

DianNao: A Small-Footprint High-Throughput Accelerator for Ubiquitous Machine-Learning

Tianshi Chen
SKLCA, ICT, China

Zidong Du
SKLCA, ICT, China

Ninghui Sun
SKLCA, ICT, China

Jia Wang
SKLCA, ICT, China

Chengyong Wu
SKLCA, ICT, China

Yunji Chen
SKLCA, ICT, China

Olivier Temam
Inria, France

Abstract

Machine-Learning tasks are becoming pervasive in a broad range of domains, and in a broad range of systems (from embedded systems to data centers). At the same time, a small set of machine-learning algorithms (especially Convolutional and Deep Neural Networks, i.e., CNNs and DNNs) are proving to be state-of-the-art across many applications. As architectures evolve towards heterogeneous multi-cores composed of a mix of cores and accelerators, a machine-learning accelerator can achieve the rare combination of efficiency (due to the small number of target algorithms) and broad application scope.

Until now, most machine-learning accelerator designs have focused on efficiently implementing the computational part of the algorithms. However, recent state-of-the-art CNNs and DNNs are characterized by their large size. In this study, we design an accelerator for large-scale CNNs and DNNs, with a special emphasis on the impact of memory on accelerator design, performance and energy.

We show that it is possible to design an accelerator with a high **throughput**, capable of performing 452 GOP/s (key NN operations such as **synaptic weight multiplications** and

neurons outputs additions) in a small **footprint** of 3.02 mm² and 485 mW; compared to a 128-bit 2GHz SIMD processor, the accelerator is 117.87x faster, and it can reduce the total energy by 21.08x. The accelerator characteristics are obtained after **layout** at 65nm. Such a high throughput in a small footprint can open up the usage of state-of-the-art machine-learning algorithms in a broad set of systems and for a broad set of applications.

1. Introduction

As architectures evolve towards **heterogeneous** multi-cores composed of a mix of cores and accelerators, designing accelerators which realize the best possible tradeoff between flexibility and efficiency is becoming a **prominent** issue.

The first question is for which category of applications one should primarily design accelerators? Together with the architecture trend towards accelerators, a second **simultaneous** and significant trend in high-performance and embedded applications is developing: many of the emerging high-performance and embedded applications, from image/video/audio recognition to automatic translation, business analytics, and all forms of robotics rely on *machine-learning techniques*. This trend even starts to **percolate** in our community where it turns out that about half of the benchmarks of PARSEC [2], a suite partly introduced to highlight the emergence of new types of applications, can be implemented using machine-learning algorithms [4]. This trend in application comes together with a third and equally remarkable trend in machine-learning where a small number of techniques, based on neural networks (especially Convolutional Neural Networks [27] and Deep Neural Networks

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ASPLOS '14, March 1–5, 2014, Salt Lake City, Utah, USA.
Copyright © 2014 ACM 978-1-4503-2305-5/14/03...\$15.00.
<http://dx.doi.org/10.1145/http://dx.doi.org/10.1145/2541940.2541967>

[16]), have been proved in the past few years to be state-of-the-art across a broad range of applications [25]. As a result, there is a unique opportunity to design accelerators which can realize the best of both worlds: significant application scope together with high performance and efficiency due to the limited number of target algorithms.

Currently, these **workloads** are mostly executed on multi-cores using **SIMD** [41], on GPUs [5], or on **FPGAs** [3]. However, the **aforementioned** trends have already been identified by a number of researchers who have proposed accelerators implementing Convolutional Neural Networks [3] or Multi-Layer **Perceptrons** [38]; accelerators focusing on other domains, such as image processing, also propose efficient implementations of some of the **primitives** used by machine-learning algorithms, such as convolutions [33]. Others have proposed **ASIC** implementations of Convolutional Neural Networks [13], or of other custom neural network algorithms [21]. However, all these works have first, and successfully, focused on efficiently implementing the computational primitives but they either **voluntarily** ignore **memory transfers** for the sake of simplicity [33, 38], or they directly plug their computational accelerator to memory via a more or less **sophisticated** DMA [3, 13, 21].

While efficient implementation of computational primitives is a first and important step with promising results, inefficient memory transfers can potentially **void** the throughput, energy or cost advantages of accelerators, i.e., an Amdahl's law effect, and thus, they should become a first-order concern, just like in processors, rather than an element factored in accelerator design on a second step. Unlike in processors though, one can **factor in** the specific nature of memory transfers in target algorithms, just like it is done for accelerating computations. This is especially important in the domain of machine-learning where there is a clear trend towards **scaling up** the size of neural networks in order to achieve better accuracy and more functionality [16, 26].

In this study, we investigate an accelerator design that can **accommodate** the most popular state-of-the-art algorithms, i.e., Convolutional Neural Networks (CNNs) and Deep Neural Networks (DNNs). We focus the design of the accelerator on memory usage, and we investigate an accelerator architecture and control both to minimize memory transfers and to perform them as efficiently as possible. We present a design at 65nm which can perform 496 16-bit fixed-point operations in parallel every 1.02ns, i.e., 452 GOP/s, in a 3.02mm², 485mW footprint (excluding **main memory accesses**). On 10 of the largest layers found in recent CNNs and DNNs, this accelerator is 117.87x faster and 21.08x more energy-efficient (including main memory accesses) on average than an 128-bit SIMD core clocked at 2GHz.

In summary, our main contributions are the following:

- A **synthesized** (place & route) accelerator design for large-scale CNNs and DNNs, the state-of-the-art machine-learning algorithms.

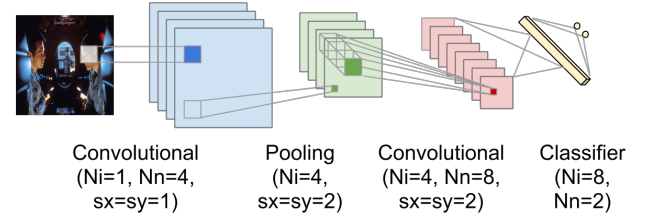


Figure 1. Neural network hierarchy containing convolutional, pooling and classifier layers.

- The accelerator achieves high throughput in a small area, power and energy footprint.
- The accelerator design focuses on memory behavior, and measurements are not **circumscribed** to computational tasks, they factor in the performance and energy impact of memory transfers.

The paper is organized as follows. In Section 2, we first provide a primer on recent machine-learning techniques and introduce the main layers composing CNNs and DNNs. In Section 3, we analyze and optimize the memory behavior of these layers, in preparation for both the baseline and the accelerator design. In section 4, we explain why an ASIC implementation of large-scale CNNs or DNNs cannot be the same as the straightforward ASIC implementation of small NNs. We introduce our accelerator design in Section 5. The methodology is presented in Section 6, the experimental results in Section 7, related work in Section 8.

2. Primer on Recent Machine-Learning Techniques

Even though the role of neural networks in the machine-learning domain has been **rocky**, i.e., initially **hyped** in the 1980s/1990s, then fading into **oblivion** with the **advent** of **Support Vector Machines** [6]. Since 2006, a subset of neural networks have emerged as achieving state-of-the-art machine-learning accuracy across a broad set of applications, partly inspired by progress in neuroscience models of computer vision, such as HMAX [37]. This subset of neural networks includes both Deep Neural Networks (DNNs) [25] and Convolutional Neural Networks (CNNs) [27]. DNNs and CNNs are strongly related, they especially differ in the presence and/or nature of convolutional layers, see later.

Processing vs. training. For now, we have implemented the fast processing of inputs (feed-forward) rather than training (backward path) on our accelerator. This derives from technical and market considerations. Technically, there is a frequent and important **misconception** that *on-line* learning is necessary for many applications. On the contrary, for many industrial applications *off-line* learning is sufficient, where the neural network is first trained on a set of data, and then shipped to the customer, e.g., trained on hand-

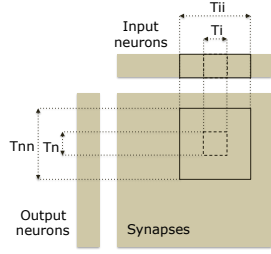


Figure 2. Classifier layer tiling.

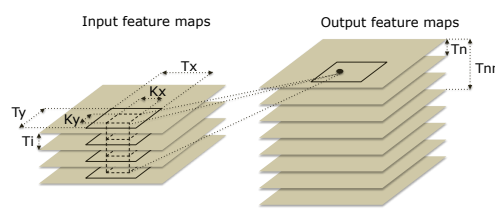


Figure 3. Convolutional layer tiling.

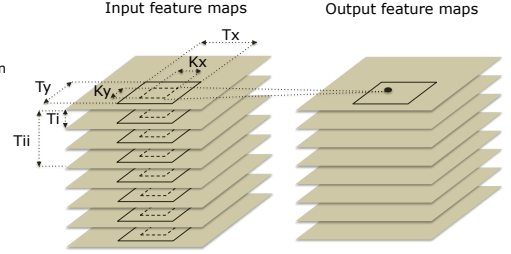


Figure 4. Pooling layer tiling.

written digits, license plate numbers, a number of faces or objects to recognize, etc; the network can be periodically taken off-line and retrained. While, today, machine-learning researchers and engineers would especially want an architecture that speeds up training, this represents a small market, and for now, we focus on the much larger market of *end users*, who need fast/efficient feed-forward networks. Interestingly, machine-learning researchers who have recently dipped into hardware accelerators [13] have made the same choice. Still, because the nature of computations and access patterns used in training (especially back-propagation) is fairly similar to that of the forward path, we plan to later augment the accelerator with the necessary features to support training.

