

BENCHIP: Benchmarking Intelligence Processors

Jinhua Tao¹, Zidong Du^{1,2}, Qi Guo^{1,2}, Huiying Lan¹, Lei Zhang¹
Shengyuan Zhou¹, Lingjie Xu³, Cong Liu⁴, Haifeng Liu⁵, Shan Tang⁶
Allen Rush⁷, Willian Chen⁷, Shaoli Liu^{1,2}, Yunji Chen¹, Tianshi Chen^{1,2}

¹ICT CAS, ²Cambricon, ³Alibaba Infrastructure Service, Alibaba Group
⁴IFLYTEK, ⁵JD, ⁶RDA Microelectronics, ⁷AMD

Abstract

The increasing attention on deep learning has tremendously spurred the design of intelligence processing hardware. The variety of emerging intelligence processors requires standard benchmarks for fair comparison and system optimization (in both software and hardware). However, existing benchmarks are unsuitable for benchmarking intelligence processors due to their non-diversity and nonrepresentativeness. Also, the lack of a standard benchmarking methodology further exacerbates this problem. In this paper, we propose BENCHIP, a benchmark suite and benchmarking methodology for intelligence processors. The benchmark suite in BENCHIP consists of two sets of benchmarks: *microbenchmarks* and *macrobenchmarks*. The microbenchmarks consist of single-layer networks. They are mainly designed for bottleneck analysis and system optimization. The macrobenchmarks contain state-of-the-art industrial networks, so as to offer a realistic comparison of different platforms. We also propose a standard benchmarking methodology built upon an industrial software stack and evaluation metrics that comprehensively reflect the various characteristics of the evaluated intelligence processors. BENCHIP is utilized for evaluating various hardware platforms, including CPUs, GPUs, and accelerators. BENCHIP will be open-sourced soon.

1 Introduction

1.1 Motivation

Recently, *deep learning* has become a *de facto* technique for intelligent processing tasks such as image classification [1, 2, 3], video captioning [4], speech recognition [5], and machine translation [6]. As deep learning architectures (i.e., artificial neural networks) are evolving towards deep topologies with complicated transformations, existing general-purpose hardware platforms such as CPUs and GPUs are not able to provide high performance and energy efficiency. As an efficient alternative for deep learning, a large variety of customized hardware architectures, ranging from specialized GPUs (e.g., NVIDIA DGX-1 [7]) to FPGAs (e.g., CNP [8] and [9]) to ASICs (e.g., DianNao [10], NeuFlow [11], and EIE [12]) have emerged. We call those specially designed architectures *intelligence processors* (IPs)¹.

Benchmarking has served long as the foundation of designing new hardware architectures. For example, the SPEC-CPU series [13] and the corresponding methodology (e.g., framework, procedure, and metrics, etc.) have been the incontestable source for evaluating and optimizing general-purpose architecture, and they have evolved over time to keep pace with the rapid advance of uncore architecture. With multicore architectures dominating the market, the PARSEC benchmark suite [14] was released in 2008 accordingly. Recently, especially in last one or two years, intelligence processors have become a hot topic. However, with the chaos of benchmarks, it is difficult to compare and quantize the various IPs, especially the rapidly increasing works emerging at top architecture conferences [12, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30]. Still, fair and comprehensive benchmarking is a necessity for assessing the progress of IPs. Moreover, with the help of benchmarking, both architects and practitioners can compare various architectures, identify their bottlenecks, and conduct corresponding system/architectural optimization. As a result, during the design and deployment of intelligence processors, one of the central tasks is benchmarking.

Benchmark suite. When benchmarking intelligence processors, designing an appropriate benchmark suite is the top-most consideration. Roughly, two main categories of benchmarks already exist: the *collective* benchmarks, which are collected and documented as a benchmark suite (e.g., BenchNN [31], DeepBench [32], and Fathom [33]), and the *personalized* benchmarks, which have been used to evaluate

¹Both CPUs and GPUs can be viewed as intelligence processors from a broader perspective, as they remain the mainstream of intelligent processing.

different intelligence processors in various computer architecture papers.

BenchNN illustrates the potential of neural networks by using them to re-implement five tasks from the well-known PARSEC benchmark suite to reveal the broad application scope of neural networks. Though BenchNN successfully demonstrated an urgent need for neural networks, it is more like a symbol rather than a practical benchmark suite. To be specific, there are three limitations: (1) BenchNN provides a very limited number of neural network applications (i.e., five in total), (2) the neural networks used in BenchNN are classic models such as simple multi-layer perceptrons that fails to characterize the latest trend reflected by state-of-the-art deep learning techniques, and (3) BenchNN does not reflect the significant diversity among neural networks—for example, both the *ferret* and *streamcluster* benchmarks are built on the Self Organizing Maps [34].

DeepBench aims at measuring the performance of basic deep learning operations across various hardware platforms. Such operations include matrix multiplication, convolution, recurrent layer, and all-reduce operations. As DeepBench only contains basic primitive operations, it cannot approximate high-level behaviors of the full-stack neural networks widely used in industry. For example, the Faster-RCNN [35] used for object detection contains far more operations than the basic operations in DeepBench. In short, due to its inability to evaluate of the full-stack neural networks used in practice, DeepBench is not suitable for benchmarking intelligence processors.

Fathom focuses on better understanding a broader variety of deep learning workloads. However, with an assembled set of only eight archetype deep learning workloads, Fathom has limited capability for benchmarking intelligence processors. First, Fathom fails on diversity, as it only contains common layers such as *Conv.* and *Fully-Connected* but lacks important layers such as *Deconv.*, *Unpooling* and *Batch normalization*. Second, Fathom fails to provide a benchmarking method for several of the most important categories of intelligence processors, e.g., customized hardware architectures. Third, Fathom cannot benchmark fine-grained performance or efficiency, which are important for intelligence processors, especially customized hardware. Thus, Fathom is also not suitable for benchmarking intelligence processors.

The *personalized* benchmarks are the neural networks that appear in computer architecture papers but not as a benchmark suite for evaluating intelligence processors. Such benchmarks include both basic operations (e.g., convolution, activation, and normalization) and full neural networks (e.g., AlexNet [1], VGG [2], and ResNet [3]). As such benchmarks vary for evaluating different architectures, they may suffer from non-diversity and nonrepresentativeness. Figure 1(a) shows that

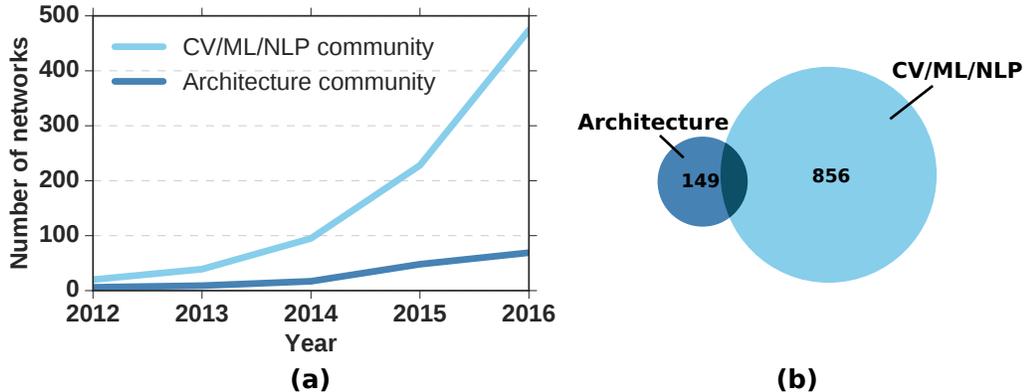


Figure 1: (a): The number of neural network models proposed in CV/ML/NLP top-tier conferences (e.g., CVPR/ICML/EMNLP) has increased rapidly, while the number of neural networks utilized in the computer architecture community (e.g., ISCA/MICRO/HPCA/ASPLOS) exhibits relatively moderate growth. (b): The neural networks used in the computer architecture community only cover a small range of the neural networks used in applications from the AI community.

there exists a large gap between the neural networks proposed in CV (computer vision)/ML(machine learning)/NLP (natural language processing) applications and those used for evaluating hardware architectures. We also notice the increasing capability of performing multiple types of neural networks in the architecture community. Interestingly, from Figure 1 (b), we observe that more than half of the personalized benchmarks utilized in computer architecture papers are *not* practically used in AI research fields (e.g., CV/ML/NLP). Even when we assume a one-year delay for transferring knowledge from the AI community to the computer architecture community, this number is still larger than 30%. This observation clearly shows a risk that conclusions made by the computer architecture community might mislead the AI community.

In short, due to their non-diversity and nonrepresentativeness, existing benchmarks are not suitable for benchmarking intelligence processors.

Benchmarking methodology. A benchmark suite only specifies *what* to measure. We also need to determine *how* to benchmark intelligence processors, i.e., develop a benchmarking methodology. The existing benchmarking methodology utilized in the computer architecture community is chaotic, mainly due to the complicated software/hardware stacks of modern computer systems. Given the prototypical nature of a newly published intelligence architecture, a full software stack may not exist to conduct a fair comparison against commodity hardware. As a result, some evaluate their architectures with a high-level programming framework [22, 26], while others are evaluated directly at the RTL level by getting rid of the software stack [12]. Since the software stack in real applications is very

crucial to overall performance, comparison at different hardware/software levels seems to be a case of comparing “apples and oranges.” To alleviate this problem, a benchmark specification for the evaluation framework, metrics, and procedure should be clearly defined.

Moreover, intelligence processors have unique features compared with general-purpose architectures that pose further challenges for the methodology. One such feature is that the prediction accuracy is also an important design tradeoff in addition to performance/power/area. While it could be a reasonable solution to slightly relax the prediction accuracy of machines learning models (via adopting, e.g., data quantization [36] or connection sparsity [37]) for performance and energy efficiency, excessive relaxation on the prediction accuracy will lead to unfair and degenerate comparisons. Therefore, it is necessary to take the accuracy, together with other traditional design tradeoffs, into consideration when designing the benchmarking methodology.

1.2 Our work

Based on the above observations, in this paper, we propose BENCHIP, a novel benchmark suite and benchmarking methodology deliberately designed for intelligence processors. The benchmark suite in BENCHIP contains two types of benchmarks: *microbenchmarks* and *macrobenchmarks*. The microbenchmarks contain 12 representative single-layer networks, such as convolution, pooling, activation, etc., that are mainly used for architectural/system optimization. The macrobenchmarks consist of 11 commonly-used full neural networks (e.g., AlexNet, VGG, and Faster-RCNN) from different intelligent processing scenarios. These are directly extracted from real applications and they are mainly used for the evaluation and comparison of hardware platforms.

Our benchmarking methodology is mainly designed for creditability, portability, and fairness. The creditability is achieved by building an industrial software stack consisting of a high-level programming model, library, and device driver. The portability is accomplished by providing a standard interface for the intermediate-level high-performance library. Thus, with a new architecture, once the corresponding library complies with the standard interface, it can be easily plugged into the software stack for evaluation. For the fairness purpose, we report evaluation metrics that take the accuracy into consideration in order to comprehensively reflect the various characteristics of intelligence processors.

