



HAL
open science

Low-complexity Approximate Convolutional Neural Networks

Renato J. Cintra, Stefan Duffner, Christophe Garcia, André Leite

► **To cite this version:**

Renato J. Cintra, Stefan Duffner, Christophe Garcia, André Leite. Low-complexity Approximate Convolutional Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems*, 2018, 10.1109/TNNLS.2018.2815435 . hal-01727219

HAL Id: hal-01727219

<https://hal.science/hal-01727219>

Submitted on 27 Aug 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Low-complexity Approximate Convolutional Neural Networks

Renato J. Cintra, *Senior Member, IEEE*, Stefan Duffner, Christophe Garcia and André Leite

Abstract—In this paper, we present an approach for minimizing the computational complexity of trained Convolutional Neural Networks (ConvNet). The idea is to approximate all elements of a given ConvNet and replace the original convolutional filters and parameters (pooling and bias coefficients; and activation function) with efficient approximations capable of extreme reductions in computational complexity. Low-complexity convolution filters are obtained through a binary (zero-one) linear programming scheme based on the Frobenius norm over sets of dyadic rationals. The resulting matrices allow for multiplication-free computations requiring only addition and bit-shifting operations. Such low-complexity structures pave the way for low-power, efficient hardware designs. We applied our approach on three use cases of different complexity: (i) a “light” but efficient ConvNet for face detection (with around 1 000 parameters); (ii) another one for hand-written digit classification (with more than 180 000 parameters); and (iii) a significantly larger ConvNet: AlexNet with ≈ 1.2 million matrices. We evaluated the overall performance on the respective tasks for different levels of approximations. In all considered applications, very low-complexity approximations have been derived maintaining an almost equal classification performance.

Index Terms—Convolutional Neural Networks, Approximation, Optimization, Numerical Computation

I. INTRODUCTION

Since their introduction in the 1990s by LeCun *et al.* [1], convolutional neural networks (ConvNets) have proven to be very powerful in many challenging computer vision tasks, such as hand-written character recognition [1], [2], embedded text detection and recognition [3], [4], automatic facial analysis [5]–[7], traffic sign recognition [8], pedestrian detection [9], vision-based navigation [10], and house numbers recognition [11], just to cite a few. Although state-of-the-art results have been reached in many different fields, ConvNets have become very popular only recently with the impressive results obtained by Krizhevsky *et al.* [12] in the recognition task, followed by Simonyan and Zisserman who won the localization challenge at the Large-Scale Visual Recognition Challenge (ImageNet) 2014 [13].

The main property of ConvNets is their capability for automatic extraction of complex and application-suitable features

from raw data (e.g., pixels in computer vision). To do so, they integrate a pipeline of convolution and pooling layers generally followed by a multi-layer perceptron, jointly performing local feature extraction and classification (or regression) in a single architecture where all parameters are learnt using the classical error back-propagation algorithm [14]. Traditionally, ConvNets are used for image processing tasks and, to this end, they are often applied on small regions of a bigger image in a sliding-window framework, for instance to detect an object. Because of weight sharing, each layer essentially performs convolution or pooling operations using small kernels inside a “retina”, and when applied to a large image, the replication of the ConvNet operation over all positions in the image can be significantly optimized by performing each convolution over the full image at once, efficiently implementing a full image transformation pipeline.

However, during training and when applying the trained ConvNet, there is still a significant number of floating-point computations (multiplications and additions). In order to accelerate these operations, most current approaches and available software rely on parallel computing using GPUs [12], [15] facilitated to some extent by the inherently parallel architectures of ConvNets. The recent trend in “deep learning” to use more and more complex models with millions of parameters requires enormous amounts of computational resources, in particular for training the models but also for applying the learnt ConvNets. Reducing the computational complexity of these models is thus of great interest to the research community as well as to industry. Moreover, such reduction in complexity of ConvNets is necessary to implement them on devices with limited resources (such as mobile devices) in order to operate at acceptable speed, for example for real-time applications, and reduce the overall power consumption.

In this paper, we focus on drastically reducing the computational cost itself by proposing a post-training approximation scheme aiming at replacing all parameters of a ConvNet with low-complexity versions. That is, for the convolution filters, only additions and bit-shifting operations are performed—no multiplication is necessary. Additionally, the activation function is sought to be replaced with low-complexity alternatives. Formulated as an optimization problem, we adopt matrix approximation techniques based on the Frobenius norm error and dyadic rational numbers represented in terms of Canonical Signed Digit (CSD) encoding.

Indeed, a sound approximation theory is a necessary step to facilitated hardware development. This is illustrated in the case of image compression where the most efficient coding schemes are based on approximate matrices realized in dedi-

Manuscript received...

Renato J. Cintra and André Leite are with the Signal Processing Group, Departamento de Estatística, Universidade Federal de Pernambuco, Recife, PE, Brazil. R. J. Cintra is also with the Department of Electrical and Computer Engineering, University of Calgary, Calgary, AB, Canada and was with the LIRIS, Institut National des Sciences Appliquées (INSA), Lyon, France (e-mail: {rjdisc, leite}@de.ufpe.br).

Stefan Duffner and Christophe Garcia are with the Université de Lyon, CNRS, INSA-Lyon, LIRIS, UMR5205, F-69621, France. (e-mail: {stefan.duffner, christophe.garcia}@liris.cnrs.fr)

This work was partially supported by the CNPq, Brazil.

cated hardware [16]–[19]. We aim at introducing an approach for approximating CNNs in order to pave the way for future efficient dedicated hardware design.

This work is organized as follows. Section II provides a literature review on the efficient numerical implementation of neural networks and, in particular, ConvNets. Section III details our approach for approximating the elements of a given ConvNet aiming at designing low-complexity structures. In Section IV we present the results of our experimental evaluation of the proposed approach for two typical ConvNet architectures. We assess the approximate ConvNets relative to their exact counterparts in terms of several figures of merit. We conclude the paper in Section V.

II. RELATED WORK

Although GPU implementations [20] allow for fast training and application of ConvNets on sufficiently equipped platforms, their integration on embedded systems for real-time applications may be more difficult due to the limited amount of available resources on these devices, which usually requires a good trade-off between performance and code size. Several previous works have tackled this problem. In early works [21]–[23] weight parameters of neural networks have been represented as power-of-two integers. Thus all multiplications can be operated as simple bit shifts. Direct training of these networks was also possible by keeping a floating-point version of each weight parameter, or otherwise use a technique called “weight dithering” [24]. Simard and Graf [25] extended this idea by encoding all other parameters as powers of two except the weights, i.e. neuron activations, gradients, and learning rates. Later, Draghici [26] conducted a broader analysis on the computational power of neural networks with reduced precision weights. Recently, Machado *et al.* [27] proposed a specific approximation scheme for sparse representation learning using only values of powers of two, and they integrated this quantization into the learning process. Finally, in the work of Courbariaux *et al.* [28], all weights are encoded as binary values (1-bit), and for the training a floating-point version is still needed. Similarly, Kim and Paris [29], proposes a completely binary neural network. However, an initial real-valued training phase is required.

More recently, special attention has been paid to hardware implementations of ConvNets, especially on FPGAs. In [30], for example, a high-level optimization methodology is applied to the implementation of the CFF face detector [6]. They propose algorithmic optimizations and advanced memory management and transform the floating-point computation into fixed-point arithmetic. Interestingly such coarse approximation could furnish very similar detection rates and low false-alarm rates on referenced datasets, for a roughly sevenfold gain of speed. Later, this work has been extended in [31], where the authors present for the first time several implementations of the CFF algorithm on FPGA, with a parallel architecture composed of a Processing Element ring and a FIFO memory, which constitutes a generic architecture capable of processing images of different sizes. Farabet and LeCun [32] also propose a scalable hardware architecture to implement

large-scale ConvNets, with a modular vision engine for large image processing, with FPGA and ASIC implementations. Chakradhar *et al.* [33] present a dynamically configurable FPGA co-processor that adapts to complex ConvNet architectures exploiting different types of parallelism. A very low-complexity ASIC design of ConvNets has been developed by Chen *et al.* [34], allowing for very high execution speeds and power consumption of state-of-the-art ConvNets. Finally, Zhang *et al.* [35] propose a FPGA design strategy and algorithmic enhancements to optimize the computational throughput and memory bandwidth for any given ConvNet architecture.

Other recent works have focused on the *algorithmic* and *memory* optimizations of large-scale ConvNets. For example, Mamalet *et al.* [36] proposed different strategies for simplifying the convolutional filters (fusion of convolutional and pooling layers, 1D separable filters), in order to modify the hypothesis space, and to speed-up learning and processing times. These convolutions can also effectively be performed by simple multiplications of the filters with the respective input images in the frequency domain [37]. However, due to the overhead of the FFT, there is only a computational gain with larger filter sizes, and if a given filter can be reused consecutively for many input images. Vanhoucke *et al.* [38] presented a set of different techniques to accelerate the computation of ConvNets on CPU, mostly for Intel and AMD CPUs, exploiting for example SIMD instructions, memory locality, and fixed-point representations. Also, many recent works [39]–[46] have focused on reducing the complexity of convolution or fully-connected layers of large-scale ConvNets by replacing the high-dimensional matrix or tensor multiplications with several low-rank matrix multiplications using different low-rank factorization methods, either at test-time or both for training and testing. Although, large gains in computational and memory resources can be obtained on complex ConvNets, these optimizations do not focus on hardware implementation and low-power constraints.