General structure. Even though Deep and Convolutional Neural Networks come in various forms, they share enough properties that a generic formulation can be defined. In general, these algorithms are made of a (possibly large) number of *layers*; these layers are executed *in sequence* so they can be considered (and optimized) independently. Each layer usually contains several sub-layers called *feature maps*; we then use the terms *input feature maps* and *output feature maps*. Overall, there are three main kinds of layers: most of the **hierarchy** is composed of convolutional and pooling (also called sub-sampling) layers, and there is a classifier at the top of the network made of one or a few layers.

Convolutional layers. The role of convolutional layers is to apply one or several *local* filters to data from the input (previous) layer. Thus, the connectivity between the input and output feature map is local instead of full. Consider the case where the input is an image, the convolution is a 2D transform between a $K_x \times K_y$ subset (window) of the input layer and a kernel of the same dimensions, see Figure 1. The kernel values are the synaptic weights between an input layer and an output (convolutional) layer. Since an input layer usually contains several input feature maps, and since an output feature map point is usually obtained by applying a convolution to the same window of *all* input feature maps, see Figure 1, the kernel is 3D, i.e., $K_x \times K_y \times N_i$, where N_i is the number of input feature maps. Note that in some cases, the connectivity is **sparse**, i.e., not all input feature maps are used for each output feature map. The typical code of a convolutional layer is shown in Figure 7, see *Original*

code. A non-linear function is applied to the convolution output, for instance $f(x) = \tanh(x)$. Convolutional layers are also characterized by the overlap between two consecutive windows (in one or two dimensions), see steps s_x, s_y for loops x, y .

In some cases, the same kernel is applied to all $K_x \times K_y$ windows of the input layer, i.e., weights are implicitly *shared* across the whole input feature map. This is characteristic of CNNs, while kernels can be specific to each point of the output feature map in DNNs [26], we then use the term *private* kernels.

Pooling layers. The role of pooling layers is to aggregate information among a set of neighbor input data. In the case of images again, it serves to retain only the salient features of an image within a given window and/or to do so at different scales, see Figure 1. An important side effect of pooling layers is to reduce the feature map dimensions. An example code of a pooling layer is shown in Figure 8 (see *Original code*). Note that each feature map is pooled separately, i.e., 2D pooling, not 3D pooling. Pooling can be done in various ways, some of the preferred techniques are the *average* and *max* operations; pooling may or may not be followed by a non-linear function.

Classifier layers. Convolution and pooling layers are interleaved within deep hierarchies, and the top of the hierarchies is usually a *classifier*. This classifier can be linear or a multi-layer (often 2-layer) perceptron, see Figure 1. An example perceptron layer is shown in Figure 5, see *Original code*. Like convolutional layers, a non-linear function is applied to the neurons output, often a sigmoid, e.g., $f(x) = \frac{1}{1+e^{-x}}$; unlike convolutional or pooling layers, classifiers usually **aggregate** (flatten) all feature maps, so there is no notion of feature maps in classifier layers.

3. Processor-Based Implementation of (Large) Neural Networks

The distinctive aspect of accelerating large-scale neural networks is the potentially high memory traffic. In this section, we analyze in details the locality properties of the different layers mentioned in Section 2, we **tune** processor-based implementations of these layers in preparation for both our baseline, and the design and utilization of the accelerator. We apply the locality analysis/optimization to all layers, and we

illustrate the bandwidth impact of these transformations with 4 of our benchmark layers (CLASS1, CONV3, CONV5, POOL3); their characteristics are later detailed in Section 6.

For the memory bandwidth measurements of this section, we use a cache simulator plugged to a virtual computational structure on which we make no assumption except that it is capable of processing T_n neurons with T_i synapses each every cycle. The cache hierarchy is inspired by Intel Core i7: L1 is 32KB, 64-byte line, 8-way; the optional L2 is 2MB, 64-byte, 8-way. Unlike the Core i7, we assume the caches have enough banks/ports to serve $T_n \times 4$ bytes for input neurons, and $T_n \times T_i \times 4$ bytes for synapses. For large T_n, T_i , the cost of such caches can be **prohibitive**, but it is only used for our limit study of locality and bandwidth; in our experiments, we use $T_n = T_i = 16$.

3.1 Classifier Layers

```
for (int nnn = 0; nnn < Nn; nnn += Tnn) { // tiling for output neurons;
  for (int iii = 0; iii < Ni; iii += Tii) { // tiling for input neurons;
    for (int nn = nnn; nn < nnn + Tnn; nn += Tn) {
      for (int n = nn; n < nn + Tn; n++)
        sum[n] = 0;
      for (int ii = iii; ii < iii + Tii; ii += Ti)
        // — Original code —
        for (int n = nn; n < nn + Tn; n++)
          for (int i = ii; i < ii + Ti; i++)
            sum[n] += synapse[n][i] * neuron[i];
      for (int n = nn; n < nn + Tn; n++)
        neuron[n] = sigmoid(sum[n]);
    }
  }
}
```

Figure 5. Pseudo-code for a classifier (here, perceptron) layer (original loop nest + locality optimization).

We consider the perceptron classifier layer, see Figures 2 and 5; the tiling loops ii and nn simply reflect that the computational structure can process T_n neurons with T_i synapses simultaneously. The total number of memory transfers is (inputs loaded + synapses loaded + outputs written): $N_i \times N_n + N_i \times N_n + N_n$. For the example layer CLASS1, the corresponding memory bandwidth is high at 120 GB/s, see CLASS1 - Original in Figure 6. We explain below how it is possible to reduce this bandwidth, sometimes drastically.

Input/Output neurons. Consider Figure 2 and the code of Figure 5 again. Input neurons are reused for each output neuron, but since the number of input neurons can range anywhere between a few tens to hundreds of thousands, they will often not fit in an L1 cache. Therefore, we tile loop ii (input neurons) with tile factor T_{ii} . A typical tradeoff of tiling is that improving one reference (here $neuron[i]$ for input neurons) increases the reuse distance of another reference ($sum[n]$ for partial sums of output neurons), so we need to tile for the second reference as well, hence loop nnn and the tile factor T_{nn} for output neurons partial sums. As expected, tiling drastically reduces the memory bandwidth requirements of input neurons, and those of output neurons increase, **albeit marginally**. The layer memory behavior is now dominated by synapses.

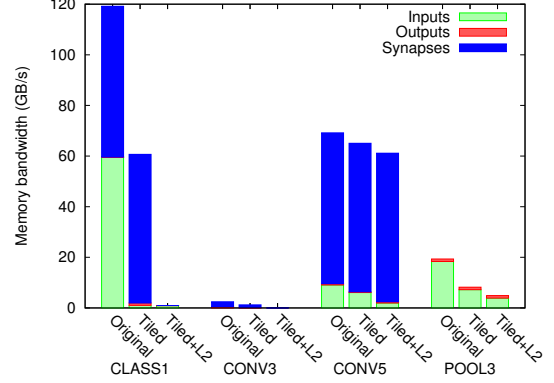


Figure 6. Memory bandwidth requirements for each layer type (CONV3 has shared kernels, CONV5 has private kernels).

Synapses. In a perceptron layer, all synapses are usually unique, and thus there is no reuse *within the layer*. On the other hand, the synapses are reused *across network invocations*, i.e., for each new input data (also called “input row”) presented to the neural network. So a sufficiently large L2 could store all network synapses and take advantage of that locality. For DNNs with private kernels, this is not possible as the total number of synapses are in the tens or hundreds of millions (the largest network to date has a billion synapses [26]). However, for both CNNs and DNNs with shared kernels, the total number of synapses range in the millions, which is within the reach of an L2 cache. In Figure 6, see CLASS1 - Tiled+L2, we **emulate** the case where reuse across network invocations is possible by considering only the perceptron layer; as a result, the total bandwidth requirements are now drastically reduced.

3.2 Convolutional Layers

We consider two-dimensional convolutional layers, see Figures 3 and 7. The two distinctive features of convolutional layers with respect to classifier layers are the presence of input and output feature maps (loops i and n) and kernels (loops k_x, k_y).

Inputs/Outputs. There are two types of reuse opportunities for inputs and outputs: the *sliding window* used to scan the (two-dimensional (x, y)) input layer, and the reuse across the N_n output feature maps, see Figure 3. The former corresponds to $\frac{K_x \times K_y}{s_x \times s_y}$ reuses at most, and the latter to N_n reuses. We tile for the former in Figure 7 (tiles T_x, T_y), but we often do not need to tile for the latter because the data to be reused, i.e., one kernel of $K_x \times K_y \times N_i$, fits in the L1 data cache since K_x, K_y are usually of the order of 10 and N_i can vary between less than 10 to a few hundreds; naturally, when this is not the case, we can tile input feature maps (ii) and introduce an second-level tiling loop iii again.

