top-N accuracy* and *mean Average Precision (mAP)* are the most frequently applied metrics to evaluate classification [61] and detection [62] tasks, respectively.
2. **Computational Complexity (#OPs).** A widely used metric to measure computational complexity is the number of floating-point operations required by a network – or, at least, those required by CONV and FC layers. The overall computational load can be determined by adding up the number of operations per layer – previously detailed in Section 3.1¹.
3. **Model Size.** It refers to the total amount of learnable parameters of a network. The memory footprint may preclude the execution of certain models on specific platforms.
4. **Memory Accesses.** In addition to network weights, relevant activations must be kept in memory during inference. The minimum number of basic memory operations for a layer forward-pass will be²:

$$\#memOPs = n(\mathbf{I}) + n(\mathbf{W}) + n(\mathbf{O}) \quad (2)$$

where $n(\mathbf{X})$ denotes the number of elements in the tensor \mathbf{X} . Thus, for example, $n(\mathbf{I})$ is equal to $H_{in}W_{in}C_{in}$.

Note that these intrinsic metrics do not directly reflect actual inference performance. However, they provide a preliminary estimate of the resources required by a network.

3.3 Inference Metrics

Relevant metrics concerning inference performance must be measured during forward-pass:

1. **Throughput.** Real-time applications rely on processing images at a prescribed frame rate. CNN inference runtime limits the maximum achievable throughput for the related computer vision algorithm.
2. **Energy Consumption.** Battery lifetime is one of the most critical constraints on embedded platforms. Therefore, a key parameter is the total energy demanded by the system during inference.

3.4 CNN implementation strategies

The way in which the network architecture relates with inference metrics is highly dependent on the CNN implementation. The software libraries underlying a particular framework implement a diversity of optimization strategies to accelerate matrix multiplication according to the available hardware resources. The most commonly implemented approach is the so-called *unrolled convolution*, in which convolutions are performed through image-to-column transformation (*im2col*) plus General Matrix-to-Matrix Multiplication (GEMM). Thus, after *im2col*, convolution *receptive fields* are unrolled into columns, whereas filters are unrolled into rows. As a result, the convolution becomes a matrix-to-matrix product that can be highly-optimized through several libraries such as ATLAS [63], OpenBLAS [64], MKL [65], and cuBLAS [66]. However, this performance optimization increases the allocated memory owing to the unrolled receptive fields. This memory overhead must be taken into account in Eq. (2), where $n(\mathbf{I})$ becomes $(k_h k_w C_{in})(H_{out} W_{out})$ for CONV layers implementing this strategy. Other approaches that have been applied to accelerate matrix multiplication include Fast Fourier Transform (FFT) [29], Winograd [28], and Strassen [30] algorithms.

¹We provide a general estimation on the minimum number of operations required for inference. Ultimately, it will be the specific interaction between hardware and software in the targeted system that will determine the actual computational complexity.

²This equation is again a plain estimation. The number of memory accesses will ultimately depend on the hardware platform – memory word size, memory hierarchy, cache size, etc. – and the computational strategy for each operation – partial matrix products, data access pattern, etc.

4 PreVIous: a Framework for Modeling and Prediction of Visual Inference Performance

4.1 General Description

The a priori evaluation of CNN performance directly on the basis of network complexity is inaccurate, even when considering only a specific hardware device [16]. For the sake of more precise and multi-platform modeling, two important aspects must be stressed:

- Network inference involves numerous types of computational operations and data access patterns. For example, FC layers include an elevated number of weights, thus requiring a great deal of memory operations, whereas ReLU layers only perform a simple operation on activations.
- Energy consumption is also highly dependent on the particular characteristics of the layers and how they are mapped into the underlying hardware resources. For instance, the energy cost of memory access varies up to two orders of magnitude depending on the considered level within the memory hierarchy [5].

Therefore, we propose to characterize the expected performance of CNN inference through *per-layer predictive models*. Not only does it make more sense according to the two aspects just mentioned, but also per-layer performance assessment is valuable for network architecture design, layer selection, and network compression.

Fig. 2 illustrates PreVIous. Basically, it comprises a first stage where a prediction model is constructed upon the characterization of PreVIousNet on the selected system, which is defined as a software framework implemented on a hardware platform. This one-time constructed model is able to predict, in a second stage, the performance of any other CNN to be run on such a system in terms of runtime and energy. Next, we describe this framework in detail.

4.2 Selected System

PreVIous is agnostic with respect to the software-hardware combination for modeling and prediction. In this study, we focus on two popular software frameworks deployed on two low-cost hardware platforms. The baseline combination integrates Caffe [67] and Raspberry Pi (RPi) 3 Model B [68]. This embedded platform (sized $85 \times 56 \times 20 \text{ mm}^3$) features a Quad Core ARM Cortex-A53 1.2 GHz CPU on a Broadcom BCM2837 System-on-Chip, 1 GB RAM LPDDR2 at 900 MHz, different network interfaces, and an external micro-SD card to provide non-volatile storage capacity. The operating system is Raspbian v9.4 Linux Kernel v4.14. All the measurements on this system were taken after booting in console mode to boost CNN inference performance and reduce energy consumption.

Once confirmed the effectiveness of PreVIous on this baseline case, we extended our analysis by changing both software and hardware in the selected system. First, we built OpenCV [69] v4.0.1 on the RPi. Second, we built Caffe on a different hardware platform, namely Odroid-XU4 [70]. This CPU-based embedded system (sized $83 \times 58 \times 20 \text{ mm}$) is more suitable for high-performance IoT applications. Its Exynos 5422 SoC implements the so-called big.LITTLE heterogeneous technology, arranging its multi-core architecture into two clusters: four “big” Cortex-A15 2 GHz cores

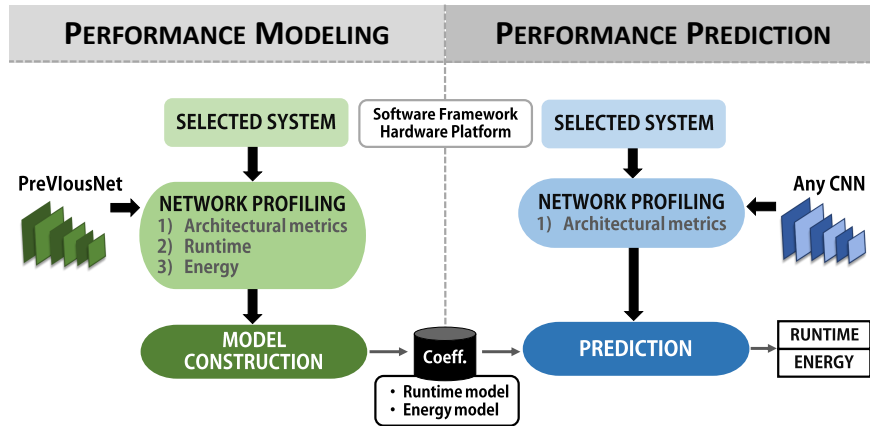


Figure 2: General overview of PreVIous. It comprises two stages: 1) performance modeling, where a prediction model is constructed for the selected system through the characterization of PreVIousNet; 2) performance prediction, where the performance of any CNN of interest to be run on the selected system is accurately predicted on the basis of the previously constructed model.

providing maximum computational performance and four power-efficient “LITTLE” Cortex-A7 1.4 GHz processors. Odroid-XU4 features 2 GB LPDDR3 RAM at 933 MHz; the operating system is Ubuntu v16.04 Linux Kernel v4.14. All the measurements were also taken in console mode while running CNNs on the cluster of cores clocked at 2 GHz. This is the configuration achieving maximum throughput.

