

Weight Block Sparsity: Training, Compilers, and Accelerators

Paolo D’Alberto Taehee Jeong Akshay Jain Shreyas Manjunath Mrinal Sarmah Samuel Hsu Nitesh Pipralia

Abstract—Inference is a synonymous for performance: for example, more and faster classifications per seconds. As faster and tailor-made chip-lets are deployed in the field also larger and larger models are domineering the public’s attention: quantization and sparsification are used to fit these large models into more reasonable sized ones. Quantization reduces the foot print per model weight and sparsification trims superfluous ones. We present the main ideas about a vertical system where convolution and matrix multiplication weights can be trained to exploit an 8x8 block sparsity, compilers recognize such a sparsity for data compaction and computation splitting into threads and cores.

Blocks present spatial locality that a vector operation can make full use and temporal locality for good reuse and cost amortization. If we take a Resnet50, we can reduce the weight by half with little accuracy loss, We can achieve speeds similar to an hypothetical Resnet25. We shall present performance estimates by accurate and complete code generation for a small and efficient set of AIE2 (Xilinx Versal FPGAs). We show block sparsity could be applied to CPU and GPUS.

Index Terms—Sparsity, Performance, FPGA, Tools

I. INTRODUCTION

We shall start by what we mean for block sparsity and for a vertical solution. This will provide a clear introduction to the subject. Block sparsity is an intuitive concept but it is also a little misunderstood because we all have a different mental picture of it. Take a matrix multiplication as in Equation 1

$$\begin{pmatrix} C_0 & C_1 \\ C_2 & C_3 \end{pmatrix} = \begin{pmatrix} A_0 & A_1 \\ A_2 & A_3 \end{pmatrix} \begin{pmatrix} 0 & B_1 \\ B_2 & 0 \end{pmatrix} \quad (1)$$

This is simply the computation

$$C_0 = A_1 B_2; C_1 = A_0 B_1; C_2 = A_3 B_2; C_3 = A_2 B_1 \quad (2)$$

and in general with proper γ_i

$$C_i = \sum_{k=0}^1 A_{2i+k} (\gamma_{2*k+i} B_{2*k+i}) \quad (3)$$

Where the matrix B is a constant and diagonal and each submatrices B_2 and B_1 can split further down and show zero blocks of even smaller sizes. In this work, we chose the basic block of being of size $B_i = 8 \times 8$. It is arbitrary but it is a great starting point for architectures based on AMD AIE2 products. For example,

$$B = \sum_i \gamma_i B_i, \quad \gamma_i \in \{0, 1\} \quad (4)$$

Each block is a constant either non zero or zero:

$$B = \begin{bmatrix} \gamma_0 B_0 & \dots & \gamma_0 B_{n-1} \\ \gamma_n B_n & \dots & \gamma_{2n-1} B_{2n-1} \\ \dots & \dots & \dots \\ \gamma_{(n-1)n} B_{(n-1)n} & \dots & \gamma_{(n-1)^2} B_{(n-1)^2} \end{bmatrix}$$

Of course, the matrix A can have any value and so C . Some applications have some constraints about the row and columns of B_i : for example, each row or column cannot be zero. In practice, we do not prune the network, we keep the same number of *channels* every where.

This is a well known data structure in the sparse computation field: We can use Compress Block Row or Column format (CBR). There are standard Matrix Sparse-Matrix Multiplication interfaces and algorithms for CPU and GPUs using this data format (where only one operand is sparse).

In the CBR format, the γ_i are not present and only non zeros elements are stored. In other architectures, we can choose to store all non zero blocks in row format and keep a matrix Γ of zeros and ones (or columns). The Γ is a bit matrix (here) but it can be represented as a short integer matrix representing the non-zero block column and in general two orders of magnitude ($8 \times 8 \times 4$) smaller than the sparse or original B matrix. It is possible the problem has inherent block sparsity. In general, we are working with *man made* block sparsity: indeed for us, the granularity of the dense block is a property of the hardware.

In classification, the model weight size determines one important aspect of the model complexity: the number of operations per single output, the relation between layers (e.g., depth), and redundancy. The last property is important for resiliency and also to leave some space to learn new features if necessary. Exploiting sparsity is another way to reduce redundancies and computations. In our context, sparsity is the zeroing of weights (convolution and fully connected): we start with dense model using full precision and we have to find a way to chose the binary matrix Γ (which is also called a mask). For a matrix multiplication the mask is $\Gamma \in \mathbb{R}^{n \times m}$. In a convolution, n is the number of output channels (divided by 8) and m is the number of input channels (as above) and each $B_i \in \mathbb{R}^{8 \times h \times w \times 8}$ where h and w are the height and width of the convolution kernel. Imagine B having a dimension going into the paper of size $h \times w$; that is, the height and width of the kernel.

We explore training techniques (PyTorch, Keras is available in github): the most successful so far is the simplest: we take a pre-trained model, we compute a Γ per layer using a function to determine the blocks more likely to be zeros (Norm) and then we train the model till convergence or accuracy achieved. We take the sparse model and we quantize to 8-bit integer computation by using Vitis-AI quantizer. The final model is a XIR quantize model (assuming addition accumulators of 32

bits). See Section III.

For FPGA accelerator using AIE2, we have a custom compiler that takes the XIR model and an abstraction of a connected set of AIE2. See Section IV. For example, A DDR, one or two memory dedicated per column (512KB each called Mem-tile), 1 to 8 columns, 1 to 8 AIE2 cores per column, and each core has 8*8KB internal memory. There are vertical connections and there are horizontal connections. Given the HW and per layer, the compiler computes the maximum sub-volume computation per core. By heuristics and following a schedule, it computes a memory allocation in mem-tile for input, outputs, and weights. It formats the weights so that to exploit spatial distribution to Mem-tiles and cores into a compact form into a single tensor of weight, bias, and Γ .

With the schedule and the DDR-MemTile allocation, we generate all the explicit communications between DDR, MemTile, and cores. Knowing the subproblem sizes per core and the computation throughput and with a clear specification of what is executed in parallel: we can estimate the execution time per layer and of the entire network with an accuracy closer to a simulation. The code is valid, the native AIE2 compiler can interpret it, with all the DMA settings in place, and we tested for a single core in HW. We shall use this code to have a detailed time estimates for all parts of the computation: we shall show estimates for two CNN models, a few different AIE designs; see Section V.

In the following Section II, we shall start with a quantitative measure about the advantages of block sparsity in general.