Synapses. For convolutional layers with *shared* kernels (see Section 2), the same kernel parameters (synaptic weights) are reused across all x_{out}, y_{out} output feature maps locations. As a result, the total bandwidth is already


```

for (int yy = 0; yy < Nyin; yy += Ty) {
  for (int xx = 0; xx < Nxin; xx += Tx) {
    for (int nnn = 0; nnn < Nn; nnn += Tnn) {
      // — Original code — (excluding nn, ii loops)
      int yout = 0;
      for (int y = yy; y < yy + Ty; y += sy) { // tiling for y;
        int xout = 0;
        for (int x = xx; x < xx + Tx; x += sx) { // tiling for x;
          for (int nn = nnn; nn < nnn + Tnn; nn += Tn) {
            for (int n = nn; n < nn + Tn; n++)
              sum[n] = 0;
            // sliding window;
            for (int ky = 0; ky < Ky; ky++)
              for (int kx = 0; kx < Kx; kx++)
                for (int ii = 0; ii < Ni; ii += Ti)
                  for (int n = nn; n < nn + Tn; n++)
                    for (int i = ii; i < ii + Ti; i++)
                      // version with shared kernels
                      sum[n] += synapse[ky][kx][n][i]
                               * neuron[ky + y][kx + x][i];
                      // version with private kernels
                      sum[n] += synapse[yout][xout][ky][kx][n][i]
                               * neuron[ky + y][kx + x][i];
            for (int n = nn; n < nn + Tn; n++)
              neuron[yout][xout][n] = non_linear_transform(sum[n]);
          } xout++; } yout++;
        } } }
      } } }

```

Figure 7. Pseudo-code for convolutional layer (original loop nest + locality optimization), both shared and private kernels versions.

low, as shown for layer CONV3 in Figure 6. However, since the total shared kernels capacity is $K_x \times K_y \times N_i \times N_o$, it can exceed the L1 cache capacity, so we tile again output feature maps (tile T_{nn}) to bring it down to $K_x \times K_y \times N_i \times T_{nn}$. As a result, the overall memory bandwidth can be further reduced, as shown in Figure 6.

For convolutional layers with *private* kernels, the synapses are all unique and there is no reuse, as for classifier layers, hence the similar synapses bandwidth of CONV5 in Figure 6. As for classifier layers, reuse is still possible across network invocations if the L2 capacity is sufficient. Even though step coefficients (s_x, s_y) and sparse input to output feature maps (see Section 2) can drastically reduce the number of private kernels synaptic weights, for very large layers such as CONV5, they still range in the hundreds of megabytes and thus will largely exceed L2 capacity, implying a high memory bandwidth, see Figure 6.

It is important to note that there is an on-going debate within the machine-learning community about shared vs. private kernels [26, 35], and the machine-learning importance of having private instead of shared kernels remains unclear. Since they can result in significantly different architecture performance, this may be a case where the architecture/performance community could weigh in on the machine-learning debate.

3.3 Pooling Layers

We now consider pooling layers, see Figures 4 and 8. Unlike convolutional layers, the number of input and output feature

```

for (int yy = 0; yy < Nyin; yy += Ty) {
  for (int xx = 0; xx < Nxin; xx += Tx) {
    for (int iii = 0; iii < Ni; iii += Tii)
      // — Original code — (excluding ii loop)
      int yout = 0;
      for (int y = yy; y < yy + Ty; y += sy) {
        int xout = 0;
        for (int x = xx; x < xx + Tx; x += sx) {
          for (int ii = iii; ii < iii + Tii; ii += Ti)
            for (int i = ii; i < ii + Ti; i++)
              value[i] = 0;
          for (int ky = 0; ky < Ky; ky++)
            for (int kx = 0; kx < Kx; kx++)
              for (int i = ii; i < ii + Ti; i++)
                // version with average pooling;
                value[i] += neuron[ky + y][kx + x][i];
                // version with max pooling;
                value[i] = max(value[i], neuron[ky + y][kx + x][i]);
          } } }
        // for average pooling;
        neuron[xout][yout][i] = value[i] / (Kx * Ky);
        xout++; } yout++;
      } } }

```

Figure 8. Pseudo-code for pooling layer (original loop nest + locality optimization).

maps is the same, and more importantly, there is no kernel, i.e., no synaptic weight to store, and an output feature map element is determined only by $K_x \times K_y$ input feature map elements, i.e., a 2D window (instead of a 3D window for convolutional layers). As a result, the only source of reuse comes from the sliding window (instead of the combined effect of sliding window and output feature maps). Since there are less reuse opportunities, the memory bandwidth of input neurons are higher than for convolutional layers, and tiling (T_x, T_y) brings less dramatic improvements, see Figure 6.

4. Accelerator for Small Neural Networks

In this section, we first evaluate a “naive” and greedy approach for implementing a hardware neural network accelerator where all neurons and synapses are laid out in hardware, memory is only used for input rows and storing results. While these neural networks can potentially achieve the best energy efficiency, we show that they are not scalable. Still, we use such networks to investigate the maximum number of neurons which can be reasonably implemented in hardware.

4.1 Hardware Neural Networks

The most natural way to map a neural network onto silicon is simply to fully lay out the neurons and synapses, so that the hardware implementation matches the conceptual representation of neural networks, see Figure 9. The neurons are each implemented as logic circuits, and the synapses are implemented as **latches** or **RAMs**. This approach has been recently used for perceptron or spike-based hardware neural networks [30, 38]. It is **compatible** with some embedded applications where the number of neurons and synapses can be small, and it can provide both high speed and low energy because the distance traveled by data is very small: from one

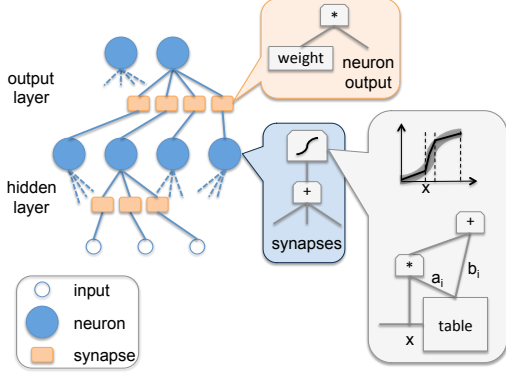


Figure 9. Full hardware implementation of neural networks.

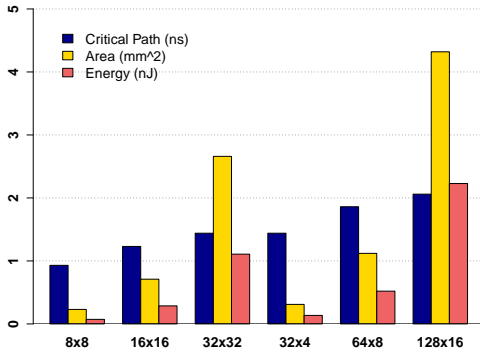


Figure 10. Energy, critical path and area of full-hardware layers.

neuron to a neuron of the next layer, and from one synaptic latch to the associated neuron. For instance, an execution time of 15ns and an energy reduction of 974x over a core has been reported for a 90-10-10 (90 inputs, 10 hidden, 10 outputs) perceptron [38].

4.2 Maximum Number of Hardware Neurons ?

However, the area, energy and delay grow quadratically with the number of neurons. We have synthesized the ASIC versions of neural network layers of various dimensions, and we report their area, critical path and energy in Figure 10. We have used Synopsys ICC for the place and route, and the TSMC 65nm GP library, standard VT. A hardware neuron performs the following operations: multiplication of inputs and synapses, addition of all such multiplications, followed by a sigmoid, see Figure 9. A $T_n \times T_i$ layer is a layer of T_n neurons with T_i synapses each. A 16x16 layer requires less than 0.71 mm², but a 32x32 layer already costs 2.66 mm². Considering the neurons are in the thousands for large-scale neural networks, a full hardware layout of just one layer would range in the hundreds or thousands of mm², and thus, this approach is not realistic for large-scale neural networks.

For such neural networks, only a fraction of neurons and synapses can be implemented in hardware. Paradoxically, this was already the case for old neural network designs

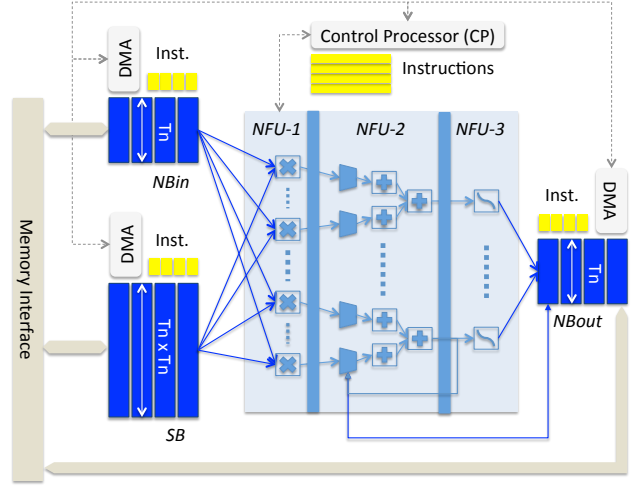


Figure 11. Accelerator.

such as the Intel ETANN [18] at the beginning of the 1990s, not because neural networks were already large at the time, but because hardware resources (number of transistors) were naturally much more scarce. The principle was to time-share the physical neurons and use the on-chip RAM to store synapses and intermediate neurons values of hidden layers. However, at that time, many neural networks were small enough that all synapses and intermediate neurons values could fit in the neural network RAM. Since this is no longer the case, one of the main challenges for large-scale neural network accelerator design has become the **interplay** between the computational and the memory hierarchy.

5. Accelerator for Large Neural Networks

In this section, we draw from the analysis of Sections 3 and 4 to design an accelerator for large-scale neural networks.