4.3 Network Profiling

This step involves the extraction of three sets of measurements during the modeling stage. The first set, i.e., architectural metrics, is also extracted for network profiling during the performance prediction stage.

1. **Architectural Metrics.** Simply by parsing the network definition it is possible to extract per-layer architectural parameters such as input/output dimensions (H, W, C), number of learnt weights, i.e., $n(\mathbf{W})$, kernel sizes (k_h, k_w), as well as estimates of computational ($\#OPs$) and memory ($\#memOPs$) requirements.
2. **Time Profiling.** Each layer composing PreVIousNet is individually run to produce per-layer runtime profiling through software functions. Specifically, we employed the `time.time()` Python method to measure the elapsed time. To ensure accurate empirical measurements, 50 layer executions were averaged.
3. **Energy Profiling.** The layers of PreVIousNet are also individually characterized in terms of energy consumption. Some embedded platforms incorporate vendor-specific power meter tools to facilitate energy profiling. Otherwise, a power analyzer must be connected to the power supply pins of the selected system. In our case, we connected the Keysight N6705C DC power analyzer to the power supply pins of the aforementioned hardware platforms. As an example, Fig. 3(a) depicts the complete power profiling of All-CNN-C [71], the simplest among the seven networks characterized to assess the prediction capacity of PreVIous. A total of 50 executions per layer were carried out. The sampling period of the power analyzer was set to the minimum possible value, i.e., 40.96 μs . For proper identification of the layers, a time interval of 300 ms was established via software to separate each set of 50 executions. In addition, we used the previously obtained time profiling to extract the portion of the power signal corresponding with the layer under characterization. For instance, Fig. 3(b) shows the extracted signal for layer ‘conv2’ of All-CNN-C. Then, the energy consumption for each layer is obtained by integrating its power signal and averaging over the 50 performed executions.

4.4 Model Construction

After network profiling, *linear regression models* per type of layer are constructed for both runtime and energy consumption. In particular, PreVIous generates regression models for the diversity of layers listed in Section 3.1, all of which are distinctively covered by PreVIousNet, as described in Section 5. Architectural metrics play the role of predictors for such models. Thus, the performance prediction stage simply consists in parsing the definition of a CNN of interest to extract its architectural metrics and apply them to the corresponding regression models for its constituent layers.

For the sake of simplicity and reducing overfitting risk, we focused on linear models. Generally speaking, a linear regression model aims at finding the best combination of a set of variables $\mathbf{x} = [x_1, x_2, \dots, x_p]$ to predict the observations

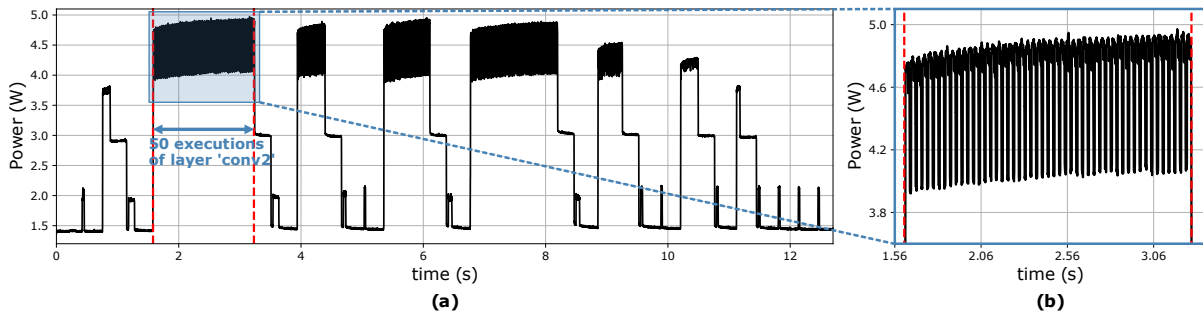


Figure 3: (a) Power signal provided by a DC power analyzer when carrying out 50 consecutive executions of each layer of All-CNN-C on the combination RPi-Caffe. A time interval of 300 ms was established via software to separate each set of executions; (b) Portion of the power signal corresponding to layer ‘conv2’ of All-CNN-C. This extracted signal is integrated and averaged over the 50 executions to obtain the energy consumption of the layer.

of a response y with minimum error. Given n observations of such variables and response, $\{\mathbf{x}_i, y_i\}_{i=1}^n$, the model can be expressed as:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon \quad (3)$$

where \mathbf{X} is the $n \times p$ matrix of predictor observations, \mathbf{w} is the $p \times 1$ vector of adjusted model coefficients, \mathbf{y} is the $n \times 1$ response vector, and ϵ is the $n \times 1$ error vector. Note that the observations for building the model are encoded by each row \mathbf{x}_i of the observation matrix \mathbf{X} .

According to this notation, PreVIous produces a regression model for each type of layer in Section 3.1 based exclusively on architectural metrics (\mathbf{x}) to estimate runtime or energy (y). We must point out that the observation set $\{\mathbf{x}_i, y_i\}_{i=1}^n$ could be obtained by collecting performance data from many different CNNs. However, two undesirable facts arise from this procedure: (1) the characterization of CNNs is time-consuming and burdensome; (2) different CNNs still share many layers with similar parameters and therefore their characterization would be redundant. Alternatively, we propose a simpler and systematic approach to produce the observation set. A single characterization is needed, i.e., that of PreVIousNet, which is a new CNN that encompasses a large variety of layers and sweeps the most usual layer parameters, features, and data dimensions. Hence, the architecture design space is comprehensively covered.

To make the most of simple linear regression models, two important points must be taken into account:

- *Dimensionality reduction.* The higher the number of predictors, the most likely the construction of an overfitted model. Thus, only variables highly correlated with the target response y must be considered for the model.
- *Model regularization.* This allows to make predictions less sensitive to a reduced set of observations. In particular, we apply standardized Ridge regularization in which the coefficients are obtained by minimizing the following expression:

$$\sum_{i=1}^n (y_i - \mathbf{x}_i \mathbf{w})^2 + \lambda \sum_{j=1}^p w_j^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2 \quad (4)$$

where λ denotes the regularization tuning parameter for controlling the strength of the Ridge penalty term. We set this penalty parameter to 1.

Considering these points, we selected $n(\mathbf{W})$, $\#OPs$, and $\#memOPs$ as the most meaningful predictors for building both runtime and energy per-layer regression models. We chose these variables because they presented the highest correlation in terms of Pearson correlation coefficient when analyzing architectural parameters vs. runtime/energy in our baseline system, i.e., *RPI-Caffe*. Indeed, this inherent linear relation supports the decision of applying linear regression models rather than more complex approaches such as Support Vector Machines, neural networks, or Gaussian process regression.

5 PreVIousNet

As mentioned in previous sections, PreVIousNet is a full-custom neural network specifically conceived for modeling the performance of a variety of layers on a selected system. Therefore, it is not applicable for vision inference. Below, we first summarize common CNN architectural parameters and layer settings contemplated in PreVIousNet. Next, we describe the designed architecture. Finally, the specific configuration employed for the performance modeling stage of PreVIous is reported.

