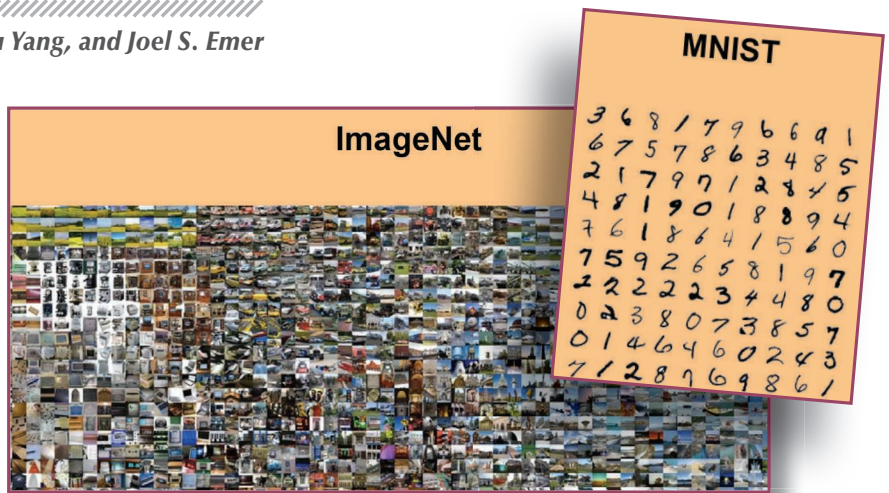


Vivienne Sze, Yu-Hsin Chen, Tien-Ju Yang, and Joel S. Emer



How to Evaluate Deep Neural Network Processors

TOPS/W (alone) Considered Harmful

A significant amount of specialized hardware has been developed for processing deep neural networks (DNNs) in both academia and industry. This article aims to highlight the key concepts required to evaluate and compare these DNN processors. We discuss existing challenges, such as the flexibility and scalability needed to support a wide range of neural networks, as well as design considerations for both the DNN processors and the DNN models themselves. We also describe specific metrics that can be used to evaluate and compare existing solutions beyond the commonly used tera-operations per second per watt (TOPS/W). This article is based on the tutorial “How to Understand and Evaluate Deep Learning Processors” that was given at the 2020 Interna-

tional Solid-State Circuits Conference, as well as excerpts from the book, *Efficient Processing of Deep Neural Networks* [36].

Motivation and Background

Over the past few years, there has been a significant amount of research on enabling the efficient processing of DNNs. The challenge of efficient DNN processing depends on balancing multiple objectives:

- high performance (including accuracy) and efficiency (including cost)
- enough flexibility to cater to a wide and rapidly changing range of workloads
- good integration with existing software frameworks.

DNN computations are composed of several processing layers (Figure 1), where, for many layers, the main computation is a weighted sum; in other words, the main computation for DNN processing is often a

multiply-accumulate (MAC) operation. The arrangement of the MAC operations within a layer is defined by the layer shape; for instance, Table 1 and Figure 2 highlight the shape parameters for layers used in convolutional neural networks (CNNs), a popular type of DNN. Because the shape parameters can vary across layers, DNNs come in a wide variety of shapes and sizes, depending on the application. (The DNN research community often refers to the shape and size of a DNN as its *network architecture*. However, to avoid confusion with the use of the word *architecture* by the hardware community, we talk about *DNN models* and their shape and size in this article.) This variety is one of the motivations for flexibility, and it causes the objectives listed previously to be highly interrelated.

Figure 3 illustrates the hardware architecture of a typical DNN processor, which is composed of an array

of processing elements (PEs), where each PE contains MAC units to perform the computation (and, optionally, some local storage) and an inter-PE communication network. The entire PE array is also connected via an on-chip network to a large global buffer, which, in turn, is connected off chip to DRAM. DNN processor designs tend to vary in terms of the number of PEs, number of levels in the memory hierarchy, amount of storage at each level, and how the PEs and memory are connected through the on-chip network.

Given the combination of such hardware and the associated DNN models, it is important to discuss the key metrics that should be considered when comparing and evaluating the strengths and weaknesses of different designs. They also can be used to evaluate proposed techniques and should be incorporated into design considerations. While efficiency is often associated with only the number of operations per second per watt [e.g., floating-point operations per second per watt (FLOPS/W) or TOPS/W], it is actually composed of many more metrics, including accuracy, throughput, latency, energy consumption, power consumption, cost, flexibility, and scalability. (Note that TOPS/W efficiency is typically reported at, and often along with, the peak performance in TOPS, which gives the maximum efficiency since it assumes maximum utilization and thus maximum amortization of overhead. However, this does not tell the complete story because processors typically do not operate at their peak TOPS and their efficiency degrades at lower utilization. It is a well-known challenge to achieve energy-proportional computing, where the efficiency stays constant across performance [37].) Reporting a comprehensive set of these metrics is important to provide a complete picture of the tradeoffs made by a proposed design or technique.

In this article, we

- discuss the importance of each of these metrics
- break down the factors that affect each metric and, when feasible,

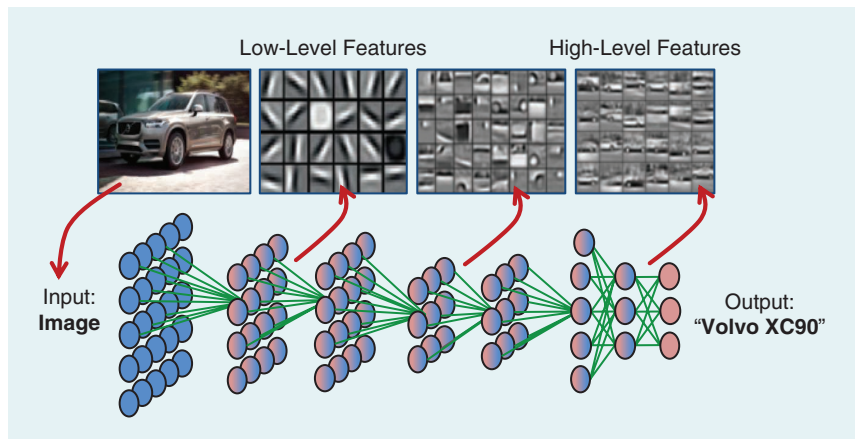


FIGURE 1: An example of image classification using a DNN. The DNN is composed of multiple layers, and the number of layers is referred to as the depth of the network. Note that the extracted features go from low level to high level as we go deeper into the network. (Source: Adapted from [1].)

TABLE 1. THE SHAPE PARAMETERS FOR THE LAYERS USED IN DNNs.

SHAPE PARAMETER	DESCRIPTION
N	Batch size of 3D feature map
M	Number of 3D filters/number of channels of output feature map (output channels)
C	Number of channels of filter/input feature map (input channels)
H/W	Spatial height/width of input feature map
R/S	Spatial height/width of filter
P/Q	Spatial height/width of output feature map

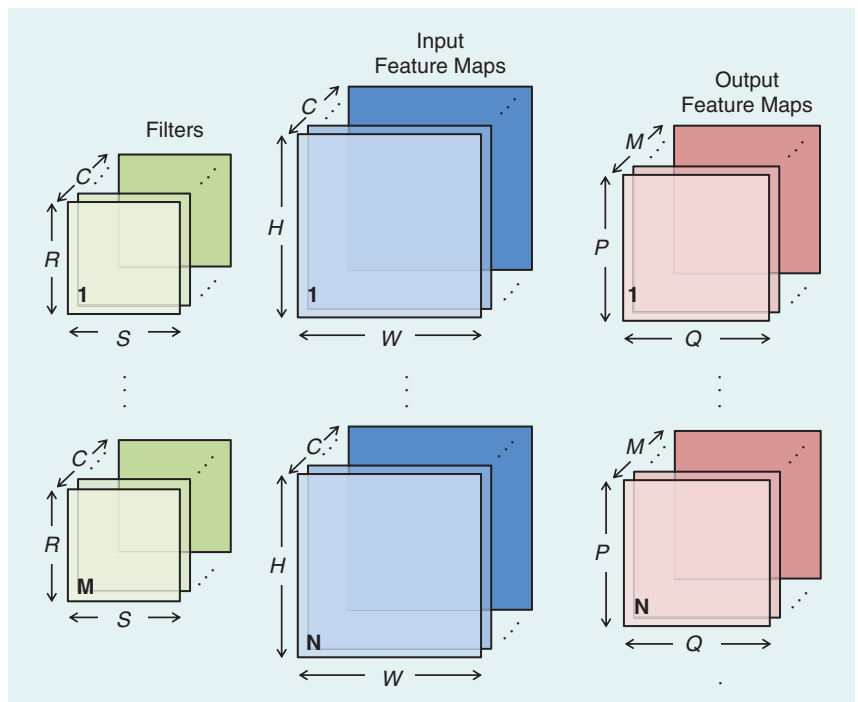


FIGURE 2: The shape parameters of the layers used in CNNs.

present equations that describe the relationship between the factors and the metrics

- describe how these metrics can be incorporated into design considerations for both the DNN hardware and the DNN model
- specify what should be reported for a given metric to enable proper evaluation.