As opposed to many previous works that integrate the approximation process into the learning [47]–[51], our approach operates on existing fully-trained models, that originally may have been aimed for standard PCs or more powerful architectures. Thus, our approximation scheme allows to integrate these models into hardware with much fewer resources.

III. APPROXIMATION APPROACH

A. General Goal

Our goal is to derive low-complexity structures capable of reducing the computational costs of a given ConvNet. Ideally, the following two conditions are simultaneously expected to be satisfied:

- (i) the computational elements of the ConvNet (convolutional filters, sub-sampling coefficients, bias values, and sigmoid function calls) are replaced by corresponding low-complexity structures;
- (ii) the performance of the ConvNet is not significantly degraded.

However, addressing both above conditions proves to be a hard task. In particular, the large number of variables, the

non-linearities, and extremely long simulation times prevents such approach. Also, to the best of our knowledge, literature furnishes no mathematical result linking the approximation of individual ConvNet elements and the final ConvNet performance. Thus, we adopt a greedy-like heuristic which consists in individually simplifying each computational structure of a ConvNet in the hope of finding a resulting structure capable of good performance [52].

In a ConvNet two main types of mathematical elements are found: (i) matrix structures and (ii) activation functions. The matrix structures are represented by convolution filter weights, sub-sampling operations, and bias values; whereas the activation function is usually a non-linear function, such as the threshold, piecewise-linear, and sigmoid functions [53].

To approximate these two classes of elements, different tools are required. For the matrix-based structures, we selected matrix approximation methods as a venue to derive low-complexity computational elements [54]–[57]. For the activation, we separate methods capable of approximating functions with efficient digital implementation [58]–[60].

B. Low-complexity Matrix Structures

In [54], [57], [61]–[64], several methods for deriving approximations of discrete transform matrices—such as the discrete cosine transform [65]—were proposed. Let \mathbf{M} be an $N \times N$ given matrix. For instance, \mathbf{M} can be a convolutional filter. In this case, a computational instantiation of \mathbf{M} applied to evaluate a single output pixel requires in principle N^2 floating point multiplications. A typical ConvNet may contain thousands of convolutional filters. For example, the classical architecture described in [12] contains 244,760 filters, which is nowadays considered a relatively small network. Therefore, to minimize such a significant computational cost, we aim at obtaining a low-complexity matrix $\hat{\mathbf{M}}$ capable of satisfying the following relation in an optimal sense: $\mathbf{M} \approx \hat{\mathbf{M}}$. The matrix $\hat{\mathbf{M}}$ is said to be an approximation for \mathbf{M} . Such approximate convolutional filters would allow the realization of computationally intensive ConvNets in limited resources architectures.

In this paper, a low-complexity matrix is a matrix of dyadic rational entries. Dyadic rational numbers are fractions of the form $m/2^n$, where n is a positive integer and m is an odd integer. Such numbers are suitable for binary arithmetic. Indeed, a multiplication by a dyadic rational consists of a multiplication by m followed by a right shift of n bits. Because m is an integer, we can take full advantage of fixed-point arithmetic. Indeed m can be given a binary representation with minimum number of adders, aiming at multiplicative irreducibility. Multiplicative irreducibility is attained whenever the minimum number of additions to implement a multiplication by m is equal to the number of ones in the binary representation of m [64]. Multiplicative irreducibility is often obtained when the CSD representation is considered [66]. Therefore, a multiplication by m can be converted into a sequence of additions and bit-shifting operations. As a consequence, low-complexity matrices are multiplierless, a very desirable property as floating-point operations are much more costly than additions and bit-shifting operations.

Standard methods for matrix approximation include: inspection [67], matrix parametrization [61], and matrix factorization [68]. However, since a typical ConvNet may contain from thousands to millions of filters, inspection-based approaches are not feasible. Methods based on the parametrization of the matrix elements are also ineffective because (i) the elements of convolutional filters are usually not clearly related, i.e., they do not satisfy identifiable mathematical relationships and (ii) the elements are not repeated. Additionally, ConvNet filters are not expected to satisfy properties, such as symmetry and orthogonality, which favors the derivation of approximations. Thus, methods based on matrix factorizations are less adequate.

C. Matrix Approximation by Linear Programming

We adopted a general approach to the problem of obtaining $\hat{\mathbf{M}}$ according to a optimization problem as described below:

$$\hat{\mathbf{M}} = \underset{\mathbf{T}}{\operatorname{arg\,min}} \operatorname{error}(\mathbf{M}, \mathbf{T}). \quad (1)$$

The above optimization problem can yield better approximate matrices if an expansion factor α is introduced [64]. By adopting the usual Frobenius norm [56] as an error measure, (1) can be recast according to the following mixed integer nonlinear programming (INLP) setup [69]:

$$(\alpha^*, \mathbf{T}^*) = \underset{\alpha, \mathbf{T}}{\operatorname{arg\,min}} \|\mathbf{M} - \alpha \cdot \mathbf{T}\|^2, \quad (2)$$

where $\alpha > 0$ is the real-valued expansion factor and $\|\cdot\|$ is the Frobenius norm [56]. The choice of the Frobenius norm is justified by the following argument. An approximate CNN must have its elements numerically ‘close’ to elements from the exact CNN. Therefore, a measure that takes into consideration distance in a energy-based manner (euclidean distance sense) emerges naturally as a means to guarantee that the approximate filtering structures (e.g., convolution kernels) are close to the exact counterpart. The Frobenius norm satisfies the above rationale. This analysis is confirmed in [39].

To ensure that the candidate matrices \mathbf{T} have low complexity, we limited the search space of the above problem to the matrices whose elements are defined over a sets of dyadic rationals \mathcal{D} . Some particular sets are [57], [67]:

$$\mathcal{D}_1 = \{-1, 0, 1\},$$

$$\mathcal{D}_2 = \{-2, -1, 0, 1, 2\},$$

$$\mathcal{D}_3 = \{-4, -3, -2, -1, 0, 1, 2, 3, 4\},$$

$$\mathcal{D}_4 = \left\{ -4, -3, -2, -1, -\frac{3}{4}, -\frac{1}{2}, -\frac{1}{4}, 0, \frac{1}{4}, \frac{1}{2}, \frac{3}{4}, 1, 2, 3, 4 \right\},$$

$$\mathcal{D}_5 = \left\{ -7, -6, -5, -4, -3, -2, -1, -\frac{3}{4}, -\frac{1}{2}, -\frac{1}{4}, 0, \right.$$

$$\left. \frac{1}{4}, \frac{1}{2}, \frac{3}{4}, 1, 2, 3, 4, 5, 6, 7 \right\},$$

$$\mathcal{D}_6 = \left\{ -4, -\frac{15}{4}, -\frac{7}{2}, -\frac{13}{4}, \dots, \frac{13}{4}, \frac{7}{2}, \frac{15}{4}, 4 \right\},$$

$$\mathcal{D}_7 = \left\{ -5, -\frac{19}{4}, -\frac{9}{2}, -\frac{17}{4}, \dots, \frac{17}{4}, \frac{9}{2}, \frac{19}{4}, 5 \right\}$$

$$\mathcal{D}_8 = \left\{ -7, -\frac{27}{4}, -\frac{13}{2}, -\frac{25}{4}, \dots, \frac{25}{4}, \frac{13}{2}, \frac{27}{4}, 7 \right\}.$$

Sets \mathcal{D}_6 , \mathcal{D}_7 , and \mathcal{D}_8 possess uniformly spaced rationals.

A straightforward way of addressing (2) is as follows. Considering a given set of dyadic rationals \mathcal{D} , for each element of $\alpha \cdot \mathbf{T}$, we simply find the closest neighbour of such element in \mathbf{D} . Such approach can be efficiently implemented by means of binary search. However, this approach is only effective as long as (2) remains an unconstrained optimization problem. Alternatively, we can consider a more flexible approach based on integer linear programming (ILP).