Contributions. The work in this paper makes the following key contributions:

Table 1: The benchmarks in BENCHIP

| Layer | # Config. | Note | NN | Dataset | Scenario |
|-----------|------------------------|----------|-----------------------|----------------------|---------------------------------|
| Conv. | 7 (3/1/3) ¹ | | LeNet-5 [38] | MNIST [38] | Hand written digits recognition |
| Pooling | 7/7 (3/1/3) | Avg./Max | RNN [39] | WSJ [40] | Speech recognition |
| FC | 7 (3/1/3) | | AlexNet [1] | ImageNet [41] | Image classification |
| ReLU | 7 (3/1/3) | | VGG [2] | ImageNet [41] | Image classification |
| Sigmoid | 7 (3/1/3) | | ResNet [42] | ImageNet [41] | Image classification |
| LRN | 7 (3/1/3) | | Faster R-CNN [35] | PASCAL VOC 2012 [43] | Object recognition |
| BN | 7 (3/1/3) | | Deep Face Recog. [44] | LFW [45] | Face recognition |
| Unpooling | 7/7 (3/1/3) | Avg./Max | DeconvNet [46] | PASCAL VOC 2012 [43] | Semantic segmentation |
| Deconv. | 7 (3/1/3) | | FCLN [47] | Visual Genome [47] | Image captioning |
| LSTM | 7 (3/1/3) | | S2VT [48] | MSVD [49] | Video captioning |
| | | | SyntaxNet [50] | English WSJ [40] | Nature language processing |

¹Total configurations (Normal configurations/Extreme small configurations/Extreme large configurations)

- We build a benchmark suite consisting of microbenchmarks and macrobenchmarks to evaluate the accuracy, performance, and energy of intelligence processors.
- We conduct a comprehensive analysis of the proposed benchmarks in order to demonstrate that the proposed benchmark suite is representative and diverse.
- We propose a benchmarking methodology that contains an industrial-level software stack and evaluation metrics for guaranteeing creditability, portability, and fairness.
- We evaluate various intelligence processors, including CPUs, GPUs, and accelerators, with BENCHIP.

2 The Benchmark Suite

In this section, we first introduce the requirements of an ideal benchmark suite for intelligence processors. We next present our benchmarks in detail. Finally, we conduct an analysis to demonstrate the rationality of the proposed benchmarks.

2.1 Design requirements

An ideal benchmark suite for intelligence processors should meet the following requirements from both the application and architecture perspectives.

- **Application perspective.** The selected benchmarks should include mainstream neural network algorithms from a broad range of application scenarios. More-

over, in order to reduce benchmarking efforts, benchmarks with similar characteristics should not be included redundantly [51].

- **Architecture perspective.** In addition to the common usage scenarios, the selected benchmarks should be able to explore the processing boundary of underlying architectures, as well as address the future trends for intelligence processors.

2.2 Benchmarks

Based on the above requirements, we construct both *microbenchmarks* and *macrobenchmarks* in BENCHIP, particularly for system optimization and platform comparison, respectively. In the context of deep learning, the microbenchmarks only include single-layer networks, while the macrobenchmarks contain entire neural networks extracted from real applications. For the microbenchmarks, individual layers are equipped with multiple configurations, including normal and extreme cases for stress-testing intelligence processors. For the macrobenchmarks, entire neural networks are used for processing end-to-end intelligent tasks where data movements between layers are also crucial.

2.2.1 Microbenchmarks

Microbenchmarks are widely used for system and architectural optimization in both industry and academia. These benchmarks measure a specific component of the system. For example, STREAM is a well-known microbenchmark for evaluating the memory bandwidth of the system [52]. The University of Edinburgh provides the OpenMP microbenchmark suite for measuring the overheads of synchronization, loop scheduling, and array operations in the OpenMP runtime library [53].

In the context of deep learning, the microbenchmarks refer to single-layer networks. Unlike DeepBench, we do not select basic operations as microbenchmarks because basic operations cannot benchmark the efficiency of data movements. We include 10 single-layer networks, as listed in Table 1. In order to fulfill the requirements mentioned in Section 2.1, the principle is threefold: the selected layers should be widely used in existing neural networks (e.g., *convolutional layer*); such layers should have significant different computational, memory, and control patterns; and such layers may have huge impacts on future designs (e.g., *deconvolution layer*).

Table 2: Statistics of the different types of layers in CVPR, EMNLP, ICML, ICRA, and NIPS (2012~2016).

| Layer Type | # Presence | Ratio(%) | Network | # Presence | Ratio(%) |
|----------------------|------------|----------|---------------------|------------|----------|
| Conv. | 452 | 19.28 | LeNet | 17 | 1.99 |
| Pooling ¹ | 440 | 18.76 | AlexNet | 159 | 18.57 |
| FC | 473 | 20.17 | VGG | 171 | 19.98 |
| LRN | 181 | 7.72 | ResNet | 5 | 0.58 |
| ReLU | 425 | 18.12 | Faster-RCNN | 34 | 3.97 |
| Sigmoid | 90 | 3.84 | CNN | 128 | 14.95 |
| Tanh | 100 | 4.26 | DNN | 82 | 9.58 |
| Deconv. | 19 | 0.81 | RNN | 260 | 30.37 |
| Unpooling | 16 | 0.68 | Others ² | 30 | 1.28 |
| LSTM | 133 | 5.67 | | | |
| BN | 16 | 0.68 | | | |
| Total | 2345 | - | Total | 856 | - |

¹Average Pooling + Max Pooling

²RBM, SNN, and other ML techniques

Conv. The convolutional layer is the most important layer in deep learning, especially for images/video tasks as it automatically extracts features from the input images/videos (2D/3D input data) by applying a set of 2D/3D filters. Regarding *representativeness*, we study network layers employed in application scenarios from the AI community (e.g., the CV/NLP/ML community). More specifically, we count the appearance time of a layer/network (# Presence) in the papers of such fields and thus the appearance ratio to the total layers/networks (# Ratio). These results are shown in Table 2. From the Table, we can see that convolutional layer appears in 19.28% of all layers.

Pooling. In deep learning, the pooling layer is required to aggregate features for later classification by downsampling 2D input data. As shown in Table 2, the pooling layer is as important as the convolutional layer, having a 18.76% appearance ratio among all layers.

FC. The fully-connected (FC) layer is the most common layer, where each output neuron is connected to all input neurons with independent synapses. According to Table 2, it has been employed in most application scenarios, with an appearance ratio as high as 20.17%.

ReLU. Nonlinear activation layer is a key component as it introduces nonlinearity into the existing deep neural networks to improve the capability of classification. The inputs and outputs can be treated as one-dimensional data with element-wise operations. Rectified Linear Units (ReLU) is becoming popular (with an appear-

ance ratio of 18.12%) as it achieves better accuracy.

Sigmoid. The Sigmoid layer is another type of activation layer. While the Sigmoid layer is becoming less popular, the appearance ratio is still about 3.84% of all employed network layers. Both ReLU and Sigmoid are typical element-wise operations but have very different functionality. Thus, we include both as microbenchmarks.

LRN. The normalization layer introduces competition between input neurons at the same position on different input feature maps, which have significantly different computation patterns from other layers. There exists two typical normalization layers: Local Response Normalization (LRN) [1] and Local Contrast Normalization (LCN) [54]. We include LRN as the representative layer for its relative high appearance ratio (7.72%) compared to many other layers (including LCN). In an LRN layer, output neurons are the results of input neurons divided by responses among neighboring feature maps in the same spatial location.

BN. Batch normalization performs normalization to each mini-batch, which allows networks to use larger learning rate and thus to accelerate neural network convergence in training. It also helps reduce the dependency on weight initialization. We select it as a microbenchmark due to the increasing trend of using BN to accelerate the training.

Deconv. The deconvolutional layer is emerging in the fields of feature extraction, image reconstruction, and image semantic segmentation as it increases the size of inputs for better matching and locating. It can be viewed as a reverse operation of convolution. We include it as a microbenchmark for its unique computational pattern and the increasing popularity of employing such a layer in emerging applications (e.g., semantic segmentation).

Unpooling. The unpooling layer has a purpose and effect similar to the deconv layer. It performs the reverse operation of pooling and reconstructs a larger size of activations. The unpooling operations have gained increasing attentions in recent years, as semantic segmentation and image reconstruction have become more popular.

LSTM. The LSTM layer is a variant of the recurrent layer, which prevents the gradient from vanishing or exploding. It contains multiple element-wise gate operations and active operations (i.e., Sigmoid and Tanh). It is selected for its high appearance ratio of 5.67% and also because of the representative of recurrent computing patterns.

Summary. In Table 2, almost all appearing layers are included as microbenchmarks except for Tanh, which is mostly used as a part of the LSTM layer. Thus, it is fair enough to say that our chosen microbenchmarks for BENCHIP has been

thorough and rational.

2.2.2 Macrobenchmarks

Macrobenchmarks indicate workloads that are directly extracted from real industrial applications for evaluating and comparing different platforms—for example, in SPEC CPU-2006 [13], all benchmarks are collected and documented from desktop software for approximating real application scenarios. The PARSEC benchmark suite comprises emerging workloads from various areas, e.g., computer vision, financial analytics, and animation physics [14] and is aimed at benchmarking dominant chip multiprocessors.

In the context of deep learning, macrobenchmarks refer to entire neural networks such as AlexNet, VGG, and Faster-RCNN. We selected the 11 entire networks, as listed in Table 1. There are three reasons for having entire networks as macrobenchmarks: (1) compared with single-layer networks in microbenchmarks, entire networks are commonly used for performing end-to-end tasks such as speech recognition; (2) using combinations of individual layers for evaluation is not sufficient because factors such as the data exchange when coupling two layers are critical to the performance and energy efficiency, especially as networks become deeper and larger with increasing processing complexity; and (3) the optimization between layers cannot be measured with microbenchmarks, which are only capable of evaluating inter-layer optimizations.

LeNet-5. LeNet-5 is a notable convolutional neural network used to classify handwritten and machine-printed characters. It is selected as the representation of handwritten digits recognition.

RNN. RNN is a network for end-to-end speech recognition. It consists of a deep bidirectional LSTM recurrent neural network and a Connectionist Temporal Classification (CTC) objective function. We select it to represent the speech recognition scenario.

AlexNet. AlexNet is a deep convolutional neural network designed for the classification task of ImageNet. We select it as it is a very popular and widely used network in both architecture and AI communities, with appearance ratios of 33.54% and 18.57%, respectively (See Table 2).

VGG. VGG is a series of networks with very deep architecture designed for image classification. It uses very small kernels (e.g., 3×3 and 1×1) for convolutional layers through the whole network. It is also a very popular network with the appearance ratio of 36.08% in architecture community and 19.98% in AI community.

ResNet. ResNet refers to the deep residual networks. The residual learning framework allows it to achieve extremely deep architecture, and it gains high accuracy from the increased depth. It uses “shortcut connections” to skip one or more layers to ease the degradation problem. We select ResNet because it is the state-of-the-art network for image classification.

Faster R-CNN. Faster R-CNN is an object detection system with high performance and detection accuracy. It is composed of two modules: a region proposal network (RPN) for generating region proposals, and a fast R-CNN detection network that uses such proposals for detection. We select it as the representation of object detection network.

Deep face recognition. Deep face recognition utilizes deep CNN networks to classify and verify face images. The networks are derived from the VGG architectures with the last FC layer replaced by a classifier or an embedding descriptor vector. We select it as the representation of face recognition scenario.