II. BLOCK-SPARSE MATRIX-MATRIX MULTIPLICATION

As mental and practical exercise, consider Γ and Ω two appropriate 0,1 matrices so that for square matrices in $\mathbb{R}^{N \times N}$

$$C = (\Gamma A) * (\Omega B) \quad (5)$$

More precisely, consider non-zero blocks of size $k \times k$ so that

$$C_{i*N+j} = \sum_k (\gamma_{i*N+k} A_{i*N+k}) (\omega_{k*N+j} B_{k*N+j}) \quad (6)$$

Thanks to the sparsity and if we store only non-zeros, then γ_{i*N+k} is at the very least contiguous but ω_{k*N+j} and the right operand accesses are far from being neither simple not contiguous.

$$\dot{\Omega} B = (\Omega B)^t = \Omega^t B^t \quad (7)$$

Although expensive, the transposition of a sparse matrix is a sorting algorithm: We start from a row order to a column order, then consider K non-zeros, $K \log_2(K)$.

$$C_{i*N+j} = \sum_k (\gamma_{i*N+k} A_{i*N+k}) (\dot{\omega}_{j*N+k} \dot{B}_{j*N+k}) \quad (8)$$

There will be a meaningful product to compute if and only if $\gamma_{i*N+k} = 1$, $\dot{\omega}_{j*N+k} = 1$, and $i*N+k = j*N+k$. Without any extra information, this is simply merge sort. If we take $\gamma_{i*N+k} = i*N+k$ when the relative block is non zero, nothing otherwise, create a vector, and do the same for $\dot{\omega}_{j*N+k}$, then

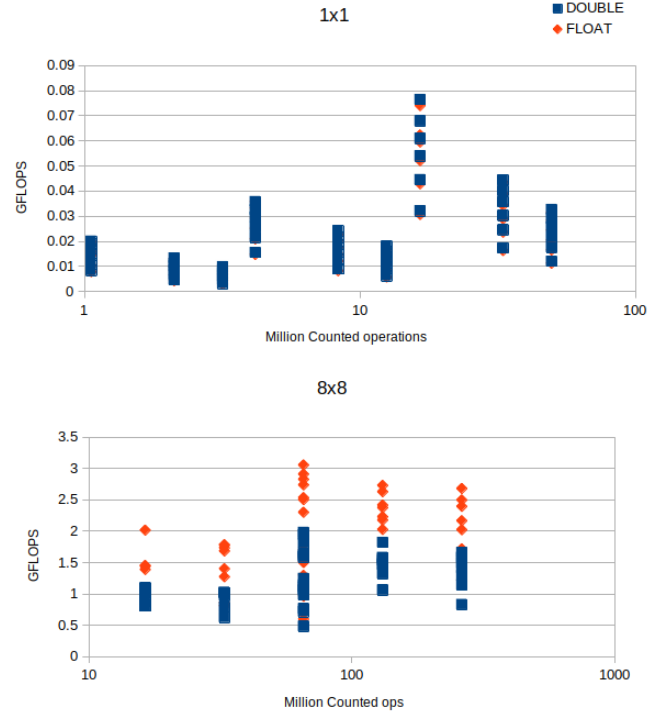


Fig. 1. Block 1x1 and 8x8 performance

we merge-sort these vectors, we do computations only on equality: we have to inspect each non zero elements M_0 and M_1 , which is $\leq O(\frac{N}{K})$. If you like to break codes yourself, see how the Sparse Sparse Matrix multiplication using Coordinate Block Structure (COO) is, please go play with [1]. We have a parallel sorting and a parallel matrix multiplication.

Now, the intuitive part, assume we want to achieve a fixed sparsity (i.e., density) of 50% for a square matrix of size N and we choose the block size $k \times k$. The number of blocks per dimension and thus the overhead for sparsity and sorting, is basically $\frac{1}{2} \frac{N}{k}$. Larger the k smaller the overhead. The relative performance of the k^3 multiplication is better as k get larger because spatial and temporal locality and optimized code for a constant/parameterized k .

In Figure 1, we present two scatter plots: on the abscissa number of effective multiplications and additions, on the ordinate the performance in GFLOPS, when the sparse matrix with dense block is 1x1 and 8x8. Given the same problem, we may use more threads and thus the Jenga of points. We can see that given the same number of effective operations, the block permits better performance and exploits better performance for each precisions: see a $2 \times$ performance for single precision computation versus double precision. This is a sparse computation and thus the poor performance (in GFLOPS) is actually expected (the peak performance in this architecture is about 500+GFLOPS).

III. BLOCK SPARSITY: TRAINING AND QUANTIZATION

In Convolutional Neural Networks, the two main operations are convolutions/correlations and fully connected layers (matrix multiplication). The block sparsity we are seeking to deploy is not naturally recurring and it has to be made. We must train the network so that we can zero blocks of data.

First, let us clarify block sparsity for convolution weights, then we clarify our training process. A convolution has a weight tensor in four dimension: $\mathbf{W} \in \mathbb{R}^{c_{out} \times h \times k \times c_{in}}$. If you can visualize the h and k dimension going into the paper: We can simplify the weight as $\tilde{\mathbf{W}} \in \mathbb{R}^{c_{out} \times c_{in}}$ and block sparsity can be simply described by a mask $\Gamma \tilde{\mathbf{W}}$. Although, we speak of a 8×8 of non zeros, this is in practice a $8 \times h \times k \times 8$ block. For the matrix multiply $h = k = 1$ and there is no difference from the previous discussions.

A. Keras

We provide a repository using Keras [2] where we implements the contents of this section: Any one can reproduce and break [3]. We target convolutions only and without quantization. The idea of the framework is simple: we take any model and we create a copy where we enhance the convolution with a (non-trainable) Γ . A convolution will have three parameters (saving the model into a different format). The forward computation is modified so that the weights used for convolution are $\Gamma \mathbf{W}$. We assume the backward computation (i.e., gradient) is done automatically from the forward definition. There is no need to change the bias. For example, we take Resnet50 from the keras application repository, we start with a $\Gamma = 1$, and we trained one epoch using imagenet repository [4]. The goal is to choose Γ so that we achieve the required sparsity and to have the minimum loss in accuracy. A little notation first:

We start from an optimized network and assume a loss function $\ell(x, w)$. In a close interval of the optimal solution w_0 we have:

$$\ell(x, w_0 + \Delta w) = \ell(x, w_0) + \nabla \ell * \Delta w + (\Delta w)^t * H * (\Delta w) + \epsilon \quad (9)$$

Where the gradient of the loss function is about zero and defined as

$$\nabla \ell * \Delta w = \sum_{w_i} \frac{\partial \ell}{\partial w_i} \Delta w_i \rightarrow 0 \quad (10)$$

The second component is the Hessian, it is symmetric and either definite positive or negative as a function of the loss.

$$(\Delta w)^t * H * (\Delta w) = \sum_{w_i} \Delta w_i \sum_{w_j} \frac{\partial^2 \ell}{\partial w_i \partial w_j} \Delta w_j \quad (11)$$

We start from w_0 , the current optimal weight, and we must choose how to reduce to zeros the weight and which ones.