Finally, we highlight tools that can be used to evaluate some of these metrics early in the design process (to enable rapid design exploration) and provide a case study on how one might bring all of these metrics together for a holistic evaluation of a given approach. First, however, we discuss each of the metrics.

Accuracy

Accuracy indicates the quality of the result for a given task. The fact that

DNNs can achieve state-of-the-art accuracy on a wide range of tasks is one of the key reasons driving their popularity and wide use today. The units used to measure accuracy depend on the task. For instance, for image classification, accuracy is reported as the percentage of correctly classified images, while, for object detection, accuracy is reported as the mean average precision, which is related to the tradeoff between true positives, false positives, and false negatives.

Factors that affect accuracy include the difficulty of the task and data set. (Ideally, robustness and fairness should be considered in conjunction with accuracy, as there is also an interplay between these factors; however, these are areas of ongoing research and beyond the scope of this article.) For instance, classification on the ImageNet data set [2] is much more dif-

ficult than on the MNIST data set [3] (Figure 4), and object detection is usually more difficult than classification. As a result, a DNN model that performs well on MNIST may not necessarily perform well on ImageNet. Achieving high accuracy on difficult tasks or data sets typically requires more complex DNN models (e.g., a larger number of MAC operations and more distinct weights, increased diversity in layer shapes, and so on), which can impact how efficiently the hardware can process the DNN model.

Accuracy should, therefore, be interpreted in the context of the difficulty of the task and data set. (As an analogy, getting nine out of 10 answers correct on a high school exam is different than nine out of 10 answers correct on a college-level exam. One must look beyond the score and consider the difficulty of the exam.) Evaluating hardware using well-studied, widely used DNN models, tasks, and data sets can allow one to better interpret the significance of the accuracy metric.

Recently, motivated by the impact of the SPEC benchmarks for general purpose computing [4], several industry and academic organizations have put together a broad suite of models, called *MLPerf*, to serve as a common set of well-studied DNN models to evaluate the performance and enable fair comparison of various software frameworks, hardware architectures, and cloud platforms for both training

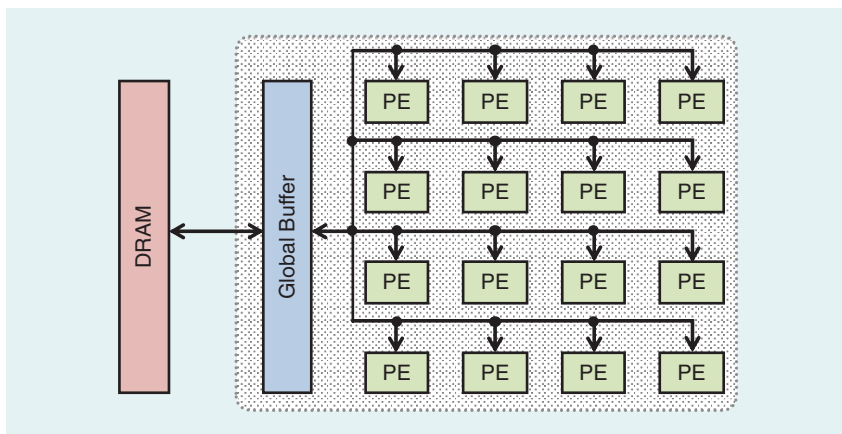


FIGURE 3: The typical hardware architecture of a DNN processor.



FIGURE 4: The (a) MNIST data set (10 classes, 60,000 training, and 10,000 testing) [3] versus the (b) ImageNet data set (1,000 classes, 1.3 million training, and 100,000 testing) [2].

and inference of DNNs [5]. (Earlier DNN benchmarking efforts, including DeepBench [6] and Fathom [7], have now been subsumed by MLPerf.) The suite comprises various types of DNNs (e.g., CNNs, recurrent neural networks, and so on) for a variety of tasks, including image classification, object identification, translation, speech to text, recommendation, sentiment analysis, and reinforcement learning.

Throughput and Latency

Throughput is used to indicate the amount of data that can be processed or the number of executions of a task that can be completed in a given time period. High throughput is often critical to an application. For instance, processing video at 30 frames/s is often necessary to deliver real-time performance. For data analytics, high throughput means that more data can be analyzed in a given amount of time. As the amount of visual data is growing exponentially, high-throughput big data analytics becomes increasingly important, particularly if an action needs to be taken based on the analysis (e.g., security or terrorist prevention, medical diagnosis, or drug discovery). Throughput is often generically reported as the number of operations per second. In the case of inference, throughput is reported as inferences per second.

Latency measures the time between the input data's arrival to a system and the generation of the result. Low latency is necessary for real-time interactive applications, such as augmented reality, autonomous navigation, and robotics. Latency is typically reported in seconds per inference.

Throughput and latency are often assumed to be directly derivable from one another. However, they are actually quite distinct. A prime example of this is the well-known approach of batching input data (e.g., batching multiple images or frames together for processing) to increase throughput since batching amortizes overhead such as loading the weights; however, batching also increases latency (e.g., at 30 frames/s and a batch of 100 frames, some frames will experience at least a 3.3-s

delay), which is not acceptable for real-time applications such as high-speed navigation, where it would reduce the time available for course correction.

Thus, achieving low latency and high throughput simultaneously can sometimes be at odds depending on the approach, and both metrics should be reported. The phenomenon described here can also be understood using Little's law [8] from queuing theory, where the average throughput and average latency are related by the average number of tasks in flight, as defined by:

$$\overline{\text{throughput}} = \frac{\text{tasks-in-flight}}{\text{latency}}$$

A DNN-centric version of Little's law would have throughput measured in inferences per second, latency measured in seconds, and inferences in flight (as the tasks-in-flight equivalent) measured in terms of the number of images in a batch being processed simultaneously. This helps to explain why increasing the number of inferences in flight to increase throughput may be counterproductive: some techniques that increase the number of inferences in flight (e.g., batching) also increase latency.

Several factors affect throughput and latency. In terms of throughput, the number of *inferences per second* is affected by

$$\frac{\text{inferences}}{\text{second}} = \frac{\text{operations}}{\text{second}} \times \frac{1}{\frac{\text{operations}}{\text{inference}}}, \quad (1)$$

where the number of *operations per second* is dictated by both the DNN hardware and DNN model, while the number of *operations per inference* is dictated by the DNN model.

When considering a system comprising multiple PEs, where a PE corresponds to a simple or primitive core that performs a single MAC operation, the number of *operations per second* can be further decomposed as follows:

$$\begin{aligned} & \frac{\text{operations}}{\text{second}} \\ &= \left(\frac{1}{\frac{\text{cycles}}{\text{operation}}} \times \frac{\text{cycles}}{\text{second}} \right) \\ & \quad \times \text{number of PEs} \\ & \quad \times \text{utilization of PEs.} \quad (2) \end{aligned}$$

The first term reflects the peak throughput of a single PE, the second term reflects the amount of parallelism, and the last term reflects degradation due to the inability of the architecture to effectively utilize the PEs. Since the main operation for processing DNNs is a MAC operation, we use the terms *number of operations* and *number of MAC operations* interchangeably.

One can increase the peak throughput of a single PE by increasing the number of *cycles per second*, which corresponds to a higher clock frequency achieved by reducing the critical path at the circuit or microarchitectural level; alternatively, one can also reduce the number of *cycles per operation*, which can be affected by the design of the MAC (e.g., a bit-serial, multicycle MAC would have more cycles per operation).

While these approaches increase the throughput of a single PE, the overall throughput can be increased by increasing the *number of PEs* and, thus, the maximum number of MAC operations that can be performed in parallel. The *number of PEs* is dictated by the area of the PE and the area cost of the system. If the area cost of the system is fixed, then increasing the *number of PEs* requires either reducing the area per PE or trading off on-chip storage area for more PEs. Reducing on-chip storage, however, can affect the *utilization of PEs*, which we discuss next.

Reducing the area per PE can also be achieved by reducing the logic associated with delivering operands to a MAC. This can be achieved by controlling multiple MAC operations with a single piece of logic. This is analogous to the situation in instruction-based systems, such as CPUs and graphics processing units (GPUs), that reduce

instruction bookkeeping overhead by using large aggregate instructions (e.g., single-instruction, multiple-data (SIMD), vector instructions, single-instruction, multiple-threads (SIMT), or tensor instructions), where a single instruction can be used to initiate multiple operations.