DeconvNet. DeconvNet is a deconvolutional network for semantic segmentation that consists of two subnetworks: convolution and deconvolution. We choose this network due to the deconvolutional and unpooling layers it contains, and as the representation of semantic segmentation scenario.

FCLN. FCLN is a learning system for image caption tasks, which is composed of a convolution network and a recurrent network. The input images are first processed by a VGG16-based convolution network and then fed into the localization layer to identify regions of interest. These regions are input into the recurrent language model to generate captions. We selected FCLN as the representative network of image caption scenario.

S2VT. S2VT is a sequence-to-sequence LSTM model trained on video-sentence pairs for generating descriptive natural language text of video clips. It uses a VGG16-based network to extract features from the raw image and then feeds the features into the LSTM network, which decodes the representation into a sequence of words. We selected S2VT as the representation of video caption scenario.

SyntaxNet. SyntaxNet is a globally normalized transition-based feed-forward network for part-of-speech tagging, dependency parsing, and sentence compression. It at core is a transition-based parser, where neural network is adopted for computing the score of decisions in certain states. We select it to represent the natural language processing scenario.

Summary. Macrobenchmarks include well-known neural networks, e.g, *LeNet-5*, *AlexNet* and *VGG*, as well as typical networks such as *RNN*. Thus, BENCHIP has both representativeness and diversity as it well covers 98.72% of the networks that appeared in application domains in the last five years. Especially, note that

most authors of this paper are from both hardware and AI industry/academia, and the selected benchmarks are commonly used in our daily production environments.

2.3 Benchmark analysis

In this section, we conduct a comprehensive analysis of the selected benchmarks to demonstrate that they can meet the requirements from both the application and architecture perspectives. The basic intuition is to characterize the benchmarks using architecture-independent characteristics. As the macrobenchmarks are composed of single-layer networks—for example, the microbenchmark *Conv.* layer provides a configuration that is exactly the second layer in the macrobenchmark *VGG*—we focus on analyzing the microbenchmarks.

In Figure 2, we still report the similarity of macrobenchmarks to better understand the variation of macrobenchmarks. We use operation amounts in different layers as feature vectors and Euclidean distances to hierarchically cluster similar networks into groups (see right part of Figure 2). The x-axis is the linkage distance between two macrobenchmarks or two clustered groups; similar networks will be grouped first. Similarly, we report the correlation of macrobenchmarks using a heatmap as shown in the left part of Figure 2. Each square block in various sizes or colors shows the correlation between the two networks which are indicated by the top and right labels. It is not surprising to observe the close relationship between the sparse and dense versions of a network as they share the same network architecture but have a different amount of operations. For example, *sparse VGG* and *VGG* merge early in the dendrogram and have cooler colors (more related) in heatmap. In addition, half the macrobenchmarks have distances larger than 3.12 (geometric mean distance, 39% of the longest distance). As well, the heatmap is almost occupied by square blocks with warmer color (less related). It clearly shows the diversity of the macrobenchmarks achieved by carefully considering the representativeness and diversity from the perspective of application and architecture (see Section 2.1).

2.3.1 Characteristics

The employed architectural-independent characteristics are listed in Table 3. They can be classified into three categories: *memory*, *computation*, and *control*, which are the main factors of overall efficiency.

Memory characteristics. As the memory wall continues to grow, we first consider three types of memory characteristics: the number of memory accesses (MemAcc), reuse distance (ReDist), and the memory footprint. The number of

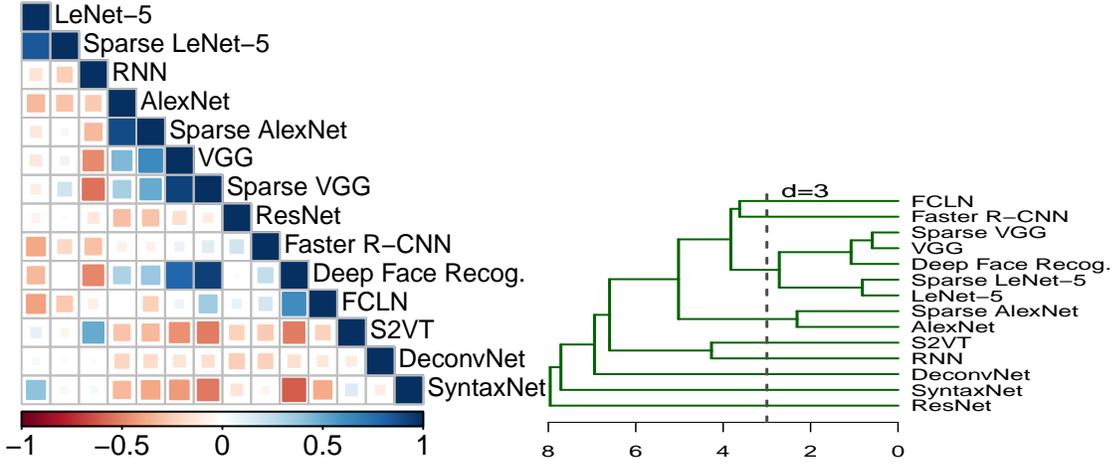


Figure 2: *Left*: Correlation heatmap of macrobenchmarks. *Right*: Similarity of macrobenchmarks.

Table 3: Architecture-independent characteristics

| Category | Name | Notes |
|-------------|--------|--|
| Memory | MemAcc | the number of total memory accesses |
| | ReDist | reuse distance |
| | InMem | memory size of input |
| | OutMem | memory size of output |
| | WghMem | memory size of weight |
| Computation | Ops | the number of operations |
| | OpMem | the ratio of operations to memory access |
| | ComPtt | computation patterns |
| Control | PR | branch prediction ratio |
| | MPR | misprediction ratio |

memory accesses is a direct metric of memory-intensiveness. Reuse distance refers to the number of different data elements accessed between two consecutive reuses of the same element [55]. It has long been used for measuring locality behaviors. Memory footprint refers to the amount of memory accessed at runtime. It is crucial to systems, which are sensitive to memory capacity such as embedded systems. In the context of neural networks, memory footprint can be further classified into three groups: input (InMem), output (OutMem), and weight (WghMem).

Computation characteristics. We define three types of computation characteristics: the number of operations (Ops), the ratio of operations to memory access (OpMem), and the computation patterns (ComPtt). The number of operations is a direct indicator of the problem size. The ratio of operations to memory access can determine whether an algorithm is compute- or memory-intensive. Also, regarding the computation patterns of deep learning algorithms, there are three types: reduction (RD), element-wise (EW) and enlargement (EL). RD refers to operations that

transform multiple input neurons into a single neuron, e.g., *Conv.*, *Pooling*, *FC*, and *LRN*. EW refers to operations that enforce element-wise transformations on the input neurons, e.g., *ReLU*, *Sigmoid*, and *BN*. EL is the operation that transforms one or a set of input neurons into multiple output neurons, e.g., *Deconv.* and *Unpooling*.

Control characteristics. The control characteristics are closely related to branch behaviors. We use two metrics to measure control characteristics: branch prediction ratio (PR) and the misprediction ratio (MPR). The branch ratio is the number of branches to that of all instructions, and the misprediction ratio is the number of mispredicted branches to the number of total branches.

2.3.2 Application perspective

Though our benchmarks include most representative neural networks from a wide range of application scenarios, their inherent characteristics should be diverse to reduce the redundancy of the benchmark suite. Based on the aforementioned architecture-independent characteristics, the diversity can be quantitatively measured. Intuitively, there are two types of diversities in the BENCHIP benchmark suite: inter-layer and intra-layer diversity. The inter-layer diversity is achieved with 12 different benchmarks, and the intra-layer diversity is achieved with seven configurations per benchmark, including three normal configurations (Cfg. A~C), one small configuration (Cfg. D), and three large (extreme) configurations (Cfg. E~G). The normal configurations are directly extracted from commonly used entire networks, while the extreme configurations can be used for stress testing hardware architectures.

Figure 3 shows the kiviatic chart of seven configurations from all 12 microbenchmarks. The axes represent characteristics listed in Table 3², and the meaning of each ring is shown in Figure 3 as well. We can observe that the characteristics of the extreme configurations are significantly different from those of the normal configurations, which well demonstrates the intra-layer diversity. For instance, for the *FC* layer, the weight size of extreme configurations (e.g., Cfg. F) is much larger than that of the normal configurations (e.g., Cfg. C). Moreover, even for the extreme configurations, as they are carefully designed to emphasize different characteristics, the diversity between them is also obvious, e.g., the input size of Cfg. F is 5.3 times smaller than that of Cfg. G for the *Pooling-Avg* layer. For the normal configurations, we can also observe sufficient diversity between different

²The ComPtt characteristic is not included in this figure since there are only three categories.

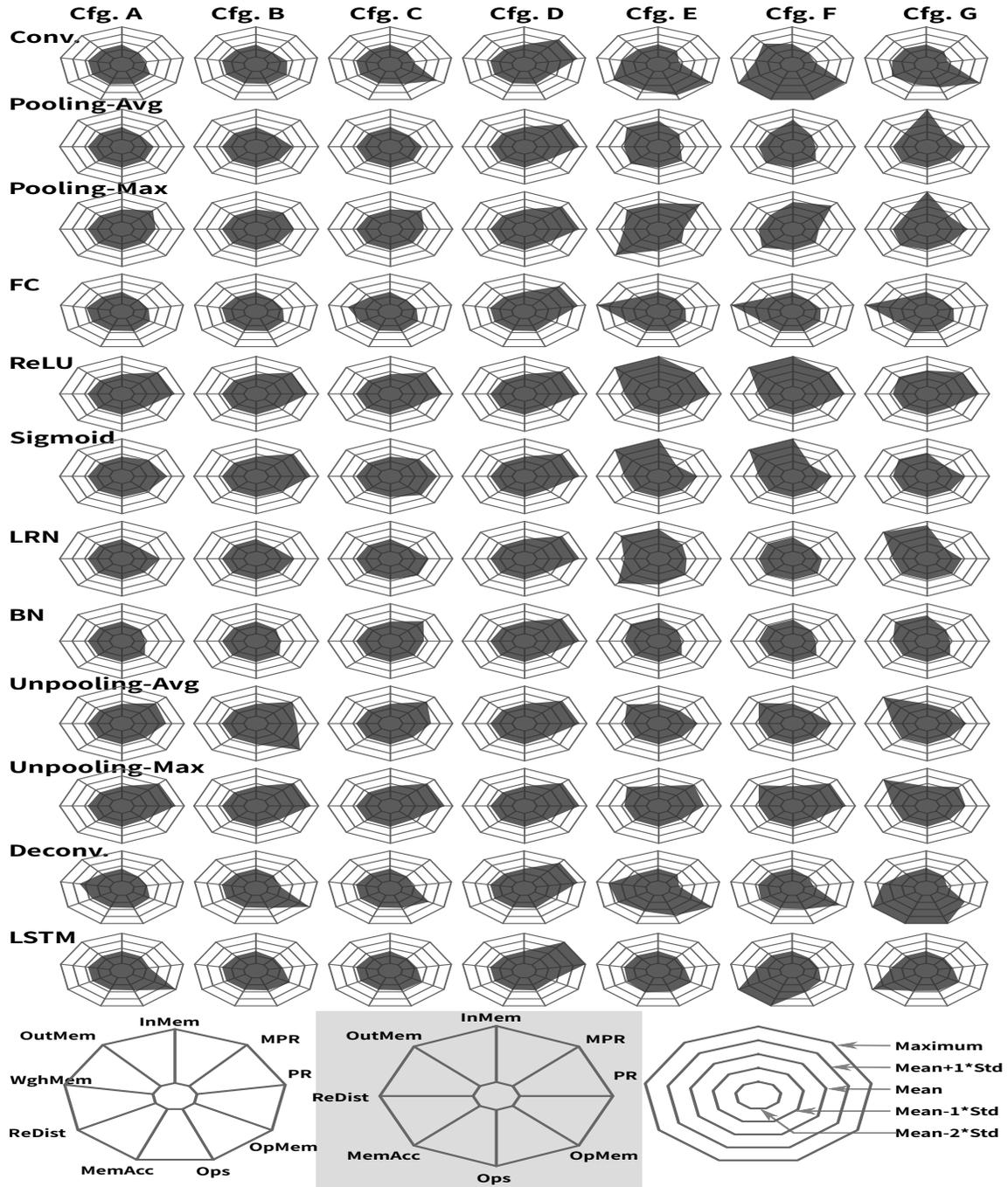


Figure 3: Kiviati chart of all configurations of microbenchmarks.