In practice, using the code available in [3], we started with a accuracy of 75% top-1 and we ended up with a 68% without accounting for quantization. The investigation using PyTorch, sparsity and quantization achieves only 3% top-1 accuracy drop (instead of 6+), making the training for sparsity quite worth while (Section III-G).

B. Γ chosen once and full training ahead

Take a convolution with $\Gamma = 1$ and weights \mathbf{W} . For each γ_i , this will be representative of a block $\mathbf{W}_i \in \mathbb{R}^{8 \times h \times w \times 8} \sim \mathbb{R}^{8 \times 8}$. We can choose the \mathbf{W}_i using a measure of importance:

- $L_2 = \sqrt{\sum_k w_k^2}$ with w_k elements of \mathbf{W}_i
- $L_1 = \sum_k |w_k|$ with
- Variance $\sigma^2 = \frac{1}{64} \sum_k (w_k - \mu)^2$ with $\mu = \frac{1}{64} \sum w_k$ or $\frac{1}{N} \sum w_k$ (the whole tensor). In signal processing σ^2 is the power of the signal.

We can then sort them in ascending order and for example choose the first 50% to set to zero. Then we start re-training. We do this for the whole network or for one convolution at a time. This is the approach used in Section III-G

C. Γ chosen in steps and small fine tuning

Let say that for every convolution, $n_i = \sum \gamma_i$, which is the number of blocks. We would like to reduce this number to $\frac{n_i}{2}$. For each epoch (say every two training epochs), we consider the current non-set-yet mask elements $\sum \gamma_i = k < \frac{n_i}{2}$. We compute our importance measure for all in ascending order. This time, we zero the first $\min(\frac{5}{100} n_i, 1)$. We keep this process until we reach 50% sparsity. At each iteration at least one block will be set to zero.

D. Γ trainable as optimization problem

We think it is worth mentioning: If we want to make Γ part of the optimization process as trainable variable we could introduce as a penalty function into the loss $\ell(x, w) + \lambda(w)$.

First let us introduce an approximation for the $\max(x)$, so when in this section you will read \max, \min , this is a log sum expontial, which is continuous, derivable everywhere, and convex:

$$\max(x) = LSE(x, \alpha) = \frac{1}{\alpha} \log \sum e^{x_i * \alpha} \quad (12)$$

With T we represent the number of non zero block in Γ

$$\lambda = -(\max(\Gamma) - \min(\Gamma)) + \beta * L2(\Gamma - T) + \iota * L1(\Gamma) \quad (13)$$

This is a simplified loss so that we minimize the valu of Γ but also we try to maximize the difference of the elements.

$$\lambda = \max(-\Gamma, 0) + \max(\Gamma - 1, 0) - (\max(\Gamma) - \min(\Gamma)) + \beta * L2(\Gamma - T) + \iota * L1(\Gamma) \quad (14)$$

This last penalty function represents our attempt to state the γ_i should be positive and in the interval $[0, 1]$ and in a such a way that we maximize the distance of the elements between 0s and 1s.

$$\lambda = \max(-\Gamma, 0) + \max(\Gamma - 1, 0) - \frac{\min(\Gamma)}{\max(\Gamma)} + \beta * L2(\Gamma - T) + \iota * L1(\Gamma) \quad (15)$$

We can exploit this functionality in Keras: the penalty function can be added quite easily and per convolution (if

you like) and it is available in the code reference. We could not use it successfully: we do not know if this is because the lack of the penalty itself, the optimizer, or the lost function, or just our own fault.

E. Hessian and Fisher Information

If we have N parameters/weights, the Hessian $H \in R^{N \times N}$ has quite the space complexity (consider even small models could have million parameters). When we are already close to the optimal solution or we are trying to move from the optimal solution without using information from the gradient, the Hessian provides the most information close by an already established solution point. There are also way to compute the Hessian and the effect of the hessian by using Fisher information [5]–[7] and this will reduce to the computation of the gradient of the loss function and the citation within. We could not reproduce the authors results.

The goal is to use the H so that we know in advance what will be the effect of zeroing weights. This will boil down to a $O(h_i * w_i)$ that we can use to estimate the importance of a block

F. Diagonal Hessian

We applied a Fisher measure and we just computed $\nabla^2 \ell$, that is we compute just the diagonal of the Hessian. Again, we use the L_2 over the normalized weight and went thought the process of training. The elements of the diagonal are not representative in general, but they are a good approximation of the contribution of a single weight. Clearly, convolution is linear and any CNN is a deep network.

The Fisher and $\nabla^2 \ell$ did not provide any main advantages. But this information is found very useful in quantization and optimizations within the same field and application. We embrace the community to help us to understand this better.

G. PyTorch

In Table I, we show the results by using one-time masks and full training: VGG-16, ResNet-50, Inceptionv3 on ImageNet20 (20 classes) and ImageNet1k (1000 classes). See Section III-B. We use three samples per class for the validation accuracy for ImageNet1k data set; instead, we use 50 samples per class for ImageNet20. Fine-tuning sparse networks on the original ImageNet data-set [4] is expensive. To reduce the training time, we chose 20 classes (from the original 1000 classes) with the least number of images per class in the training data-set and this choice will affect the accuracy because there are fewer sample for re-training.

Classification accuracy on ImageNet1k drops by only 1 - 2% after applying 50% sparsity with a 8×8 block (this is without any quantization). We experiment with different block shapes such as 16×4 and 4×16 on ResNet-50, but the accuracy is slightly worse compared to block of 8×8 . Fine-grained sparsity (1×1 block) does not sacrifice any accuracy (i.e., almost any). We use the sparsified Resnet50, we quantized using Vitis AI, and we shall use it for estimate the time advantages (i.e., Section V).