For fixed values of α , the mixed INLP problem posed in (2) can be efficiently solved by means of binary (zero-one) linear programming. In other words, we aim at converting a nonlinear problem into a linear one. Indeed, let $m_{i,j}$, $i, j = 1, 2, \dots, N$, denote the entries of \mathbf{M} and $r \in \mathcal{D}$ be a dyadic rational. We adopt the following binary decision variables:

$$x_{i,j}(r) = \begin{cases} 1, & \text{if } t_{i,j} = r, \\ 0, & \text{otherwise.} \end{cases}$$

For binary (zero-one) variables we have $y^2 = y$, where y is a dummy variable. This fact paves the way for the linearization of the above-mentioned optimization problem. Therefore, (2) can be re-written according to the following binary linear programming problem [70]–[72]:

$$\min_{x_{i,j}(r)} \sum_{i=1}^N \sum_{j=1}^N \sum_{r \in \mathcal{D}} (m_{i,j} - \alpha \cdot r)^2 \cdot x_{i,j}(r), \quad (3)$$

subject to

$$\sum_{r \in \mathcal{D}} x_{i,j}(r) = 1, \quad i, j = 1, 2, \dots, N.$$

The above constraint is to ensure that each element $m_{i,j}$ is approximated by a unique dyadic rational in \mathcal{D} . The solution of the above problem is denoted as $x_{i,j}^{(\alpha)}$, $i, j = 1, 2, \dots, N$, $r \in \mathcal{D}$, being linked to the choice of α . Such binary (zero-one) solution can be employed to compute the actual entries $t_{i,j}^{(\alpha)}$ of the low-complexity matrix associated to the considered α according to:

$$t_{i,j}^{(\alpha)} = \sum_{r \in \mathcal{D}} r \cdot x_{i,j}^{(\alpha)}(r). \quad (4)$$

The resulting low-complexity matrix is denoted by \mathbf{T}_α . The approximation error is implied by (2) and can be computed according to:

$$\text{Error}(\alpha) = \|\mathbf{M} - \alpha \cdot \mathbf{T}_\alpha\|^2.$$

Because a sequence of values for α is selected, the above problem is solved for each instantiation; furnishing the sequence of errors indexed by α : $\text{Error}(\alpha)$. Being a linear programming problem, each instantiation can be solved efficiently and very quickly by contemporary computational packages [73], [74]. State-of-the-art solvers can obtain solutions for ILP problems at an average computation complexity in $\mathcal{O}(N)$ [75] or $\mathcal{O}(N \log N)$ [70]–[72]. Finally, we determine the global optimum value α^* according to:

$$\alpha^* = \arg \min_{\alpha} \text{Error}(\alpha), \quad (5)$$

which can be solved by simple minimization over a vector

of values. Associated to α^* , we also obtain $\mathbf{T}^* \triangleq \mathbf{T}_{\alpha^*}$, which is the global optimal low-complexity matrix. Therefore, the sought approximation $\hat{\mathbf{M}}$ is given by:

$$\hat{\mathbf{M}} = \alpha^* \cdot \mathbf{T}^*. \quad (6)$$

The above ILP approach allows the user to easily include constraints to the optimization problem. This is relevant for further investigation in this topic; in particular when specific mathematical properties are expected to be enforced on the resulting low-complexity matrices (for instance, 2D filter normalization [76, p. 115]).

We emphasize that the solving method for (2) is only required to be efficient enough to cope with the time constraints at the design phase of the approximate neural network. In other words, solvers available in contemporary optimization packages are suitable; and the choice of the particular method for solving (2) is not a critical for our approach. Additionally, we note that the optimization solver is simply a step for obtaining the final neural network. Once the approximate structures are found, optimization solver are clearly not required anymore.

D. Example

To illustrate the procedure, we selected \mathcal{D}_8 and considered the search space of the expansion factor to be the interval $[0.25, 1]$ with a step of 10^{-3} . Additionally, we consider the following particular convolutional filter employed in the ConvNet described in Section IV-B:

$$\mathbf{M}_0 = \begin{bmatrix} 1.5200701 & 1.0317051 & 0.7906240 & -0.2153791 & -0.2340538 \\ 1.3982610 & 2.1860176 & 2.0152923 & 1.5620477 & 0.8270900 \\ -0.6848867 & 0.7470516 & 1.6923728 & 1.2537112 & 1.1946758 \\ -1.2387477 & -0.5483563 & 0.1261987 & 0.8677799 & 0.7742613 \\ -1.4691808 & -1.2178997 & -0.2924347 & 0.2172496 & 0.1325074 \end{bmatrix}.$$

Solving (2) for the above matrix, we obtain:

$$\alpha^* = 0.30931,$$

$$\mathbf{T}^* = \begin{bmatrix} 5 & 3.25 & 2.5 & -0.75 & -0.75 \\ 4.5 & 7 & 6.5 & 5 & 2.75 \\ -2.25 & 2.5 & 5.5 & 4 & 3.75 \\ -4 & -1.75 & 0.5 & 2.75 & 2.5 \\ -4.75 & -4 & -1 & 0.75 & 0.5 \end{bmatrix}$$

$$= \frac{1}{4} \cdot \begin{bmatrix} 20 & 13 & 10 & -3 & -3 \\ 18 & 28 & 26 & 20 & 11 \\ -9 & 10 & 22 & 16 & 15 \\ -16 & -7 & 2 & 11 & 10 \\ -19 & -16 & -4 & 3 & 2 \end{bmatrix}.$$

Fig. 1(a) depicts the Frobenius norm error for varying values of α (cf. 2). For very small α , the values of $\alpha \cdot \mathbf{M}$ are close to zero. So the discrete entries of the candidate matrices \mathbf{T} are unable to provide a good approximation. As α increases, a similar effect happens. However, for intermediate values the minimum can be found. Fig. 1(b) shows details in the vicinity of the optimum. The curves shown in Fig. 1 are piecewise concatenations of parabolae. This is due to the quadratic nature of the coefficients $(r - \alpha \cdot m_{i,j})^2$ of the linear programming problem in (3). Each parabola is linked to a particular approximate candidate \mathbf{T} .

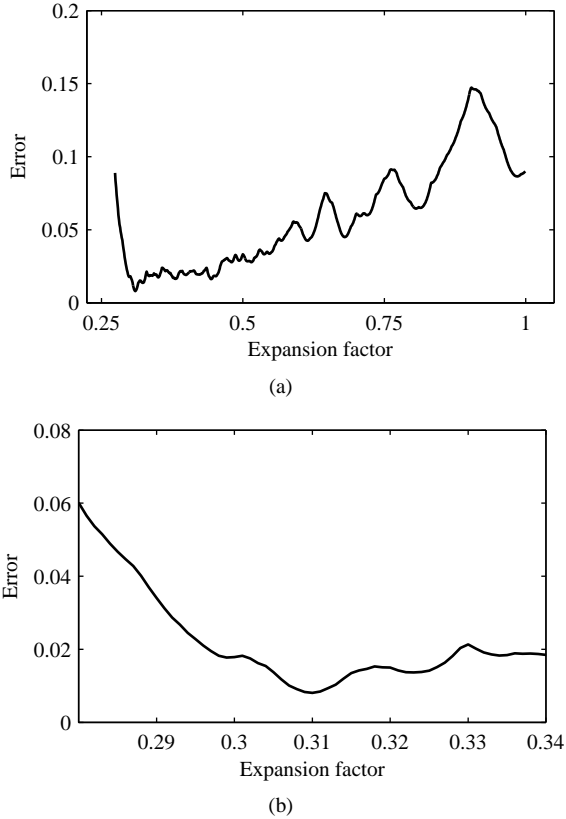


Fig. 1: Approximation error for the particular matrix \mathbf{M}_0 : (a) error curve over the considered search interval for α and (b) detailed view around the optimum value of α .

Notice that the low-complexity matrix \mathbf{T}^* is expressed in terms of small integers, which can be given simple binary expansions (e.g., $22 = 2^6 - 2^4 - 2^2$). Similarly, we have that $\alpha^* = 0.30931 \approx 2^{-2} + 2^{-4} - 2^{-8} = 0.30859375$.

Therefore, considering (6), the actual fully multiplierless approximation is furnished by:

$$\hat{\mathbf{M}} = (2^{-4} + 2^{-6} - 2^{-10}) \cdot \begin{bmatrix} 20 & 13 & 10 & -3 & -3 \\ 18 & 28 & 26 & 20 & 11 \\ -9 & 10 & 22 & 16 & 15 \\ -16 & -7 & 2 & 11 & 10 \\ -19 & -16 & -4 & 3 & 2 \end{bmatrix}.$$

E. Activation Function Approximations

Although there are several types of activation functions, we focus our analyses on the continuous tanh-sigmoid function, which is defined according to the hyperbolic tangent function [53]. As indicated in [53], [77], the mathematical expression for the tanh-sigmoid function is given by:

$$\phi(x) = a \cdot \tanh(b \cdot x), \quad (7)$$

where $a = 1.7159$ and $b = 2/3$. This particular activation function has been originally proposed by LeCun [78] and adopted in several working models as the Convolutional Face Finder (CFF) [6].

In [58], [60], [79]–[82], approximations for the related sigmoid function given by $y = 1/(1 + e^{-x})$ were examined,

including the Alippi and Storti-Gajani (ASG) approximation [83], the piecewise linear approximation of a non-linear function (PLAN) [84], and simple linear [58] and quadratic [79] approximations. Based on these approximations, we derived expressions for the tanh-sigmoid approximations as shown in Table I. We adopted an 8-bit representation for a resulting in the following approximate value: $\hat{a} = 7/4$.