5.1 Layer Parameters

Concerning convolutions, PreVIousNet includes both typical settings and special cases of CONV layers implemented in embedded CNNs:

- *Standard CONV.* Adjustable settings include:
 - Kernel size (k_h, k_w): conventionally, an odd value is set for both dimensions. Besides, state-of-the-art CNNs feature small kernel sizes to reduce the computational load.
 - Number of kernels N : to expand the channel dimension, $N > C_{in}$ kernels are normally applied.
 - Stride s : in case of strided convolutions, the most common value is $s = 2$.
- *Depthwise CONV.* In this type of layer, computation is saved by applying one kernel filter to each input channel.
- *Pointwise CONV.* This non-spatial convolution uses 1×1 kernels. Two variants are possible:

- *Bottleneck*. It reduces the computational load of subsequent layers by shrinking the channel dimension, i.e., $C_{out} < C_{in}$.
- *General*. It increases the channel dimension without performing spatial operations. This type of layer is applied either to revert the bottleneck channel-shrinking effect [72] or to build separable convolutions [10].

Concerning other types of layers, only *Pooling* layers have adjustable (k_h, k_w, s) parameters. In these layers, the maximum operation is commonly performed over 2×2 patches. In addition, some networks employ so-called global average pooling to replace memory-intensive FC layers. In this particular case, an average operation is performed on the entire input feature map.

5.2 Architecture

For each layer in a CNN, both *input data dimensions* and *layer parameters* determine the computational load and memory requirements, thus affecting the execution performance. This is the fact that inspired the design of the main architecture of PreVIOU, denoted as PreVIOU-01 in Fig. 4(a). We aimed at covering a wide range of possibilities within the architecture design space. In this regard, note that:

1. Data dimensions and computational load progressively increase as the network goes deeper – i.e., moving rightwards in Fig. 4(a) through the *levels* of the network. The network input dimensions at the first level – H , W , and C – are adjustable variables of PreVIOU.
2. Various layers and parameters are contemplated in parallel branches inserted at each level – vertically displayed in Fig. 4(a).

Based on these characteristics, PreVIOU-01 was designed as follows. At each level, a cluster of parallel *CONV* layers encompasses the aforementioned strategies: standard convolutions, pointwise, depthwise, and bottlenecks. Concerning activation layers (*BN*, *Scale*, *ReLU*), they perform their operation on diversely shaped intermediate *fmaps* of the network – i.e., at different levels and network branches. Likewise, *Pooling* layers with varied configurations are introduced at different levels. Finally, *Eltwise* and *Concat* layers operate on equally sized pairs of tensors coming from previous branches of the network. Overall, PreVIOU-01 comprises 52 layers: 15 *CONV*, 7 *BN*, 7 *Scale*, 7 *ReLU*, 6 *Pooling*, 5 *Eltwise*, and 5 *Concat*.

Note that *FC* and *Softmax* layers are not included in PreVIOU-01. These layers deal with a special case of data structure: 1D vectors instead of 3D tensors. An additional observation is that FC layers consume a notable amount

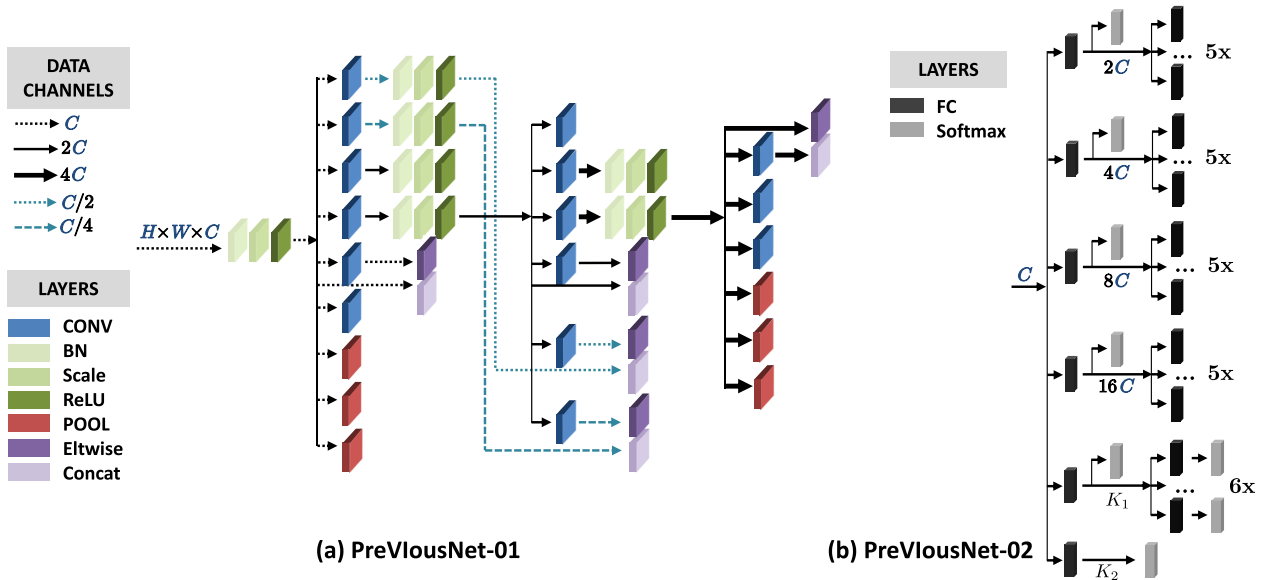


Figure 4: Macroarchitecture of PreVIOU [73] used in the modeling stage of PreVIOU. Seven types of layers with different configurations are contemplated by PreVIOU-01 (a), whereas *FC* and *Softmax* layers are covered by PreVIOU-02 (b). The network input dimensions at the first level – H , W , and C – are adjustable variables of PreVIOU.

of memory. Therefore, these layers are characterized through a different architecture, i.e., PreVIousNet-02, depicted in Fig. 4(b). Thus, PreVIousNet comprises two compact specialized networks. In PreVIousNet-02, various clusters of parallel FC layers deal with input and output vectors featuring diverse sizes. The network input is a $1 \times 1 \times C$ vector, where C can be adjusted. Then, various sizes of data are processed by the network layers: either resulting from applying an expansion factor to the input ($2C$, $4C$, etc.), or using customized vector lengths K_i . For instance, $K_1 = 10$ and $K_2 = 1000$ can be used as they are common output sizes in classification networks trained on ImageNet [74], CIFAR [75], and MNIST [76]. Similarly, *Softmax* layers operate on vectors with varied dimensions. As a whole, PreVIousNet-02 includes 44 layers: 32 *FC* and 12 *Softmax*.

PreVIousNet is publicly available in Caffe format [73]. Further details about the network can be observed from this model definition. For instance, although no weight file is provided, loading the network in Caffe will automatically initialize the weights according to the ‘‘MSRA’’ initialization scheme [77]. Remarkably, the proposed architecture is not unique. It can be adjusted according to the specificity of the networks to be characterized during the performance prediction stage of PreVIous.