TABLE I
ACCURACIES OF THE SPARSITY MODELS

Model	Dataset	Baseline Acc.(%)	Sparsity		
			block	ratio (%)	Acc.(%)
Inception-v3	ImageNet1k	77.2	8x8	50	75.5
ResNet-50	ImageNet1k	76.7	8x8	50	74.6
VGG-16	ImageNet1k	70.6	8x8	50	69.7
ResNet-50	ImageNet20	96.1	8x8	25	95.1
ResNet-50	ImageNet20	96.1	8x8	50	92.0
ResNet-50	ImageNet20	96.1	8x8	75	87.1
ResNet-50	ImageNet20	96.1	1x1	25	96.0
ResNet-50	ImageNet20	96.1	1x1	50	95.6
ResNet-50	ImageNet20	96.1	1x1	75	93.5
VGG-16	ImageNet20	92.0	8x8	50	89.6
VGG-16	ImageNet20	92.0	1x1	50	92.3
VGG-16	ImageNet20	92.0	1x1	75	91.7

IV. COMPILER AND CODE GENERATION AIE

Within our Xilinx/AMD effort, we can take a Py-Torch/Keras, quantize it using Vitis AI, and create an intermediate representation by using what we call Xilinx Intermediate Representation (XIR). This is a graph of operations that read and write tensors. A convolution has one quantized input: we use the layout format BHWC, the tensors are represented in INT8 with a position where the fraction starts (power of two scale). It computes a tensor using the same layout as the input and with a proper scale. The weights and bias are properties of the convolutions, as such they can be tailored.

For this project, we are at the third generation of a compiler for custom architectures (previously DPUV1 and DPUV3int8 [8], [9]). The main common hardware feature for this new architecture is the availability of two main memory levels: DDR and MemTile, and for such memory we import the philosophy of memory data layout and memory allocation (of the previous compilers).

All weights are statically prepared into DDR and we move them explicitly towards the inner levels. Inputs and outputs have designated space DDR (to communicate externally with PCI connection or RAM). DDR can and it will be used for swap tensors in case the inner levels do not allow allocation. The compiler can explore different schedules: the memory allocation to Mem-Tile is basically coloring algorithm plus some heuristics. In this architecture, we do not allow *streaming* of neither data or weights (because they share space in Mem-Tile). In previous architecture, we could.

A. AIE Hardware abstraction

For this presentation, see Figure 2, we work with a mesh of 4x4 AIE2 cores connected by 4 horizontal and 4 vertical interconnections. Each core has 8 banks memories for a total 64KB. About 6 banks are used as input/output/weight buffers and 2 banks are used as temporary space for kernels. Each core can request and send data to its direct neighbors (if aware of connection and control).

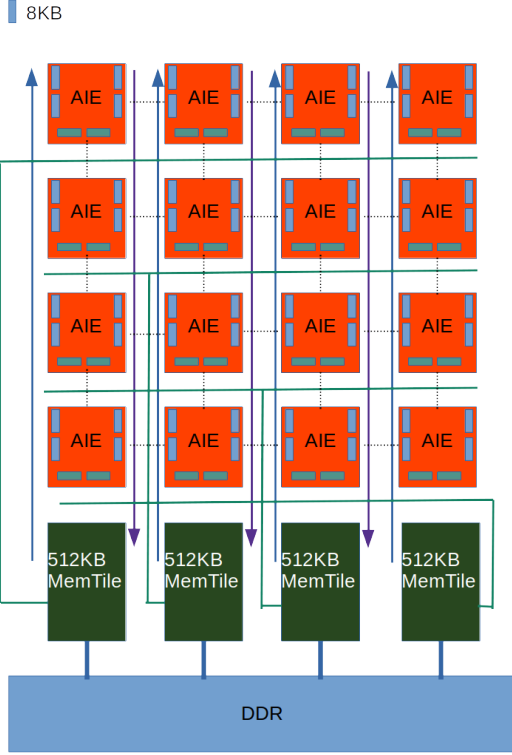


Fig. 2. 4x4 AIE representation

There are four mem-tiles: each 512KB and each can either be connected to one columns and its direct neighbor column, or it can be connected to a row and its neighbor. The total amount of space is 2MB.

A Memtile can broad-cast data per column or per row. We can dedicate a memtile for weights, one for activations, or we can share. To maximize the computation parallelism, the core will write data per column into memtile: there are different designs that can be used. We shall explore the case where Inputs/Outputs/Weights are evenly distributed across mem-tiles although the distribution is different per tensor.

The DDR is connected with two channels to write into each mem-tile and each mem-tile can use two channels to write into DDR. DDR and Mem-tile communications are parallel. The abstraction can be extended to more complex connections (from 1×1 , 8×2 , to 8×8 and more rows of memtiles designed for larger chips). But 4×4 is representative for the AIE engines you will find in the next generation of CPU+FPGA chips.

B. Sub Volume, Data compression and Data Movement

There are a few assumptions. The computation is split by column: each output tensor is split evenly by width, so each memtile will store different columns by width, each core will compute different output channels, and the computation streams computing the height of the tensor by using core input and output ping/pong. As often as possible, we store the weights in the core and we reuse them as much as possible (unless we need to stream the weight instead). The set of

cores is a cohort and symmetric computations are always selected. This means that we do not merge two operations like convolution and max-pool and give three columns to convolution and one column to max-pool. Other approach in the literature will do this more advanced methodology.

Before we start describing the memory allocation, we have to explain and compute the size of the sub-volume computation at AIE core. That is, if we have the inputs, output, and weights in memtile, what is the largest computation we can do in the AIE? The minimum computation is at least one output channel and one row by height. If this is not possible, we try to reduce the size of the width (e.g., shaping the tensor in memtile by using DDR spills) and we can manage to split the input channels (this will require to split the weights accordingly and accumulate). We call W-Split the distribution of tensor by columns in the AIE mesh. We can COUT-split, this will require the partial transfer of weights (but we can stream by height). We can CIN-split when we need to split by input channel, this is the last resort because it is also the most expensive (accumulation of the outputs). The split of the computation is also called tiling. Once we have the subvolume, we know if input and output tensor need W-split or CIN-split. This in general, means that the tensor cannot fit completely in Mem-tile and part of it will be moved. Or the computation will need to read from mem-tile the same inputs or weights multiple times.

At this stage, the subvolume describes the smallest shape of the weights we need to manage. We compress the weight accordingly to such a subvolume so that any data movement will always be a multiple of the subvolume and can be a single load. Such a compress data will have the same properties whether is sparse or not. Of course, sparse data is smaller and we compute fewer operations.