The *number of PEs* and the peak throughput of a single PE indicate only the theoretical maximum throughput (i.e., peak performance) when all PEs are performing computation (100% utilization). In reality, the achievable throughput depends on the actual utilization of those PEs, which is affected by several factors as follows:

$$\begin{aligned} & \text{utilization of PEs} \\ &= \frac{\text{number of active PEs}}{\text{number of PEs}} \\ & \times \text{utilization of active PEs.} \quad (3) \end{aligned}$$

The first term reflects the ability to distribute the workload to PEs, while the second term reflects how efficiently those active PEs are processing the workload. The *number of active PEs* is the number of PEs that receive work (the ratio of active PEs to the total number of PEs can be referred to as the *active PE percentage*); therefore, it is desirable to distribute the workload to as many PEs as possi-

ble. The ability to distribute the workload is determined by the flexibility of the architecture, for instance, the on-chip network, to support the different layer shapes in a DNN model as explored in [9] and [10].

Within the constraints of the on-chip network, the *number of active PEs* is also determined by the specific allocation of work to PEs by the mapping process. The mapping process involves the placement and scheduling in space and time of every MAC operation (including the delivery of the appropriate operands) onto the PEs. The mapper can be thought of as a compiler for the DNN processor [11]. The mapping process, on a layer-by-layer basis, is explored in detail in [12]–[14]. Additional challenges regarding the flexibility of mapping are discussed in the “Energy Efficiency and Power Consumption” section.

The *utilization of active PEs* is largely dictated by the timely delivery of work to the PEs such that the active PEs do not become idle while waiting for the data to arrive. This can be affected by the bandwidth (BW) and latency of the (on-chip and off-chip) memory and network. The BW requirements can be affected by the amount of data reuse available in the DNN model and the amount of data reuse that can be exploited

by the memory hierarchy and data-flow. The dataflow determines the order of operations and where data are stored and reused. The amount of data reuse can also be increased using a larger batch size, which is one of the reasons that increasing batch size can increase throughput. The challenges of data delivery and memory BW are discussed in [14] and [15]. The *utilization of active PEs* can also be affected by the imbalance of work allocated across PEs, which may occur when exploiting sparsity (i.e., avoiding unnecessary work associated with multiplications by zero); PEs with less work become idle and, thus, have lower utilization.

There is also an interplay between the *number of PEs* and the *utilization of PEs*. For instance, one way to reduce the likelihood that a PE needs to wait for data is to store some data locally near or within the PE. However, this requires increasing the chip area allocated to on-chip storage, which, given a fixed chip area, would reduce the number of PEs. Therefore, a key design consideration is the area allocation between compute (which increases the number of PEs) versus on-chip storage (which increases the utilization of PEs).

The impact of these factors can be captured using *Eyexam*, a systematic way of understanding the performance limits for DNN processors as a function of specific characteristics of the DNN model and processor design. *Eyexam* includes and extends the well-known roofline model [16]. The roofline model, as illustrated in Figure 5, relates average BW demand and peak computational ability to performance.

The goal of *Eyexam* is to provide a fine-grained performance profile for a DNN processor. It is a sequential analysis process that involves seven major steps, as shown in Figure 6. The process starts with the assumption that the architecture has infinite processing parallelism, storage capacity, and data BW. Therefore, it has infinite performance (as measured in MAC operations per cycle).

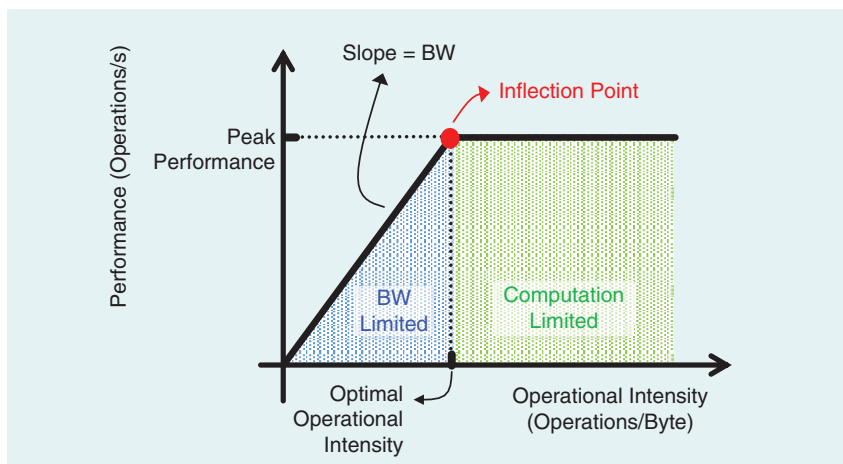


FIGURE 5: The roofline model. The peak operations per second is indicated by the bold line; when the operation intensity (dictated by the amount of compute per byte of data) is low, the operations per second is limited by the data delivery. The design goal is to operate as closely as possible to the peak operations per second for the operation intensity of a given workload.

For each of the following steps, certain constraints are added to reflect changes in the assumptions on the DNN processor or workload. The associated performance loss can, therefore, be attributed to that change, and the final performance at one step becomes the upper bound for the next step.

- **Step 1 (layer shape and size):** In this first step, we look at the impact of the workload constraint so that there is all spatial (i.e., parallel) processing and no temporal (i.e., serial) processing. Therefore, the performance upper bound is determined by the finite size of the workload (i.e., the number of MAC operations in the layer).
- **Step 2 (dataflow):** In this step, we specify the dataflow, which determines the order of operations and where data are stored and reused, and examine the impact of this architectural constraint. Imposing a dataflow forces a serialization of processing and reduces the performance upper bound, which is the maximum parallelism of the dataflow.
- **Step 3 (number of PEs):** In this step, we restrict the system to a finite number of PEs and look at the impact of this architectural constraint. A finite number of PEs can degrade performance whenever there is more parallel work to do than that number of PEs. In addition,

some of the PEs will be idle (i.e., reducing the number of active PEs) if the amount of work is not an integer multiple of the number of PEs (i.e., the work cannot be equally divided among the PEs).

- **Step 4 (physical dimensions of the PE array):** In this step, we consider the physical dimensions of the PE array and data delivery network (e.g., arranging 12 PEs as 3×4 , 2×6 , 4×3 , and so on). The spatial partitioning and associated on-chip network are often constrained per data type (e.g., input activation or filter weight), which can cause additional performance loss because the required data cannot be delivered to the PEs.
- **Step 5 (storage capacity):** In this step, we consider the impact of making the buffer storage have finite capacity. Lack of storage can limit parallelism when there is insufficient storage to hold intermediate results and, thus, degrade performance.
- **Step 6 (data BW):** In this step, we consider the impact of a finite BW for delivering data across the different levels of the memory hierarchy. The amount of data that needs to be transferred between each level of the memory hierarchy for each step of computation and the available data BW determine whether the PEs can be kept busy.

- **Step 7 (varying data access patterns):** In this step, we consider the impact of BW varying across time due to the dynamically changing data access patterns (step 6 addresses only average BW). This includes ramp-up time to initially load values and ramp-down time to drain values after completion. Many common solutions are available to address this issue, including using double buffering, but these can increase the area or reduce the amount of reuse.

Table 2 summarizes the constraints applied at each step of the Eyexam process.

Up until this point, we have discussed how hardware design decisions impact performance (i.e., throughput and latency). We now consider how the choice of DNN model can also have an effect. Specifically, while the number of *operations per inference* in (1) depends on the DNN model, the number of *operations per second* depends on both the DNN model and the hardware. Thus, designing DNN models with efficient layer shapes (also referred to as *efficient network architectures*, such as MobileNet [17]) can reduce the number of MAC operations in the DNN model and, consequently, the number of *operations per inference*. However, such DNN models can result in a wide range of layer shapes, some of which

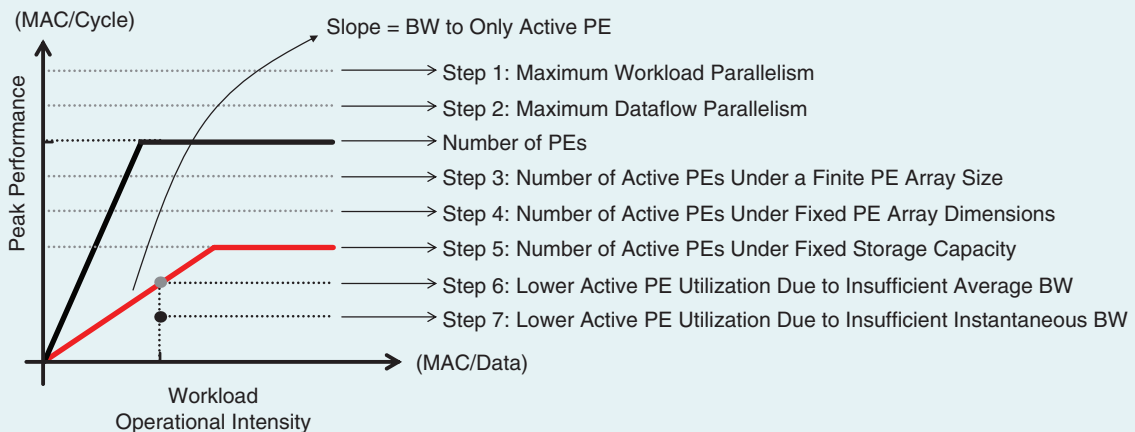


FIGURE 6: The impact of the Eyexam steps on the roofline model.

