# Modeling Deep Learning Accelerator Enabled GPUs

Md Aamir Raihan [*], Negar Goli [*] , and Tor M. Aamodt

Electrical and Computer Engineering
University of British Columbia

{araihan, negargoli93, aamodt}@ece.ubc.ca

*Abstract*—**The efficacy of deep learning has resulted in its use in a growing number of applications. The Volta graphics processor unit (GPU) architecture from NVIDIA introduced a specialized functional unit, the "tensor core", that helps meet the growing demand for higher performance for deep learning. In this paper we study the design of the tensor cores in NVIDIA's Volta and Turing architectures. We further propose an architectural model for the tensor cores in Volta. When implemented a GPU simulator, GPGPU-Sim, our tensor core model achieves 99.6% correlation versus an NVIDIA Titan V GPU in terms of average instructions per cycle when running tensor core enabled GEMM workloads. We also describe support added to enable GPGPU-Sim to run CUTLASS, an open-source CUDA C++ template library providing customizable GEMM templates that utilize tensor cores.**

*Keywords*-**Tensor Core, Tesla Titan V, Turing RTX 2080, CUTLASS library, GPGPU-Sim**

## I. INTRODUCTION

Deep neural networks (DNNs) are having impact in a growing number of areas but the benefits of DNNs come at the expense of high computational cost. Deep learning based data analytics has recently emerged as an important technique [1]. DNNs have enabled breakthroughs in speech recognition [2], [3], image recognition [4], [5] and computer vision [6], [7]. DNNs require performing a large number of multi-dimensional matrix (or tensor) computations. Recent research has explored how to accelerate these operations [8]–[15] and many companies are developing custom hardware for these workloads [16]–[18].

GPUs are commonly used for deep learning, especially during training, as they provide an order of magnitude higher performance versus a comparable investment in CPUs [19]. Specific effort has been directed at optimizing GPU hardware and software for accelerating tensor operations found in DNNs. On the hardware side, in the Volta architecture NVIDIA introduced a specialized function unit called a Tensor Core for this purpose. Tensor cores are also found on NVIDIA's more recent Turing architecture [20] and NVIDIA's T4 Turing-base GPUs are further optimized for inference tasks [21]. NVIDIA claims [22] tensor cores provide a speedup of $3\times$ on the Tesla V100 GPU when

running mixed precision training. Five out of six 2018 Gordon Bell Award Finalists employed tensor cores to improve application performance and three did so specifically by accelerating machine learning [23].

However, to the best of our knowledge, the underlying design of tensor cores has not been publicly described by NVIDIA. Thus, we investigated the NVIDIA tensor cores found in both Volta and Turing architectures. Informed by our analysis we extended GPGPU-Sim [24] to include a model for tensor cores.

This paper makes the following contributions:

- It shows how different threads cooperate in transferring an input matrix to each tensor core.
- It gives an in-depth analysis of the execution of the tensor operation on the tensor cores and describes the microbenchmarks we used to perform our analysis.
- It proposed a microarchitectural model of tensor cores consistent with the characteristics revealed through our microbenchmarks.
- It describes our functional and timing model changes for modeling tensor cores in GPGPU-Sim.
- It describes support we added to enable applications built with NVIDIA's CUTLASS library to run on GPGPU-Sim.
- It quantifies the accuracy of our modified GPGPU-Sim by running tensor core enabled kernels generated with CUTLASS and thereby demonstrating an IPC correlation of 99.6%.

We believe the observations made in this paper will provide useful guidance to those wishing to explore how to incorporate deep learning accelerators within GPUs. The corresponding changes to model tensor cores in GPGPU-Sim should provide the academic community a helpful baseline for comparing alternative approaches. The changes to enable CUTLASS to run on GPGPU-Sim should ease study of architectural characteristics of custom kernels on frameworks such as PyTorch (which was recently enabled to run on GPGPU-Sim [25]).

* equal contribution

## II. Background

This section briefly summarizes, at a high-level, relevant aspects of the Volta GPU architecture as documented by NVIDIA, the source-code and instruction-level interfaces for programming Tensor Cores on NVIDIA GPUs before finally describing what NVIDIA has disclosed about the design of their Tensor Cores.

### A. Volta Microarchitecture

The first GPU to include accelerators for machine learning was NVIDIA's Volta [26]. Recent NVIDIA GPUs including Volta are generally composed of multiple Streaming Multiprocessors (SM) connected by an on-chip network to multiple memory partitions. Each memory partition contains a portion of the last-level cache and connects the GPU to off-chip DRAM. As described by NVIDIA, multiple tensor cores are included inside each SM. The SM design in Volta is partitioned into four processing blocks which NVIDIA refers to as Sub-Cores. As shown in Figure 1, each sub-core in Volta has two tensor cores, one Warp scheduler, one dispatch unit, and a 64 KB register file.