F. Complexity

As a consequence of the above approximations, we have substantial savings in computation costs. Indeed, a single call of the original $N \times N$ matrix \mathbf{M} requires N^2 multiplications of floating-point entries per pixel. On the other hand, the proposed approximation $\hat{\mathbf{M}}$ contains only small integers that can be very efficiently encoded with minimal number of adders [66]. Similarly, the expansion factor α^* can be given a truncated rational approximation in the form of dyadic rationals. The same rationale also applies to the remaining computational structures of the original ConvNet. Thus, the final resulting structure is fully *multiplierless*—only additions and bit-shifting operations are required. In terms of hardware realization, the number of arithmetic operations translate into chip area and power consumption [85], [86]. Thus, in limited resource scenarios (e.g., embedded systems and wireless sensors), approximations may provide an effective way of porting large ConvNets into physical realization.

To summarize, the proposed approximation approach consists of:

- (i) finding approximate convolutional filter by solving (1) for each exact convolutional filter from a given ConvNet;
- (ii) converting scaling factors, sub-sampling coefficients, and bias values into CSD representation aiming at the minimization of computation costs and multiplicative irreducibility;
- (iii) approximating the activation function to a simple function.

IV. EXPERIMENTS

We studied the effectiveness of the proposed approximation approach on two classical computer vision problems: (i) a binary and (ii) a multi-class classification problem. The first application is face detection, where the ConvNet classifies image regions as face or non-face. The second one is handwritten digit recognition, where the trained model is used to classify a given image patch into one of the ten digits “0” to “9”. For each of the two applications, we trained a ConvNet in a classical way and evaluated its performance in terms of precision and recall, for the given application.

The first ConvNet is relatively small, whereas the second model (for digit recognition) contains much more parameters. We aim at demonstrating that our proposed approach is able to effectively process larger networks.

After approximating the parameters of the models, we compared their performance with their respective original, exact versions. Note that we do not aim at improving the state-of-the-art in face detection or hand-written digit recognition. Indeed, current literature presents concrete example of

TABLE I: Approximations for the tanh-sigmoid

ASG-based	$\sigma_1(x) = \hat{a} \cdot \begin{cases} -1 + \frac{1 + \frac{\ x\ - x}{2^{ x }}}{2}, & x < 0, \\ 1 - \frac{1 + \frac{\ x\ - x}{2^{ x }}}{2}, & x \geq 0. \end{cases}$	PLAN-based	$\sigma_2(x) = \hat{a} \cdot \begin{cases} -1, & x < -5, \\ \frac{x}{16} - \frac{89}{128}, & -5 \leq x < -\frac{19}{8}, \\ \frac{x}{4} - \frac{1}{4}, & -\frac{19}{8} \leq x < -1, \\ \frac{x}{2}, & -1 \leq x < 1, \\ \frac{x}{4} + \frac{1}{4}, & 1 \leq x < \frac{19}{8}, \\ \frac{x}{16} + \frac{11}{16}, & \frac{19}{8} \leq x < 5, \\ 1, & x \geq 5. \end{cases}$
Linear I	$\sigma_3(x) = \hat{a} \cdot \begin{cases} -1, & x < -4, \\ \frac{x}{4}, & -4 \leq x < 4, \\ 1, & x \geq 4. \end{cases}$	Linear II	$\sigma_4(x) = \hat{a} \cdot \begin{cases} -1, & x < -2, \\ \frac{x}{2}, & -2 \leq x < 2, \\ 1, & x \geq 2. \end{cases}$
Quadratic I	$\sigma_6(x) = \hat{a} \cdot \begin{cases} -1, & x < -4, \\ (\frac{x}{4} + 1)^2 - 1, & -4 \leq x < 0, \\ 1 - (\frac{x}{4} + 1)^2, & 0 \leq x < 4, \\ 1, & x \geq 4. \end{cases}$	Quadratic II	$\sigma_6(x) = \hat{a} \cdot \begin{cases} -1, & x < -2, \\ (\frac{x}{2} + 1)^2 - 1, & -2 \leq x < 0, \\ 1 - (\frac{x}{2} + 1)^2, & 0 \leq x < 2, \\ 1, & x \geq 2. \end{cases}$

complex models, such as multi-view or part-based detectors for face detection [87] and huge ensemble classifiers for digit recognition [88]. Our goal is to demonstrate—based on common representative models—that the complexity of a given trained ConvNet model can be reduced significantly by approximating its parameters while maintaining a very similar performance.

Hereafter an approximate network based on dyadic set \mathcal{D}_i is referred to as A_i .

A. Binary Classification

Our first set of experiments employs a ConvNet that was trained for face detection in grey-scale images. Thus, such network is a binary classifier that decides whether the given input image is a face or not. As a working model, we selected the classical face detector called Convolutional Face Finder (CFF) proposed by Garcia and Delakis [6]. This model is a relatively “light” ConvNet with an input size of 32×36 and six layers: four layers alternating convolution and average pooling operations, with 4, 4, 14, and 14 maps, respectively, followed by 14 neurons and one single final output neuron. The first convolution layer contains four filters of size 5×5 , the second one contains 20 filters of size 3×3 , and the 14 neurons of the first neuron layer are treated as convolutions of size 6×7 , each neuron being connected to only one map. Pooling maps contain a single coefficient, and all maps and neurons have an additional bias. The entire ConvNet has 951 trainable parameters in total. A thorough description of this particular ConvNet is supplied in [6]. The employed activation function is the exact continuous tanh-sigmoidal function as detailed in (7).

After training the ConvNet as described in [6], we approximated all the convolution filter matrices with low-complexity versions. We created several approximations using the different sets of dyadic rationals described in the previous section: $\mathcal{D}_1, \mathcal{D}_1, \dots, \mathcal{D}_8$. We also replaced all average pooling coeffi-

TABLE II: Arithmetic cost for CFF-based models

Model	Operation			
	Mult.	Add.	CSD Add.	Bit-shifting
Exact	882	843	-	-
A_1	0	843	235	346
A_2	0	843	251	362
A_3	0	843	377	488
A_4	0	843	457	568
A_5	0	843	506	617
A_6	0	843	756	867
A_7	0	843	842	953
A_8	0	843	1028	1139

cients and bias terms with their closest CSD representation using 8 bits, being 7 bits for the fractional part.

Table II lists the arithmetic costs of the exact ConvNet compared to its approximations. Floating-point multiplications, direct additions, additions due to the CSD expansion, and bit-shifting operations were counted. The exact structure requires both floating-point multiplications and additions. In contrast, the approximate methods completely eliminates the need for multiplications at the expense of much simpler operations: additions and bit-shifting operations. Because the approximate quantities can be easily represented in fixed-point arithmetic representation, it is suitable for hardware implementation. Additionally, the hardware implementation of bit-shifting operations require virtually no cost, because it can be implemented by simple physical wiring. As a result, we have a very favorable trade-off: multiplications are exchanged for additions.

In order to analyse the effect of the approximation on the actual performance of the ConvNet, we evaluated the different versions on three standard face detection benchmarks: FDDB [89] (2845 images), AFW [90] (205 images), and

Pascal Faces dataset [91] (851 images); and we used the improved annotation and evaluation protocol proposed by Mathias *et al.* [87].

Tables III-V show the average precision rates for different combinations of sigmoid and weight matrix approximations relative to the exact model for the three datasets.

Overall, the “Linear II” sigmoid approximation provides the best results, followed by “ASG”, “Quadratic I”, and “PLAN”. In terms of weight matrix approximations, A_1 – A_5 generally give unsatisfactory results, and A_7 performs best. Also, it is interesting to note that the finer approximation A_8 gave worse results than A_7 . We further evaluated some variants, where different layers of the ConvNet have been approximated with different sets of dyadic rationals. For such mixed approximate structures, we have denoted them by $A_{i,j,k,l}$, where the subscripts indicate the selected dyadic set for each layer. In other words, indices i, j indicate that the dyadic sets \mathcal{D}_i and \mathcal{D}_j , respectively, are employed in the first two convolution layers; and similarly, indices k, l correspond to the adoption of the dyadic sets \mathcal{D}_k and \mathcal{D}_l , respectively, for the two final fully-connected layers. We found that the first layer requires a finer approximation than the other layers. This allowed us to maintain a good performance with very low-complexity approximations (e.g., A_3 , A_4) for these later layers. This can be explained according to the following: (i) by its own very nature, the layers have different degrees of importance; (ii) errors in initial layers tend to propagate through the succeeding layers; and (iii) error propagation can potentially be amplified along the layers. Fig. 2 shows the receiver operating characteristic (ROC) curves of the best performing approximations for the three datasets.

These results are quite impressive given the fact that we considerably reduced the precision of each parameter of the ConvNet, and given the highly non-linear classification problem where the frontier between the face and non-face classes can be very thin and complex.

Fig. 3 shows some face detection results from the exact model (top) and the approximation $A_{7,3,3,3}$ (bottom), i.e. a finer approximation for the first layer and a very coarse one for the rest of the layers. The results are almost identical.

B. Multi-class Classification

We studied a second case where a ConvNet has been trained for a classical multi-class classification problem: the MNIST hand-written digit recognition dataset [92]. To show that the proposed approximations can also be applied to larger networks we trained a ConvNet with a different architecture containing again six layers but much more maps and around 180000 parameters and more than 5300 matrices in total. The input is a 32×32 grey-scale image, and the network is composed of five convolution maps (5×5 kernels) followed by five average pooling maps (connected one-to-one), 50 convolution maps (3×3 , fully connected), 50 average pooling maps (connected one-to-one), 100 neurons (6×6 matrices, fully connected), and the 10 final output neurons corresponding to the 10 digits to classify.