5.3 Network Configuration

The input size of PreVIousNet-01 ($H \times W \times C$) can be properly set according to the most common tensor sizes handled by CNNs. Let us consider SqueezeNet [7] as an example of embedded CNN. In this network, the number of input channels ranges from 3 to 512, whereas height and width of *fmaps* decrease following the sequence 227, 113, 56, 28, 14. According to this example, a characterization of PreVIousNet-01 with varied input tensor sizes is required to collect as much information as possible to build accurate prediction models. In our experiments, we empirically set the following four input dimensions: (1) $56 \times 56 \times 32$, (2) $28 \times 28 \times 64$, (3) $14 \times 14 \times 64$, and (4) $7 \times 7 \times 64$. Thus, we are particularly sampling³ CONV layers with $H_{in} = W_{in} = \{56, 28, 14, 7\}$ and $C_{in} = \{32, 64, 128, 256\}$. Out of these ranges, the predicted performance values must be extrapolated from the corresponding regression models. However, this extrapolation is precise, as will be shown in Section 6.

Concerning PreVIousNet-02, we specifically run this network with an input vector sized $1 \times 1 \times 256$. Consequently, the following common vector lengths are considered $C_{in} = \{256, 512, 1024, 2048, 4096\}$, plus customized values $\{K_1, K_2\} = \{10, 1000\}$. Other input configurations could be used, according to architectures of interest or device limitations.

6 Experimental Results

As a first step, we completed the *performance modeling* stage of PreVIous to create the models for both runtime and energy consumption on the IoT devices described in Section 4.2. In particular, prediction models were built upon the performance profiling of PreVIousNet-01 and PreVIousNet-02 under the 5 configurations specified in Section 5.3 — four for PreVIousNet-01 and one for PreVIousNet-02. Matrix \mathbf{X} in Eq. (3) was obtained for each type of layer from the architectural parameters of all the corresponding layers in the aforementioned configurations of PreVIousNet. Likewise, response vector \mathbf{y} in Eq. (3) was built for both runtime and energy consumption from the corresponding profilings described in Section 4.3.

Then, we conducted the *performance prediction* stage of PreVIous on seven popular CNNs, most of them suitable for embedded devices: AlexNet [78], All-CNN-C [71], MobileNet [10], ResNet-18 [72], SimpleNet [79], SqueezeNet [7], and Tiny YOLO [80]. These networks were trained on ImageNet dataset [74] for 1000-category classification, except for All-CNN-C and SimpleNet, which perform classification on CIFAR-10 and CIFAR-100 [75], respectively, and Tiny YOLO, trained on COCO dataset [81] for object detection. As a whole, 399 CNN layers were assessed in this extensive study.

6.1 Layerwise Predictions

To evaluate the precision of the per-layer prediction models resulting from PreVIous, we compared layerwise predictions with actual profiling measurements of the corresponding layers in all the considered CNNs. As an example, Fig. 5 illustrates the high accuracy of the runtime predictions from PreVIous when compared to the empirical measurements in our baseline system, i.e., *RPi-Caffe*.

³Note that in PreVIousNet-01, the number of channels increases over network levels, whereas the input *fmap* resolution remains constant. The motivation is to build a simplified network to be evaluated under different configurations.

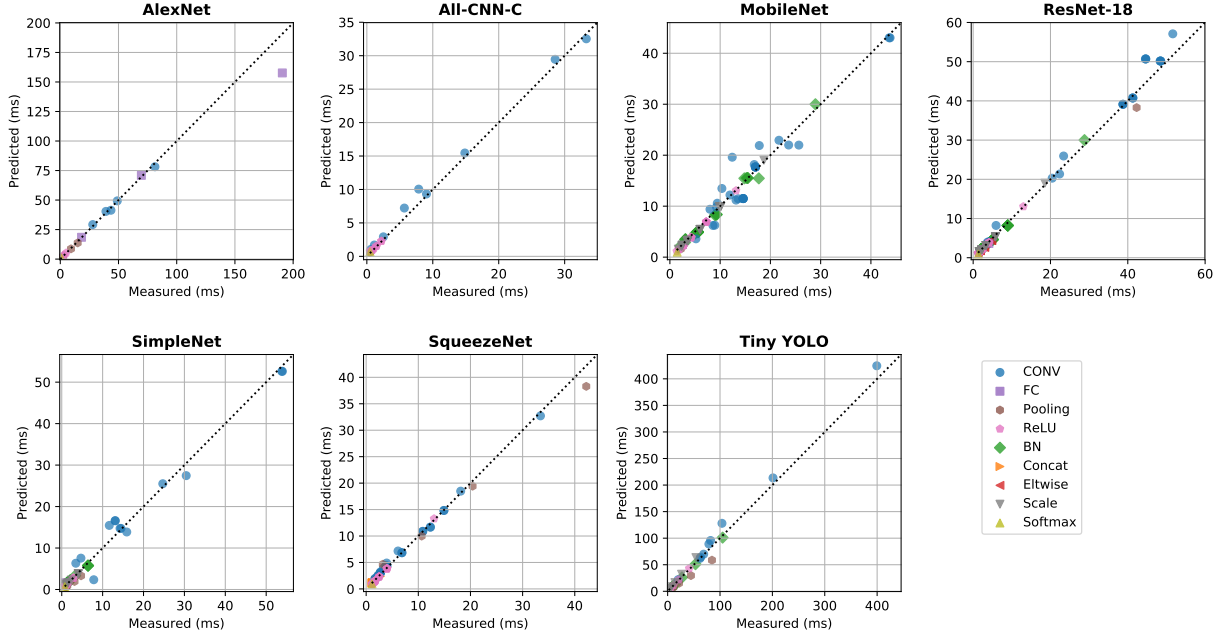


Figure 5: Comparison between *per-layer inference runtime* predictions (y-axis) and actual measurements (x-axis) for the baseline combination *RPi-Caffe*. Similar results are obtained for *RPi-OpenCV* and *XU4-Caffe*. The dashed line depicts an ideal estimation in which predictions exactly match actual measurements.

A summary of the results in Fig. 5 is presented in Table 1. Each row reports the total time obtained by adding up the runtime of all the layers composing the corresponding network, that is:

$$t = \sum_{l=1}^{N_l} t_l \quad \hat{t} = \sum_{l=1}^{N_l} \hat{t}_l \quad (5)$$

where t_l denotes per-layer measurements, \hat{t}_l denotes per-layer predictions, N_l is the number of layers in the CNN, t is the total measured time, and \hat{t} is the total predicted time.

As an example of fine-grained performance assessment of a neural network, Fig. 6 shows the time actually required to complete each layer of All-CNN-C [71] in *RPi-Caffe* compared with model estimations. Note that the prediction model from PreVIOUS correctly identifies ‘conv2’ and ‘conv5’ as the layers demanding the majority of the inference time. This identification is extremely useful to boost automatic optimization algorithms or NAS engines.

⁴The well-known AlexNet architecture also includes two *Local Response Normalization* (LRN) layers. These layers have not been contemplated in our layerwise prediction because they have been superseded by BN.

Table 1: Per-layer runtime predictions on RPi-Caffe system. Detailed profiling is shown in Fig. 5.