TABLE 2. A SUMMARY OF THE STEPS IN EYEXAM.

STEP	CONSTRAINT	TYPE	NEW PERFORMANCE BOUND	REASON FOR PERFORMANCE LOSS
1	Layer size and shape	Workload	Maximum workload parallelism	Finite workload size
2	Dataflow loop nest	Architectural	Maximum dataflow parallelism	Restricted dataflows defined by loop nest
3	Number of PEs	Architectural	Maximum PE parallelism	Additional restriction to mappings due to shape fragmentation
4	Physical dimensions of PE array	Architectural	Number of active PEs	Additional restriction to mappings due to shape fragmentation for each dimension
5	Fixed storage capacity	Architectural	Number of active PEs	Additional restriction to mappings due to storage of intermediate data (depends on dataflow)
6	Fixed data BW	Microarchitectural	Maximum data BW to active PEs	Insufficient average BW to active PEs
7	Varying data access patterns	Microarchitectural	Actual measured performance	Insufficient instant BW to active PEs

may have poor *utilization of PEs*, and thus reduce the overall *operations per second*, as shown in (2).

A deeper consideration of the *operations per second* is that all operations are not created equal, so *cycles per operation* may not be constant. For example, if we consider the fact that anything multiplied by zero is zero, some MAC operations are ineffectual (i.e., they do not change the accumulated value). The number of ineffectual operations is a function of both the weights in the DNN model and the input data. These ineffectual MAC operations can require fewer cycles or no cycles at all. Conversely, we need to process only effectual (or nonzero) MAC operations, where both inputs are nonzero; this is referred to as *exploiting sparsity*. Various hardware architectures have been proposed to exploit sparsity [18]–[20].

Processing only effectual MAC operations can increase the *(total) operations per second* by increasing the *(total) operations per cycle*. (By *total operations*, we mean both effectual and ineffectual operations.) Ideally, the hardware would skip all ineffectual operations; however, in practice, designing hardware to skip all ineffectual operations can be challenging and result in increased hardware complexity and overhead. For instance, it might be easier to design hardware that recognizes zeros in only one of the operands (e.g., weights) rather than both. There-

fore, the ineffectual operations can be further divided into those that are exploited by the hardware (i.e., skipped) and those that are unexploited by the hardware (i.e., not skipped). The number of operations actually performed by the hardware is, therefore, the *effectual operations plus unexploited ineffectual operations*.

Equation (4) at the bottom of the page shows how *operations per cycle* can be decomposed into

- the number of *effectual operations plus unexploited ineffectual operations per cycle*, which remains somewhat constant for a given hardware architecture design
- the ratio of *effectual operations over effectual operations plus unexploited ineffectual operations*, which refers to the ability of the hardware to exploit ineffectual operations (ideally, unexploited ineffectual operations should be zero, and this ratio should be one)
- the number of *effectual operations out of (total) operations*, which is related to the amount of sparsity and depends on the DNN model.

As the amount of sparsity increases [i.e., the number of *effectual operations out of (total) operations* decreases], the *operations per cycle* increases, as shown in (4); this subsequently increases *operations per second*, as shown in (2).

However, exploiting sparsity requires additional hardware to identify when inputs to the MAC are zero to avoid performing unnecessary MAC operations. The additional hardware can increase the critical path, which decreases *cycles per second*, and can also increase the area of the PE, which reduces the *number of PEs* for a given area. Both of these factors can reduce the *operations per second*, as shown in (2). Therefore, the complexity of the additional hardware can result in a tradeoff between reducing the number of *unexploited ineffectual operations* and increasing the critical path or reducing the number of PEs.

Finally, designing hardware and DNN models that support reduced precision (i.e., fewer bits per operand and per operation) can also increase the number of *operations per second*. Fewer bits per operand means that

$$\begin{aligned}
 \frac{\text{operations}}{\text{cycle}} &= \frac{\text{effectual operations} + \text{unexploited ineffectual operations}}{\text{cycle}} \\
 &\times \frac{\text{effectual operations}}{\text{effectual operations} + \text{unexploited ineffectual operations}} \\
 &\times \frac{1}{\frac{\text{effectual operations}}{\text{operations}}} \quad (4)
 \end{aligned}$$

TABLE 3. THE CLASSIFICATION OF FACTORS THAT AFFECT INFERENCE PER SECOND.

FACTOR	HARDWARE	DNN MODEL	INPUT DATA
Operations per inference		✓	
Operations per cycle	✓		
Cycles per second	✓		
Number of PEs	✓		
Number of active PEs	✓	✓	
Utilization of active PEs	✓	✓	
Effectual operations out of (total) operations		✓	✓
Effectual operations plus unexploited ineffectual operations per cycle	✓		

the memory BW required to support a given operation is reduced, which can increase the *utilization of the active PEs* since they are less likely to be starved for data. In addition, the area of each PE can be reduced, which can increase the *number of PEs* for a given area. Both of these factors can increase the *operations per second*, as shown in (2). Note, however, that if *multiple* levels of precision need to be supported, additional hardware is required [21]; this can, once again, increase the critical path and also increase the area of the PE, both of which can reduce the *operations per second*, as shown in (2).

In this section, we discussed multiple factors that affect the number of *inferences per second*. Table 3 classifies whether the factors are dictated by the hardware, the DNN model, or both.

In summary, the number of MAC operations in the DNN model alone is not sufficient for evaluating the throughput and latency. While the DNN model can affect the number of MAC operations per inference based on the network architecture (i.e., layer shapes) and the sparsity of the weights and activations, the overall impact that the DNN model has on throughput and latency depends on the ability of the hardware to add support to recognize these approaches without significantly reducing the *utilization of PEs, number of PEs, or cycles per second*. This is why the number of MAC operations is not necessarily a good proxy for throughput and latency (see Figure 7),

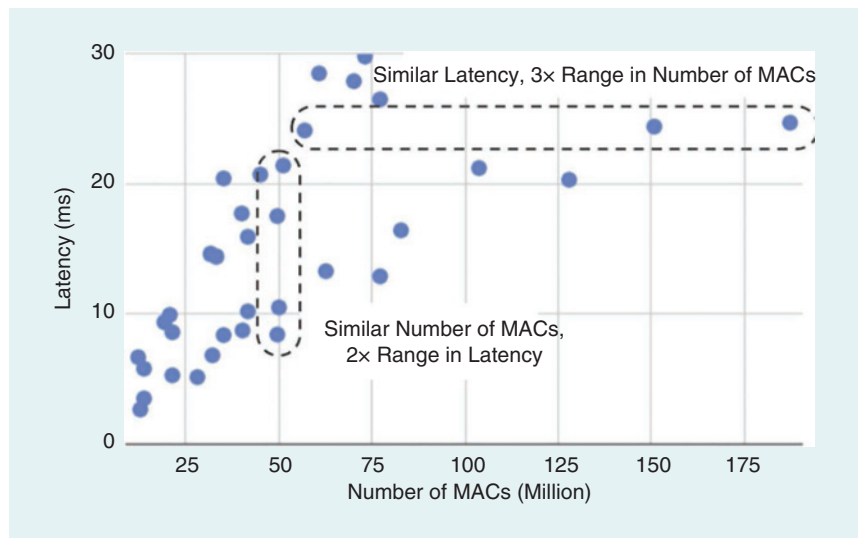


FIGURE 7: The number of MAC operations in various DNN models versus latency measured on a Google Pixel phone. Clearly, the number of MAC operations is not a good predictor of latency. (Source: From [26].)

and it is often more effective to design efficient DNN models with hardware in the loop. Various works have proposed techniques for designing DNN models with hardware in the loop [22]–[25].