After having trained this ConvNet model on the MNIST dataset, we approximated all the convolution filters, the fully-

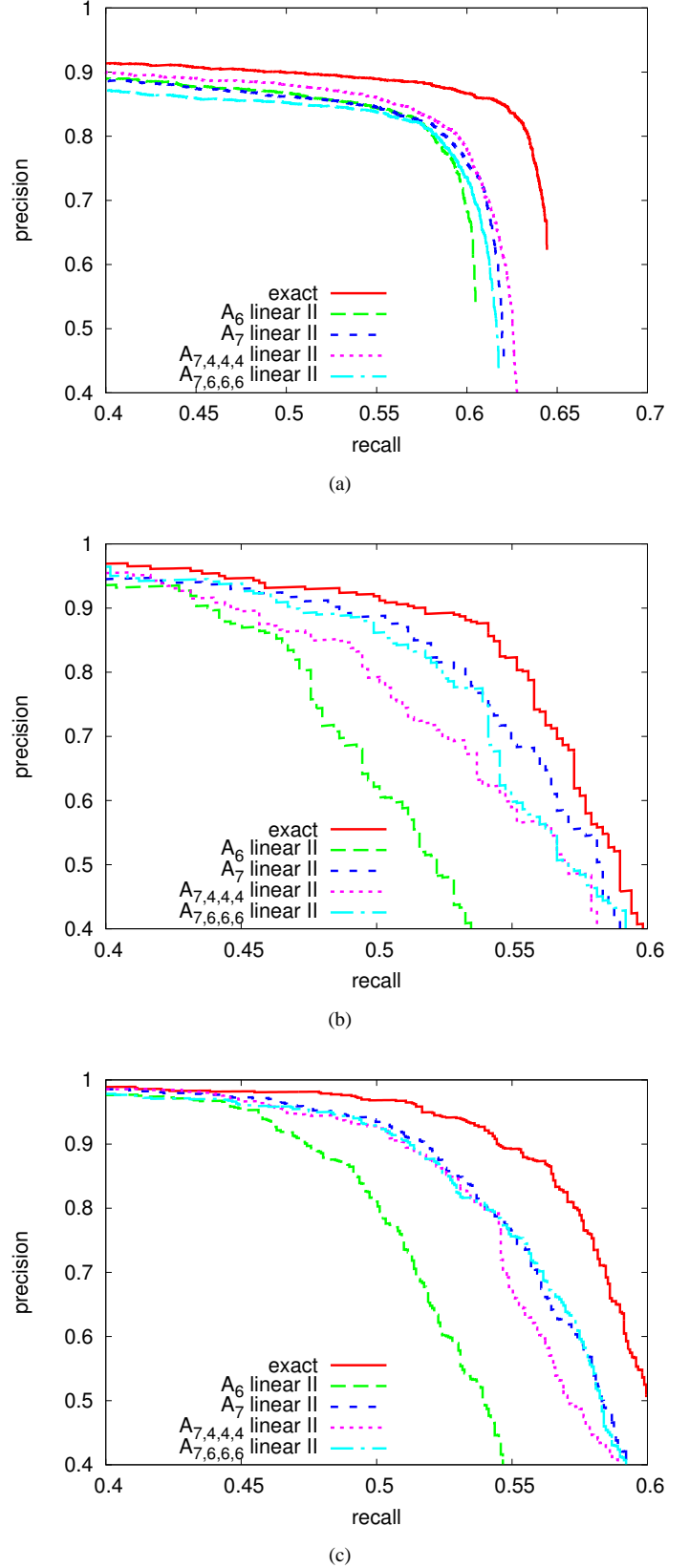


Fig. 2: ROC curves for the (a) FDDB, (b) AFW, and (c) Pascal datasets comparing the face detection performance of the original (exact) ConvNet model with different approximations.

TABLE III: Average precision for CFF with the FDDB test set and different approximations relative to the exact model

	Exact	ASG	PLAN	Linear I	Linear II	Quadratic I	Quadratic II
Exact	1.000	0.953	0.894	0.002	0.988	0.887	0.938
A ₁	0.000	0.000	0.000	0.000	0.000	0.000	0.000
A ₂	0.001	0.001	0.000	0.000	0.001	0.000	0.002
A ₃	0.549	0.311	0.432	0.005	0.510	0.426	0.180
A ₄	0.523	0.602	0.365	0.001	0.558	0.383	0.602
A ₅	0.490	0.098	0.517	0.000	0.657	0.510	0.073
A ₆	0.938	0.914	0.817	0.002	0.888	0.797	0.949
A ₇	0.960	0.943	0.847	0.002	0.963	0.848	0.935
A ₈	0.821	0.792	0.648	0.001	0.810	0.650	0.793
A _{7,3,3,3}	0.917	0.902	0.886	0.003	0.947	0.888	0.851
A _{7,4,4,4}	0.967	0.931	0.855	0.000	0.976	0.861	0.928
A _{7,6,6,6}	0.959	0.907	0.862	0.004	0.959	0.863	0.933

TABLE IV: Average precision for CFF with the AFW test set and different approximations relative to the exact model

	Exact	ASG	PLAN	Linear I	Linear II	Quadratic I	Quadratic II
Exact	1.000	0.839	0.829	0.000	1.041	0.797	0.383
A ₁	0.000	0.000	0.000	0.000	0.000	0.000	0.000
A ₂	0.001	0.001	0.000	0.000	0.000	0.000	0.000
A ₃	0.220	0.041	0.260	0.004	0.199	0.248	0.014
A ₄	0.422	0.356	0.317	0.000	0.457	0.336	0.275
A ₅	0.220	0.015	0.314	0.000	0.217	0.251	0.002
A ₆	0.864	0.794	0.715	0.000	0.893	0.694	0.700
A ₇	0.978	0.844	0.753	0.000	0.985	0.755	0.614
A ₈	0.698	0.553	0.576	0.000	0.565	0.544	0.169
A _{7,3,3,3}	0.955	0.525	0.835	0.004	0.914	0.816	0.184
A _{7,4,4,4}	1.020	0.827	0.755	0.000	0.954	0.761	0.576
A _{7,6,6,6}	0.967	0.881	0.787	0.000	0.979	0.788	0.827

TABLE V: Average precision for CFF with the ‘‘Pascal faces’’ test set and different approximations relative to the exact model

	Exact	ASG	PLAN	Linear I	Linear II	Quadratic I	Quadratic II
Exact	1.000	0.933	0.774	0.001	1.039	0.747	0.893
A ₁	0.001	0.001	0.000	0.000	0.001	0.001	0.001
A ₂	0.006	0.006	0.004	0.000	0.006	0.003	0.009
A ₃	0.234	0.128	0.239	0.000	0.234	0.230	0.098
A ₄	0.357	0.375	0.259	0.001	0.361	0.258	0.370
A ₅	0.433	0.143	0.386	0.000	0.449	0.345	0.101
A ₆	0.862	0.844	0.656	0.000	0.894	0.638	0.847
A ₇	0.917	0.879	0.707	0.001	0.971	0.702	0.889
A ₈	0.725	0.700	0.520	0.001	0.697	0.517	0.617
A _{7,3,3,3}	0.899	0.752	0.755	0.004	0.899	0.752	0.645
A _{7,4,4,4}	0.930	0.889	0.710	0.000	0.967	0.708	0.906
A _{7,6,6,6}	0.931	0.906	0.725	0.004	0.970	0.722	0.916



Fig. 3: Some CFF face detection results on the AFW dataset. *Top*: exact model; *bottom*: approximation $A_{7,3,3,3}$. Despite the very coarse approximation, the results are very close. In the second last image, the approximate model even detects an additional face, missed by the original CFF. However, in the last example a false detection is produced.

TABLE VI: Arithmetic cost for MNIST-based models

Model	Operation			
	Mult.	Add.	CSD Add.	Bit-shifting
Exact	183375	178110	-	-
A_1	0	178110	12740	23325
A_2	0	178110	12722	23307
A_3	0	178110	49127	59712
A_4	0	178110	61228	71813
A_5	0	178110	65211	75796
A_6	0	178110	141401	151986
A_7	0	178110	158595	169180
A_8	0	178110	188417	199002

connected layer matrices and the activation functions which is based on the tanh-sigmoid function. Again, all pooling coefficients and bias terms were replaced by their closest CSD representation using 8 bits. The computation cost of the exact and approximate structures is shown in Table VI. Similarly to the previous experiment, the approximate models have totally eliminated the multiplicative costs. Floating-point arithmetic is not required; being fixed-point arithmetic adequate. The cost of the extra additions due to the CSD representation is very low compared to the multiplicative cost required by the exact model. The cost of bit-shifting operations is negligible.

Table VII shows the relative classification rates on the MNIST test set for the different approximations, and Fig. 4 depicts the respective ROC curves of the best-performing approximations (combined for the 10 classes).