The main components of the accelerator are the following: an input buffer for input neurons (NBin), an output buffer for output neurons (NBout), and a third buffer for synaptic weights (SB), connected to a computational block (performing both synapses and neurons computations) which we call the Neural Functional Unit (NFU), and the control logic (CP), see Figure 11. We first describe the NFU below, and then we focus on and explain the **rationale** for the storage elements of the accelerator.

5.1 Computations: Neural Functional Unit (NFU)

The spirit of the NFU is to reflect the **decomposition** of a layer into computational blocks of T_i inputs/synapses and T_n output neurons. This corresponds to loops i and n for both classifier and convolutional layers, see Figures 5 and Figure 7, and loop i for pooling layers, see Figure 8.

Arithmetic operators. The computations of each layer type can be decomposed in either 2 or 3 stages. For classifier layers: multiplication of synapses \times inputs, additions of all

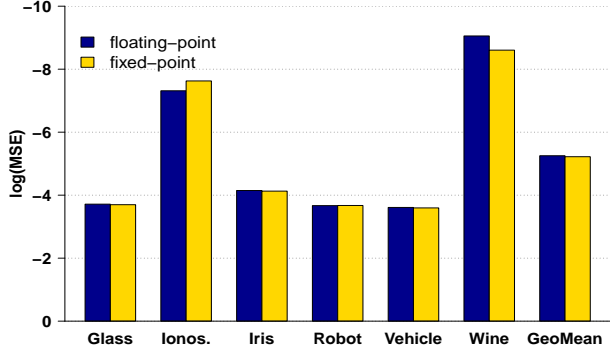


Figure 12. 32-bit floating-point vs. 16-bit fixed-point accuracy for UCI data sets (metric: $\log(\text{Mean Squared Error})$).

Type	Error Rate
32-bit floating-point	0.0311
16-bit fixed-point	0.0337

Table 1. 32-bit floating-point vs. 16-bit fixed-point accuracy for MNIST (metric: error rate).

multiplications, sigmoid. For convolutional layers, the stages are the same; the nature of the last stage (sigmoid or another non-linear function) can vary. For pooling layers, there is no multiplication (no synapse), and the pooling operations can be average or max. Note that the adders have multiple inputs, they are in fact *adder trees*, see Figure 11; the second stage also contains *shifters* and max operators for pooling layers.

Staggered pipeline. We can pipeline all 2 or 3 operations, but the pipeline must be staggered: the first or first two stages (respectively for pooling, and for classifier and convolutional layers) operate as normal pipeline stages, but the third stage is only active after all additions have been performed (for classifier and convolutional layers; for pooling layers, there is no operation in the third stage). From now on, we refer to stage n of the NFU pipeline as NFU- n .

NFU-3 function implementation. As previously proposed in the literature [23, 38], the sigmoid of NFU-3 (for classifier and convolutional layers) can be efficiently implemented using *piecewise linear interpolation* ($f(x) = a_i \times x + b_i, x \in [x_i, x_{i+1}]$) with *negligible* loss of accuracy (16 segments are sufficient) [24], see Figure 9. In terms of operators, it corresponds to two 16×1 16-bit multiplexers (for segment boundaries selection, i.e., x_i, x_{i+1}), one 16-bit multiplier (16-bit output) and one 16-bit adder to perform the interpolation. The 16-segment coefficients (a_i, b_i) are stored in a small RAM; this allows to implement *any* function, not just a sigmoid (e.g., hyperbolic tangent, linear functions, etc) by just changing the RAM segment coefficients a_i, b_i ; the segment boundaries (x_i, x_{i+1}) are *hardwired*.

16-bit fixed-point arithmetic operators. We use 16-bit fixed-point arithmetic operators instead of word-size (e.g., 32-bit) floating-point operators. While it may seem surpris-

Type	Area (μm^2)	Power (μW)
16-bit truncated fixed-point multiplier	1309.32	576.90
32-bit floating-point multiplier	7997.76	4229.60

Table 2. Characteristics of multipliers.

ing, there is *ample* evidence in the literature that even smaller operators (e.g., 8 bits or even less) have almost no impact on the accuracy of neural networks [8, 17, 24]. To illustrate and further confirm that notion, we trained and tested multi-layer perceptrons on data sets from the UC Irvine Machine-Learning repository, see Figure 12, and on the standard MNIST machine-learning benchmark (handwritten digits) [27], see Table 1, using both 16-bit fixed-point and 32-bit floating-point operators; we used 10-fold cross-validation for testing. For the fixed-point operators, we use 6 bits for the integer part, 10 bits for the fractional part (we use this fixed-point configuration throughout the paper). The results are shown in Figure 12 and confirm the very small accuracy impact of that tradeoff. We *conservatively* use 16-bit fixed-point for now, but we will explore smaller, or variable-size, operators in the future. Note that the arithmetic operators are *truncated*, i.e., their output is 16 bits; we use a standard n -bit *truncated multiplier* with *correction constant* [22]. As shown in Table 2, its area is 6.10x smaller and its power 7.33x lower than a 32-bit floating-point multiplier at 65nm, see Section 6 for the CAD tools methodology.

5.2 Storage: NBin, NBout, SB and NFU-2 Registers

The different storage structures of the accelerator can be *construed* as modified buffers of scratchpads. While a cache is an excellent storage structure for a general-purpose processor, it is a sub-optimal way to exploit reuse because of the cache access overhead (tag check, associativity, line size, speculative read, etc) and cache conflicts [39]. The efficient alternative, scratchpad, is used in *VLIW* processors but it is known to be very difficult to compile for. However a scratchpad in a dedicated accelerator realizes the best of both worlds: efficient storage, and both efficient and easy exploitation of locality because only a few algorithms have to be manually adapted. In this case, we can almost directly translate the locality transformations introduced in Section 3 into mapping commands for the buffers, mostly *modulating* the tile factors. A code mapping example is provided in Section 5.3.2

We explain below how the storage part of the accelerator is organized, and which limitations of cache architectures it overcomes.

5.2.1 Split buffers.

As explained before, we have split storage into three structures: an input buffer (NBin), an output buffer (NBout) and a synapse buffer (SB).

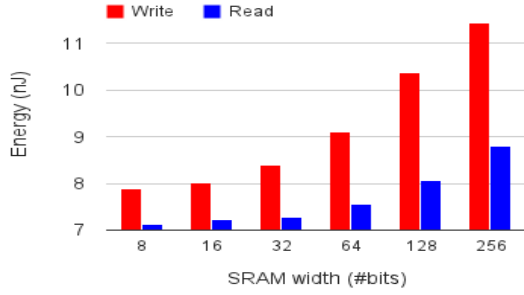


Figure 13. Read energy vs. SRAM width.

Width. The first benefit of splitting structures is to tailor the SRAMs to the appropriate read/write width. The width of both NBin and NBout is $T_n \times 2$ bytes, while the width of SB is $T_n \times T_n \times 2$ bytes. A single read width size, e.g., as with a cache line size, would be a poor tradeoff. If it's adjusted to synapses, i.e., if the line size is $T_n \times T_n \times 2$, then there is a significant energy penalty for reading $T_n \times 2$ bytes out of a $T_n \times T_n \times 2$ -wide data bank, see Figure 13 which indicates the SRAM read energy as a function of bank width for the TSMC process at 65nm. If the line size is adjusted to neurons, i.e., if the line size is $T_n \times 2$, there is a significant time penalty for reading $T_n \times T_n \times 2$ bytes out. Splitting storage into dedicated structures allows to achieve the best time and energy for each read request.

Conflicts. The second benefit of splitting storage structures is to avoid conflicts, as would occur in a cache. It is especially important as we want to keep the size of the storage structures small for cost and energy (leakage) reasons. The alternative solution is to use a highly associative cache. Consider the constraints: the cache line (or the number of ports) needs to be large ($T_n \times T_n \times 2$) in order to serve the synapses at a high rate; since we would want to keep the cache size small, the only alternative to tolerate such a long cache line is high associativity. However, in an n -way cache, a fast read is implemented by speculatively reading all n ways/banks in parallel; as a result, the energy cost of an associative cache increases quickly. Even a 64-byte read from an 8-way associative 32KB cache costs 3.15x more energy than a 32-byte read from a direct mapped cache, at 65nm; measurements done using CACTI [40]. And even with a 64-byte line only, the first-level 32KB data cache of Core i7 is already 8-way associative, so we need an even larger associativity with a very large line (for $T_n = 16$, the line size would be 512-byte long). In other words, a highly associative cache would be a costly energy solution in our case. Split storage and precise knowledge of locality behavior allows to entirely remove data conflicts.

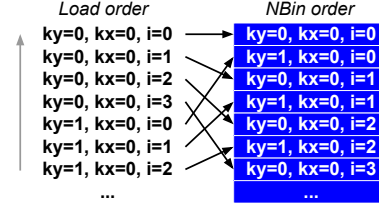


Figure 14. Local transpose ($K_y = 2, K_x = 1, N_i = 4$).