Similarly, the number of PEs in the hardware and their peak throughput are not sufficient for evaluating the throughput and latency. It is critical to report the actual runtime of the DNN models on the hardware to account for other effects, such as utilization of the PEs, as highlighted in (2). Ideally, this evaluation should be performed on clearly specified DNN models, for instance, those that are part of the MLPerf benchmarking suite. In addition, batch size should be reported in conjunction with the throughput to evaluate latency.

Energy Efficiency and Power Consumption

Energy efficiency indicates the amount of data that can be processed or the number of executions of a task that can be completed for a given unit of energy. High energy efficiency is important when processing DNNs at the edge in embedded devices with limited battery capacity (e.g., smartphones, smart sensors, robots, and wearables). Edge processing may be preferred over the cloud for certain applications due to latency, privacy, or communication BW limitations. Energy efficiency is often generically reported as the number of operations per joule. In the case of inference, energy efficiency is reported as inferences per joule and energy

consumption is reported as joules per inference.

Power consumption is used to indicate the amount of energy consumed per unit time. Increased power consumption results in increased heat dissipation; accordingly, the maximum power consumption is dictated by a design criterion typically called the *thermal design power (TDP)*, which is the power that the cooling system is designed to dissipate. Power consumption is important when processing DNNs in the cloud, as data centers have stringent power ceilings due to cooling costs; similarly, handheld and wearable devices also have tight power constraints since the user is often quite sensitive to heat and the form factor of the device limits the cooling mechanisms (e.g., no fans). Power consumption is typically reported in watts or joules per second.

Power consumption in conjunction with energy efficiency limits the throughput as follows:

$$\frac{\text{inferences}}{\text{second}} \leq \text{Max} \left(\frac{\text{joules}}{\text{second}} \right) \times \frac{\text{inferences}}{\text{joule}} \quad (5)$$

Therefore, if we can improve energy efficiency by increasing the number of *inferences per joule*, we can increase the number of *inferences per second* and, thus, the throughput of the system. Several factors affect energy efficiency. The number of *inferences per joule* can be decomposed into

$$\frac{\text{inferences}}{\text{joule}} = \frac{\text{operations}}{\text{joule}} \times \frac{1}{\frac{\text{operations}}{\text{inference}}} \quad (6)$$

where the number of *operations per joule* is dictated by both the hardware and DNN model, while the number of *operations per inference* is dictated by the DNN model.

Various design considerations for the hardware will affect the energy per operation (i.e., joules per operation). The energy per operation can be broken down into the energy required to move the input and output data and the energy required to perform the MAC computation:

$$\text{Energy}_{\text{total}} = \text{Energy}_{\text{data}} + \text{Energy}_{\text{MAC}} \quad (7)$$

For each component, the joules per operation (here, an operation can be a MAC operation or a data movement) is computed as

$$\frac{\text{joules}}{\text{operation}} = \alpha \times C \times V_{\text{DD}}^2 \quad (8)$$

where C is the total switching capacitance, V_{DD} is the supply voltage, and α is the switching activity, which indicates how often the capacitance is charged.

The energy consumption is dominated by the data movement, as the capacitance for data movement tends to be much higher than the capacitance for arithmetic operations, such as a MAC (Figure 8). Furthermore, the switching capacitance increases with the distance the data need to travel to reach the PE, which consists of the distance to get out of the memory where the data are stored and the distance to cross the network between the memory and the PE. Accordingly, larger memories and longer interconnects (e.g., off chip) tend to consume more energy than smaller and closer memories due to the capacitance of the long wires employed. To reduce the energy consumption of data movement, we can exploit data reuse, where the data are moved once from a distant large memory (e.g., off-chip DRAM) and reused for multiple operations from a local smaller memory (e.g., an on-chip buffer or scratchpad within the PE). Optimizing data movement is a major consideration in the design of DNN processors, as explored in [15] and [27]. In addition, advanced device and memory technologies can be used to reduce the switching capacitance between compute and memory, for instance, by enabling in-memory computing [28], [29].

This raises the issue of the appropriate scope over which energy efficiency and power consumption should be reported. Including the entire system (out to the fans and power supplies) is beyond the scope of this article. Conversely, ignoring off-chip memory accesses, which can vary greatly among chip designs, can

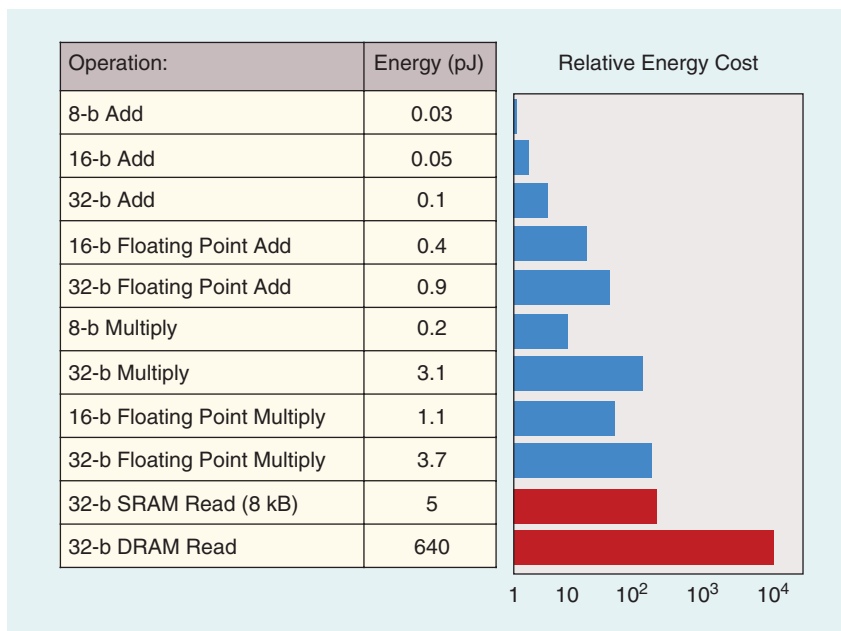


FIGURE 8: The energy consumption for various arithmetic operations and memory accesses in a 45-nm process. The relative energy cost (computed relative to the 8-b add) is shown on a log scale. The energy consumption of data movement (red) is significantly higher than arithmetic operations (blue). SRAM: static random-access memory. (Source: Adapted from [30].)

easily result in a misleading perception of the efficiency of the system. Therefore, it is critical to report not only the energy efficiency and power consumption of the chip but also the energy efficiency and power consumption of the off-chip memory (e.g., DRAM) or the amount of off-chip accesses (e.g., DRAM accesses) if no specific memory technology is specified; for the latter, it can be reported in terms of the total amount of data that is read and written off chip per inference.

Reducing the joules per MAC operation can be achieved by reducing the switching activity and/or capacitance at the circuit or microarchitecture level. This can also be achieved by reducing precision (e.g., reducing the bit width of the MAC operation), as shown in Figure 8. Note that the impact of reducing precision on accuracy must also be considered.

For instruction-based systems, such as CPUs and GPUs, this can also be achieved by reducing instruction book-keeping overhead. For example, using large aggregate instructions (e.g., SIMD, vector instructions, SIMT, or tensor instructions), a single instruction can be used to initiate multiple operations.

Similar to the throughput metric discussed in the “Throughput and Latency” section, the number of *operations per inference* depends on the DNN model; however, the *operations per joules* may be a function of the ability of the hardware to exploit sparsity to avoid performing ineffectual MAC operations. Equation (9) at the bottom of the page shows how *operations per joule* can be decomposed into

- the number of *effectual operations plus unexploited ineffectual operations per joule*, which remains somewhat constant for a given hardware architecture design
- the ratio of *effectual operations over effectual operations plus unexploited ineffectual operations*, which refers to the ability of the hardware to exploit ineffectual operations (ideally, unexploited ineffectual operations should be zero, and this ratio should be one)

- the number of *effectual operations out of (total) operations*, which is related to the amount of sparsity and depends on the DNN model.

For hardware that can exploit sparsity, increasing the amount of sparsity [i.e., decreasing the number of *effectual operations out of (total) operations*] can increase the number of *operations per joule*, which, subsequently, increases *inferences per joule*, as shown in (6). While exploiting sparsity has the potential of increasing the number of *(total) operations per joule*, the additional hardware will decrease the *effectual operations plus unexploited ineffectual operations per joule*. To achieve a net benefit, the decrease in *effectual operations plus unexploited ineffectual operations per joule* must be more than offset by the decrease of *effectual operations out of (total) operations*.

In summary, we want to emphasize that the number of MAC operations and weights in the DNN model are not sufficient for evaluating energy efficiency. From an energy perspective, all MAC operations or weights are not created equal. This is because the number of MAC operations and weights do not reflect where the data are accessed and how much the data are reused, both of which have a significant impact on the operations per joule. Therefore, the number of MAC operations and weights are not necessarily a good proxy for energy consumption, and it is often more effective to design efficient DNN models with hardware in the loop.