Besides the addition of tensor cores, Volta includes other enhancements relevant to performance of machine learning workloads: In comparison to Pascal, NVIDIA's prior GPU architecture, each streaming multiprocessor (SM) in Volta has twice as many scheduling units along with separate integer and 32-bit floating point (FP32) cores. In addition, handling of divergent threads is different in Volta versus prior GPUs in that both paths following a branch can be executed by threads within a single warp in an interleaved fashion.

NVIDIA typically releases several GPUs with the same underlying architecture but different amounts of on-chip resources. For Volta, we focus in this paper on the Titan V GPU. The SM inside the Titan V has the same number of registers as Pascal. However, the Titan V GPU has 24 more SMs and thus can support more threads, warps, and thread blocks compared to prior generation GPUs.

### B. Warp Matrix Function (WMMA) API

CUDA 9.0 [28] introduced a "warp matrix function" C++ language API to enable programmers to use the tensor cores on supported GPUs. This interface is also referred to as the CUDA C++ "warp-level matrix multiply and accumulate" (WMMA) API [29], [30]. It is well known that tiling can improve memory locality for dense matrix operations [31]. The WMMA API exposes tensor cores to the GPU programmer as warp-wide operations for performing the computation $D = A \times B + C$, where $A$, $B$, $C$ and $D$ can be tiles of larger matrices. Using the WMMA API, all threads in a warp cooperatively work together to perform a matrix-multiply and accumulate operation on these tiles. NVIDIA's WMMA API currently specifies a limited set of tile sizes. The sizes for tiles $A$, $B$, $C$ and $D$ are represented
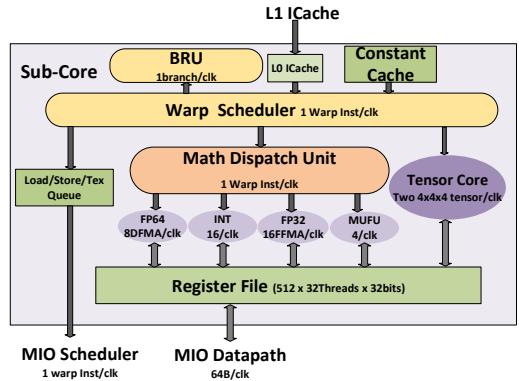


Figure 1: Votla SM Sub-Core (reproduced from [27])

using the notation $M \times N \times K$, where $M \times K$ is the dimension of Tile A, $K \times N$ is the dimension of Tile B and thus $C$ and $D$ have dimension $M \times N$. CUDA 9.0 supports only one tile sizes, $16 \times 16 \times 16$, while later versions allow additional flexibility.

Using NVIDIA's terminology, each tile is further divided into "fragments" where a fragment is a set of tile elements that are mapped into the registers of a single thread. Thus, input matrices are distributed across different threads and each thread contains only a portion of a tile. NVIDIA specifically states [28] the mapping of tile elements to registers is unspecified. Naively, considering a $16 \times 16$ tile contains 256 elements, one possibility would be that each thread in a warp with 32 threads would store an $\frac{256}{32} = 8$ element fragment in eight separate general-purpose registers. In Section III we show that current GPU's do something more sophisticated.

The CUDA WMMA API provides three new functions: `load_matrix_sync`, `store_matrix_sync` and `mma_sync`. All three functions perform an implicit warp-wide barrier synchronization before computing a result. The `load_matrix_sync` and `store_matrix_sync` functions are used for loading and storing a portion of the input matrices in the general-purpose registers accessible to each thread. The `mma_sync` function performs a warp synchronous matrix multiply-accumulate operation producing an $M \times N$ (e.g., $16 \times 16$) result in the general-purpose registers associated with the tile for the D matrix.

NVIDIA provides four high-level programming interfaces for using tensor cores: the WMMA API described above and three CUDA libraries: cuBLAS [32], cuDNN [33], [34] and CUTLASS [35], [36]. In addition, many deep learning frameworks have included support for tensor cores [37], [38].

### C. PTX Instruction Set

NVIDIA's toolchain compiles CUDA into host code that runs on the CPU and device code that runs on the GPU.

```
wmma.load.a.sync.layout.shape.type ra, [pa] {stride};
wmma.load.b.sync.layout.shape.type rb, [pb] {stride};
wmma.load.c.sync.layout.shape.type rc, [pc] {stride};
wmma.mma.sync.alayout.blayout.shape.dtype.ctype rd, ra, rb, rc;
wmma.store.d.sync.layout.shape.type rd, [pd] {stride};
```

Figure 2: Tensor Core PTX instructions

The device code is first compiled into a device-independent machine-language instruction set architecture known as Parallel Thread eXecution (PTX) before being compiled into device-specific machine code (SASS).