5.2.2 Exploiting the locality of inputs and synapses.

DMAs. For spatial locality exploitation, we implement three DMAs, one for each buffer (two load DMAs, one store DMA for outputs). DMA requests are issued to NBin in the form of instructions, later described in Section 5.3.2. These requests are buffered in a separate FIFO associated with each buffer, see Figure 11, and they are issued as soon as the DMA has sent all the memory requests for the previous instruction. These DMA requests FIFOs enable to *decouple* the requests issued to all buffers and the NFU from the current buffer and NFU operations. As a result, DMA requests can be preloaded far in advance for tolerating long latencies, as long as there is enough buffer capacity; this *preloading* is akin to prefetching, albeit without speculation. Due to the combined role of NBin (and SB) as both scratchpads for reuse and preload buffers, we use a dual-port SRAM; the TSMC 65nm library rates the read energy overhead of dual port SRAMs for a 64-entry NB at 24%.

Rotating NBin buffer for temporal reuse of input neurons. The inputs of all layers are split into chunks which fit in NBin, and they are reused by implementing NBin as a *circular buffer*. In practice, the rotation is naturally implemented by changing a register index, much like in a software implementation, there is no physical (and costly) movement of buffer entries.

Local transpose in NBin for pooling layers. There is a tension between convolutional and pooling layers for the data structure organization of (input) neurons. As mentioned before, K_x, K_y are usually small (often less than 10), and N_i is about an order of magnitude larger. So memory fetches are more efficient (long stride-1 accesses) with the input feature maps as the innermost index of the three-dimensional neurons data structure. However, this is inconvenient for pooling layers because one output is computed *per input feature map*, i.e., using only $K_x \times K_y$ data (while in convolutional layers, all $K_x \times K_y \times N_i$ data are required to compute one output data). As a result, for pooling layers, the logical data structure organization is to have k_x, k_y as the innermost dimensions so that all inputs required to compute one output are consecutively stored in the NBin buffer. We resolve this tension by introducing a mapping function in NBin which has the effect of *locally transposing* loops k_y, k_x and loop i so that data is loaded along loop i , but it is stored in NBin and thus sent to NFU along loops k_y, k_x first; this is accom-

published by interleaving the data in NBin as it is loaded, see Figure 14.

For synapses and SB, as mentioned in Section 3, there is either no reuse (classifier layers, convolutional layers with private kernels and pooling layers), or reuse of shared kernels in convolutional layers. For outputs and NBout, we need to reuse the partial sums, i.e., see reference $sum[n]$ in Figure 5. This reuse requires additional hardware modifications explained in the next section.

5.2.3 Exploiting the locality of outputs.

In both classifier and convolutional layers, the partial output sum of T_n output neurons is computed for a chunk of input neurons contained in NBin. Then, the input neurons are used for another chunk of T_n output neurons, etc. This creates two issues.

Dedicated registers. First, while the chunk of input neurons is loaded from NBin and used to compute the partial sums, it would be inefficient to let the partial sum exit the NFU pipeline and then re-load it into the pipeline for each entry of the NBin buffer, since data transfers are a major source of energy expense [14]. So we introduce *dedicated registers* in NFU-2, which store the partial sums.

Circular buffer. Second, a more complicated issue is what to do with the T_n partial sums when the input neurons in NBin are reused for a new set of T_n output neurons. Instead of sending these T_n partial sums back to memory (and to later reload them when the next chunk of input neurons is loaded into NBin), we temporarily rotate them out to NBout. A priori, this is a conflicting role for NBout which is also used to store the final output neurons to be written back to memory (write buffer). In practice though, as long as all input neurons have not been integrated in the partial sums, NBout is idle. So we can use it as a temporary storage buffer by rotating the T_n partial sums out to NBout, see Figure 11. Naturally, the loop iterating over output neurons must be tiled so that no more output neurons are computing their partial sums simultaneously than the capacity of NBout, but that is implemented through a second-level tiling similar to loop *nnn* in Figure 5 and Figure 7. As a result, NBout is not only connected to NFU-3 and memory, but also to NFU-2: one entry of NBout can be loaded into the dedicated registers of NFU-2, and these registers can be stored in NBout.

5.3 Control and Code

5.3.1 CP.

In this section, we describe the control of the accelerator. One approach to control would be to hardwire the three target layers. While this remains an option for the future, for now, we have decided to use *control instructions* in order to explore different implementations (e.g., partitioning and scheduling) of layers, and to provide machine-learning researchers with the flexibility to try out different layer implementations.

CP	SB			NBin			NBout			NFU								
END	READ OP	REUSE	ADDRESS	SIZE	READ OP	REUSE	STRIDE	STRIDE BEGIN	STRIDE END	ADDRESS	SIZE	NFU-1 OP	NFU-2 OP	NFU-2 IN	NFU-2 OUT	NFU-3 OP	OUTPUT BEGIN	OUTPUT END

Table 3. Control instruction format.

[illegible]

Table 4. *Subset of classifier/perceptron code ($N_i = 8192$, $N_o = 256$, $T_n = 16$, 64-entry buffers).*

A layer execution is broken down into a set of instructions. Roughly, one instruction corresponds to the loops ii, i, n for classifier and convolutional layers, see Figures 5 and 7, and to the loops ii, i in pooling layers (using the interleaving mechanism described in Section 5.2.3), see Figure 8. The instructions are stored in an SRAM associated with the *Control Processor* (CP), see Figure 11. The CP drives the execution of the DMAs of the three buffers and the NFU. The term “processor” only relates to the aforementioned “instructions”, later described in Section 5.3.2, but it has very few of the traditional features of a processor (mostly a PC and an adder for loop index and address computations); from a hardware perspective, it is more like a configurable FSM.

5.3.2 Layer Code.

Every instruction has five slots, corresponding to the CP itself, the three buffers and the NFU, see Table 3.

Because of the CP instructions, there is a need for code generation, but a compiler would be overkill in our case as only three main types of codes must be generated. So we have implemented three dedicated code generators for the three layers. In Table 4, we give an example of the code generated for a classifier/perceptron layer. Since $T_n = 16$ (16×16 -bit data per buffer row) and NBin has 64 rows, its capacity is 2KB, so it cannot contain all the input neurons ($N_i = 8192$, so 16KB). As a result, the code is broken down to operate on **chunks** of 2KB; note that the first instruction of NBin is a **LOAD** (data fetched from memory), and that it is marked as *reused* (flag immediately after load); the next instruction is a *read*, because these input neurons are rotated in the buffer for the next chunk of T_n neurons, and the read is also marked as *reused* because there are 8 such rotations ($\frac{16KB}{2KB}$); at the same time, notice that the output of NFU-2 for the first (and next) instruction is NBout, i.e., the partial output neurons sums are rotated to NBout, as explained in Section 5.2.3, which is why the NBout instruction is

WRITE; notice also that the input of NFU-2 is RESET (first chunk of input neurons, registers reset). Finally, when the last chunk of input neurons are sent (last instruction in table), the (store) DMA of NBout is set for writing 512 bytes (256 outputs), and the NBout instruction is STORE; the NBout write operation for the next instructions will be NOP (DMA set at first chunk and automatically storing data back to memory until DMA elapses).

Note that the architecture can implement either per-image or batch processing [41], only the generated layer control code would change.

6. Experimental Methodology

Measurements. We use three different tools for performance/energy measurements.

Accelerator simulator. We implemented a custom cycle-accurate, bit-accurate C++ simulator of the accelerator fabric, which was initially used for architecture exploration, and which later served as the specification for the Verilog implementation. This simulator is also used to measure time in number of cycles. It is plugged to a main memory model allowing a bandwidth of up to 250 GB/s.

CAD tools. For area, energy and critical path delay (cycle time) measurements, we implemented a Verilog version of the accelerator, which we first synthesized using the Synopsys Design Compiler using the TSMC 65nm GP standard VT library, and which we then placed and routed using the Synopsys ICC compiler. We then simulated the design using Synopsys VCS, and we estimated the power using PrimeTime PX.

SIMD. For the SIMD baseline, we use the GEM5+McPAT [28] combination. We use a 4-issue superscalar x86 core with a 128-bit (8×16 -bit) SIMD unit (SSE/SSE2), clocked at 2GHz. The core has a 192-entry ROB, and a 64-entry load/store queue. The L1 data (and instruction) cache is 32KB and the L2 cache is 2MB; both caches are 8-way associative and use a 64-byte line; these cache characteristics correspond to those of the Intel Core i7. The L1 miss latency to the L2 is 10 cycles, and the L2 miss latency to memory is 250 cycles; the memory bus width is 256 bits. We have aligned the energy cost of main memory accesses of our accelerator and the simulator by using those provided by McPAT (e.g., 17.6nJ for a 256-bit read memory access).

We implemented a SIMD version of the different layer codes, which we manually tuned for locality as explained in Section 3 (for each layer, we perform a stochastic exploration to find good tile factors); we compiled these programs using the default -O optimization level but the inner loops were written in assembly to make the best possible use of the SIMD unit. In order to assess the performance of the SIMD core, we also implemented a standard C++ version of the different benchmark layers presented below, and on average (geometric mean), we observed that the SIMD core provides

Layer	N_x	N_y	K_x	K_y	N_i	N_o	Description
CONV1	500	375	9	9	32	48	Street scene parsing (CNN) [13], (e.g., identifying “building”, “vehicle”, etc)
POOL1	492	367	2	2	12	-	
CLASS1	-	-	-	-	960	20	
CONV2*	200	200	18	18	8	8	Detection of faces in YouTube videos (DNN) [26], largest NN to date (Google)
CONV3	32	32	4	4	108	200	Traffic sign identification for car navigation (CNN) [36]
POOL3	32	32	4	4	100	-	
CLASS3	-	-	-	-	200	100	
CONV4	32	32	7	7	16	512	Google Street View house numbers (CNN) [35]
CONV5*	256	256	11	11	256	384	Multi-Object recognition in natural images (DNN) [16], winner 2012 ImageNet competition
POOL5	256	256	2	2	256	-	

Table 5. Benchmark layers (CONV=convolutional, POOL=pooling, CLASS=classifier; CONVx* indicates private kernels).

a 3.92x improvement in execution time and 3.74x in energy over the x86 core.