	Measured (ms) t	Predicted (ms) \hat{t}	Error (%)
AlexNet ⁴	561.64	526.75	-6.21%
All-CNN-C	115.68	123.22	6.52%
MobileNet	943.73	908.74	-3.71%
ResNet-18	1032.84	1049.10	1.57%
SimpleNet	347.59	349.22	0.47%
SqueezeNet	348.15	343.54	-1.32%
Tiny YOLO	1691.37	1740.03	2.88%
	<i>Average (absolute values)</i>		3.24%

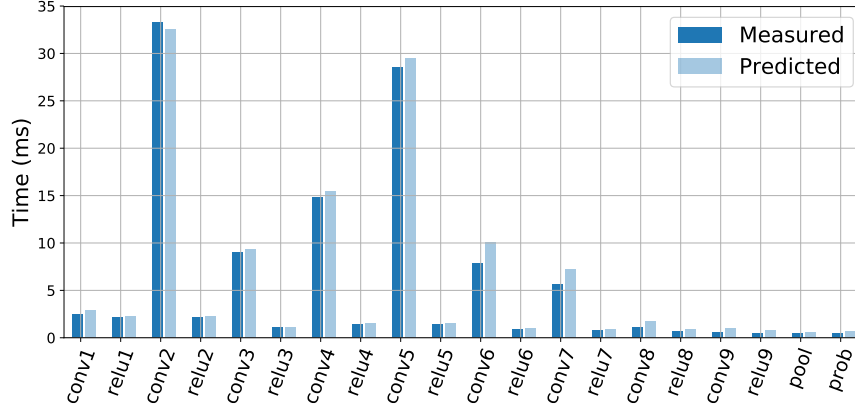


Figure 6: Layerwise runtime on All-CNN-C [71], trained for classification on CIFAR10 dataset [75]. Network profiling measurements were taken on the baseline combination *RPi-Caffe*

Finally, Table 2 reports the average absolute error of per-layer performance prediction from PreVIous for the seven considered CNNs on the three studied systems. Remarkably, high prediction accuracy is maintained in all the cases for both runtime and energy, thereby validating PreVIous as a general framework for CNN performance modeling.

Table 2: Average per-layer absolute prediction error of PreVIous for the seven considered CNNs on the three studied systems.

	Runtime	Energy
RPi3 - Caffe	3.24%	5.30%
RPi3 - OpenCV	4.08%	3.63%
XU4 - Caffe	3.82%	5.01%

6.2 Network Predictions

Note that Section 6.1 is focused on aggregated per-layer measurements. Indeed, a usual procedure followed in previous works on network optimization or NAS [35, 37, 40, 43] consists in estimating the global forward-pass performance of a network by adding up per-layer metrics, as expressed in Eq. (5) for runtime. In principle, this approach should be valid given that layers are sequentially executed during CNN inference in many realizations. However, in practice, when per-layer measurements have been independently taken, their direct addition may not coincide with the actual network inference performance [48, 49, 53]. This mismatch arises from aspects such as software optimizations (e.g., layer fusion or constant folding) and processor strategies (e.g., data prefetching or data re-utilization in the memory hierarchy). To take this fact into account, we measured the performance of the complete forward-pass of PreVIousNet in terms of runtime and energy. This forward-pass characterization allows us to write the following expression:

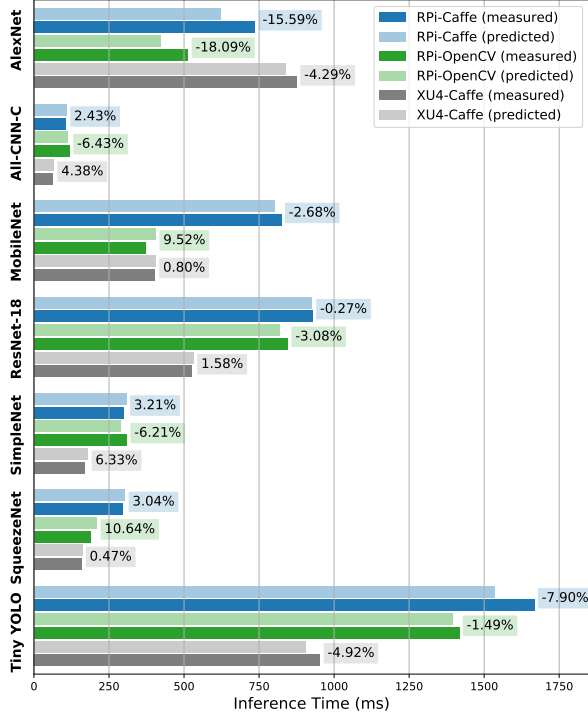
$$\hat{y} = c \sum_{l=1}^{N_l} \hat{y}_l \quad (6)$$

where \hat{y} represents the total predicted runtime or energy, \hat{y}_l denotes the per-layer predictions either for runtime – \hat{t}_l in Eq. (5) – or energy, and c is a coefficient resulting from linear regression between the direct addition of predictions and the corresponding actual measurement for the complete forward-pass of the 5 configurations of PreVIousNet on each software-hardware combination. The values of c for each case are reported in Table 3. Eq. (6) enables the comparison of predictions from PreVIous for *complete network inference* against the corresponding experimental measurements. This comparison is depicted in Figs. 7a and 7b for runtime and energy, respectively⁵. PreVIous also provides good estimates of complete forward-pass inference, with deviations below 10% in 18 out of 21 studied cases for runtime and 15 out of 21 cases for energy.

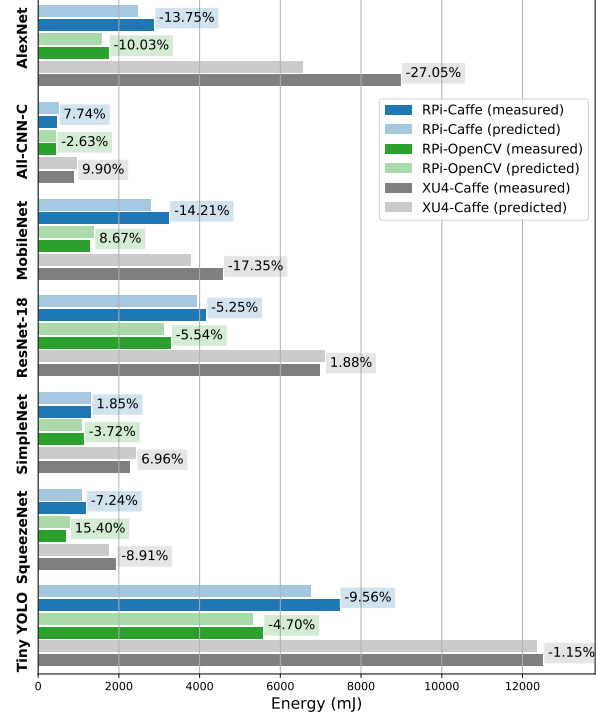
⁵For AlexNet, profiling measurements of the deprecated LRN layers were added in the summation in Eq. (6) for the sake of fair comparison.

Table 3: Value of the empirical coefficient c in Eq. (6) for the three studied systems.

	Runtime	Energy
RPi3 - Caffe	0.88	1.08
RPi3 - OpenCV	0.85	0.89
XU4 - Caffe	0.93	1.09



(a) Network runtime predictions vs. actual measurements for all the assessed combinations of CNN model, hardware platform, and software framework.