To evaluate the energy efficiency and power consumption of the entire system, it is critical to report not only the energy efficiency and

power consumption of the chip but also the energy efficiency and power consumption of the off-chip memory (e.g., DRAM) or the amount of off-chip accesses (e.g., DRAM accesses) if no specific memory technology is specified; for the latter, it can be reported in terms of the total amount of data that is read and written off-chip per inference. As with throughput and latency, the evaluation should be performed on clearly specified and, ideally, widely used DNN models. Various tools can be used to help with energy estimation, as shown in Figure 9.

Hardware Cost

To evaluate the desirability of a given architecture or technique, it is also important to consider the hardware cost of the design. Hardware cost is used to indicate the monetary cost to build a system.

There is also the cost associated with operating a system, such as the electricity bill and the cooling cost, which are primarily dictated by the energy efficiency and power consumption, respectively. In addition, there is the cost associated with designing the system. The operating cost is covered by the “Energy Efficiency and Power Consumption” section, and we limited our coverage of the design cost to the fact that custom DNN processors have a higher design cost (after amortization) than off-the-shelf CPUs and GPUs. We consider anything beyond this, e.g., the economics of the semiconductor business, including how to price platforms, to be outside the scope of this article. Considering the hardware cost of the design is important from both an industry and a research perspective as it dictates whether a system is financially viable.

$$\frac{\text{operations}}{\text{joule}} = \frac{\text{effectual operations} + \text{unexploited ineffectual operations}}{\text{joule}} \times \frac{\text{effectual operations}}{\text{effectual operations} + \text{unexploited ineffectual operations}} \times \frac{1}{\frac{\text{effectual operations}}{\text{operations}}} \quad (9)$$

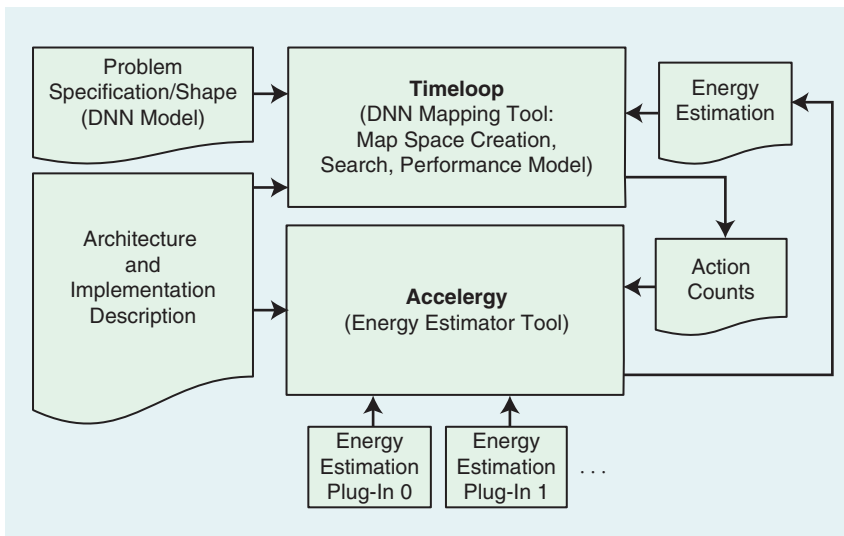


FIGURE 9: Timeloop [13] with the integration of Accelegy [34] as an energy estimation model. Timeloop sends projected action counts for a mapping to Accelegy and receives an energy estimation to guide its search. Accelegy plug-ins allow for the customization of component energy estimation. These tools are available at <http://accelegy.mit.edu/tutorial.html>.

From an industry perspective, the cost constraints are related to volume and market; for instance, embedded processors have much more stringent cost limitations than processors in the cloud. One of the key factors that affects cost is the chip area (e.g., square millimeters) in conjunction with the process technology (e.g., 45-nm CMOS), which constrains the amount of on-chip storage as well as amount of compute (e.g., the number of PEs for DNN processors, the number of cores for CPUs and GPUs, the number of digital signal processing (DSP) engines for field-programmable gate arrays (FPGAs), and so on). To report information related to area without specifying a specific process technology, one can report the amount of on-chip memory (e.g., the storage capacity of the global buffer) and compute (e.g., the number of PEs) as a proxy for area.

Another important factor is the amount of off-chip BW, which dictates the cost and complexity of the packaging and printed circuit board (PCB) design [e.g., High Bandwidth Memory (HBM) [31] to connect to off-chip DRAM, NVLink [38] to connect to other GPUs, and so on] as well as whether additional chip area is required for a transceiver to handle

signal integrity at high speeds. The off-chip BW, which is typically reported in gigabits per second and the number of ports, can be used as a proxy for packaging and PCB cost.

There is also an interplay between the costs attributable to the chip area and off-chip BW. For instance, increasing on-chip storage, which increases chip area, can reduce off-chip BW. Accordingly, both metrics should be reported to provide perspective on the total cost of the system.

Of course, reducing cost alone is not the only objective. The design objective is, invariably, to maximize the throughput or energy efficiency for a given cost, specifically to maximize inferences per second per cost (e.g., U.S. dollars) and/or inferences per joule per cost. This is closely related to the previously discussed property of utilization; to be cost efficient, the design should aim to utilize every PE to increase inferences per second, since each PE increases the area and, thus, the cost of the chip; similarly, the design should aim to effectively utilize all of the on-chip storage to reduce off-chip BW or increase operations per off-chip memory access as expressed by the roofline model (see Figure 5), as each byte of on-chip memory also increases cost.

Flexibility

The merit of a DNN processor is also a function of its flexibility, which refers to the range of DNN models that can be supported on the DNN processor and the ability of the software environment (e.g., the mapper) to maximally exploit the capabilities of the hardware for any desired DNN model. Given the fast-moving pace of DNN research and deployment, it is increasingly important that DNN processors support a wide range of DNN models and tasks.

We can define *support* in two tiers. The first tier requires only that the hardware needs to be able to functionally support different DNN models (i.e., the DNN model can run on the hardware). The second tier requires that the hardware also maintain efficiency (i.e., high throughput and energy efficiency) across different DNN models.

To maintain efficiency, the hardware should not rely on certain properties of the DNN models to achieve efficiency, as the properties of DNN models are diverse and evolving rapidly. For instance, a DNN processor that can efficiently support the case where the entire DNN model (i.e., all of the weights) fits on chip may perform extremely poorly when the DNN model grows larger, which is likely, given that the size of DNN models continues to increase over time; a more flexible processor would be able to efficiently handle a wide range of DNN models, even those that exceed on-chip memory.

The degree of flexibility provided by a DNN processor presents a complex tradeoff with processor cost. Specifically, additional hardware usually needs to be added to flexibly support a wider range of workloads and/or improve their throughput and energy efficiency. Thus, the design objective is to reduce the overhead (e.g., area cost and energy consumption) of supporting flexibility while maintaining efficiency across the wide range of DNN models. Therefore, evaluating flexibility would entail ensuring that the extra

hardware is a net benefit across multiple workloads.

Flexibility has become increasingly important when we factor in the many techniques being applied to the DNN models with the promise to make them more efficient, since they increase the diversity of workloads that need to be supported. These techniques include DNNs with different network architectures (i.e., different layer shapes, which impact the amount of required storage and compute and the available data reuse that can be exploited), different levels of precision (i.e., different numbers of bits across layers and data types), and different degrees of sparsity (i.e., the number of zeros in the data). There are also different types of DNN layers and computations beyond MAC operations (e.g., activation functions) that need to be supported.

Actually getting a performance or efficiency benefit from these techniques invariably requires additional hardware. Again, it is important that the overhead of the additional hardware does not exceed the benefits of these techniques. This encourages a hardware and DNN model code-sign approach.

To date, exploiting the flexibility of DNN hardware has relied on mapping processes that act like static per-layer compilers. As the field moves to DNN models that change dynamically, mapping processes will need to dynamically adapt at runtime to changes in the DNN model or input data while still maximally exploiting the flexibility of the hardware to improve efficiency.

In summary, to assess the flexibility of DNN processors, their efficiency (e.g., inferences per second, inferences per joule) should be evaluated on a wide range of DNN models. The MLPerf benchmarking workloads are a good start; however, additional workloads may be needed to represent efficient techniques, such as efficient network architectures, as well as reduced precision and sparsity. The workloads should match the desired application. Ideally, since there

can be many possible combinations, it would also be beneficial to define the range and limits of DNN models that can be efficiently supported on a given platform (e.g., the maximum number of weights per filter or DNN model; minimum amount of sparsity; required structure of the sparsity; levels of precision, such as 8, 4, 2, or 1 b; types of layers and activation functions; and so on).