Benchmarks. For benchmarks, we have selected the largest convolutional, pooling and/or classifier layers of several recent and large neural network structures. The characteristics of these 10 layers plus a description of the associated neural network and task are shown in Table 5.

7. Experimental Results

7.1 Accelerator Characteristics after Layout

The current version uses $T_n = 16$ (16 hardware neurons with 16 synapses each), so that the design contains 256 16-bit truncated multipliers in NFU-1 (for classifier and convolutional layers), 16 adder trees of 15 adders each in NFU-2 (for the same layers, plus pooling layer if average is used), as well as a 16-input shifter and max in NFU-2 (for pooling layers), and 16 16-bit truncated multipliers plus 16 adders in NFU-3 (for classifier and convolutional layers, and optionally for pooling layers). For classifier and convolutional layers, NFU-1 and NFU-2 are active every cycle, achieving $256 + 16 \times 15 = 496$ fixed-point operations every cycle; at 0.98GHz, this amounts to 452 GOP/s (Giga fixed-point Operations per second). At the end of a layer, NFU-3 would be active as well while NFU-1 and NFU-2 process the remaining data, reaching a peak activity of $496 + 2 \times 16 = 528$ operations per cycle (482 GOP/s) for a short period.

We have done the synthesis and layout of the accelerator with $T_n = 16$ and 64-entry buffers at 65nm using Synopsys tools, see Figure 15. The main characteristics and power/area breakdown by component type and functional block are shown in Table 6. We brought the critical path delay down to 1.02ns by introducing 3 pipeline stages in NFU-1 (multipliers), 2 stages in NFU-2 (adder trees), and 3 stages in NFU-3 (piecewise linear function approximation) for a total of 8 pipeline stages. Currently, the critical path is in the issue

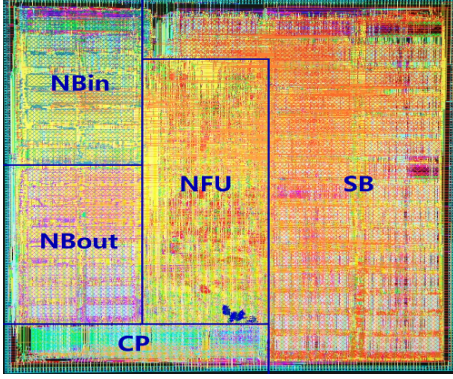


Figure 15. Layout (65nm).

Component or Block	Area in μm^2	Power (%)	Power in mW	Critical path in ns (%)
ACCELERATOR	3,023,077		485	1.02
Combinational	608,842	(20.14%)	89	(18.41%)
Memory	1,158,000	(38.31%)	177	(36.59%)
Registers	375,882	(12.43%)	86	(17.84%)
Clock network	68,721	(2.27%)	132	(27.16%)
Filler cell	811,632	(26.85%)		
SB	1,153,814	(38.17%)	105	(22.65%)
NBin	427,992	(14.16%)	91	(19.76%)
NBout	433,906	(14.35%)	92	(19.97%)
NFU	846,563	(28.00%)	132	(27.22%)
CP	141,809	(5.69%)	31	(6.39%)
AXIMUX	9,767	(0.32%)	8	(2.65%)
Other	9,226	(0.31%)	26	(5.36%)

Table 6. Characteristics of accelerator and breakdown by component type (first 5 lines), and functional block (last 7 lines).

logic which is in charge of reading data out of NBin/NBout; next versions will focus on how to reduce or pipeline this critical path. The total RAM capacity (NBin + NBout + SB + CP instructions) is 44KB (8KB for the CP RAM). The area and power are dominated by the buffers (NBin/NBout/SB) at respectively 56% and 60%, with the NFU being a close second at 28% and 27%. The percentage of the total cell power is 59.47%, but the routing network (included in the different components of the table breakdown) accounts for a significant share of the total power at 38.77%. At 65nm, due to the high toggle rate of the accelerator, the leakage power is almost negligible at 1.73%.

Finally, we have also evaluated a design with $T_n = 8$, and thus 64 multipliers in NFU-1. The total area for that design is 0.85 mm^2 , i.e., 3.59x smaller than for $T_n = 16$ due to the reduced buffer width and the fewer number of arithmetic operators. We plan to investigate larger designs with $T_n = 32$ or 64 in the near future.

7.2 Time and Throughput

In Figure 16, we report the speedup of the accelerator over SIMD, see $SIMD/Acc$. Recall that we use a 128-bit SIMD processor, so capable of performing up to 8 16-bit operations

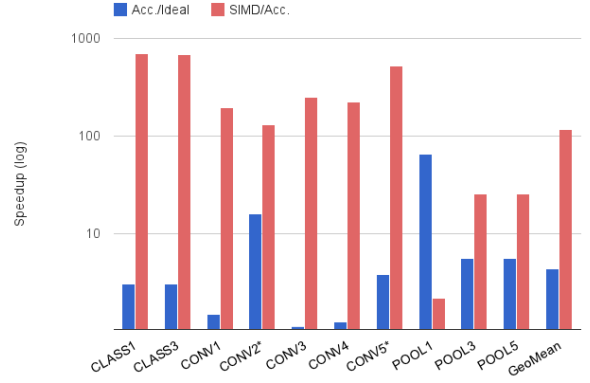


Figure 16. Speedup of accelerator over SIMD, and of ideal accelerator over accelerator.

every cycle (we naturally use 16-bit fixed-point operations in the SIMD as well). As mentioned in Section 7.1, the accelerator performs 496 16-bit operations every cycle for both classifier and convolutional layers, i.e., $62x$ more ($\frac{496}{8}$) than the SIMD core. We empirically observe that on these two types of layers, the accelerator is on average 117.87x faster than the SIMD core, so about 2x above the ratio of computational operators ($62x$). We measured that, for classifier and convolutional layers, the SIMD core performs 2.01 16-bit operations per cycle on average, instead of the upper bound of 8 operations per cycle. We traced this back to two major reasons.

First, better latency tolerance due to an appropriate combination of preloading and reuse in NBin and SB buffers; note that we did not implement a prefetcher in the SIMD core, which would partly bridge that gap. This explains the high performance gap for layers CLASS1, CLASS3 and CONV5 which have the largest feature maps sizes, thus the most spatial locality, and which then benefit most from preloading, giving them a performance boost, i.e., 629.92x on average, about 3x more than other convolutional layers; we expect that a prefetcher in the SIMD core would cancel that performance boost. The spatial locality in NBin is exploited along the input feature map dimension by the DMA, and with a small N_i , the DMA has to issue many short memory requests, which is less efficient. The rest of the convolutional layers (CONV1 to CONV4) have an average speedup of 195.15x; CONV2 has a lesser performance (130.64x) due to private kernels and less spatial locality. Pooling layers have less performance overall because only the adder tree in NFU-2 is used (240 operators out of 496 operators), 25.73x for POOL3 and 25.52x for POOL5.

In order to further analyze the relatively poor behavior of POOL1 (only 2.17x over SIMD), we have tested a configuration of the accelerator where all operands (inputs and synapses) are ready for the NFU, i.e., ideal behavior

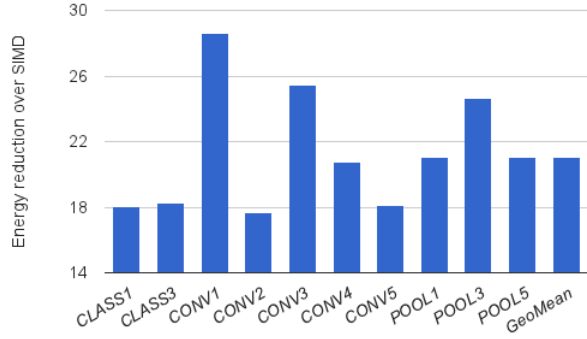


Figure 17. Energy reduction of accelerator over SIMD.

of NBin, SB and NBout; we call this version “Ideal”, see Figure 16. We see that the accelerator is significantly slower on POOL1 and CONV2 than the ideal configuration (respectively 66.00x and 16.14x). This is due to the small size of their input/output feature maps (e.g., $N_i = 12$ for POOL1), combined with the fewer operators used for POOL1. So far, the accelerator has been geared towards large layers, but we can address this weakness by implementing a 2D or 3D DMA (DMA requests over i, k_x, k_y loops); we leave this optimization for future work.

The second reason for the speedup over SIMD beyond 62x lays in control and scheduling overhead. In the accelerator, we have tried to minimize lost cycles. For instance, when output neurons partial sums are rotated to NBout (before being sent back to NFU-2), the oldest buffer row (T_n partial sums) is eagerly rotated out to the NBout/NFU-2 input latch, and a multiplexer in NFU-2 ensures that either this latch or the NFU-2 registers are used as input for the NFU-2 stage computations; this allows a rotation without any pipeline stall. Several such design optimizations help achieve a slowdown of only 4.36x over the ideal accelerator, see Figure 16, and in fact, 2.64x only if we exclude CONV2 and POOL1.