To perform operations on Tensor Cores at the PTX level, NVIDIA introduced three PTX instructions in PTX version 6.0 [30] with the syntax shown in Figure 2. In this figure the "sync" qualifier indicates that the instruction waits for all threads in the warp to synchronize before beginning execution. The PTX manual uses the term "operand matrix" to refer to a tile. The "layout" qualifier specifies whether an operand matrix is stored in memory with a row-major or column-major layout. The "shape" qualifier represents the fragment size of the operand matrices (e.g., $16 \times 16 \times 16$ is specified by setting shape to m16n16k16). The "type" qualifier indicates the precision of the operand matrices, i.e. FP16 or FP32. For Volta, the A and B matrices must be FP16 but the C operand matrix can be either FP16 or FP32. NVIDIA's Turing architecture supports additional integer arithmetic modes initially targeted for inference. In these, the operand matrices A and B can be 8, 4, or 1-bit signed or unsigned integers and operand matrices C and D are kept in higher-precision INT32 format to avoid overflow during accumulation [39].

The operand matrices A, B and C must be loaded from memory to the register-file prior to initiating a matrix-multiply operation. This data movement is accomplished via three wmma.load PTX instructions. Specifically, wmma.load.a, wmma.load.b and wmma.load.c load the matrices A, B and C respectively into registers ra, rb and rc where ra, rb and rc represent sets of general-purpose registers distributed across the threads of a warp corresponding with the notion of a fragment. pa, pb, pc are the memory address where operand matrices A, B and C are stored in memory.

Typically, input tiles loaded from memory are a portion of a larger matrix. To help accessing tiles of a larger matrix, wmma.load and wmma.store support strided-memory access. The "stride" operand specifies the beginning of each row (or column).

The wmma.mma PTX instruction performs a warp-level matrix-multiply with accumulate operation. This instruction computes $D = A \times B + C$ using registers a, b and c which contain the matrix A, B and C respectively. The computed results are stored in general-purpose registers d in each



Figure 3: Tensor cores complete one $4 \times 4$ MACC operation per cycle ($D = A * B + C$). Reproduces Figure 8 in [26].

thread.

### D. Tensor Core

Each tensor core is a programmable compute unit specialized for accelerating machine learning workloads. The Tesla Titan V GPU contains 640 tensor cores distributed across 80 SMs, with eight tensor cores per SM, providing a theoretical performance of 125 TFLOPS at an operational frequency of 1530 MHz. According to NVIDIA [26], each tensor core can complete a single $4 \times 4$ matrix-multiply-and-accumulation (MACC) each clock cycle, i.e. $D = A \times B + C$, where $A, B, C$ are $4 \times 4$ matrices as shown in Figure 3. While individual tensor cores operate on $4 \times 4$ matrices at any one time, as noted earlier, the WMMA API exposes the tensor cores on tile-sizes which are much larger. Naively, a multiply of two $16 \times 16$ matrices decomposes into a blocked matrix-multiply involving four $4 \times 4$ matrix-multiply accumulates for each of the sixteen $4 \times 4$ submatrices of the result matrix. Thus, each mma_sync at the CUDA C++ WMMA level or each wmma.mma operation at the PTX level may be implemented with 64 separate tensor core operations. The tensor cores have two modes of operation: FP16 and mixed-precision. In FP16 mode, the tensor core reads three $4 \times 4$ 16-bit floating-point matrices as source operands whereas in mixed-precision mode it reads two $4 \times 4$ 16-bit floating-point matrices along with a third $4 \times 4$ 32-bit floating-point accumulation matrix.

### III. DEMYSTIFYING NVIDIA'S TENSOR CORES

In this section we describe the results of our attempt to better understand the low-level implementation details of tensor cores on recent GPUs. Our analysis extends and refines that of Jia et al. [40] who examined the distribution of matrix operand elements to registers for mixed precision mode in column-major layout. In their work, Jia et al. [40] refer to a group of four consecutive threads within a warp as

```
<FRAGMENT_DECLARATION> a_frag;
wmma::load_matrix_sync(a_frag, mem_addr, stride );
for(int i=0; i < a_frag.num_elements; i++) {
  float t=static_cast<float>(a_frag.x[i]);
  printf("THREAD%d CONTAINS %.2f\n",threadIdx.x,t);
}
```
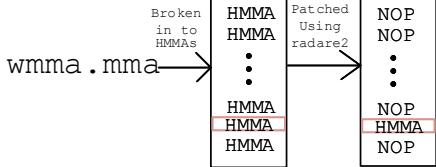
Figure 4: Microbenchmark for decoding thread fragments