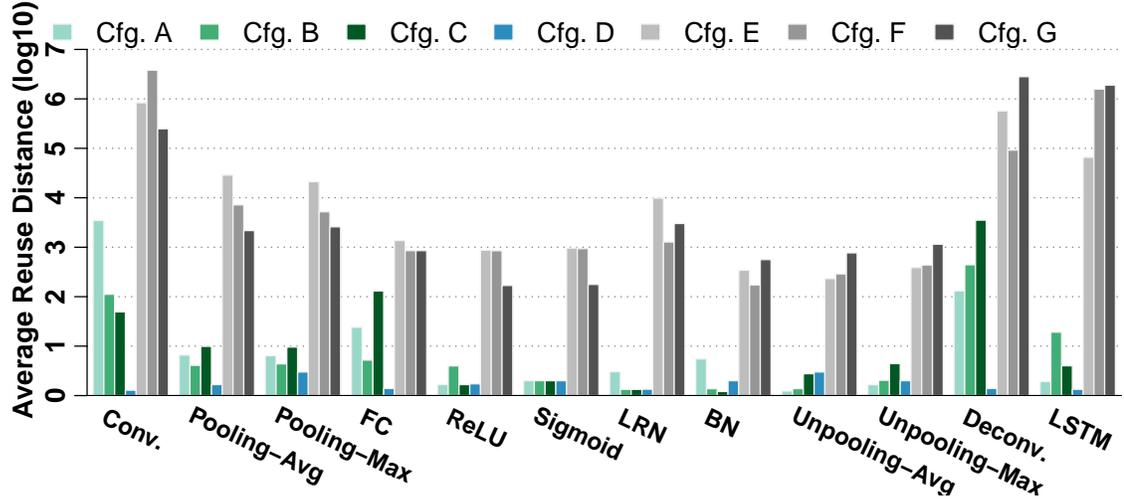


Figure 4: Average reuse distances of microbenchmarks.

benchmarks, where the inter-layer diversity can be observed. Taking the normal configuration Cfg. A as an example, the average reuse distance of *Conv.* (i.e., 3468) is two orders of magnitude larger than that of *FC* (i.e., 24), see Figure 4.

Memory characteristics. Figure 4 shows the average reuse distances of all microbenchmarks, where the reuse distance ranges from 1.2 to $3.8E+7$. For the normal configurations of most benchmarks, the reuse distances are less than 10 (e.g., *Pooling*, *ReLU*, and *LRN*), and thus designing an architecture with on-chip local memory (i.e., cache or scratchpad memory) of $10 \times element_size$ can significantly reduce off-chip memory accesses (except for the *Conv.*, *FC*, and *Deconv.*). In an extreme case such as Cfg. F of *Conv.*, the reuse distance is more than 10^6 , and an on-chip memory buffer would be more than 10MB for holding all the data, which is prohibitively costly with existing SRAM technology. In this case, the on-chip eDRAM employed in DaDianNao [56] and on-die 3D-stacked DRAM [57] would be a potential solution for alleviating the penalty of vast off-chip memory accesses caused by large reuse distances.

Figure 5 further shows the reuse distances of all configurations of *Conv.* and *FC*, which have a relatively large average reuse distance. For the *Conv.* layer, the reuse distance ranges from 2^1 to 2^{29} for the normal configurations. While for the extreme configurations, the reuse distance would exceed 2^{30} , posing a higher challenge on the underlying memory hierarchy. For the *FC* layer, the reuse distances can be roughly classified into three groups: 2^2 , 2^{17} , and 2^{22} . Thus, for the *FC* layer, it is intuitive to customize the memory hierarchy with three levels to accommodate the reuse distances listed above.

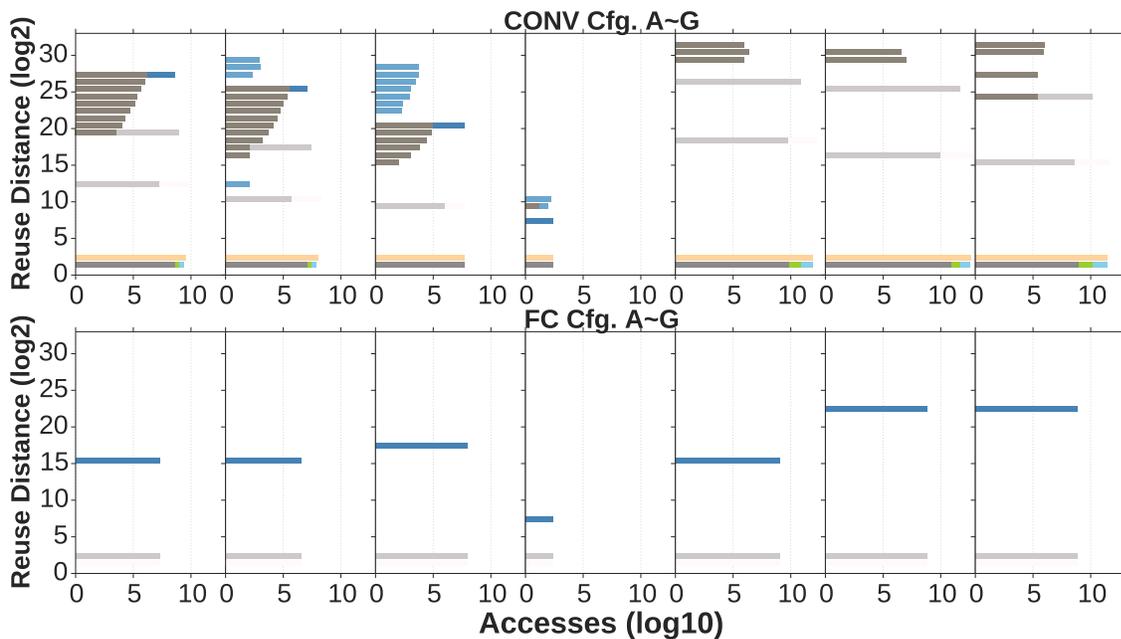


Figure 5: Reuse distances of all configurations of *Conv.* and *FC*.

Computation characteristics. Figure 7 shows the number of operations (Ops) of all benchmarks, where the maximal number of operations could be $8E+12$ (Cfg. F of *Conv.*). For the normal configurations, most Ops are less than 10^8 , while the Ops reach about 10^{10} in some cases. Given a specific benchmark, the Ops also varies significantly to well demonstrate the diversity. A more detailed analysis is the ratio of operations to memory access (i.e., OpMem). The layer with the highest value is the Cfg. F of the *Conv.* layer (> 33), while the Cfg. F of the *FC* layer is only 0.039. This observation is in accordance with the intuition that the *Conv.* layer is compute intensive and the *FC* layer is relatively memory intensive.

Control characteristics. Figure 6 illustrates the MPR of all benchmarks. An interesting observation is that the smallest configuration always has the highest MPR compared with other configurations, as the number of computations is much less than that of others. Another observation is that for the *Pooling* layer, the MPR values of the average pooling and max pooling are completely different. For example, for the normal configuration Cfg. A, the MPR of average pooling is less than 0.5%, while the MPR of max pooling is about 4%. The underlying reason is that max pooling consists of comparison operations that do not exist in average pooling.

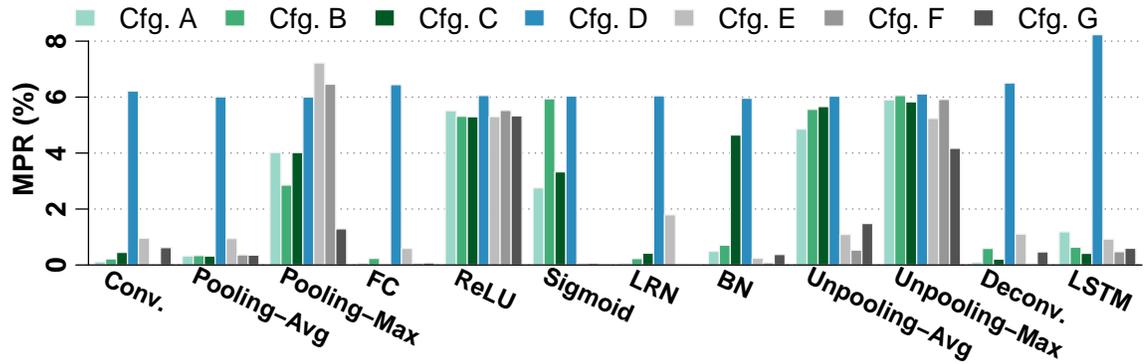


Figure 6: MPR of microbenchmarks.

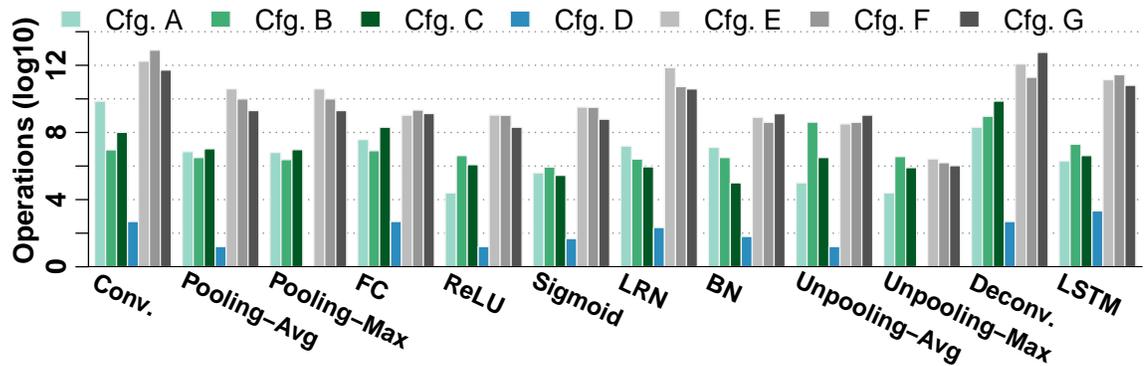


Figure 7: The number of operations of microbenchmarks.

2.3.3 Architecture perspective

From the architecture perspective, the benchmarks should be able to explore the boundary of the processing ability of intelligence processors. This is achieved by using the extreme configurations (Cfg. E~G) of each benchmark. In addition, the benchmarks should be timely and reflect future trends in the intelligence processors.