7.3 Energy

In Figure 17, we provide the energy ratio between the SIMD core and the accelerator. While high at 21.08x, the average energy ratio is actually more than an order of magnitude smaller than previously reported energy ratios between processors and accelerators; for instance Hameed et al. [14] report an energy ratio of about 500x, and 974x has been reported for a small Multi-Layer Perceptron [38]. The smaller ratio is largely due to the energy spent in memory accesses, which was voluntarily not factored in the two aforementioned studies. Like in these two accelerators and others, the energy cost of computations has been considerably reduced by a combination of more efficient computational operators (especially a massive number of small 16-bit fixed-point

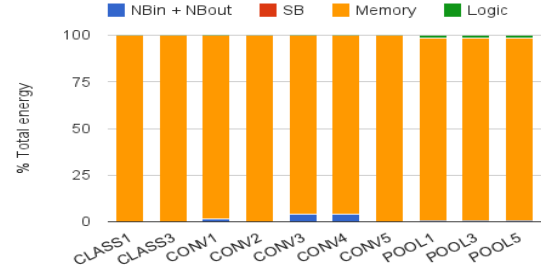


Figure 18. Breakdown of accelerator energy.

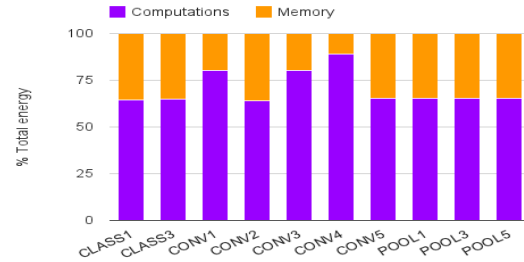


Figure 19. Breakdown of SIMD energy.

truncated multipliers in our case), and small custom storage located close to the operators (64-entry NBin, NBout, SB and the NFU-2 registers). As a result, there is now an Amdahl’s law effect for energy, where any further improvement can only be achieved by bringing down the energy cost of main memory accesses. We tried to artificially set the energy cost of the main memory accesses in both the SIMD and accelerator to 0, and we observed that the average energy reduction of the accelerator increases by more than one order of magnitude, in line with previous results.

This is further illustrated by the breakdown of the energy consumed by the accelerator in Figure 18 where the energy of main memory accesses obviously dominates. A distant second is the energy of NBin/NBout for the convolutional layers with shared kernels (CONV1, CONV3, CONV4). In this case, a set of shared kernels are kept in SB so the memory traffic due to synapses becomes very low, as explained in Section 3 (shared kernels + tiling), but the input neurons must still be reloaded for each new set of shared kernels, hence the still noticeable energy expense. The energy of the computational logic in pooling layers (POOL1, POOL3, POOL5) is similarly a distant second expense, this time because there is no synapse to load. The slightly higher energy reduction of pooling layers (22.17x on average), see Figure 17, is due to the fact the SB buffer is not used (no synapse), and the accesses to NBin alone are relatively cheap due to its small width, see Figure 13.

The SIMD energy breakdown is in sharp contrast, as shown in Figure 19, with about two thirds of the energy spent in computations, and only one third in memory accesses.

While finding a computationally more efficient approach to SIMD made sense, future work for the accelerator should focus on reducing the energy spent in memory accesses.

8. Related Work

Due to stringent energy constraints, such as Dark Silicon [10, 32], there is a growing consensus that future high-performance micro-architectures will take the form of heterogeneous multi-cores, i.e., combinations of cores and accelerators. Accelerators can range from processors tuned for certain tasks, to ASIC-like circuits such as H264 [14], or more flexible accelerators capable of targeting a broad range of, but not all, tasks [12, 44] such as QsCores [42], or accelerators for image processing [33].

The accelerator proposed in this article follows this spirit of targeting a specific, but broad, domain, i.e., machine-learning tasks here. Due to recent progress in machine-learning, certain types of neural networks, especially Deep Neural Networks [25] and Convolutional Neural Networks [27] have become state-of-the-art machine-learning techniques [26] across a broad range of applications such as web search [19], image analysis [31] or speech recognition [7].

While many implementations of hardware neurons and neural networks have been investigated in the past two decades [18], the main purpose of hardware neural networks has been fast modeling of biological neural networks [20, 34] for implementing neurons with thousands of connections. While several of these neuromorphic architectures have been applied to computational tasks [30, 43], the specific bio-inspired information representation (spiking neural networks) they rely on may not be competitive with state-of-the-art neural networks, though this remains an open debate at the threshold between neuroscience and machine-learning.

However, recently, due to simultaneous trends in applications, machine-learning and technology constraints, hardware neural networks have been increasingly considered as potential accelerators, either for very dedicated functionalities within a processor, such as branch prediction [1], or for their fault-tolerance properties [15, 38]. The latter property has also been leveraged to trade application accuracy for energy efficiency through hardware neural processing units [9, 11].

The focus of our accelerator is on large-scale machine-learning tasks, with layers of thousands of neurons and millions of synapses, and for that reason, there is a special emphasis on interactions with memory. Our study not only confirms previous observations that dedicated storage is key for achieving good performance and power [14], but it also highlights that, beyond exploiting locality at the level of registers located close to computational operators [33, 38], considering memory as a prime-order concern can profoundly affect accelerator design.

Many of the aforementioned studies stem from the architecture community. A symmetric effort has started in the machine-learning community where a few researchers have been investigating hardware designs for speeding up neural network processing, especially for real-time applications. Neuflow [13] is an accelerator for fast and low-power implementation of the feed-forward paths of CNNs for vision systems. It organizes computations and register-level storage according to the sliding window property of convolutional and pooling layers; but in that respect, it also ignores much of the first-order locality coming from input and output feature maps. Its interplay with memory remains limited to a DMA, there is no significant on-chip storage, though the DMA is capable of performing complex access patterns. A more complex architecture, albeit with similar performance as Neuflow, has been proposed by Kim et al. [21] and consists of 128 SIMD processors of 16 PEs each; the architecture is significantly larger and implements a specific neural vision model (neither CNNs nor DNNs), but it can achieve 60 frame/sec (real-time) multi-object recognition for up to 10 different objects. Maashri et al. [29] have also investigated the implementation of another neural network model, the bio-inspired HMAX for vision processing, using a set of custom accelerators arranged around a switch fabric; in the article, the authors allude to locality optimizations across different orientations, which are roughly the HMAX equivalent of feature maps. Closer to our community again, but solely focusing on CNNs, Chakradhar et al. [3] have also investigated the implementation of CNNs on reconfigurable circuits; though there is little emphasis on locality exploitation, they pay special attention to properly mapping a CNN in order to improve bandwidth utilization.

9. Conclusions

In this article we focus on accelerators for machine-learning because of the broad set of applications and the few key state-of-the-art algorithms offer the rare opportunity to combine high efficiency and broad application scope. Since state-of-the-art CNNs and DNNs mean very large networks, we specifically focus on the implementation of large-scale layers. By carefully exploiting the locality properties of such layers, and by introducing storage structures custom designed to take advantage of these properties, we show that it is possible to design a machine-learning accelerator capable of high performance in a very small area footprint. Our measurements are not circumscribed to the accelerator fabric, they factor in the performance and energy overhead of main memory transfers; still, we show that it is possible to achieve a speedup of 117.87x and an energy reduction of 21.08x over a 128-bit 2GHz SIMD core with a normal cache hierarchy. We have obtained a layout of the design at 65nm.

Besides a planned tape-out, future work includes improving the accelerator behavior for short layers, slightly altering the NFU to include some of the latest algorithmic im-

provements such as Local Response Normalization, further reducing the impact of main memory transfers, investigating scalability (especially increasing T_n), and implementing training in hardware.

Acknowledgments

This work is supported by a Google Faculty Research Award, the Intel Collaborative Research Institute for Computational Intelligence (ICRI-CI), the French ANR MHANN and NEMESIS grants, the NSF of China (under Grants 61003064, 61100163, 61133004, 61222204, 61221062, 61303158), the 863 Program of China (under Grant 2012AA012202), the Strategic Priority Research Program of the CAS (under Grant XDA06010403), the 10,000 and 1,000 talent programs.