(b) Network energy predictions vs. actual measurements for all the assessed combinations of CNN model, hardware platform, and software framework.

Additionally, we further analyzed the less accurate cases such as the predictions on AlexNet – specially for energy modeling on Odroid XU4. Regarding this network, both runtime and energy consumption were underestimated by PreVIous in all cases, as shown by the negative errors reported in Figs. 7a and 7b. The reason of this behavior lies in the divergence between per-layer and complete-CNN inference performance. This divergence is due to the fact that, as reflected in Table 3, we apply a single value for the c term in Eq. (6) per hardware-software combination for the sake of simplicity and generalization. For instance, in the worst prediction case – AlexNet⁶ energy estimation on *XU4-Caffe* –, the ratio between per-layer and complete-CNN forward-pass measurements is 1.60, which is notably higher than the value of c resulting from linear regression for the combination *XU4-Caffe*, i.e., 1.09. The same situation was identified in *RPi-Caffe* and *XU4-Caffe* for MobileNet energy consumption. By contrast, per-layer predictions of PreVIous were certainly accurate, with absolute network errors below 5% in the vast majority of the 42 studied cases – see summaries in Table 2.

6.3 Discussion

Some key points must be stressed about the results presented in this section:

⁶Note that the well-known AlexNet model can be considered the starting point in the evolution leading to current CNN models. Therefore, its architecture somehow differs from modern ones, featuring deprecated layers and an elevated amount of learnt parameters.

1. We have intentionally assessed a diversity of network layers, as opposed to previous approaches, mostly focused on CONV layers. Indeed, Fig. 5 highlights the non-negligible – even dominant in some cases – contribution of certain layers, e.g., Pooling in SqueezeNet or BN in MobileNet, to the total inference time. This proves the importance of their consideration for performance modeling [13, 48, 82].
2. The proposed methodology has been validated on various systems suitable for IoT applications. This verification includes the process of model construction and the prediction capacity of PreVIous. The remarkable aspect here is that only five systematic network characterizations – i.e., the considered configurations of PreVIousNet – suffice to build accurate prediction models for a particular system.
3. Our study is, to the best of our knowledge, the most comprehensive in the literature in joint terms of number of CNNs, types of layers, performance metrics, and hardware-software combinations. The prediction accuracy is also, in global terms, the highest among similar reported works. Table 4 presents a comparison of our study vs. such similar works⁷. Note that most of them made use of high-end GPUs; in our case, we focused on low-cost low-power small-sized IoT devices. The last column summarizes the prediction accuracy for each case in terms of Mean Absolute Percentage Error (MAPE), which is defined over the considered CNNs as the average of the absolute value of the difference between the complete network prediction and the corresponding actual measurement divided by the measurement. Thus, the values of this column in our case are the average of the absolute values in Figs. 7a and 7b for each selected system. Concerning related works, we calculated the MAPE according to the individual errors reported for each characterized CNN. Note that we have covered a much wider spectrum of layers than the other studies, i.e., 9 types of layers vs. 3 types at most. This is the basis for achieving better predictions over a larger set of CNNs following a common procedure for both runtime and energy.

7 Conclusions

This study demonstrates that it is possible to predict the performance of CNNs on embedded vision devices with high accuracy through a simple procedure. Taking into account the growing and ever-changing zoo of CNN models, such

⁷Among the related studies described in Sec. 2, this table contains those which addressed global CNN performance modeling — as opposed to single-layer characterization — and reported numerical prediction results.

Table 4: Comparison of CNN modeling accuracy between this study and related works in the literature. MAPE stands for Mean Absolute Percentage Error associated with complete network inference.

Network Runtime			
Ref.	System	CNNs	MAPE
[48] [†]	TK1 CPU – Caffe	NIN, VGG19M	4.71%
	TK1 GPU – Caffe		23.70%
	TX1 CPU – Caffe		39.91%
	TX1 GPU – Caffe	NIN, VGG19M, SqueezeNet, MobileNet	31.51%
[49] [‡]	Titan X GPU – TensorFlow	VGG16, AlexNet, NIN, Overfeat, CIFAR10-6conv	7.96%
	GTX1070 GPU – TensorFlow	AlexNet, NIN	12.32%
	GTX1070 GPU – Caffe		16.17%
This study [⌘]	RPi3 CPU – Caffe		5.02%
	RPi3 CPU – OpenCV	AlexNet, All-CNN-C, MobileNet, ResNet-18, SimpleNet, SqueezeNet, Tiny YOLO	7.92%
	XU4 CPU – Caffe		3.25%
Network Energy			
Ref.	System	CNNs	MAPE
[48] [†]	TX1 CPU – Caffe		39.08%
	TX1 GPU – Caffe	NIN, VGG19M, SqueezeNet, MobileNet	15.30%
[49] [‡]	Titan X GPU – TensorFlow	VGG16, AlexNet, NIN, Overfeat, CIFAR10-6conv	2.25%
	GTX1070 GPU – TensorFlow	AlexNet, NIN	8.40%
	GTX1070 GPU – Caffe		21.99%
[50] [§]	TX1 CPU – Caffe	AlexNet, ResNet-50, SqueezeNet, GoogLeNet, SqueezeNetRes, VGG-small, Places-CNDS-8s, All-CNN-C, Inception-BN, MobileNet	12.26%
This study [⌘]	RPi3 CPU – Caffe		8.52%
	RPi3 CPU – OpenCV	AlexNet, All-CNN-C, MobileNet, ResNet-18, SimpleNet, SqueezeNet, Tiny YOLO	7.24%
	XU4 CPU – Caffe		10.46%

Types of assessed layers: [†] CONV; [‡] CONV, FC, and Pooling; [§] CONV; [⌘] CONV, FC, Pooling, ReLU, BN, Concat, Eltwise, Scale, and Softmax.

a priori prediction is key for rapid exploration and optimal implementation of visual inference. The utility of the proposed methodology, i.e., PreVIous, is two-fold. First, fine-grained layer performance prediction facilitates network architecture design and optimization. Second, network performance estimation can assist in CNN selection to fulfill prescribed IoT requirements such as latency and battery lifetime.

Simplicity is indeed a major asset of PreVIous. Only the characterization of a single architecture is required for performance modeling. We also make use of linear regression to reduce model complexity. In addition, the procedure does not rely on any specific measurement tool, being agnostic with respect to the selected hardware-software combination.

Future work will address the design of further versions of PreVIousNet in order to consider new types of layers or even entire building blocks. For instance, recurrent building blocks of highly optimized architectures can be characterized as a whole, e.g., the *Fire* module of SqueezeNet or *separable convolutions* of MobileNets. This approach can also be exploited by automatic algorithms to explore new architectures optimally adapted to specific embedded systems.

8 Acknowledgments

This work was supported by Spanish Government MICINN (European Region Development Fund, ERDF/FEDER) through project RTI2018-097088-B-C31, European Union H2020 MSCA through project ACHIEVE-ITN (Grant No. 765866), and by the US Office of Naval Research through Grant No. N00014-19-1-2156.