Scalability

Scalability has become increasingly important due to the wide use cases for DNNs. This is demonstrated by emerging technologies employed not just for scaling up the size of the chip but also for building systems with multiple chips (often referred to as *chiplets*) [32] or even wafer-scale chips [33]. *Scalability* refers to how well a design can be scaled up to achieve higher performance (i.e., latency and throughput) and energy efficiency when increasing the amount of resources (e.g., the number of PEs and on-chip storage). This evaluation is done under the assumption that the system does not have to be significantly redesigned (e.g., the design only needs to be replicated), since major design changes can be expensive in terms of time and cost. Ideally, a scalable design can be used for low-cost embedded devices and high-performance devices in the cloud simply by scaling up the resources.

Ideally, the performance would scale linearly and proportionally with the number of PEs. When the problem size (e.g., the batch size) is held constant, this is referred to as strong scaling and is the more challenging type of scaling. On the other hand, scaling performance while allowing the problem size to increase (e.g., by increasing batch size) is called weak scaling and is also an important objective in some situations.

Similarly, the energy efficiency would also improve with more on-chip storage; however, this would likely be nonlinear (e.g., increasing the on-chip storage such that the en-

tire DNN model fits on chip would result in an abrupt improvement in energy efficiency). In practice, this is often challenging due to factors, such as the reduced utilization of PEs and the increased cost of data movement due to long-distance interconnects.

Scalability can be connected with cost efficiency by considering how inferences per second per cost and inferences per joule per cost change with scale. For instance, if throughput increases linearly with the number of PEs (with proportional scaling of all storage), then the inferences per second per cost could be constant. It is also possible for the inferences per second per cost to improve superlinearly with an increasing number of PEs due to increased sharing of data across PEs. On the other hand, inferences per joule per cost might remain constant or even improve as a consequence of more sharing of data by multiple PEs.

In summary, to understand the scalability of a DNN processor design, it is important to report its performance and efficiency metrics as the number of PEs and storage capacity increase. This may include how well the design might handle technologies used for scaling up, such as interchip interconnect.

Interplay Among Different Metrics

It is important that all metrics be accounted for to fairly evaluate the design tradeoffs. For instance, without the accuracy given for a specific data set and task, one could run a simple DNN model and easily claim low power, high throughput, and low cost—however, the processor might not be usable for a meaningful task. Alternatively, without reporting the off-chip BW, one could build a processor with only MACs and easily claim low cost, high throughput, high accuracy, and low *chip* power—however, when evaluating *system* power, the off-chip memory access would be substantial. Finally, the test setup should also be reported, including whether the results are

measured or obtained from simulation and how many images were tested. If obtained from simulation, it should be clarified whether it was after synthesis or place-and-route and which library corner (e.g., process corner, supply voltage, temperature) was used.

Clearly, there are many important metrics to consider when designing DNN processors. At the same time, the design space for DNN processors is very large. As a result, it would be helpful to have the ability to rapidly explore the design space early in the design process and accurately estimate the various metrics for the proposed designs. An accurate estimation requires proper consideration of how the properties of the hardware, such as mapping, and properties of the workload, such as DNN model shape, precision, and sparsity, impact metrics, such as throughput and energy efficiency. One example toolset that allows one to perform rapid exploration and evaluation is the combination of Timeloop [13] and Accelergy [34], as depicted in Figure 9. Both tools accept a template-based specification of a proposed architecture, and, given a DNN model description, Timeloop searches for an optimal mapping using an analytical performance model and generates activity counts that allow Accelergy to generate architecture-level energy estimates (all before a detailed register-transfer-level (RTL) description of the design is available). Accelergy also accepts component energy costs, which is especially useful for understanding the impact of new technologies [35].

Summary

In this article, we presented the various key metrics for evaluating DNN processors, discussed the importance of each metric and their interrelationships, and, where appropriate, included equations that can be used to tease apart the factors that contribute to those metrics. We also showed how those metrics are related to both the hardware design and the DNN models and highlight-

ed why hardware/model codesign is important. Finally, given those metrics, the evaluation process for whether a DNN system is a viable solution for a given application might go as follows.

- The accuracy determines if the system can perform the given task.
- The latency and throughput determine whether it can run fast enough and in real time.
- The energy and power consumption primarily dictate the form factor of the device where the processing can operate.
- The cost, which is primarily dictated by the chip area and external memory BW requirements, determines how much one would pay for the solution.
- The flexibility determines the range of tasks it can support.
- The scalability determines whether the same design effort can be amortized for deployment in multiple domains (e.g., in the cloud and at the edge) and if the system can efficiently be scaled with DNN model size.

Portions of this article are based on our book, *Efficient Processing of Deep Neural Network* [36]. This excerpt has described various metrics that are important for evaluating DNN processors. The remainder of the book expands on how to design DNN processors and DNN models that optimize for these metrics.

Acknowledgments

This work was funded in part by DARPA YFA, DARPA contract HR0011-18-3-0007, the MIT Center for Integrated Circuits and Systems (CICS), the MIT-IBM Watson AI Lab, the MIT Quest for Intelligence, NSF E2CDA 1639921, and gifts/faculty awards from Nvidia, Facebook, Google, Intel, and Qualcomm.

References

- [1] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng, "Unsupervised learning of hierarchical representations with convolutional deep belief networks," *Commun. ACM*, vol. 54, no. 10, pp. 95–103, 2011. doi: 10.1145/2001269.2001295.
- [2] O. Russakovsky et al., "ImageNet large scale visual recognition challenge," *Int. J.*

- Comput. Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- [3] Y. LeCun, C. Cortes, and C. J. C. Burges, "The MNIST database of handwritten digits." Accessed on: July 2020. [Online]. Available: <http://yann.lecun.com/exdb/mnist/>
- [4] Standard Performance Evaluation Corporation (SPEC). Accessed on: July 2020. [Online]. Available: <https://www.spec.org/>
- [5] MLPerf. Accessed on: July 2020. [Online]. Available: <https://mlperf.org/>
- [6] "DeepBench," GitHub. [Online]. Accessed on: July 2020. Available: <https://github.com/baidu-research/DeepBench>
- [7] R. Adolf, S. Rama, B. Reagen, G.-Y. Wei, and D. Brooks, "Fathom: Reference workloads for modern deep learning methods," in *Proc. Int. Symp. Workload Characterization (IISWC)*, 2016. doi: 10.1109/IISWC.2016.7581275.
- [8] J. D. Little, "A proof for the queuing formula: $L = \lambda w$," *Oper. Res.*, vol. 9, no. 3, pp. 383–387, 1961. doi: 10.1287/opre.9.3.383.
- [9] Y.-H. Chen, T.-J. Yang, J. Emer, and V. Sze, "Eyeriss v2: A flexible accelerator for emerging deep neural networks on mobile devices," *IEEE J. Emerg. Sel. Topics Circuits Syst.*, vol. 9, no. 2, pp. 292–308, 2019. doi: 10.1109/JETCAS.2019.2910232.
- [10] H. Kwon, A. Samajdar, and T. Krishna, "MAERI: Enabling flexible dataflow mapping over DNN accelerators via reconfigurable interconnects," in *Proc. Architectural Support Programming Languages and Operating Systems (ASPLOS)*, 2018, pp. 461–475. doi: 10.1145/3173162.3173176.
- [11] Y.-H. Chen, J. Emer, and V. Sze, "Using dataflow to optimize energy efficiency of deep neural network accelerators," *IEEE Micro.*, vol. 37, no. 3, pp. 12–21, May–June 2017. doi: 10.1109/MM.2017.54.
- [12] Y.-H. Chen, T. Krishna, J. Emer, and V. Sze, "Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks," *IEEE J. Solid-State Circuits*, vol. 51, no. 1, pp. 127–138, 2017. doi: 10.1109/JSSC.2016.2616357.
- [13] A. Parashar et al., "Timeloop: A systematic approach to DNN accelerator evaluation," in *Proc. Int. Symp. Performance Analysis Systems and Software (ISPASS)*, 2019, pp. 304–315. doi: 10.1109/ISPASS.2019.00042.
- [14] H. Kwon, P. Chararasi, M. Pellauer, A. Parashar, V. Sarkar, and T. Krishna, "Understanding reuse, performance, and hardware cost of DNN dataflow: A data-centric approach," in *Proc. Int. Symp. Microarchitecture (MICRO)*, 2019, pp. 754–768. doi: 10.1145/3352460.3352522.
- [15] Y.-H. Chen, J. Emer, and V. Sze, "Eyeriss: A spatial architecture for energy-efficient dataflow for convolutional neural networks," in *Proc. Int. Symp. Computer Architecture (ISCA)*, 2016, pp. 367–379. doi: 10.1109/ISCA.2016.40.
- [16] S. Williams, A. Waterman, and D. Patterson, "Roofline: An insightful visual performance model for multicore architectures," *Commun. ACM*, vol. 52, no. 4, pp. 65–76, Apr. 2009. doi: 10.1145/1498765.1498785.
- [17] A. G. Howard et al., "Mobilenets: Efficient convolutional neural networks for mobile vision applications. 2017. [Online]. Available: [arXiv:1704.04861](https://arxiv.org/abs/1704.04861)
- [18] J. Albericchio, P. Judd, T. Hetherington, T. Aamodt, N. E. Jerger, and A. Moshovos, "Cnvlutin: Ineffectual-neuron-free deep neural network computing," in *Proc. Int. Symp. Computer Architecture (ISCA)*, 2016, pp. 1–13. doi: 10.1109/ISCA.2016.11.