C. Schedule and memory allocation

This is a two level memory hierarchy. During the scheduling of each layer, we evaluate what can fit in mem tile. Here, Activations and weights share the space. At each step the memory allocation will check if we can allocate a tensor in memtile. If we cannot, we evict all tensors into DDR and then split/time the computation. If we can, we do, and we will evict the tensors no longer needed accordingly.

At the end of this stage, every tensor will be associated to an address in memtile or DDR (or both). If there are only DDR addresses, the compiler will take the basic computation and, by heuristics, will split the computation (thus the output tensor) by width, output channel, height, and input channel (non necessarily in this order). Every computation will have the DDR to and from MemTile and its data movements for the first two levels of the memory hierarchy. The heuristics have a simple objective: find the largest problem fitting the memory level, with the assumption that the output tensor will be build by row and streaming using ping/pong double buffering (considering padding, strides and kernels sizes).

We use a implementation that takes a recursive and iterative approach in tiling [9]: For clarity, \sum is a sum or reduction and $\dot{\sum}$ is a parallel loop and a W-split can be written as

$$Y = \text{Conv}(X, W) = \sum_w \text{Conv}(X_w, W) \quad (16)$$

The split is pre-computed as function of the foot print, before and after each convolution there will be an explicit data movement. At this stage each input, output, and weights will have an address associated with each sub-computation. Consider a subcomputation as data movement, computation, and data movement, and data movements are optional or decorators. Then the code generation of each $\text{Conv}(X_w, W)$ is independent and, recursively and as needed, there will be specific splits of the computation accordingly. This could be simply described using loops but it is actually a tree (i.e, a root with k children). Sub convolutions so determined may have different memory requirements, we try as much as possible to obtain as symmetric computation as possible. If the convolution has strides and large kernel, each sub-convolution may have overlap data but defined addresses and data movement if necessary. For example, computing the output by rows and the weights are reused.

$$\forall w, \text{Conv}(X_w, W) \rightarrow \sum_h \text{Conv}(X_{w,h}, W) \quad (17)$$

D. Code Generation

The compiler recursively creates a list of operations smaller and smaller that can actually be executed Mem-Tile to Mem-Tile. In practice, there is a further decomposition using only AIE cores but it is completely determined by the subvolume computation as we explained previously. Then, we can explicitly determine the data movements from mem-tile to core and back.

A peculiar feature of the current system is the concept of *iteration*. Using locks and chaining them (locks like in semaphores): We can create a single iteration instruction from the prospective of a single core (as a SIMD instruction) but driving all cores at once (ASM-like code)

```
LOAD Lock k_0 mem-tile add core add iter i
CONV iteration i
SAVE Lock k_1 mem-tile add core add iter i
```

There is an implicit lock (say k_x) that is used for the pong and the system cycles in between locks (k_x and k_0)

These three operations will be execute a number of iterations i and using a ping/pong they will load different slices of data and compute different slices of data. The key ingredient for this to work is the volumes of input and output tensors are all the same.

In our environment padding is common and we can manage by a custom load (from memtile to core) and this requires a custom load that will not be repeated

```
## Head top padding
LOAD Lock k_0 mem-tile add_0 core add iter 1
CONV iteration 1
SAVE Lock k_1 mem-tile add_1 core add iter 1
## Body iteration
## k_0 -> k_2 -> k_4 Lock Chain
LOAD Lock k_2 mem-tile add_2 core add iter 9
```

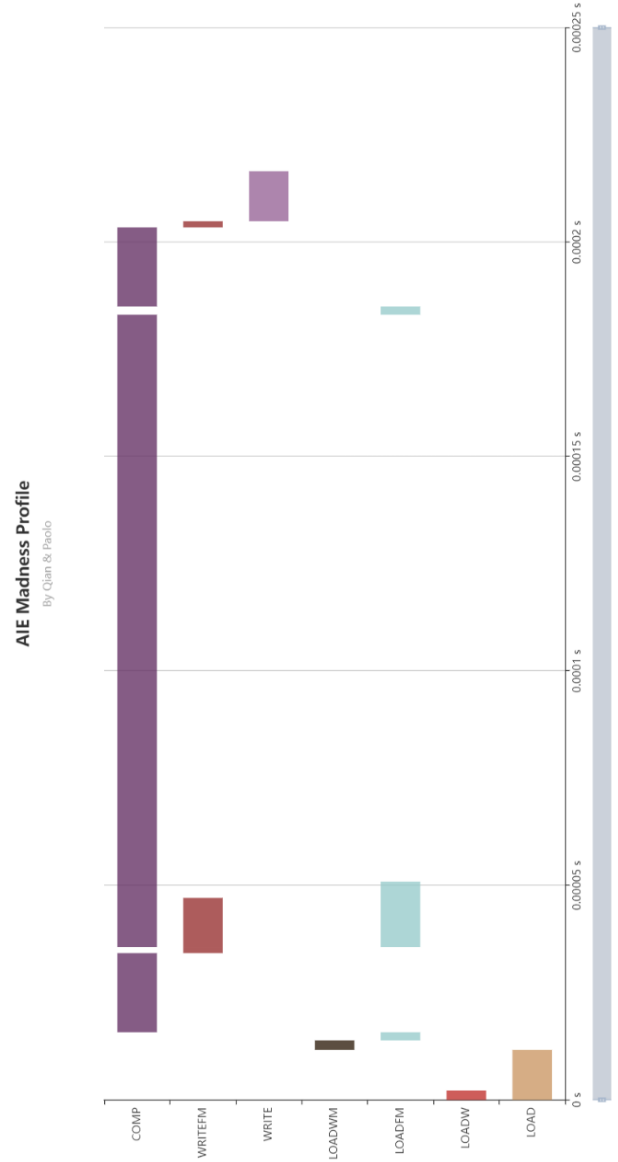


Fig. 3. Resnet single convolution with padding for 4x4 AIE

```
CONV iteration 7
SAVE Lock k_3 mem-tile add_3 core add iter 9
## tail
LOAD Lock k_4 mem-tile add_4 core add iter 1
CONV iteration 1
SAVE Lock k_5 mem-tile add_5 core add iter 1
```

See Figure 3 for how this code will play out in practice. At this stage we have all the information and the code generation per layer is a two pass process: first, we generate code for the all load/store and then we combine them into chain having dependency so that to be logically correct and as fast as

possible. Take the code above and we introduce a chain information.