Stress testing. According to Figure 4, the average reuse distances of extreme configurations are much larger than those of the normal configurations, exhibiting different behaviors on conventional architectures. Given the *Conv.* layer as an example, on a specific CPU (i.e., Intel i5-3470), the L2 cache miss rate of Cfg. F (67%) is 8.4x larger than that of the normal configuration as Cfg. C (8%). Moreover, the average reuse distances of extreme configurations in *Conv.*, *Deconv.*, and *LSTM* are larger than those of others, which puts extreme pressure on the memory hierarchy for these benchmarks. For example, the L3 cache miss rate of Cfg. C in *Conv.* is 1.26x larger than that of Cfg. C in *ReLU*.

In Figure 7, the number of operations of the most extreme configurations is

larger than that of the normal configurations, except for *Unpooling-Max*. The reason is that an extremely large unpooling window size leads to few operations. Regarding the ratio of operations to memory access (OpMem), the extreme configurations are quite different from the normal configurations. Taking the *Conv.* layer as an example, the OpMem of Cfg. F is 33.1, while the OpMem of Cfg. A is only 0.36. In this case, the memory system is more crucial for extreme configurations than the normal ones.

Regarding the control characteristics in Figure 6, the extreme configurations do not exhibit significantly different behaviors. For the *ReLU* layer, the MPR of extreme configurations is almost the same as that of the normal configurations. A potential reason is that the number of comparison operations increases in proportion to the input size, leading to a comparable branch and misprediction ratio regardless of the input size.

Future trends. We expect that there are at least two trends for deep learning. The first is that more algorithms and network structures will be proposed for achieving *higher* accuracy on solving *more* generalized problems. As our benchmarks already have sufficient diversity in terms of computation, memory, and control behaviors, the problem can be alleviated to a certain extent. Besides, we will also update the BENCHIP benchmark suite frequently to keep pace with the rapid advance of deep learning by either adding new benchmarks or eliminating out-of-date benchmarks. The second trend is that the network model tends to be more compact by sacrificing accuracy for performance/energy efficiency. The most notable techniques include the *sparse* model and *low-precision* models (e.g., models with an 8-bit fixpoint and even 1-bit values). To reflect this trend, for the microbenchmarks, we provide sparse and 16-bit models for *Conv.* and *FC* that requires weight values. For the macrobenchmarks, we provide sparse models for LeNet-5, AlexNet, and VGG, as they are the basis of many others networks (e.g., Faster-RCNN and S2VT). The research community is very active in considering both accuracy and architectural improvement, including pruning and compression. In the future, we will update the benchmarks by providing more sparse networks and models.

3 Benchmarking Methodology

In this section, we first present the guidelines and overview of the proposed benchmarking methodology in BENCHIP. Then, we detail the benchmarking framework.

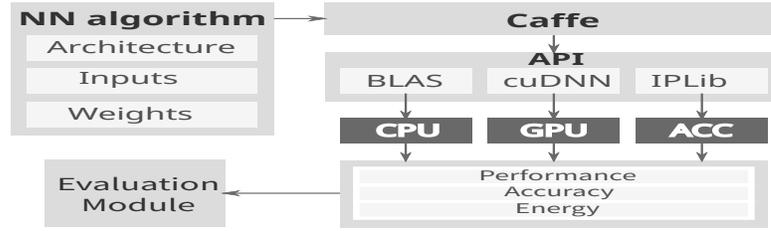


Figure 8: The overall framework of BENCHIP’s benchmarking methodology.

3.1 Guidelines and overview

BENCHIP has two main objectives. The first is to provide a reference toolset for comparing emerging intelligence processors in a fair and easy-to-use fashion. The second is to facilitate identifying performance/energy bottlenecks for potential optimization. The benchmarking methodology should follow at least the following guidelines:

- **Creditability.** The benchmarking methodology should replicate the real application scenario.
- **Portability.** The benchmarking framework should be portable across different platforms, i.e., CPUs, GPUs, and accelerators.
- **Fairness.** The benchmarking methodology should provide fair comparison through specifying the rules and metrics.

To meet the above requirements, we offer a general software stack with a portable interface for different intelligence processors, and the software stack is a replicate of the most widely used industrial software stack, i.e., the stack consisting of high-level programming model, high-performance library, and device driver. The fairness is mainly guaranteed by specifying several benchmarking rules and specially designed evaluation metrics.

3.2 Benchmarking framework

In Figure 8, the benchmarking framework is built upon Caffe [58], where all necessary layers and networks are implemented. Currently, we support Caffe firstly because it is a widely used production-quality framework. We will support other software frameworks progressively. The descriptions of network architecture, input datasets, and pre-trained synaptic weights are treated as standard inputs. The Caffe framework interfaces with the underlying high-performance libraries for different architectures—for example, BLAS for CPUs and cuDNN for GPUs. Given a new hardware accelerator, the corresponding library should be offered, so as to be integrated into the Caffe framework. The platform should be able to report the

performance and energy. Also, the entire framework should report the accuracy of the neural networks. Finally, these individual results are sent to the evaluation module for producing final evaluation results.

Benchmark specification. In BENCHIP, the following files are provided. (1) Configuration file. Since the benchmarking framework is built upon Caffe, the same configure file is used for describing the network architecture. In more detail, we provide the Caffe-compatible *.prototxt* file to characterize both single-layer and entire networks, specifying parameters of each layer and connections between layers. Note that single layers are also considered as independent networks composed of a data layer and the tested layer. (2) Learned parameters (e.g., weight and bias for Conv and FC, γ for BN). For each single layer and entire network, the *.caffemodel* file, where the parameters of every weight layer in that network are stored, is offered. (3) Reference output. BENCHIP offers reference output to measure the precision of tested intelligence processors. Regarding microbenchmarks, we provide input data and reference output data, and the difference between the reference output and the computed output is measured using Mean Squared Error (MSE). Regarding macrobenchmarks, since we select networks from different application scenarios, the evaluated metrics vary as well. Intuitively, their original metrics are employed. For example, we employ the *Top-5* error rate for measuring the networks that aim at image classification (e.g., AlexNet) and METEOR [59] for networks targeting video captioning (e.g., S2VT).

Library interface. We provide a standard library interface, allowing convenient evaluation of different architectures with the Caffe-style framework. It comprises a set of high-level function interfaces implemented in C/C++, each of which corresponds to a basic operation in neural networks. Given a new intelligence processor, the standard interface can be overridden with its own library. Listing 1 presents examples of such interfaces and how to implement them with user-provided library. The `accConvolutionLayer` and virtual function `Forward_acc()` are provided by our benchmarking framework. The virtual function is then implemented by the user-provided library function `accConvolutionForward()`. We also provide an interface for fusing multiple layers—`accFusionLayer::Forward_acc()`—since some processors could be able to obtain better performance by eliminating boundaries between consecutive layers. As shown in Listing 1, the standard interface is implemented with the user-provided `accMultiLayerForward()` so as to optimize the combined execution of the convolutional and pooling layers.

Evaluation metrics. There are four important design considerations for intelligence processors: performance, energy, area, and accuracy. These are deployed as the evaluation metrics. Performance and energy should be reported when run-

Listing 1: Example of Caffe implementation

```
// standard interface for a single layer
void accConvolutionLayer::Forward_acc(...) {
    ...
    accConvolutionForward(...); // user-provided lib
    ...}
// standard interface for fused layers
void accFusionLayer::Forward_acc(...) {
    ...
    accMultiLayerForward(Conv, Pool, ...); // user-provided lib
    ...}

```

ning the benchmarks. On CPUs and GPUs, such metrics can be collected using PAPI [60] and nvsim, respectively. On the accelerator, if it is a simulation prototype, the corresponding simulator should be able to report performance, energy, and area. The collection of the accuracy metric is built into the framework as well.

For microbenchmarks, such metrics already cover various design tradeoffs, and they are sufficient for optimization. For macrobenchmarks, we also provide several synthesized scores for better evaluating design tradeoffs balancing, from the *efficiency* perspective. Such efficiency includes *operations per Joule* as energy efficiency, *operations per second* as computation efficiency and *accuracy affected by area savings* as silicon efficiency.

With such direct and synthesized metrics, it is relatively easy to conduct platform comparisons and optimizations in multiple dimensions, e.g., performance, energy, and various efficiencies. This allows architects to consider potential designs for flexibly balancing different tradeoffs. Given two platforms, macrobenchmarks can be deployed on both platforms to obtain the overall evaluation results. With these results, end-users can perform an overall comparison of these two platforms within their interests. For the weak platform, the architects can leverage microbenchmarks to conduct a problem diagnosis. For the better one, microbenchmarks can be used for further improvements.

4 Benchmarking with BenchIP

In this section, we show how to benchmark IPs by using BENCHIP. We first measure eight IPs (three CPUs, three GPUs and two accelerators) with macrobenchmarks to obtain the overall scores. Then, we delve into the details of the execution efficiency of IPs using the fine-grained microbenchmarks.

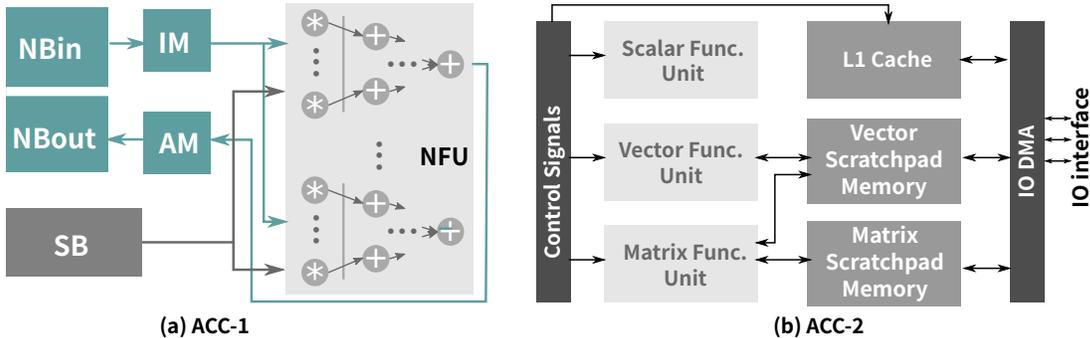


Figure 9: The frameworks of ACC-1 and ACC-2.

4.1 Evaluated IPs

CPU. We select three types of CPUs—CPU-E, CPU-D and CPU-S—for different scenarios. CPU-E is an embedded processor, i.e., 4xA57 Atlas/2MB L2. CPU-D is a desktop processor, i.e., Intel i5-3470@3.2GHz with 16GB memory. CPU-S is a server processor, i.e., Intel E5-2620 v2@2.1GHz with 64GB memory.

GPU. Similarly, we select three GPUs with different levels of capabilities: GPU-E, GPU-D, and GPU-S. GPU-E is an embedded GPU, i.e., Nvidia Maxwell GeForce [61] (256 CUDA cores, 512 GFlops with FP32). GPU-D is a desktop GPU, i.e., GeForce GTX950 [62] (640 CUDA cores, 2GB GDDR5, 1.57 TFlops). GPU-S is a server GPU, i.e., Nvidia Tesla K40 [63] (2880 processors, 12GB GDDR5, 4.29 TFlops).