The results show that the approximations, although very coarse, have a very small effect on the overall performance of the ConvNet. In a multi-class setting, the trained ConvNet, at least in this particular case, is much more robust to the loss in precision of the weights induced by our approximation scheme compared to the binary classifier. For example, a very coarse approximation like $A_{3,3,1,1}$ leads to a relative performance decrease of less than 1%.

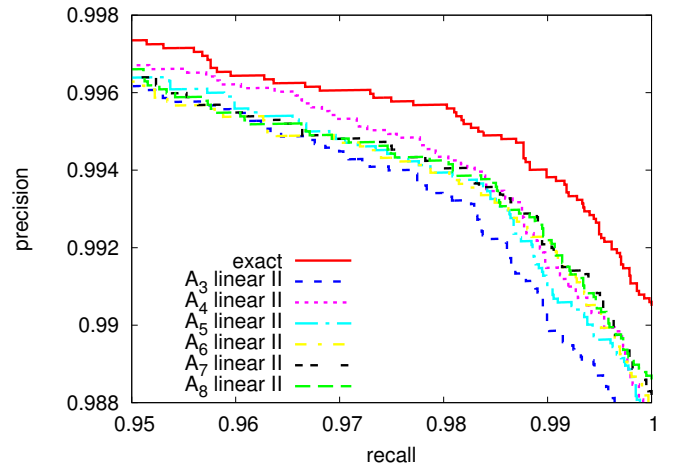


Fig. 4: ROC curves for the MNIST hand-written digit classification test set comparing the performance of the original (exact) ConvNet model with its approximations.

C. Large-Scale Deep Neural Networks

Finally, we applied our approximation approach to a deeper and more complex network architecture, the well-known AlexNet proposed by Krizhevsky *et al.* [12], and the ImageNet dataset [13] for image classification with 1000 classes. This model contains more than 1.2 million matrices and 5096 vectors. We approximated all convolution filter matrices of the fully-trained 8-layer ConvNet using two different sets of dyadic rationals for different layers, a very coarse set \mathcal{D}_9 and a slightly finer set \mathcal{D}_{10} :

$$\mathcal{D}_9 = \left\{ -2, -1, -\frac{1}{2}, -\frac{1}{8}, 0, \frac{1}{8}, \frac{1}{2}, 1, 2 \right\},$$

$$\mathcal{D}_{10} = \left\{ -2, -1, -\frac{1}{2}, -\frac{1}{4}, -\frac{1}{8}, 0, \frac{1}{8}, \frac{1}{4}, \frac{1}{2}, 1, 2 \right\}.$$

Again, all other coefficients are approximated by their closest 8-bit CSD representation. The pooling layers do not have any coefficient here, and only linear and Rectified Linear

TABLE VII: Mean classification rates for the MNIST test set and different approximations relative to the exact model.

	Exact	ASG	PLAN	Linear I	Linear II	Quadratic I	Quadratic II
Exact	1.0000	1.0000	0.9847	0.9680	0.9978	1.0000	1.0000
A_1	0.9684	0.9684	0.9588	0.9260	0.9615	0.9684	0.9684
A_2	0.9643	0.9643	0.9627	0.8805	0.9573	0.9643	0.9643
A_3	0.9961	0.9961	0.9848	0.9655	0.9944	0.9961	0.9961
A_4	0.9973	0.9973	0.9863	0.9700	0.9969	0.9973	0.9973
A_5	0.9976	0.9976	0.9866	0.9666	0.9969	0.9976	0.9976
A_6	0.9991	0.9991	0.9868	0.9701	0.9973	0.9991	0.9991
A_7	0.9992	0.9992	0.9846	0.9680	0.9977	0.9992	0.9992
A_8	0.9994	0.9994	0.9848	0.9675	0.9981	0.9994	0.9994
$A_{3,3,1,1}$	0.9931	0.9931	0.9749	0.9625	0.9924	0.9931	0.9931
$A_{3,1,1,1}$	0.9891	0.9891	0.9684	0.9580	0.9866	0.9891	0.9891
$A_{4,4,1,1}$	0.9937	0.9937	0.9780	0.9618	0.9943	0.9937	0.9937
$A_{4,1,1,1}$	0.9885	0.9885	0.9655	0.9572	0.9872	0.9885	0.9885

TABLE VIII: Classification accuracy and top-5 accuracy for ImageNet and different approximations relative to the exact AlexNet model

	Absolute		Relative	
	Accuracy	Top-5	Accuracy	Top-5
Exact	0.5682	0.7995	1.0000	1.0000
A_9	0.4862	0.7288	0.8558	0.9117
A_{10}	0.5463	0.7820	0.9616	0.9782
$A_{10,9,9,9,9}$	0.5423	0.7794	0.9544	0.9750
$A_{10,10,9,9,9,9}$	0.5442	0.7796	0.9578	0.9751

Units (ReLU) are used as activation function, which are already of very low complexity and thus do not require any approximation.

We used the ImageNet 2012 validation set to evaluate our different approximations. And, as usual in the literature, we compute the classification accuracy as well as the top-5 accuracy for the 50000 test images. Table VIII shows the results. The approximation A_{10} with the set \mathcal{D}_{10} gives the best performance, with a relative decrease in accuracy of only 3.84% and 2.18% on the top-5 accuracy. However, as the following line shows, we can achieve almost the same performance using the coarser set \mathcal{D}_9 for all convolution layers except the first one. This again suggests that a finer approximation of the first layer is required to prevent a drastic performance drop.

V. CONCLUSION

We presented a novel scheme for approximating the parameters of a trained ConvNet, notably the convolution filters, neuron weights, as well as pooling and bias coefficients. Activation functions were also approximated. The particularity of the matrix approximations is that they allow for an extremely efficient implementation—software or hardware—using only additions and bit-shifts, and no multiplication. We thoroughly evaluated the impact of this parameter approximation measuring the overall performance of ConvNets on three different

use cases: one smaller ConvNet for face detection, a larger ConvNet for hand-written digit classification, and a much more complex, deep ConvNet for large-scale image classification.

For all three models, our proposed scheme was able to produce low-complexity approximations without a significant loss in performance.

These results suggest that huge reductions in computational complexity of trained ConvNet models can be obtained, and extremely efficient hardware implementations can be realized. Further studies need to be undertaken to analyse the impact of this type of approximations for more use cases and different architectures.

REFERENCES

- [1] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceeding of the IEEE*, Nov. 1998.
- [2] K. Chellapilla, S. Puri, and P. Simard, “High performance convolutional neural networks for document processing,” in *Proc. of the Int. Workshop on Frontiers in Handwriting Recognition (IWFHR’06)*, 2006.
- [3] M. Delakis and C. Garcia, “Text detection with convolutional neural networks,” in *VISAPP 2008: Proceedings of the Third International Conference on Computer Vision Theory and Applications*, vol. 2, (Funchal, Madeira, Portugal), pp. 290–294, Jan. 2008.
- [4] K. Elagouni, C. Garcia, F. Mamalet, and P. Sébillot, “Text recognition in multimedia documents: a study of two neural-based OCRs using and avoiding character segmentation,” *IJDAR*, vol. 17, no. 1, pp. 19–31, 2014.
- [5] R. Vaillant, C. Monroq, and Y. Le Cun, “An original approach for the localization of objects in images,” in *IEE Proc on Vision, Image, and Signal Processing*, pp. 141(4):245–250, 1994.
- [6] C. Garcia and M. Delakis, “Convolutional face finder: A neural architecture for fast and robust face detection,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 26, no. 11, pp. 1408–1423, 2004.
- [7] M. Osadchy, Y. LeCun, M. L. Miller, and P. Perona, “Synergistic face detection and pose estimation with energy-based model,” in *Proc. of Advances in Neural Information Processing Systems (NIPS’05)*, 2005.
- [8] D. C. Ciresan, U. Meier, J. Masci, and J. Schmidhuber, “Multi-column deep neural network for traffic sign classification,” *Neural Networks*, vol. 32, pp. 333–338, 2012.
- [9] P. Sermanet, K. Kavackuoglu, S. Chintala, and Y. LeCun, “Pedestrian detection with unsupervised multi-stage feature learning,” in *2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, June 23-28, 2013*, pp. 3626–3633, 2013.
- [10] R. Hadsell, P. Sermanet, M. Scoffier, A. Erkan, K. Kavackuoglu, U. Muller, and Y. LeCun, “Learning long-range vision for autonomous off-road driving,” *Journal of Field Robotics*, Feb. 2009.