References

- [1] ITU's Telecommunication Standardization Sector, "Overview of the Internet of Things," International Telecommunication Union, Tech. Rep., 2012. [Online]. Available: <https://www.itu.int/ITU-T/recommendations/rec.aspx?rec=11559>
- [2] J. Lin, W. Yu, N. Zhang, X. Yang, H. Zhang, and W. Zhao, "A survey on Internet of Things: Architecture, enabling technologies, security and privacy, and applications," *IEEE Internet of Things Journal*, vol. 4, no. 5, pp. 1125–1142, 2017.
- [3] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [4] M. Verhelst and B. Moons, "Embedded deep neural network processing," *IEEE Solid-State Circuits Magazine*, vol. 9, no. 4, pp. 55–65, 2017.
- [5] V. Sze, Y. Chen, T. Yang, and J. Emer, "Efficient processing of deep neural networks: A tutorial and survey," *Proceedings of the IEEE*, vol. 105, no. 12, pp. 2295–2329, 2017.
- [6] D. Velasco-Montero, J. Fernández-Berni, R. Carmona-Galán, and A. Rodríguez-Vázquez, "Optimum selection of DNN model and framework for edge inference," *IEEE Access*, vol. 6, pp. 51 680–51 692, 2018.
- [7] 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 <1MB model size," *arXiv*, no. 1602.07360, 2016.
- [8] A. Gholami, K. Kwon, B. Wu, Z. Tai, X. Yue, P. H. Jin, S. Zhao, and K. Keutzer, "SqueezeNext: Hardware-aware neural network design," *arXiv*, no. 1803.10615, 2018.
- [9] B. Wu, F. N. Iandola, P. H. Jin, and K. Keutzer, "SqueezeDet: Unified, small, low power fully convolutional neural networks for real-time object detection for autonomous driving," *arXiv*, no. 1612.01051, 2016.
- [10] 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*, no. 1704.04861, 2017.
- [11] M. Sandler, A. G. Howard, M. Zhu, A. Zhmoginov, and L. Chen, "Inverted residuals and linear bottlenecks: Mobile networks for classification, detection and segmentation," *arXiv*, no. 1801.04381, 2018.
- [12] X. Zhang, X. Zhou, M. Lin, and J. Sun, "ShuffleNet: An extremely efficient convolutional neural network for mobile devices," *arXiv*, no. 1707.01083, 2017.
- [13] N. Ma, X. Zhang, H. Zheng, and J. Sun, "ShuffleNet V2: practical guidelines for efficient CNN architecture design," *arXiv*, no. 1807.11164, 2018.
- [14] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "XNOR-Net: ImageNet classification using binary convolutional neural networks," in *European Conference on Computer Vision*, Oct 2016, pp. 525–542.

- [15] D. Li, X. Chen, M. Becchi, and Z. Zong, "Evaluating the energy efficiency of deep convolutional neural networks on CPUs and GPUs," in *2016 IEEE Int. Conferences on Big Data and Cloud Computing (BDCloud), Social Computing and Networking (SocialCom), Sustainable Computing and Communications (SustainCom) (BDCloud-SocialCom-SustainCom)*, Oct 2016, pp. 477–484.
- [16] A. Canziani, A. Paszke, and E. Culurciello, "An analysis of deep neural network models for practical applications," *arXiv*, vol. 1605.07678, 2016.
- [17] A. Ignatov, R. Timofte, W. Chou, K. Wang, M. Wu, T. Hartley, and L. V. Gool, "AI benchmark: Running deep neural networks on android smartphones," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [18] S. Bianco, R. Cadène, L. Celona, and P. Napolitano, "Benchmark analysis of representative deep neural network architectures," *IEEE Access*, vol. 6, pp. 64 270–64 277, 10 2018.
- [19] Y. Wang, G. Wei, and D. Brooks, "Benchmarking TPU, GPU, and CPU platforms for deep learning," *arXiv*, no. abs/1907.10701, 2019.
- [20] Z. Cai, X. He, J. Sun, and N. Vasconcelos, "Deep learning with low precision by half-wave gaussian quantization," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5406–5414, 2017.
- [21] S. Zhou, Y. Wu, Z. Ni, X. Zhou, H. Wen, and Y. Zou, "DoReFa-Net: Training low bitwidth convolutional neural networks with low bitwidth gradients," *arXiv*, no. 1606.06160, 2016.
- [22] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1," *arXiv*, no. 1602.02830, 2016.
- [23] A. Zhou, A. Yao, Y. Guo, L. Xu, and Y. Chen, "Incremental network quantization: Towards lossless cnns with low-precision weights," in *International Conference on Learning Representations (ICLR 2017)*, 2017.
- [24] P. Molchanov, S. Tyree, T. Karras, T. Aila, and J. Kautz, "Pruning convolutional neural networks for resource efficient transfer learning," in *International Conference on Learning Representations*, Nov 2017.
- [25] H. Hu, R. Peng, Y. Tai, and C. Tang, "Network trimming: A data-driven neuron pruning approach towards efficient deep architectures," *arXiv*, no. 1607.03250, 2016.
- [26] J. Zou, T. Rui, Y. Zhou, C. Yang, and S. Zhang, "Convolutional neural network simplification via feature map pruning," *Computers & Electrical Engineering*, vol. 70, 2018.
- [27] Y. He, X. Zhang, and J. Sun, "Channel pruning for accelerating very deep neural networks," in *2017 IEEE International Conference on Computer Vision (ICCV)*, Oct 2017, pp. 1398–1406.
- [28] A. Lavin and S. Gray, "Fast algorithms for convolutional neural networks," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 4013–4021.
- [29] M. Mathieu, M. Henaff, and Y. LeCun, "Fast training of convolutional networks through FFTs," in *International Conference on Learning Representations (ICLR2014)*, 2014.
- [30] J. Cong and B. Xiao, "Minimizing computation in convolutional neural networks," in *Artificial Neural Networks and Machine Learning - ICANN 2014*. Springer International Publishing, 2014, pp. 281–290.
- [31] I. Freeman, L. Roese-Koerner, and A. Kummert, "EffNet: An efficient structure for convolutional neural networks," in *2018 25th IEEE International Conference on Image Processing (ICIP)*, 2018, pp. 6–10.
- [32] Y.-D. Kim, E. Park, S. Yoo, T. Choi, L. Yang, and D. Shin, "Compression of deep convolutional neural networks for fast and low power mobile applications," *arXiv*, no. 1511.06530, 2015.
- [33] V. Lebedev, Y. Ganin, M. Rakhuba, I. V. Oseledets, and V. S. Lempitsky, "Speeding-up convolutional neural networks using fine-tuned cp-decomposition," in *International Conference on Learning Representations (ICLR 2015)*, 2015.
- [34] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8697–8710.
- [35] T.-J. Yang, A. Howard, B. Chen, X. Zhang, A. Go, M. Sandler, V. Sze, and H. Adam, "NetAdapt: Platform-aware neural network adaptation for mobile applications," in *European Conference on Computer Vision (ECCV)*, September 2018.
- [36] Y. He, J. Lin, Z. Liu, H. Wang, L.-J. Li, and S. Han, "AMC: AutoML for model compression and acceleration on mobile devices," in *European Conference on Computer Vision (ECCV)*, Sept 2018, pp. 815–832.
- [37] S. Han, H. Cai, L. Zhu, J. Lin, K. Wang, Z. Liu, and Y. Lin, "Design automation for efficient deep learning computing," *arXiv*, no. 1904.10616, 2019.