ACC. We select two prototypes of ASIC IPs—ACC-1 and ACC-2—by carefully reimplementing DianNao [10] and Cambricon [25], respectively, with exactly the same parameters in the papers. The reason for selecting DianNao and Cambricon as prototypes of ACC-1 and ACC-2 is twofold. First, DianNao is a notable accelerator that has moderate performance and cost that supports basic matrix/vector operations in neural networks. Second, Cambricon is one of the few accelerators that supports various kinds of neural network algorithms but still has performance comparable with the state-of-the-art high-performance accelerator, DaDianNao [56]. Thus, ACC-1 and ACC-2 can cover various scenarios.

In Figure 9, we present the framework of both ACC-1 and ACC-2. ACC-1 consists of a neural functional unit (NFU), an input/output neuron buffer (NBin/N-Bout), a synaptic buffer (SB), an indexing module (IM), and an assemble module (AM). Note that we modify the architecture of DianNao to produce ACC-1, by adding several features such as sparsity (i.e., IM). ACC-2 consists of a matrix function unit (MFU), a vector function unit (VFU), a scalar function unit (SFU), an address generation unit (AGU), and accompanied scratchpad memory.

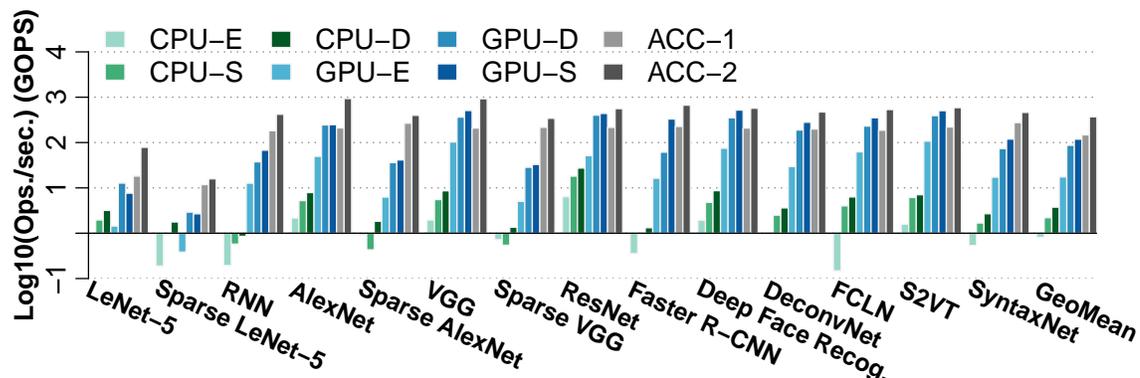


Figure 10: Giga operations per second.

4.2 Benchmarking with macrobenchmarks

In Figure 10 and 11, we report the synthesized scores from the efficiency perspective, i.e., energy efficiency using *giga operations per Joule* (GOPJ) and performance efficiency using *giga operations per second* (GOPS). Further, we report the fine-grained metrics, i.e., *performance* and *energy*, in Figure 12 and 13, respectively. We focus on performance and energy because area and accuracy basically do not vary across different macrobenchmarks and IPs, respectively. Thus, we have the following observations. Here also note that the overhead of the Caffe software stack is small, less than 2% of the total runtime on CPUs/GPUs. And for all the customized intelligence processors, we fetch end-to-end measured data for benchmarking. Thus, the results would not be affected by software stack.

Observation #1. *Embedded CPUs trade performance for silicon cost reduction better than energy efficiency.* Lightweight embedded CPUs are well known for being suitable for power limited scenarios such as embedded systems, mainly because of their low energy costs. Interestingly, as we observed, an embedded CPU, i.e., CPU-E, improves silicon area utilization better than energy efficiency. CPU-E achieves 9.22x and 11.11x better silicon area utilization but only 6.91x and 2.23x better efficiency in utilizing every joule for computation (GOPJ) than CPU-S and CPU-D, 5.03x and 3.12x slower in performance, respectively.

Observation #2. *GPUs still improve energy efficiency as well as provide high performance efficiency.* GPUs let applications exploit parallelism with tremendous computation units and memory bandwidth, which motivated their widely usage in intelligence processing. As observed, we confirmed that GPUs still leverage the energy for more operations [64]. On average, GPU-E, GPU-D and GPU-S achieve 4.65x, 23.14x, and 31.55x better computation efficiency than CPU-D

while still achieving certain improvements in energy efficiency, i.e., 5.74x, 7.02x, and 10.15x, respectively. Interestingly, GPU-D has the best utilization of such a tradeoff. Regarding the efficiencies, GPU-D achieves relatively high performance and energy efficiency, i.e., 69.17% of GPU-E (122.42% of GPU-S) in GOPJ and 73.33% of GPU-S (497.99% of GPU-E) in GOPS. When detailed in energy costs and execution times, compared to GPU-S, GPU-D has 0.84x energy consumption and 1.41x slower performance with 40.64% area cost. Moreover, GPU-D achieves 4.91x speedup over GPU-E and 1.47x energy costs, with 1.43x more area cost.

Observation #3. *Customized accelerators provide significantly better energy efficiency with a relatively lower performance efficiency improvement and silicon area cost reduction.* We observed that ACC-1 and ACC-2 achieve 1.69x and 4.26x performance efficiency but 220.10x and 170.80x energy efficiency when compared against GPU-D, and 39.21x and 98.55x performance efficiency but 1545.78x and 1199.56x energy efficiency when compared against CPU-S. Regarding both performance and energy efficiency, ACC-1 and ACC-2 have higher utilization in silicon area, i.e., 44.31%/4.00% and 390.06%/35.15% of smallest/largest CPU (CPU-E/CPU-S), 18.77%/165.42% and 1.14%/10.03% of smallest/largest GPU (GPU-E/GPU-S).

Observation #4. *ACC-1 is more efficient than ACC-2, which complies with the intuition that with sparse support the accelerator benefits from such sparsity.* Interestingly, ACC-1 is more efficient than ACC-2 in GOPJ (1.29x) but 39.79% inefficient in GOPS. Moreover, due to the sparsity, ACC-1 is slightly more efficient in GOPJ, where ACC-1 has a 2.99x speedup on average with sparse support and ACC-2 gains nothing as it processes no differently on sparse and dense networks. Furthermore, note that despite the very small network LeNet-5, both CPUs and GPUs cannot leverage the benefit of computations and data amount reduction due to sparsity, i.e., 2.09x and 4.15x slower when compared against dense version. When detailed in energy costs and execution times, ACC-1 performs 2.51x slower than ACC-2 but costs 1.29x less in energy and has far less area cost, averaging on all macrobenchmarks.

In short, while using BENCHIP, we made the following observations from the efficiency perspective:(1) embedded CPUs trades performance efficiently for small area better than low power, (2) GPUs focus on performance improvement but still improve the energy efficiency, and desktop GPUs have the best tradeoff, (3)customized accelerators are suitable for almost all scenarios, and (4) accelerators with sparse support are much better than their dense counterparts.

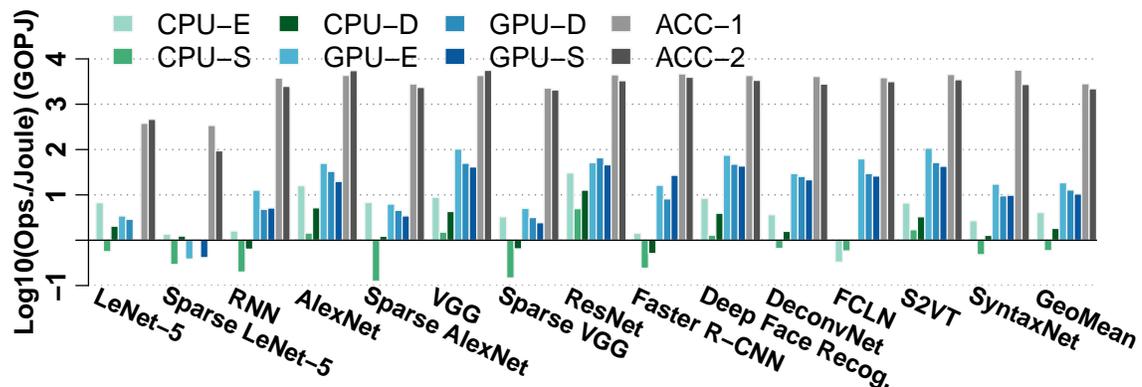


Figure 11: Giga operation per joule.

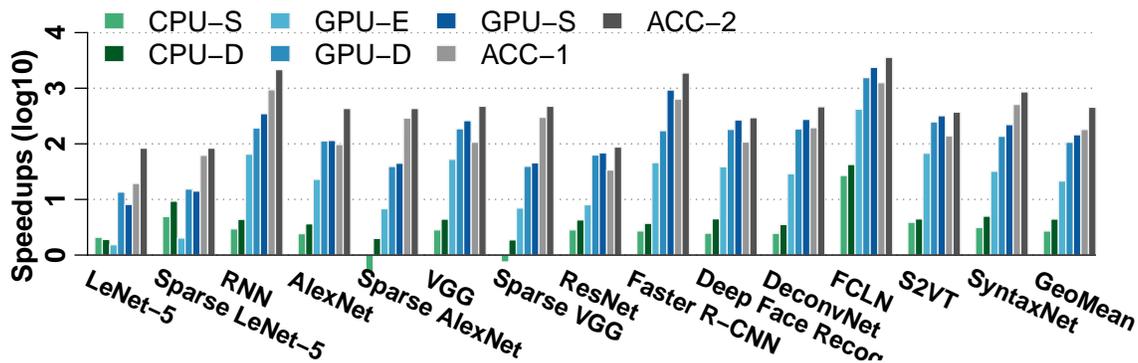


Figure 12: Speedups of evaluated IPs with macrobenchmarks (Normalized to CPU-E).

4.3 Benchmarking with microbenchmarks

For the platforms with relatively low/high scores, we can study the reasons for the inefficiency/efficiency with microbenchmarks. As illustrative examples, we further evaluate especially CPU-D, GPU-E, and ACC-2 with microbenchmarks, as they perform worse on macrobenchmarks. In Figure 14, we shown the performance of IPs on certain microbenchmarks, i.e., *Conv.*, *FC*, *ReLU*, and *Deconv.*, which cover all three computation patterns (corresponding to RD, RD, EW, and EL).

Regarding observation #1, the desktop-level CPU-D achieves better performance than both CPU-E and the server end CPU-S, i.e., 5.03x and 1.61x on average, respectively, as shown in Figure 12. Averaging on microbenchmarks with all the configurations, CPU-D achieves results consistent with those on macrobenchmarks, with 3.41x and 1.48x speedups over CPU-E and CPU-S. Also CPUs perform uni-

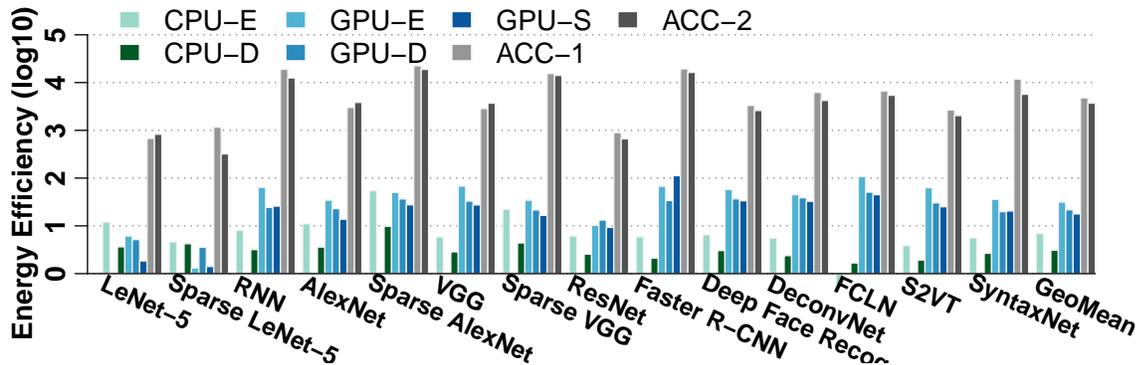


Figure 13: Normalized energy consumption of evaluated IPs with macrobenchmarks (Normalized to CPU-S).

formly on normal or larger configurations (Cfg. A~Cfg. C, Cfg. D~Cfg. G), where the speedups vary from 3.36x/1.50x to 4.07x/1.61x.