- [11] P. Sermanet, S. Chintala, and Y. LeCun, "Convolutional neural networks applied to house numbers digit classification," in *Proceedings of the 21st International Conference on Pattern Recognition, ICPR 2012, Tsukuba, Japan, November 11-15, 2012*, pp. 3288–3291, 2012.
- [12] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States.*, pp. 1106–1114, 2012.
- [13] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.
- [14] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, pp. 533–536, 1986.
- [15] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, "OverFeat: Integrated recognition, localization and detection using convolutional networks," *CoRR*, vol. abs/1312.6229, 2013.
- [16] P. K. Meher, S. Y. Park, B. K. Mohanty, K. S. Lim, and C. Yeo, "Efficient integer DCT architectures for HEVC," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 24, pp. 168–178, Jan 2014.
- [17] Y.-J. Chen, S. Orantara, and T. Nguyen, "Video compression using integer DCT," in *Proceedings 2000 International Conference on Image Processing (Cat. No.00CH37101)*, vol. 2, pp. 844–845 vol.2, Sept. 2000.
- [18] A. M. Joshi, V. Mishra, and R. M. Patrikar, "Design of real-time video watermarking based on integer dct for H.264 encoder," *International Journal of Electronics*, vol. 102, no. 1, pp. 141–155, 2015.
- [19] V. A. Coutinho, R. J. Cintra, F. M. Bayer, P. A. M. Oliveira, R. S. Oliveira, and A. Madanayake, "Pruned discrete Tchebichef transform approximation for image compression," *Circuits, Systems, and Signal Processing*, 2018.
- [20] D. Scherer, H. Schulz, and S. Behnke, "Accelerating large-scale convolutional neural networks with parallel graphics multiprocessors," in *Proceedings of the International Conference on Artificial Neural Networks*, pp. 82–91, 2010.
- [21] B. White and M. Elmasry, "The digi-neocognitron: A digital neocognitron neural network model for VLSI," *IEEE Transactions on Neural Networks*, vol. 3, no. 1, pp. 73–85, 1992.
- [22] M. Marchesi, G. Orlando, F. Piazza, and A. Uncini, "Fast neural networks without multipliers," *IEEE Transactions on Neural Networks*, vol. 4, no. 1, pp. 53–62, 1993.
- [23] H. Kwan and C. Tang, "Multiplierless multilayer feedforward neural network design suitable for continuous input-output mapping," *Electronic Letters*, vol. 29, no. 14, pp. 1259–1260, 1993.
- [24] J. Vincent and D. Myers, "Weight dithering and wordlength selection for digital backpropagation networks," *BT Technology Journal*, vol. 10, no. 3, pp. 124–133, 1992.
- [25] P. Simard and H. P. Graf, "Backpropagation without multiplication," in *Proceedings of the Annual Conference on Neural Information Processing Systems*, pp. 232–239, 1994.
- [26] S. Draghici, "On the capabilities of neural networks using limited precision weights," *Neural Networks*, vol. 15, no. 3, pp. 395–414, 2002.
- [27] E. L. Machado, C. J. Miosso, R. von Borries, M. Coutinho, P. de Azevedo Berger, T. Marques, and R. P. Jacobi, "Computational cost reduction in learned transform classifications," *CoRR*, 2015.
- [28] M. Courbariaux, Y. Bengio, and J.-P. David, "BinaryConnect: Training deep neural networks with binary weights during propagations," in *Proceedings of the Annual Conference on Neural Information Processing Systems*, 2015.
- [29] M. Kim and S. Paris, "Bitwise neural networks," in *ICML Workshop on Resource-Efficient Machine Learning*, 2015.
- [30] F. Mamalet, S. Roux, and C. Garcia, "Real-time video convolutional face finder on embedded platforms," *EURASIP J. Emb. Sys.*, vol. 2007, 2007.
- [31] N. Farrugia, F. Mamalet, S. Roux, F. Yang, and M. Paindavoine, "Fast and robust face detection on a parallel optimized architecture implemented on FPGA," *IEEE Trans. Circuits Syst. Video Techn.*, vol. 19, no. 4, pp. 597–602, 2009.
- [32] C. Farabet, B. Martini, P. Akselrod, S. Talay, Y. LeCun, and E. Culurciello, "Hardware accelerated convolutional neural networks for synthetic vision systems," in *International Symposium on Circuits and Systems (ISCAS 2010), May 30 - June 2, 2010, Paris, France*, pp. 257–260, 2010.
- [33] S. Chakradhar, M. Sankaradas, V. Jakkula, and S. Cadambi, "A dynamically configurable coprocessor for convolutional neural networks," in *Proceedings of the International Symposium on Computer Architecture*, pp. 247–257, 2010.
- [34] T. Chen, Z. Du, N. Sun, J. Wang, C. Wu, Y. Chen, and O. Temam, "DianNao: A small-footprint high-throughput accelerator for ubiquitous machine-learning," in *Proceedings of the International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, pp. 269–284, 2014.
- [35] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, "Optimizing FPGA-based accelerator design for deep convolutional neural networks," in *Proceedings of the ACM/SIGDA International Symposium on Field-Programmable Gate Arrays - FPGA*, pp. 161–170, 2015.
- [36] F. Mamalet and C. Garcia, "Simplifying convnets for fast learning," in *Artificial Neural Networks and Machine Learning - ICANN 2012 - 22nd International Conference on Artificial Neural Networks, Lausanne, Switzerland, September 11-14, 2012, Proceedings, Part II*, pp. 58–65, 2012.
- [37] M. Mathieu, M. Henaff, and Y. LeCun, "Fast training of convolutional networks through FFTs," in *Proceedings of the International Conference on Learning Representations*, 2014.
- [38] V. Vanhoucke, A. Senior, and M. Mao, "Improving the speed of neural networks on CPUs," in *Proceedings of the Deep Learning and Unsupervised Feature Learning NIPS Workshop*, 2011.
- [39] K. Osawa and R. Yokota, "Evaluating the compression efficiency of the filters in convolutional neural networks," in *Artificial Neural Networks and Machine Learning - ICANN 2017 (A. Lintas, S. Rovetta, P. V. P., and A. Villa, eds.)*, vol. 10614 of *Lecture Notes in Computer Science*, pp. 459–466, Springer, 2017.
- [40] J. Xue, J. Li, and Y. Gong, "Restructuring of deep neural network acoustic models with singular value decomposition," in *Interspeech*, 2013.
- [41] T. N. Sainath, B. Kingsbury, V. Sindhwani, E. Arisoy, and B. Ramabhadran, "Low-rank matrix factorization for deep neural network training with high-dimensional output targets," in *ICASSP*, pp. 6655–6659, 2013.
- [42] M. Denil, B. Shakibi, L. Dinh, and N. de Freitas, "Predicting parameters in deep learning," in *Proceedings of the Annual Conference on Neural Information Processing Systems*, 2013.
- [43] E. Denton, W. Zaremba, and J. Bruna, "Exploiting linear structure within convolutional networks for efficient evaluation," in *Proceedings of the Annual Conference on Neural Information Processing Systems*, 2014.
- [44] Z. Yang, M. Moczulski, M. Denil, N. D. Freitas, A. Smola, L. Song, and Z. Wang, "Deep Fried Convnets," in *Proceedings of the International Conference on Computer Vision*, 2014.
- [45] M. Jaderberg, A. Vedaldi, and A. Zisserman, "Speeding up convolutional neural networks with low rank expansions," in *Proceedings of the British Machine Vision Conference*, 2014.
- [46] V. Lebedev, Y. Ganin, M. Rakhuba, I. Oseledets, and V. Lempitsky, "Speeding up convolutional neural networks using fine-tuned CP decomposition," in *Proceedings of the International Conference on Learning Representations*, 2015.
- [47] M. Courbariaux, Y. Bengio, and J.-P. David, "Binaryconnect: Training deep neural networks with binary weights during propagations," in *Advances in Neural Information Processing Systems*, pp. 3123–3131, 2015.
- [48] M. Courbariaux and Y. Bengio, "Binarynet: Training deep neural networks with weights and activations constrained to +1 or -1," *CoRR*, vol. abs/1602.02830, 2016.
- [49] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "XNOR-Net: Imagenet classification using binary convolutional neural networks," *CoRR*, vol. abs/1603.05279, 2016.
- [50] M. Courbariaux, Y. Bengio, and J. David, "Low precision arithmetic for deep learning," *CoRR*, vol. abs/1412.7024, 2014.
- [51] D. Miyashita, E. H. Lee, and B. Murmann, "Convolutional neural networks using logarithmic data representation," *CoRR*, vol. abs/1603.01025, 2016.
- [52] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*. The MIT Press, 3rd ed., July 2009.
- [53] S. Haykin, *Neural Networks*. Upper Saddle River, NJ: Prentice-Hall, Inc., 1999.
- [54] R. J. Cintra, "An integer approximation method for discrete sinusoidal transforms," *Circuits, Systems, and Signal Processing*, vol. 30, no. 6, pp. 1481–1501, 2011.
- [55] P. A. M. Oliveira, R. J. Cintra, F. M. Bayer, S. Kulasekera, and A. Madanayake, "A discrete Tchebichef transform approximation for image and video coding," *IEEE Signal Processing Letters*, vol. 22, pp. 1137–1141, Jan. 2015.
- [56] G. A. F. Seber, *A Matrix Handbook for Statisticians*. Wiley-Interscience, 2007.