References

- [1] R. S. Amant, D. A. Jimenez, and D. Burger. Low-power, high-performance analog neural branch prediction. In *International Symposium on Microarchitecture*, Como, 2008.
- [2] C. Bienia, S. Kumar, J. P. Singh, and K. Li. The PARSEC benchmark suite: Characterization and architectural implications. In *International Conference on Parallel Architectures and Compilation Techniques*, New York, New York, USA, 2008. ACM Press.
- [3] S. Chakradhar, M. Sankaradas, V. Jakkula, and S. Cadambi. A dynamically configurable coprocessor for convolutional neural networks. In *International symposium on Computer Architecture*, page 247, Saint Malo, France, June 2010. ACM Press.
- [4] T. Chen, Y. Chen, M. Duranton, Q. Guo, A. Hashmi, M. Lipasti, A. Nere, S. Qiu, M. Sebag, and O. Temam. BenchNN: On the Broad Potential Application Scope of Hardware Neural Network Accelerators. In *International Symposium on Workload Characterization*, 2012.
- [5] A. Coates, B. Huval, T. Wang, D. J. Wu, and A. Y. Ng. Deep learning with cots hpc systems. In *International Conference on Machine Learning*, 2013.
- [6] C. Cortes and V. Vapnik. Support-Vector Networks. In *Machine Learning*, pages 273–297, 1995.
- [7] G. Dahl, T. Sainath, and G. Hinton. Improving Deep Neural Networks for LVCSR using Rectified Linear Units and Dropout. In *International Conference on Acoustics, Speech and Signal Processing*, 2013.
- [8] S. Draghici. On the capabilities of neural networks using limited precision weights. *Neural Netw.*, 15(3):395–414, 2002.
- [9] Z. Du, A. Lingamneni, Y. Chen, K. V. Palem, O. Temam, and C. Wu. Leveraging the Error Resilience of Machine-Learning Applications for Designing Highly Energy Efficient Accelerators. In *Asia and South Pacific Design Automation Conference*, 2014.
- [10] H. Esmaeilzadeh, E. Blem, R. S. Amant, K. Sankaralingam, and D. Burger. Dark Silicon and the End of Multicore Scaling. In *Proceedings of the 38th International Symposium on Computer Architecture (ISCA)*, June 2011.
- [11] H. Esmaeilzadeh, A. Sampson, L. Ceze, and D. Burger. Neural Acceleration for General-Purpose Approximate Programs. In *International Symposium on Microarchitecture*, number 3, pages 1–6, 2012.
- [12] K. Fan, M. Kudlur, G. S. Dasika, and S. A. Mahlke. Bridging the computation gap between programmable processors and hardwired accelerators. In *HPCA*, pages 313–322. IEEE Computer Society, 2009.
- [13] C. Farabet, B. Martini, B. Corda, P. Akselrod, E. Culurciello, and Y. LeCun. NeuFlow: A runtime reconfigurable dataflow processor for vision. In *CVPR Workshop*, pages 109–116. Ieee, June 2011.
- [14] R. Hameed, W. Qadeer, M. Wachs, O. Azizi, A. Solomatnikov, B. C. Lee, S. Richardson, C. Kozyrakis, and M. Horowitz. Understanding sources of inefficiency in general-purpose chips. In *International Symposium on Computer Architecture*, page 37, New York, New York, USA, 2010. ACM Press.
- [15] A. Hashmi, A. Nere, J. J. Thomas, and M. Lipasti. A case for neuromorphic ISAs. In *International Conference on Architectural Support for Programming Languages and Operating Systems*, New York, NY, 2011. ACM.
- [16] G. Hinton and N. Srivastava. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv: ...*, pages 1–18, 2012.
- [17] J. L. Holli and J.-N. Hwang. Finite Precision Error Analysis of Neural Network Hardware Implementations. *IEEE Transactions on Computers*, 42(3):281–290, 1993.
- [18] M. Holler, S. Tam, H. Castro, and R. Benson. An electrically trainable artificial neural network (ETANN) with 10240 floating gate synapses. In *Artificial neural networks*, pages 50–55, Piscataway, NJ, USA, 1990. IEEE Press.
- [19] P. Huang, X. He, J. Gao, and L. Deng. Learning deep structured semantic models for web search using clickthrough data. In *International Conference on Information and Knowledge Management*, 2013.
- [20] M. M. Khan, D. R. Lester, L. A. Plana, A. Rast, X. Jin, E. Painkras, and S. B. Furber. SpiNNaker: Mapping neural networks onto a massively-parallel chip multiprocessor. In *IEEE International Joint Conference on Neural Networks (IJCNN)*, pages 2849–2856. Ieee, 2008.
- [21] J.-y. Kim, S. Member, M. Kim, S. Lee, J. Oh, K. Kim, and H.-j. Yoo. A 201.4 GOPS 496 mW Real-Time Multi-Object Recognition Processor With Bio-Inspired Neural Perception Engine. *IEEE Journal of Solid-State Circuits*, 45(1):32–45, Jan. 2010.
- [22] E. J. King and E. E. Swartzlander Jr. Data-dependent truncation scheme for parallel multipliers. In *Signals, Systems & Computers, 1997. Conference Record of the Thirty-First Asilomar Conference on*, volume 2, pages 1178–1182. IEEE, 1997.
- [23] D. Larkin, A. Kinane, V. Muresan, and N. E. O’Connor. An Efficient Hardware Architecture for a Neural Network Activation Function Generator. In J. Wang, Z. Yi, J. M. Zurada, B.-L. Lu, and H. Yin, editors, *ISNN (2)*, volume 3973 of *Lecture Notes in Computer Science*, pages 1319–1327. Springer, 2006.

- [24] D. Larkin, A. Kinane, and N. E. O'Connor. Towards Hardware Acceleration of Neuroevolution for Multimedia Processing Applications on Mobile Devices. In *ICONIP (3)*, pages 1178–1188, 2006.
- [25] H. Larochelle, D. Erhan, A. Courville, J. Bergstra, and Y. Bengio. An empirical evaluation of deep architectures on problems with many factors of variation. In *International Conference on Machine Learning*, pages 473–480, New York, New York, USA, 2007. ACM Press.
- [26] Q. V. Le, M. A. Ranzato, R. Monga, M. Devin, K. Chen, G. S. Corrado, J. Dean, and A. Y. Ng. Building High-level Features Using Large Scale Unsupervised Learning. In *International Conference on Machine Learning*, June 2012.
- [27] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86, 1998.
- [28] S. Li, J. H. Ahn, R. D. Strong, J. B. Brockman, D. M. Tullsen, and N. P. Jouppi. McPAT: an integrated power, area, and timing modeling framework for multicore and manycore architectures. In *Proceedings of the 42nd Annual IEEE/ACM International Symposium on Microarchitecture, MICRO 42*, pages 469–480, New York, NY, USA, 2009. ACM.
- [29] A. A. Maashri, M. Debole, M. Cotter, N. Chandramoorthy, Y. Xiao, V. Narayanan, and C. Chakrabarti. Accelerating neuromorphic vision algorithms for recognition. *Proceedings of the 49th Annual Design Automation Conference on - DAC '12*, page 579, 2012.
- [30] P. Merolla, J. Arthur, F. Akopyan, N. Imam, R. Manohar, and D. Modha. A digital neurosynaptic core using embedded crossbar memory with 45pJ per spike in 45nm. In *IEEE Custom Integrated Circuits Conference*, pages 1–4. IEEE, Sept. 2011.
- [31] V. Mnih and G. Hinton. Learning to Label Aerial Images from Noisy Data. In *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, pages 567–574, 2012.
- [32] M. Muller. Dark Silicon and the Internet. In *EE Times "Designing with ARM" virtual conference*, 2010.
- [33] W. Qadeer, R. Hameed, O. Shacham, P. Venkatesan, C. Kozyrakis, and M. A. Horowitz. Convolution engine: balancing efficiency & flexibility in specialized computing. In *International Symposium on Computer Architecture*, 2013.
- [34] J. Schemmel, J. Fieres, and K. Meier. Wafer-scale integration of analog neural networks. In *International Joint Conference on Neural Networks*, pages 431–438. Ieee, June 2008.
- [35] P. Sermanet, S. Chintala, and Y. LeCun. Convolutional Neural Networks Applied to House Numbers Digit Classification. In *Pattern Recognition (ICPR)*, . . . , 2012.
- [36] P. Sermanet and Y. LeCun. Traffic sign recognition with multi-scale Convolutional Networks. In *International Joint Conference on Neural Networks*, pages 2809–2813. Ieee, July 2011.
- [37] T. Serre, L. Wolf, S. Bileschi, M. Riesenhuber, and T. Poggio. Robust object recognition with cortex-like mechanisms. *IEEE transactions on pattern analysis and machine intelligence*, 29(3):411–26, Mar. 2007.
- [38] O. Temam. A Defect-Tolerant Accelerator for Emerging High-Performance Applications. In *International Symposium on Computer Architecture*, Portland, Oregon, 2012.
- [39] O. Temam and N. Drach. Software assistance for data caches. *Future Generation Computer Systems*, 11(6):519–536, 1995.
- [40] S. Thoziyoor, N. Muralimanohar, and J. Ahn. CACTI 5.1. *HP Labs, Palo Alto, Tech*, 2008.
- [41] V. Vanhoucke, A. Senior, and M. Z. Mao. Improving the speed of neural networks on CPUs. In *Deep Learning and Unsupervised Feature Learning Workshop, NIPS 2011*, 2011.
- [42] G. Venkatesh, J. Sampson, N. Goulding-hotta, S. K. Venkata, M. B. Taylor, and S. Swanson. QsCORES : Trading Dark Silicon for Scalable Energy Efficiency with Quasi-Specific Cores Categories and Subject Descriptors. In *International Symposium on Microarchitecture*, 2011.
- [43] R. J. Vogelstein, U. Mallik, J. T. Vogelstein, and G. Cauwenberghs. Dynamically reconfigurable silicon array of spiking neurons with conductance-based synapses. *IEEE Transactions on Neural Networks*, 18(1):253–265, 2007.
- [44] S. Yehia, S. Girbal, H. Berry, and O. Temam. Reconciling specialization and flexibility through compound circuits. In *International Symposium on High Performance Computer Architecture*, pages 277–288, Raleigh, North Carolina, Feb. 2009. Ieee.