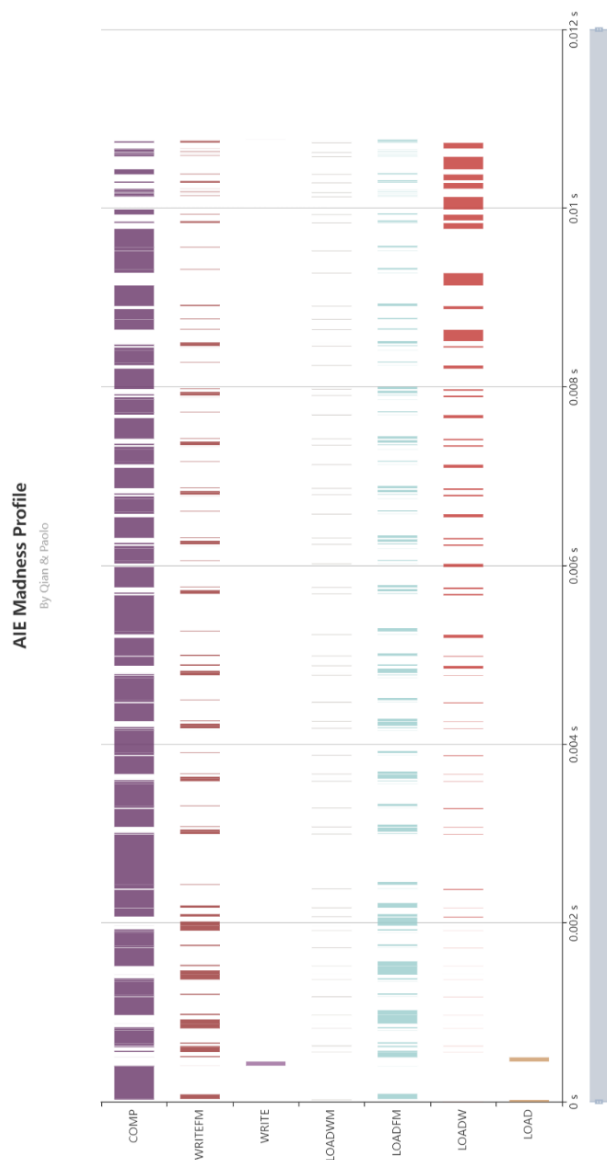


Fig. 4. Resnet50 for 4x4 AIE with dense weights

E. Time Estimation

At this stage, we need to explain how we can capture the execution time and visualize it as in Figure 4.

F. Time estimates for DDR to Mem-Tile

We have two communication types: activations and weights. Per mem-tile there are two dedicated channels.

- if we share activations and weights in the same mem-tile (we use the moniker WAWAWAWA), we can use one channel for activations and one for weights. Thus the

loads from DDR to MEM-tile (LOAD and LOADW) are parallel with a bandwidth of 4GB/s. Writes from mem-tile to DDR (WRITE) can use both channels (8GB/s). We shall expand this.

- If activations and weights go to different mem-tiles (2 and 2), each load is parallel and 8GB/s (we use the moniker WAAW). Writes are identical.

Although, mem-tile are circular buffers to improve the memory allocation, the classic streaming is not really applicable for two main reasons: First, when we share weights they must stay still to be reused, if weights are not shared then in general inputs and outputs have different rates of consumption and thus eventually one will overwrite the other. Ping/pong techniques will decrease the space available but it will hide latency: at this time we try to increase the reuse.

G. Time estimates for Mem-Tile to AIEs

The Memtile connections with AIE cores can be designed differently as a function of the purpose. We assume here a few channels with again 4GB/s bandwidth. One memtile can broadcast inputs to columns of cores (and to the nearest neighbors). These communications are for activations (LOADFM). One Memtile can broadcast to rows of cores (or the nearest neighbor), these are for weights (LOADWM). We assume that the column and row communications are parallel and each memtile core connection is parallel.

Every communication with iteration 1 is synchronous and binding: that is sequential, the load, convolution, and store is one after the other and every core is independent. For synchronization and for bookkeeping, we assume that AIE2 weights communication (from memtile) to core are synchronous and halting.

Every communication with iteration larger than one, we assume that load, computation (COMP), and store (WRITEFM) are executed in parallel and the overall execution time is the maximum of the estimated time multiplied by the number of iterations.

We estimate the execution time of a subvolume by the number of operations divided by the maximum number of operations per cycle which is in our scenario: $4 \times 8 \times 8 = 256$ operations per cycle and 1GHz frequency. This is obviously a very optimistic validation and not only for convolutions (asking for a ratio of 1 in efficiency) but also for operations like average pools and element wise additions. The execution time is a feature of the analysis but for us the estimate of the communications is more compelling and we can easily mute the computation contribution.

We do not account the wait and synchronization necessary to reprogram the fabric.

During and especially once we generated the code for all data movements, data volumes and destinations, the trigger of the computation, and the sub volume per computation we can estimate quite accurately the execution time for each operations and account their parallel execution.

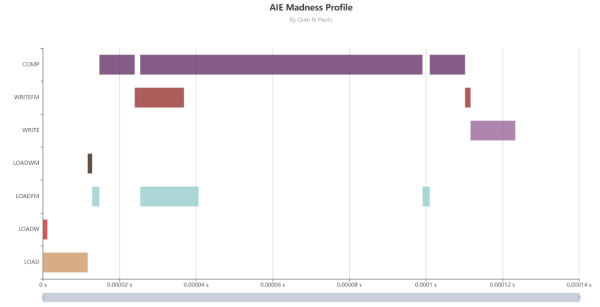


Fig. 5. Resnet single convolution with padding and sparsity for 4x4 AIE

H. Convolution example

Here and in Figure 3 and 5, we give a full convolution example with and without sparsity and with padding. In this way we can explain how the time is really estimated and also how the hardware works in principle.

It is clear from the Figures above that there are three computations (COMP). We load the weight and activations once in memtiles. There are actually one load per memtile for a total of 4 loads per activations and 4 loads for weights. Because each load is to a different memtile, they are parallel. The activation and weights communications are using two different channels and then are in parallel with 4GB/s bandwidth.

There is a single load of the weights from mem-tiles to each core. This is done once and it is blocking (LOADWM) but it can be as easily made non blocking and parallel to the activations. There is a computation using padding (top-padding) and you can see the sequential execution of load to cores (LOADFM), computation (COMP), and write to memtile (WRITEFM). This computation has iteration 1. There are 9 iterations for three instructions: we can see the load, the computation, and write in parallel. This is obviously a simplification; there will be a little load poking out at the beginning and a writing poking out at the end. Then we conclude with the final computation with padding at the bottom of the computation.