Regarding observation #2, the embedded system GPU, GPU-E, does not always have a shorter execution time than the CPUs, especially on small networks (*LeNet-5* and *sparse LeNet-5*). With the microbenchmarks, we observe that GPU-E has 0.02x speedup on average compared against CPU-D on Cfg. D which is the extremely small cases. GPUs cannot take full advantage of their high computational power because the small computational kernels map poorly on their hundreds and thousands of threads.

Regarding observation #4, ACC-2 achieves the best performance on every macrobenchmark and it is the same for all the microbenchmarks, i.e., a 2.2x speedup averagely. Moreover, compared against ACC-1, ACC-2 performs better on large configurations (3.19x speedup) than normal configurations (1.80x speedup) where its high parallelism and throughput (peak performance of 2114 GOPS) can be well leveraged.

Additionally, ACC-1 achieves better performance than all the CPUs and GPUs on most CNNs and DNNs but fails on complex and irregular networks such as *ResNet* and *DeconvNet*. It is because ACC-1 is designed to efficiently process CNNs/DNNs. Thus, ACC-1 obtains a performance benefit over GPU-S on layers commonly used in CNNs and DNNs, i.e., *Conv.* (1.43x) and *FC* (16.25x).

ACC-1 can be very effective on small networks, as it performs 217.78x faster than GPU-S on Cfg. D. On normal configurations such as Cfg. A~C, ACC-1 can still keep its efficiency with a 1.71x speedup. However, on heavy workloads, GPUs can well leverage their high parallelism. Thus, ACC-1 cannot exceed GPU-S (0.65x speedup on Cfg. D~G).

5 Discussion

Overall Score. In current version, we do not provide an overall score for different IPs. The reason is two-fold. First, most architecture designs involve a multi-dimension optimization process where architects make decisions on a series of tradeoffs that leads to a final design for some specific targets. Thus, in such case detailed evaluation results can be used directly for design improvement and an overall score would be meaningless. Second, overall scores appropriate for various platforms, e.g., CPU, GPU, FPGA, DSP, and customized accelerators do not exist. One possible definition of an overall score could be the overall efficiency, i.e., synthesizing energy efficiency, performance efficiency, and silicon efficiency. For example, the overall efficiency equation can be defined as $GeoMean(f(\frac{Ops}{Energy}) \times g(\frac{Ops}{Times}) \times h(\frac{Acc/Ref_acc}{Area}))$, where the $f(\cdot)$, $g(\cdot)$ and $h(\cdot)$ are potential scaling/mapping functions. However, due to the huge gaps among various platforms, those functions can only be user-defined or fitting equations learning from collected data after BENCHIP is opened for public. Therefore, at present we only provide direct measured metrics and synthesized efficiency scores.

6 Related Work

We review related work from three aspects: neural network models, accelerators, and benchmarks.

Models. In recent years, deep learning has been applied to various scenarios and achieved great success. For image recognition, variants of CNNs, as powerful feature extractors, are applied to different tasks (e.g., image classification and face recognition) [1, 2, 42, 44]. Image/video captioning adopts both CNN and LSTM for extracting features and generating sentences, respectively [4, 47]. In semantic segmentation, deconvolution and unpooling are used for reconstructing features [46]. In addition to image processing, deep learning also achieves great performance on speech recognition [65], and natural language processing [50], where LSTM/RNN are the core learning models. Such diverse scenarios have led to a number of variants of neural network architectures and models that have seldom been considered in previous studies of neural network accelerators.

Accelerators. With the growing consensus that CPUs and GPUs, as the traditional platforms for running neural network models, are not able to provide high energy-efficiency for specific workloads like CNN, researchers began to seek the possibility of implementing neural networks on other platforms (e.g., FPGA [9]

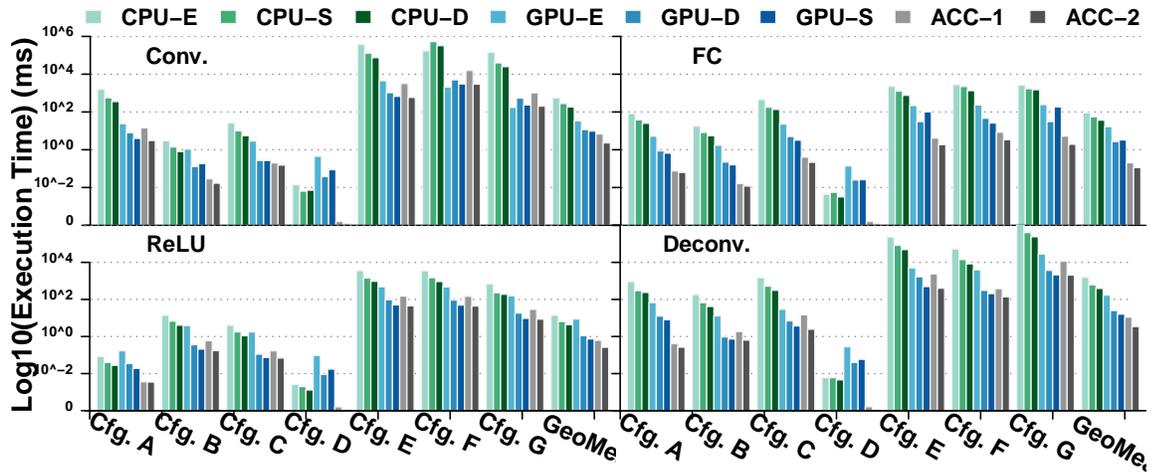


Figure 14: Performance of IPs on microbenchmarks (*Conv.*, *FC*, *ReLU* and *Deconv*).

and ASIC [10, 12, 23, 24, 54, 56]). The DianNao family [10, 54, 56] contains a series of neural network processors to support most operations of CNN and DNN.

Han *et al.* [12] designed EIE, which is a specially optimized architecture for sparse networks. Chi *et al.* [23] designed PRIME, a PIM architecture which accelerates NN algorithms in ReRAM-based main memory. The above studies evaluate their processors on different sets of deep learning operations and configurations, making it impossible to compare two processors in a fair fashion.

Benchmarks. Though the exploration of intelligence processors continues to grow, a suitable benchmark for the evaluation and optimization of such hardware is still absent. BenchNN and DeepBench are two benchmarks that partially address this problem. However, due to their non-diversity and nonrepresentativeness, they cannot be used for designing and optimizing state-of-the-art intelligence processors. In the computer architecture community, several benchmarks have already been used. Chen *et al.* [56] used 10 layers all extracted from AlexNet and two customized RBMs for evaluation. Chi *et al.* [23] employed six networks, five extracted from LeNet-5 and one from VGG. Han *et al.* employed nine DNNs extracted from AlexNet, VGG, and Neural Talk [66]. Shafiee *et al.* [24] employed seven CNNs (four from VGG, three from MSRA [67]), and two DNNs (one from DeepFace [68] and the other one from [69]). Due to the lack of a standard benchmark suite, employing such personalized benchmarks in the design of existing neural network accelerators is relatively casual, which is unsuitable for benchmarking intelligence processors.

7 Conclusions

This paper proposes BENCHIP, a benchmark suite and methodology for benchmarking intelligence processors. The benchmark suite of BENCHIP consists of two sets of benchmarks, i.e., microbenchmarks and macrobenchmarks—for fair comparison and system optimization. The benchmark methodology is built upon an industrial software stack consisting of a high-level programming model, library, and device driver. We also provide evaluation metrics for comprehensively reflecting the various characteristics of the intelligence processors being evaluated. BENCHIP is employed to evaluate various hardware platforms, including CPUs, GPUs, and neural network accelerators. BENCHIP, which is already being used in several companies, will be open-sourced soon to facilitate the design and evaluation of intelligence processors in a broad sense.

References

- [1] A. Krizhevsky, G. E. Hinton, I. Sutskever, and G. E. Hinton, “ImageNet Classification with Deep Convolutional Neural Networks,” *Advances In Neural Information Processing Systems*, pp. 1–9, 2012.
- [2] K. Simonyan and A. Zisserman, “Very Deep Convolutional Networks for Large-Scale Image Recognition,” pp. 1–14, 2015.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *arXiv preprint arXiv:1506.01497*, 2015.
- [4] S. Venugopalan, M. Rohrbach, J. Donahue, R. J. Mooney, T. Darrell, and K. Saenko, “Sequence to sequence – video to text,” in *Proceedings of the 2015 International Conference on Computer Vision (ICCV’15)*, 2015.
- [5] O. Abdel-Hamid, A. rahman Mohamed, H. Jiang, L. Deng, G. Penn, and D. Yu, “Convolutional neural networks for speech recognition,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 22, pp. 1533–1545, 2014.
- [6] A. Eriguchi, K. Hashimoto, and Y. Tsuruok, “Tree-to-sequence attentional neural machine translation,” in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL’16)*, 2016.

- [7] “The nvidia dgx-1 deep learning system.” NVIDIA DGX-1.
- [8] C. Farabet, C. Poulet, J. Y. Han, and Y. LeCun, “Cnp: An fpga-based processor for convolutional networks,” in *Proceedings of the 2009 International Conference on Field Programmable Logic and Applications (FPL’09)*, pp. 32–37, 2009.
- [9] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, “Optimizing fpga-based accelerator design for deep convolutional neural networks,” in *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA’15)*, pp. 161–170, 2015.
- [10] T. Chen, Z. Du, N. Sun, J. Wang, and C. Wu, “DianNao: a small-footprint high-throughput accelerator for ubiquitous machine-learning,” in *Proceedings of the 19th international conference on Architectural support for programming languages and operating systems (ASPLOS)*, (Salt Lake City, UT, USA), pp. 269–284, 2014.
- [11] C. Farabet, B. Martini, B. Corda, P. Akselrod, E. Culurciello, and Y. LeCun, “NeuFlow: A runtime reconfigurable dataflow processor for vision,” in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 109–116, Ieee, jun 2011.
- [12] S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. A. Horowitz, and W. J. Dally, “EIE: Efficient Inference Engine on Compressed Deep Neural Network,” in *Proceedings of the 43th Annual International Symposium on Computer Architecture (ISCA’16)*, vol. 16, 2016.
- [13] “The standard performance evaluation corporation (spec).” SPEC CPU.
- [14] C. Bienia, S. Kumar, J. P. Singh, and K. Li, “The PARSEC benchmark suite: Characterization and architectural implications,” in *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques*, pp. 72–81, 2008.
- [15] M. Alwani, H. Chen, M. Ferdman, and P. Milder, “Fused-layer CNN accelerators,” in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pp. 1–12, 2016.