- [38] A. Gordon, E. Eban, O. Nachum, B. Chen, H. Wu, T.-J. Yang, and E. Choi, “MorphNet: Fast & simple resource-constrained structure learning of deep networks,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [39] M. Tan, B. Chen, R. Pang, V. Vasudevan, M. Sandler, A. Howard, and Q. V. Le, “MnasNet: Platform-aware neural architecture search for mobile,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [40] B. Wu, X. Dai, P. Zhang, Y. Wang, F. Sun, Y. Wu, Y. Tian, P. Vajda, Y. Jia, and K. Keutzer, “FBNet: Hardware-aware efficient ConvNet design via differentiable neural architecture search,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [41] D. Stamoulis, R. Ding, D. Wang, D. Lymberopoulos, B. Priyantha, J. Liu, and D. Marculescu, “Single-Path NAS: Designing hardware-efficient convnets in less than 4 hours,” *ArXiv*, no. 1904.02877, 2019.
- [42] H. Cai, L. Zhu, and S. Han, “ProxylessNAS: Direct neural architecture search on target task and hardware,” in *International Conference on Learning Representations*, 2019.
- [43] X. Dai, P. Zhang, B. Wu, H. Yin, F. Sun, Y. Wang, M. Dukhan, Y. Hu, Y. Wu, Y. Jia, P. Vajda, M. Uyttendaele, and N. K. Jha, “ChamNet: Towards efficient network design through platform-aware model adaptation,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [44] D. Stamoulis, E. Cai, D.-C. Juan, and D. Marculescu, “Hyperpower: Power- and memory-constrained hyperparameter optimization for neural networks,” in *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2017, pp. 19–24.
- [45] Y. Chen, T. Krishna, J. S. Emer, and V. Sze, “Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks,” *IEEE Journal of Solid-State Circuits*, vol. 52, no. 1, pp. 127–138, Jan 2017.
- [46] T.-J. Yang, Y.-H. Chen, and V. Sze, “Designing energy-efficient convolutional neural networks using energy-aware pruning,” *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6071–6079, 2016.
- [47] Y. Wu and V. Sze, “Accelerigy: An architecture-level energy estimation methodology for accelerator designs,” in *Int. Conference On Computer Aided Design*, 2019.
- [48] Z. Lu, S. Rallapalli, K. S. Chan, S. Pu, and T. La Porta, “Augur: Modeling the resource requirements of ConvNets on mobile devices,” *IEEE Transactions on Mobile Computing*, 2019.
- [49] E. Cai, D. Juan, D. Stamoulis, and D. Marculescu, “NeuralPower: Predict and deploy energy-efficient convolutional neural networks,” *arXiv*, no. 1710.05420, 2017.
- [50] C. F. Rodrigues, G. D. Riley, and M. Luján, “Fine-grained energy profiling for deep convolutional neural networks on the Jetson TX1,” *arXiv*, no. 1803.11151, 2018.
- [51] E. Li, Z. Zhou, and X. Chen, “Edge intelligence: On-demand deep learning model co-inference with device-edge synergy,” in *Proceedings of the 2018 Workshop on Mobile Edge Communications (MECOMM’18)*, 2018.
- [52] Y. Kang, J. Hauswald, C. Gao, A. Rovinski, T. Mudge, J. Mars, and L. Tang, “Neurosurgeon: Collaborative intelligence between the cloud and mobile edge,” in *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems*, 2017, pp. 615–629.
- [53] A. E. Eshratifar and M. Pedram, “Energy and performance efficient computation offloading for deep neural networks in a mobile cloud computing environment,” in *Proceedings of the 2018 on Great Lakes Symposium on VLSI*, 2018, p. 111–116.
- [54] E. Gianniti and L. Zhang, “Performance prediction of gpu-based deep learning applications,” in *International Symposium on Computer Architecture and High Performance Computing SBAC-PAD*, Sept 2018.
- [55] H. Qi, E. R. Sparks, and A. Talwalkar, “Paleo: A performance model for deep neural networks,” in *International Conference on Learning Representations (ICLR)*, 2017.
- [56] D. Mishkin, N. Sergievskiy, and J. Matas, “Systematic evaluation of convolution neural network advances on the imagenet,” *Computer Vision and Image Understanding*, vol. 161, pp. 11 – 19, 2017.
- [57] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *Proceedings of the 32nd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 37, 2015, pp. 448–456.
- [58] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015, pp. 1–9.

- [59] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 2818–2826.
- [60] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi, "Inception-v4, Inception-ResNet and the impact of residual connections on learning," *AAAI Conference on Artificial Intelligence*, 2016.
- [61] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet large scale visual recognition challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
- [62] M. Everingham, L. V. Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The Pascal visual object classes (VOC) challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [63] "Automatically tuned linear algebra software (ATLAS)," <http://math-atlas.sourceforge.net>.
- [64] K. Goto and R. A. van de Geijn, "Anatomy of high-performance matrix multiplication," *ACM Trans. Math. Softw.*, vol. 34, no. 3, pp. 12:1–12:25, May 2008.
- [65] "Intel math kernel library," <https://software.intel.com/en-us/mkl>.
- [66] "Dense linear algebra on GPUs," <https://developer.nvidia.com/cublas>.
- [67] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proceedings of 22nd ACM International Conference on Multimedia*, 2014, pp. 675–678.
- [68] "Raspberry Pi 3 Model B," <https://www.raspberrypi.org/products/raspberry-pi-3-model-b/>.
- [69] "OpenCV," <https://opencv.org/>.
- [70] "Odroid XU4," <https://wiki.odroid.com/odroid-xu4/odroid-xu4>.
- [71] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for simplicity: The All Convolutional Net," *arXiv*, no. 1412.6806, 2014.
- [72] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [73] D. Velasco-Montero, J. Fernández-Berni, R. Carmona-Galán, and A. Rodríguez-Vázquez, "PreVIous. Github Repository." [Online]. Available: <https://github.com/DVM000/PreVIous.git>
- [74] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *IEEE Conference on Computer Vision and Pattern Recognition*, June 2009, pp. 248–255.
- [75] A. Krizhevsky, V. Nair, and G. Hinton, "The CIFAR dataset." [Online]. Available: <https://www.cs.toronto.edu/~kriz/cifar.html>
- [76] Y. LeCun, C. Cortes, and C. J. Burges, "The MNIST database of handwritten digits." [Online]. Available: <http://yann.lecun.com/exdb/mnist/>
- [77] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification," in *International Conference on Computer Vision (ICCV)*, 2015.
- [78] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems, NIPS*, 2012, pp. 1097–1105.
- [79] S. H. Hasanpour, M. Rouhani, M. Fayyaz, and M. Sabokrou, "Lets keep it simple, using simple architectures to outperform deeper and more complex architectures," *arXiv*, no. 1608.06037, 2016.
- [80] J. Redmon and A. Farhadi, "YOLO9000: better, faster, stronger," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [81] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common Objects in Context," in *European Conference on Computer Vision*, C. Springer, Ed., 2014, pp. 740–755.
- [82] D. Li, X. Wang, and D. Kong, "DeepRebirth: Accelerating deep neural network execution on mobile devices," in *Proc. of 32 AAAI Conference On Artificial Intelligence*, 2018, pp. 2322–2330.