[19] S. Han et al., "EIE: Efficient inference engine on compressed deep neural network," in *Proc. Int. Symp. Computer Architecture (ISCA)*, 2016. doi: 10.1145/3007787.3001163.

[20] A. Parashar et al., "SCNN: An accelerator for compressed-sparse convolutional neural networks," in *Proc. Int. Symp. Computer Architecture (ISCA)*, 2017, pp. 27–40. doi: 10.1145/3079856.3080254.

[21] V. Camus, L. Mei, C. Enz, and M. Verhelst, "Review and benchmarking of precision-scalable multiply-accumulate unit architectures for embedded neural-network processing," *IEEE J. Emerg. Select. Topics Circuits Syst.*, vol. 9, no. 4, pp. 697–711, 2019. doi: 10.1109/JETCAS.2019.2950386.

[22] T.-J. Yang, Y.-H. Chen, and V. Sze, "Designing energy-efficient convolutional neural networks using energy-aware pruning," in *Proc. Conf. Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6071–6079. doi: 10.1109/CVPR.2017.643.

[23] T.-J. Yang et al., "NetAdapt: Platform-aware neural network adaptation for mobile applications," in *Proc. European Conf. Computer Vision (ECCV)*, 2018, pp. 289–304. doi: 10.1007/978-3-030-01249-6_18.

[24] M. Tan et al., "MnasNet: Platform-aware neural architecture search for mobile," in *Proc. Conf. Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 2815–2823. doi: 10.1109/CVPR.2019.00293.

[25] T.-J. Yang and V. Sze, "Design considerations for efficient deep neural networks on processing-in-memory accelerators," in *Proc. Int. Electron Devices Meeting (IEDM)*, 2019, pp. 22.1–22.4. doi: 10.1109/IEDM19573.2019.8993662.

[26] B. Chen and J. M. Gilbert, "Introducing the CVPR 2018 On-device Visual Intelligence Challenge Google AI blog," Apr. 20, 2018. [Online]. Available: <https://ai.googleblog.com/2018/04/introducing-cvpr-2018-on-device-visual.html>

[27] M. Gao, J. Pu, X. Yang, M. Horowitz, and C. Kozyrakis, "TETRIS: Scalable and efficient neural network acceleration with 3D memory," in *Proc. Architectural Support Programming Languages and Operating Systems (ASPLOS)*, 2017. doi: 10.1145/3037697.3037702.

[28] S. Yu, "Neuro-inspired computing with emerging nonvolatile memory," *Proc. IEEE*, vol. 106, no. 2, pp. 260–285, 2018. doi: 10.1109/JPROC.2018.2790840.

[29] N. Verma et al., "In-memory computing: Advances and prospects," *IEEE Solid State Circuits Mag.*, vol. 11, no. 3, pp. 43–55, 2019. doi: 10.1109/MSSC.2019.2922889.

[30] M. Horowitz, "Computing's energy problem (and what we can do about it)," in *Proc. Int. Solid-State Circuits Conf. (ISSCC)*, 2014, pp. 1–4. doi: 10.1109/ISSCC.2014.6757323.

[31] *High Bandwidth Memory (HBM) DRAM*, JESD235, 2013.

[32] Y. S. Shao et al., "Simba: Scaling deep-learning inference with multi-chip-module-based architecture," in *Proc. IEEE/ACM Int. Symp. Microarchitecture (MICRO)*, 2019, pp. 14–27. doi: 10.1145/3352460.3358302.

[33] S. Lie, "Wafer scale deep learning," in *Proc. IEEE Hot Chips 31 Symp. (HCS)*, 2019.

[34] Y. N. Wu, J. S. Emer, and V. Sze, "Accelergy: An architecture-level energy estimation methodology for accelerator designs," in *Proc. IEEE/ACM Int. Conf. Computer Aided Design (ICCAD)*, 2019, pp. 1–8. doi: 10.1109/ICCAD45719.2019.8942149.

[35] Y. N. Wu, V. Sze, and J. S. Emer, "An architecture-level energy and area estimator for processing-in-memory accelerator designs," in *Proc. Int. Symp. Performance*

Analysis Systems and Software (ISPASS), 2020.

[36] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, *Efficient Processing of Deep Neural Networks* (Synthesis Lectures in Computer Architecture series). San Rafael, CA: Morgan & Claypool, 2020.

[37] L. A. Barroso and U. Hölzle, "The case for energy-proportional computing," *Computer*, vol. 40, no. 12, pp. 33–37, Dec. 2007. doi: 10.1109/MC.2007.443.

[38] "NVIDIA NVLink: High-speed GPU interconnect," NVIDIA, Santa Clara, CA. [Online]. Available: <https://www.nvidia.com/en-us/design-visualization/nvlink-bridges/>

About the Authors

Vivienne Sze (sze@mit.edu) is an associate professor in the Electrical Engineering and Computer Science Department, Massachusetts Institute of Technology. She is a recipient or corecipient of various awards, including the AFOSR and DARPA Young Faculty Award; the Edgerton Faculty Award; faculty awards from Google, Facebook, and Qualcomm; the Symposium on VLSI Circuits Best Student Paper Award; the IEEE Custom Integrated Circuits Conference Outstanding Invited Paper Award; and the IEEE Micro Top Picks Award. Her research interests include computing systems that enable energy-efficient machine learning, computer vision, and video compression/processing for a wide range of applications, including autonomous navigation, digital health, and the internet of things. She is a Senior Member of the IEEE.

Yu-Hsin Chen (allenhsin@alum.mit.edu) received his B.S. degree in electrical engineering from National Taiwan University, Taipei, in 2009 and his M.S. and Ph.D. degrees in electrical engineering and computer science from the Massachusetts Institute of Technology (MIT), Cambridge, in 2013 and 2018, respectively. He is currently a research scientist at Facebook, focusing on hardware/software codesign to enable on-device artificial intelligence for augmented/virtual reality systems. Previously, he was a research scientist in Nvidia's Architecture Research Group. He was the recipient of the 2018 Jin-Au Kong Outstanding Doctoral Thesis Prize in Electrical Engineering at MIT, 2015 Nvidia Graduate Fellowship, 2015 ADI Outstanding Student Designer Award, and 2017 IEEE

SSCS Predoctoral Achievement Award. His work on the dataflows for CNN accelerators was selected as one of the Top Picks in Computer Architecture in 2016. He is a Member of the IEEE.

Tien-Ju Yang (tjy@mit.edu) received his B.S. degree in electrical engineering in 2010 and his M.S. degree in electronics engineering in 2012 from National Taiwan University. He is currently a Ph.D. degree candidate in electrical engineering and computer science at the Massachusetts Institute of Technology, working on efficient deep neural network design. His research interests include the areas of deep learning, computer vision, machine learning, image/video processing, and very-large-scale integration system design. He won first place in the 2011 National Taiwan University Innovation Contest. He also cotaught a tutorial, Efficient Image Processing With Deep Neural Networks, at the 2019 IEEE International Conference on Image Processing. He is a Student Member of the IEEE.

Joel S. Emer (jsemer@mit.edu) is a senior distinguished research scientist with Nvidia's Architecture Research Group and a Professor of the Practice at the Massachusetts Institute of Technology. He was with Intel, where he was an Intel fellow and the director of microarchitecture research. At Intel, he led the VSSAD Group, of which he had previously been a member at Compaq and Digital Equipment Corporation. He has made architectural contributions to a number of VAX, Alpha, and X86 processors and has contributed to the quantitative approach to processor performance evaluation, simultaneous multithreading technology, processor reliability analysis, cache organization, and spatial architectures for deep learning. He is a Fellow of ACM and member of the National Academy of Engineering. He received the Eckert–Mauchly Award and ECE alumni awards from Purdue University and the University of Illinois. He has had six papers selected for the IEEE Micro's Top Picks in Computer Architecture. He is a Fellow of the IEEE.