In this particular scenario, the computation dominates the execution time and compression basically cut the execution time by half: from 200 μs to 130. There are convolutions in Resnet that realize exactly $2 \times$ performance but also there are convolution that are dominated by the read or by the writing, and where sparsity help only in space saving.

V. RESULTS

In this section, we present the performance of sparsity applied to all the convolutions (but the first one) for Resnet 50 Figure 4 and 6, and Inception V3 Figure 7.

When we generate the code for each instruction, we compute the execution time. So if we inspect the assembly code we will find time information in the context whether or not each instruction contribute directly. For data movement to and from DDR and mem-tile, we reduce the contribution (sum directly).

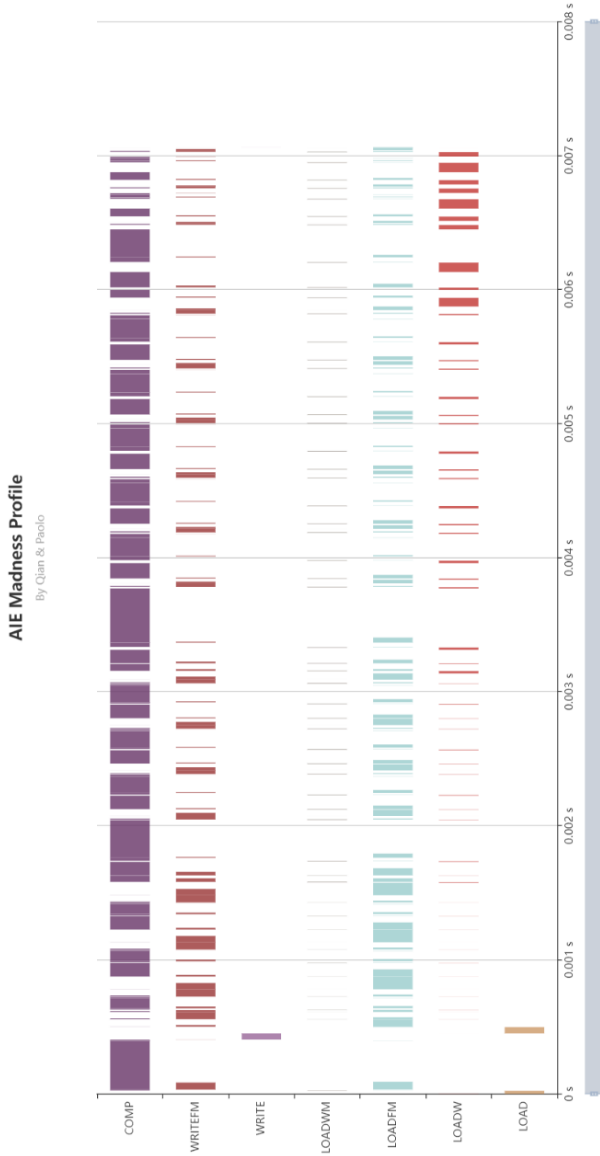


Fig. 6. Resnet50 for 4x4 AIE with 50% sparse weights

There is no streaming or communication overlapping, thus the sum.

For mem-tile to and from core communications and core computations, we create a time series. We explain in the previous section how we account for the execution time for instruction with and without iterations. All of this will be an attribute of the layer (node in the graph computation). To create a complete time estimate, we just need to take the schedule of the computation and the graph, visit each node accordingly to the schedule, write a *json* file describing

the time series, then by *javascript* we can visualize the time series using a browser. The Figure in this paper are generated directly.

For a 4×4 AIE set up, Resnet 50 fits in memtile from the beginning to the last operation (beside the first convolution and this is by construction). The estimates advantage by sparsity is almost completely achieved. See Figure 4 and 6.

For the estimate of Inception V3, we use a 6×6 AIE set up. This allows to have a limited spill into DDR and thus the sparsity delivery once again an overall $2\times$ speed up. See Figure 7. Inception V3 is mostly convolutions of different sizes.

VI. CONCLUSIONS AND CONTEXT

This is a multifaceted problem and we present a complete solution from training techniques, compilers, code generation, HW definition, and time estimations.

This could be seen as a quantization and sparsification problem: How to reduce the footprint of a CNN network. There are post training techniques that are targeting quantization and unstructured sparsity [10] and all the reference within and its citations. For the type of sparsity we are aiming at, we need to be more aggressive and training for it (as starting point [11]) but our sparsity is not really a property of the model, it is a sparsity that the software can describe and the hardware can take advantage without having any hardware support.

The difficulty of generate code for complex architectures can be described and solved in different ways. There are recent attempts of introducing SW/HW solution for spatial processors like ours [12]–[14]. The main difference:

- We use software developed for FPGA custom IPs because we can reuse the 2-level of memory hierarchy abstraction.
- We use heuristics to explore scheduling.
- We use heuristics to memory allocation.
 - There is a memory allocation so temporal locality is exploited and different (layer) schedules can be investigated.
 - All passes know all tensor shapes at compile time (no run time support).
- We write *all* the data movements and the code can and is used by a AIE2 compiler (binary compiler if you like) to create executable codes.
 - Code generated for 1×1 have been simulated, run and validate in HW.
 - we assume at this time that core to core communication by columns (larger kernel and stride require overlapping column data) or communication by row (reduction of partial sums) are inherently taken care by the kernel computation
 - Kernel computation (the basic computation in the AIE2) has effects on the shapes and layout of the operands, the HW architect instruct the compiler about the tensor requirements.
 - The compiler does not generate fully correct addresses but all the tensor shapes are valid and correct for any $K \times K$.