- [57] C. J. Tablada, F. M. Bayer, and R. J. Cintra, "A class of DCT approximations based on the Feig-Winograd algorithm," *Signal Processing*, vol. 113, pp. 38–51, 2015.
- [58] M. T. Tommiska, "Efficient digital implementation of the sigmoid function for reprogrammable logic," *IEEE Proceedings - Computers and Digital Techniques*, vol. 150, no. 6, pp. 403–411, 2003.
- [59] M. Zhang, S. Vassiliadis, and J. G. Delgado-Frias, "Sigmoid generators for neural computing using piecewise approximations," *IEEE Transactions on Computers*, vol. 45, pp. 1045–1049, Sept. 1996.
- [60] K. Basterretxea, J. M. Tarela, and I. del Campo, "Approximation of sigmoid function and the derivative for hardware implementation of artificial neurons," *IEE Proceedings - Circuits, Devices and Systems*, vol. 151, no. 1, pp. 18–24, 2004.
- [61] S. Bouguezel, M. O. Ahmad, and M. N. S. Swamy, "A low-complexity parametric transform for image compression," in *IEEE International Symposium on Circuits and Systems*, (Rio de Janeiro, BR), pp. 2145–2148, May 2011.
- [62] S. Bouguezel, M. O. Ahmad, and M. N. S. Swamy, "A multiplication-free transform for image compression," in *2nd International Conference on Signals, Circuits and Systems*, (Monastir, TN), pp. 1–4, 2008.
- [63] S. Bouguezel, M. O. Ahmad, and M. N. S. Swamy, "Low-complexity 8x8 transform for image compression," *Electronics Letters*, vol. 44, pp. 1249–1250, Oct. 2008.
- [64] V. Britanak, P. Yip, and K. R. Rao, *Discrete Cosine and Sine Transforms*. Academic Press, 2007.
- [65] N. Ahmed, T. Natarajan, and K. R. Rao, "Discrete cosine transform," *IEEE Transactions on Computers*, vol. C-23, pp. 90–93, Jan. 1974.
- [66] A. G. Dempster and M. D. Macleod, "Constant integer multiplication using minimum adders," *IEE Proceedings - Circuits, Devices and Systems*, vol. 141, pp. 407–413, Oct. 1994.
- [67] R. J. Cintra, F. M. Bayer, and C. J. Tablada, "Low-complexity 8-point DCT approximations based on integer functions," *Signal Processing*, vol. 99, pp. 201–214, June 2014.
- [68] F. M. Bayer and R. J. Cintra, "DCT-like transform for image compression requires 14 additions only," *Electronics Letters*, vol. 48, pp. 919–921, July 2012.
- [69] J. Lee and S. Leyffer, *Mixed Integer Nonlinear Programming*. The IMA Volumes in Mathematics and its Applications, Springer New York, 2011.
- [70] C. H. Papadimitriou and K. Steiglitz, *Combinatorial Optimization: Algorithms and Complexity*. Dover Publications, 1998.
- [71] D. Bienstock and G. Nemhauser, *Integer Programming and Combinatorial Optimization*. Lecture Notes in Computer Science, Springer, 2004.
- [72] J. K. Lenstra and A. H. G. R. Kan, "Computational complexity of discrete optimization problems," *Annals of Discrete Mathematics*, vol. 4, pp. 121–140, 1979.
- [73] MATLAB, *version 7.10.0 (R2010a)*. Natick, Massachusetts: The MathWorks Inc., 2010.
- [74] R Core Team, *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2013.
- [75] N. Megiddo, "Linear programming in linear time when the dimension is fixed," *Journal of the ACM (JACM)*, vol. 31, no. 1, pp. 114–127, 1984.
- [76] W. Burger and M. Burge, *Digital Image Processing: An Algorithmic Introduction Using Java*. Texts in Computer Science, Springer London, 2016.
- [77] Y. LeCun, "Generalization and network designs strategies," Technical Report CRG-TR-89-4, Department of Computer Science, 1989.
- [78] Y. LeCun, "Generalization and network design strategies," in *Proceedings of the International Conference Connectionism in Perspective*, (University of Zurich), pp. 10–13, Oct. 1988.
- [79] M. Zhang, S. Vassiliadis, and J. G. Delgado-Frias, "Sigmoid generators for neural computing using piecewise approximations," *IEEE Transactions on Computers*, vol. 45, pp. 1045–1049, Sept. 1996.
- [80] D. Larkin, A. Kinane, V. Muresan, and N. O'Connor, "An efficient hardware architecture for a neural network activation function generator," in *Advances in Neural Networks - ISNN*, pp. 1319–1327, 2006.
- [81] O. Temam, "A defect-tolerant accelerator for emerging high-performance applications," in *Proceedings of the International Symposium on Computer Architecture*, pp. 356–367, 2012.
- [82] J. Schlessman, "Approximation of the sigmoid function and its derivative using a minimax approach," Master's thesis, Lehigh University, Aug. 2002.
- [83] C. Alippi and G. Storti-Gajani, "Simple approximation of sigmoidal functions: realistic design of digital neural networks capable of learning," in *Proc. of IEEE Int. Symp. on Circuits and Systems*, (Singapore), pp. 1505–1508, June 1991.
- [84] H. Amin, K. M. Curtis, and B. R. Hayes-Gill, "Piecewise linear approximation applied to nonlinear function of a neural network," *IEE Proc. Circuits, Devices Sys.*, vol. 144, no. 6, pp. 313–317, 1997.
- [85] R. E. Blahut, *Fast Algorithms for Digital Signal Processing*. Cambridge University Press, 2010.
- [86] U. S. Potluri, A. Madanayake, R. J. Cintra, F. M. Bayer, S. Kulasekera, and A. Edirisuriya, "Improved 8-point approximate DCT for image and video compression requiring only 14 additions," *IEEE Transactions on Circuits and Systems I*, vol. 61, pp. 1727–1740, June 2014.
- [87] M. Mathias, R. Benenson, M. Pedersoli, and L. Van Gool, "Face detection without bells and whistles," in *Proceedings of the European Conference on Computer Vision*, pp. 720–735, 2014.
- [88] D. C. Cireşan, U. Meier, and J. Schmidhuber, "Multi-column deep neural networks for image classification," in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, 2012.
- [89] V. Jain and E. Learned-Miller, "FDDB: A benchmark for face detection in unconstrained settings," Tech. Rep. UM-CS-2010-009, University of Massachusetts, Amherst, 2010.
- [90] X. Zhu and D. Ramanan, "Face detection, pose estimation, and landmark localization in the wild," in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, 2012.
- [91] J. Yan, X. Zhang, Z. Lei, and S. Li, "Face detection by structural models," *Image and Vision Computing*, vol. 32, no. 10, pp. 790–799, 2014.
- [92] Y. LeCun, C. Cortes, and C. J. Burges, "The MNIST database of handwritten digits." <http://yann.lecun.com/exdb/mnist/>, 2015.



Renato J. Cintra (SM'2010) received the B.Sc., M.Sc., and D.Sc. degrees in electrical engineering from the Universidade Federal de Pernambuco (UFPE), Recife, Brazil, in 1999, 2001, and 2005, respectively. He is an associate professor at the *Departamento de Estatística*, UFPE. In 2014–2015 he was a visiting researcher at the INSA, Lyon, France. Currently he is a visiting professor at the University of Calgary, AB, Canada. He is a SIAM member and serves as an Associate Editor for the *IEEE Geoscience and Remote Sensing Letters*; the *Circuits, Systems, and Signal Processing* (Springer); the *IET Circuits, Devices & Systems*; and the *Journal of Communication and Information Systems* (SBrT/Comsoc). His long-term topics of research include theory and methods for digital signal and image processing, numerical analysis, and applied probability.



Stefan Duffner received a Bachelor's degree in Computer Science from the University of Applied Sciences Konstanz, Germany in 2002 and a Master's degree in Applied Computer Science from the University of Freiburg, Germany in 2004. He performed his dissertation research at Orange Labs in Rennes, France, on face image analysis with statistical machine learning methods, notably convolutional neural networks, and in 2008, he obtained a Ph.D. degree in Computer Science from the University of Freiburg. He then worked for 4 years as a post-doctoral researcher at the Idiap Research Institute in Martigny, Switzerland, in the field of computer vision and video analysis for visual object tracking applications. As of today, Stefan Duffner is an associate professor in the IMAGINE team of the LIRIS research lab at INSA Lyon.



Christophe Garcia is Full Professor at INSA Lyon. From 2011 to 2015, he has been head of the IMAGINE research team of the LIRIS laboratory. He is now Deputy Director of the LIRIS laboratory and Vice-President for Research in Information and Digital Society at Insa Lyon. His current technical and research activities are in the areas of ConvNets, neural networks, pattern recognition and computer vision. He holds 17 industrial patents and has published more than 140 articles in international conferences and journals. He has served as an associate

editor of Int. Journal of Visual Communication and Image Representation (Elsevier), Image and Video Processing (Hindawi) and Pattern Analysis and Application (Springer-Verlag).



André Leite received the B.Sc., M.Sc., and D.Sc. degrees in electrical engineering from Universidade Federal de Pernambuco (UFPE), Brazil, in 2003, 2005, and 2009, respectively. Since 2010, he has been with the Department of Statistics, UFPE. His areas of interest include dynamical systems, control theory and optimization, fundamentals of probability and decision theory.