- [16] P. Judd, J. Albericio, and A. Moshovos, “Stripes: Bit-Serial Deep Neural Network Computing,” in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, vol. 6056, pp. 1–1, 2016.
- [17] M. Rhu, N. Gimelshein, J. Clemons, A. Zulfiqar, and S. W. Keckler, “Virtualizing Deep Neural Networks for Memory-Efficient Neural Network Design,” in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2016.
- [18] S. Zhang, Z. Du, L. Zhang, H. Lan, S. Liu, L. Li, Q. Guo, T. Chen, and Y. Chen, “Cambricon-X : An Accelerator for Sparse Neural Networks,” in *Proceedings of the 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO-49)*, 2016.
- [19] Y. Ji, Y. Zhang, S. Li, P. Chi, C. Jiang, P. Qu, Y. Xie, and W. Chen, “NEUTRAMS : Neural Network Transformation and Co-design under Neuromorphic Hardware Constraints,” in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, no. October, 2016.
- [20] D. Kim, J. Kung, S. Chai, S. Yalamanchili, and S. Mukhopadhyay, “Neurocube: A Programmable Digital Neuromorphic Architecture with High-Density 3D Memory,” in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 380–392, 2016.
- [21] R. LiKamWa, Y. Hou, Y. Gao, M. Polansky, and L. Zhong, “RedEye: Analog ConvNet Image Sensor Architecture for Continuous Mobile Vision,” in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 255–266, 2016.
- [22] J. Albericio, P. Judd, T. Hetherington, T. Aamodt, N. E. Jerger, and A. Moshovos, “Cnvlutin: Ineffectual-Neuron-Free Deep Neural Network Computing,” in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 1–13, 2016.
- [23] P. Chi, S. Li, C. Xu, T. Zhang, J. Zhao, Y. Liu, Y. Wang, and Y. Xie, “PRIME: A Novel Processing-in-Memory Architecture for Neural Network Computation in ReRAM-Based Main Memory,” in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 27–39, 2016.
- [24] A. Shafiee, A. Nag, N. Muralimanohar, R. Balasubramonian, J. P. Strachan, M. Hu, R. S. Williams, and V. Srikumar, “Isaac: A convolutional neural

network accelerator with in-situ analog arithmetic in crossbars,” in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 14–26, June 2016.

- [25] S. Liu, Z. Du, J. Tao, D. Han, T. Luo, Y. Xie, Y. Chen, and T. Chen, “Cambricon: An Instruction Set Architecture for Neural Networks,” in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 393–405, 2016.
- [26] Y.-H. Chen, J. Emer, and V. Sze, “Eyeriss: A Spatial Architecture for Energy-Efficient Dataflow for Convolutional Neural Networks,” in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 367–379, 2016.
- [27] B. Reagen, P. Whatmough, R. Adolf, S. Rama, H. Lee, S. K. Lee, J. M. Hernandez-Lobato, G.-Y. Wei, and D. Brooks, “Minerva: Enabling Low-Power, Highly-Accurate Deep Neural Network Accelerators,” in *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pp. 267–278, 2016.
- [28] L. Song, X. Qian, H. Li, and Y. Chen, “PipeLayer : A Pipelined ReRAM-Based Accelerator for Deep Learning Basics of Deep Neural Network,” in *Proceedings of The 23rd IEEE Symposium on High Performance Computer Architecture (HPCA)*, 2017.
- [29] W. Lu, G. Yan, J. Li, S. Gong, Y. Han, and X. Li, “FlexFlow: A Flexible Dataflow Accelerator Architecture for Convolutional Neural Networks,” in *Proceedings of The 23rd IEEE Symposium on High Performance Computer Architecture (HPCA)*, pp. 553–564, 2017.
- [30] W. Wzr, V. Surfhhv, L. V. Ghsor, H. G. Rq, K. Hqg, D. Rq, P. D. Q. Fkdoohqjlqj, P. Ohduqlqj, H. J. L. W. Wdnhv, W. Zhhnv, W. R. Wudlq, R. Q. Irxu, and K. Hqg, “Towards Pervasive and User Satisfactory CNN across GPU Microarchitecture,” in *Proceedings of The 23rd IEEE Symposium on High Performance Computer Architecture (HPCA)*, 2017.
- [31] T. Chen, Y. Chen, M. Durantou, 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,” *2012 IEEE International Symposium on Workload Characterization (IISWC)*, pp. 36–45, nov 2012.

- [32] B. Research, “Deepbench,” 2016.
- [33] R. Adolf, S. Rama, B. Reagen, G.-y. Wei, and D. Brooks, “Fathom: Reference Workloads for Modern Deep Learning Methods,” in *2016 IEEE International Symposium on Workload Characterization (IISWC)*, 2016.
- [34] F. Murtagh and M. Hernández-Pajares, “The kohonen self-organizing map method: An assessment,” *Journal of Classification*, vol. 12, no. 2, pp. 165–190, 1995.
- [35] S. Ren, K. He, and R. Girshick, “Faster-R-Cnn-Towards-Real-Time-Object-Detection-With-Region-Proposal-Networks,” in *Advances in neural information processing systems*, pp. 1–9, 2015.
- [36] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, “XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks,” *arXiv preprint*, pp. 1–17, 2016.
- [37] H. Song, J. Pool, J. Tran, and W. J. Dally, “Learning Both Weights and Connections for Efficient Neural Networks,” in *Advances in Neural Information Processing Systems (NIPS’15)*, pp. 1135—1143, 2015.
- [38] Y. LeCun, L. L. Bottou, Y. Bengio, and P. Haffner, “Gradient-Based Learning Applied to Document Recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [39] A. Graves and N. Jaitly, “Towards End-To-End Speech Recognition with Recurrent Neural Networks,” *JMLR Workshop and Conference Proceedings*, vol. 32, no. 1, pp. 1764–1772, 2014.
- [40] M. P. Marcus, B. Santorini, and M. A. Marcinkiewicz, “Building a large annotated corpus of English: The Penn Treebank.,” *Computational Linguistics*, vol. 19, no. 2, pp. 313–330, 1993.
- [41] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, “ImageNet Large Scale Visual Recognition Challenge,” p. 37, sep 2014.
- [42] K. He, “Deep Residual Learning for Image Recognition,”

- [43] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, “The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results.” <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>.
- [44] O. M. Parkhi, A. Vedaldi, and A. Zisserman, “Deep Face Recognition,” in *Proceedings of the British Machine Vision Conference 2015*, no. Section 3, pp. 41.1–41.12, 2015.
- [45] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, “Labeled faces in the wild: A database for studying face recognition in unconstrained environments,” Tech. Rep. 07-49, University of Massachusetts, Amherst, October 2007.
- [46] H. Noh, S. Hong, and B. Han, “Learning Deconvolution Network for Semantic Segmentation,” in *The IEEE International Conference on Computer Vision (ICCV)*, vol. 1, 2015.
- [47] J. Johnson, A. Karpathy, and L. Fei-Fei, “DenseCap: Fully Convolutional Localization Networks for Dense Captioning,” *arXiv preprint*, 2015.
- [48] S. Venugopalan, M. Rohrbach, T. Darrell, J. Donahue, K. Saenko, and R. Mooney, “Sequence to Sequence - Video to Text,” in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4534–4542, 2015.
- [49] D. L. Chen and W. B. Dolan, “Collecting Highly Parallel Data for Paraphrase Evaluation,” *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics*, pp. 190–200, 2011.
- [50] D. Andor, C. Alberti, D. Weiss, A. Severyn, A. Presta, K. Ganchev, S. Petrov, and M. Collins, “Globally Normalized Transition-Based Neural Networks,” in *arXiv preprint*, pp. 2442–2452, 2016.
- [51] A. Phansalkar, A. Joshi, and L. K. John, “Analysis of redundancy and application balance in the spec cpu2006 benchmark suite,” in *Proceedings of the 34th Annual International Symposium on Computer Architecture (ISCA’07)*, pp. 412–423, 2007.

- [52] J. D. McCalpin, “Memory bandwidth and machine balance in current high performance computers,” *IEEE Computer Society Technical Committee on Computer Architecture (TCCA) Newsletter*, pp. 19–25, 1995.
- [53] J. M. Bull and D. O’Neill, “A microbenchmark suite for openmp 2.0,” *SIGARCH Comput. Archit. News*, vol. 29, no. 5, pp. 41–48, 2001.
- [54] Z. Du, R. Fasthuber, T. Chen, P. Ienne, L. Li, X. Feng, Y. Chen, and O. Temam, “ShiDianNao: Shifting Vision Processing Closer to the Sensor,” in *Proceedings of the 42nd Annual International Symposium on Computer Architecture*, pp. 92–104, 2015.
- [55] C. Ding and Y. Zhong, “Predicting whole-program locality through reuse distance analysis,” in *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pp. 245–257, 2003.
- [56] Y. Chen, T. Luo, S. Liu, S. Zhang, L. He, J. Wang, L. Li, T. Chen, Z. Xu, N. Sun, and O. Temam, “DaDianNao: A Machine-Learning Supercomputer,” in *Proceedings of the 47th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO-47)*, pp. 609–622, 2015.
- [57] J. T. Pawlowski, “Hybrid memory cube: breakthrough dram performance with a fundamentally re-architected dram subsystem,” in *Proceedings of the 23rd Hot Chips Symposium (HotChips’11)*, 2011.
- [58] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, “Caffe: Convolutional architecture for fast feature embedding,” *arXiv preprint arXiv:1408.5093*, 2014.
- [59] M. Denkowski and A. Lavie, “Meteor Universal: Language Specific Translation Evaluation for Any Target Language,” pp. 376–380, 2014.
- [60] P. J. Mucci, S. Browne, C. Deane, and G. Ho, “Papi: A portable interface to hardware performance counters,” in *In Proceedings of the Department of Defense HPCMP Users Group Conference*, pp. 7–10, 1999.
- [61] NVIDIA, “Nvidia tegra x1: Nvidia’s new mobile superchip.”
- [62] NVIDIA, “Geforce gtx 950m: specifications.”

- [63] NVIDIA, “Nvidia tesla gpu accelerators.”
- [64] S. W. Keckler, W. J. Dally, B. Khailany, M. Garland, and D. Glasco, “GPUs and the future of parallel computing,” *IEEE Micro*, vol. 31, no. 5, pp. 7–17, 2011.
- [65] A. Graves, A.-r. Mohamed, and G. Hinton, “Speech recognition with deep recurrent neural networks,” in *2013 IEEE international conference on acoustics, speech and signal processing*, pp. 6645–6649, IEEE, 2013.
- [66] A. Karpathy and F. Li, “Deep visual-semantic alignments for generating image descriptions,” in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pp. 3128–3137, 2015.
- [67] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV), ICCV ’15*, (Washington, DC, USA), pp. 1026–1034, IEEE Computer Society, 2015.
- [68] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf, “Deepface: Closing the gap to human-level performance in face verification,” in *Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR ’14*, (Washington, DC, USA), pp. 1701–1708, IEEE Computer Society, 2014.
- [69] Q. V. Le, M. Ranzato, R. Monga, M. Devin, G. Corrado, K. Chen, J. Dean, and A. Y. Ng, “Building high-level features using large scale unsupervised learning,” in *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*, 2012.