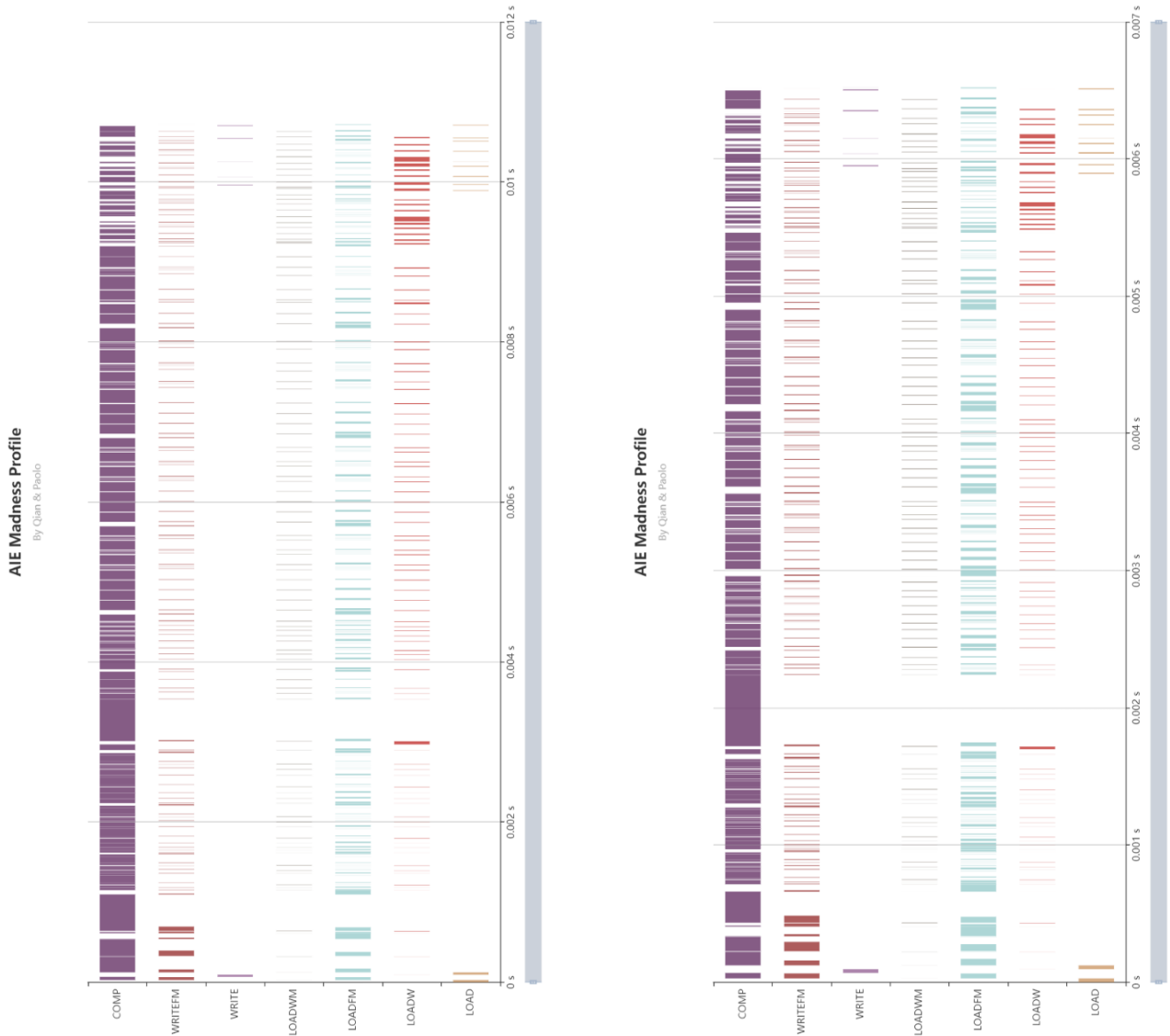


Fig. 7. Inception V3 for 6x6 AIE with dense (left) weights and sparse (right)

- No full or partial simulation is needed, and now time estimates can be used to quantify scheduling and HW configurations choices: yes there is further arbitrary connection and hardware configurations that can be changed and the effect can be only validate by writing code for it..

REFERENCES

- [1] P. D'Alberto, <https://github.com/paolodalberto/SparseFastMM/>, 2013.
- [2] F. Chollet *et al.*, "Keras," <https://keras.io>, 2015.
- [3] P. D'Alberto, "Block sparsity and training," <https://github.com/paolodalberto/BlockSparsityTraning>, 2020.
- [4] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [5] Z. Yao, A. Gholami, S. Shen, K. Keutzer, and M. W. Mahoney, "Adahessian: An adaptive second order optimizer for machine learning," *AAAI (Accepted)*, 2021.
- [6] S. Yu, Z. Yao, A. Gholami, Z. Dong, M. W. Mahoney, and K. Keutzer, "Hessian-aware pruning and optimal neural implant," *CoRR*, vol. abs/2101.08940, 2021. [Online]. Available: <https://arxiv.org/abs/2101.08940>
- [7] B. Zandonati, A. A. Pol, M. Pierini, O. Sirkin, and T. Kopetz, "Fit: A metric for model sensitivity," 2022.
- [8] P. D'Alberto, V. Wu, A. Ng, R. Nimaiyar, E. Delaye, and A. Sirasao, "Xdnn: Inference for deep convolutional neural networks," *ACM Trans. Reconfigurable Technol. Syst.*, vol. 15, no. 2, jan 2022. [Online].

Available: <https://doi.org/10.1145/3473334>

- [9] P. D'Alberto, J. Ma, J. Li, Y. Hu, M. Bollavaram, and S. Fang, "DPUV3INT8: A compiler view to programmable FPGA inference engines," *CoRR*, vol. abs/2110.04327, 2021. [Online]. Available: <https://arxiv.org/abs/2110.04327>
- [10] E. Frantar, S. Ashkboos, T. Hoefler, and D. Alistarh, "Gptq: Accurate post-training quantization for generative pre-trained transformers," 2023.
- [11] B. Hawks, J. M. Duarte, N. J. Fraser, A. Pappalardo, N. Tran, and Y. Umuroglu, "Ps and qs: Quantization-aware pruning for efficient low latency neural network inference," *CoRR*, vol. abs/2102.11289, 2021. [Online]. Available: <https://arxiv.org/abs/2102.11289>
- [12] Q. Huang, M. Kang, G. Dinh, T. Norell, A. Kalaiah, J. Demmel, J. Wawrzynnek, and Y. S. Shao, "Cosa: Scheduling by constrained optimization for spatial accelerators," *2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA)*, pp. 554–566, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:233740109>
- [13] E. Russo, M. Palesi, G. Ascia, D. Patti, S. Monteleone, and V. Catania, "Memory-aware dnn algorithm-hardware mapping via integer linear programming," *Proceedings of the 20th ACM International Conference on Computing Frontiers*, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:260442361>
- [14] J. Cai, Y. Wei, Z. Wu, S. Peng, and K. Ma, "Inter-layer scheduling space definition and exploration for tiled accelerators," *Proceedings of the 50th Annual International Symposium on Computer Architecture*, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:259177591>