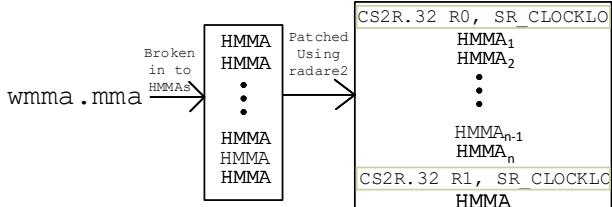Figure 5: Analyzing data accessed by tensor cores



Figure 6: Analyzing tensor core timing

a "thread group". We find it more convenient to shorten this to *threadgroup*, which we do in the remainder of the paper. As there are 32 threads in a warp, there are 8 *threadgroups* in a warp. We will refer to the *threadgroup id*[1] of a given thread, which is given by $\lfloor \frac{threadIdx}{4} \rfloor$.

### A. Microbenchmarks

In this section we discuss the microbenchmarks[2] we used for analyzing the implementation of tensor cores. We employ two types of microbenchmarks: Ones designed to determine how data move into and out of the tensor cores and others used to determine how long the tensor cores take to perform operations.

*1) Fragment to thread mapping:* Figure 4 contains a portion of the CUDA code employed in Section III-B to determine the mapping between operand matrix elements and threads. This code is part of a larger general matrix multiplication (GEMM) kernel operating on a $16 \times 16$ matrices. Each thread loads a segment of the input matrix and prints it to the output console. By initializing each element of the input matrix with different values it is straightforward to uncover the mapping from operand matrix element to thread with a warp.

*2) Analyzing machine instructions:* As described in detail in Section III-C wmma.mma PTX instructions are mapped into multiple HMMA SASS instructions. Figure 5 illustrates, at a

[1]Similar to "group id" in Jia et al. [40].
[2]https://github.com/gpgpu-sim/gpgpu-sim_simulations/tree/master/benchmarks/src/cuda/tensorcore-microbenchmarks

high level, the operation of our microbenchmarks used for analyzing how data is accessed by HMMA instructions. We use radare2 [41] to replace all HMMA operations except one with "no operation" (NOP) instructions. Figure 6 illustrates at a high-level the approach used by our microbenchmarks for analyzing the timing of low level operations on tensor cores. To develop these microbenchmarks we used radare2 to add code that reads the clock register before the $1^{st}$ and after the $n^{th}$ HMMA instruction.

### B. Operand matrix element mapping

In this section we summarize the results of our analysis of the distribution of matrix elements to threads.

*1) Volta Tensor Cores:* Figures 7a and 7b summarize how the elements of matrix operands are mapped to the registers of individual threads within a warp. The large rectangle (❶) represents $16 \times 16$ operand matrix A or B for both FP16 and mixed-precision modes of operation. Smaller squares are individual elements of the operand matrix and elements in the same row are stored contiguously in memory. Each *threadgroup* loads a different $4 \times 16$ sub-matrix, which we will refer to as a segment. The four segments that make up the operand matrix are highlighted with different shading.

The upper right-hand portion of Figure 7a (❷,❸) shows how the elements within a segment are distributed among the threads of a threadgroup. Our analysis found that on Volta, each segment is loaded by two different *threadgroups*. Thus, each element of the A and B operand matrices are loaded by two different threads in a warp on Volta. The bottom portion of Figure 7a (❹) combined with the top-left portion (❶) summarize the exact mapping. For example, we found the first four consecutive rows of operand matrix A are loaded by *threadgroup* 0 and 2.

The distribution of matrix elements to threads for operand matrix A stored in row-major layout is the same as the distribution of operand matrix B stored in column-major layout and vice-versa. For the operand matrix A in row-major layout, each thread inside the *threadgroup* loads 16 consecutive elements using two coalesced 128-bit wide load instructions (❷) whereas in column major layout each thread inside the *threadgroup* loads four blocks of four consecutive elements via four coalesced 64-bit wide load instructions, each with a stride distance of 64 elements (❸).

As illustrated in Figure 7b, the distribution of matrix elements to threads is different for operand matrix C. Specifically, for operand matrix C each *threadgroup* loads a $8 \times 4$ segment of the matrix C. Also, the specific distribution within the threadgroup now depends on whether the matrix C is FP16 or FP32 and is independent of the layout. 32-bit wide (partially coalesced) load instructions are used to access elements of matrix C in both modes of operation.

*2) Turing Tensor Cores:* Figure 8 summarizes the distribution of operand matrix elements to threads for tensor cores in NVIDIA's Turing architecture. Turing's tensor cores

(a) Operand matrices A and B.
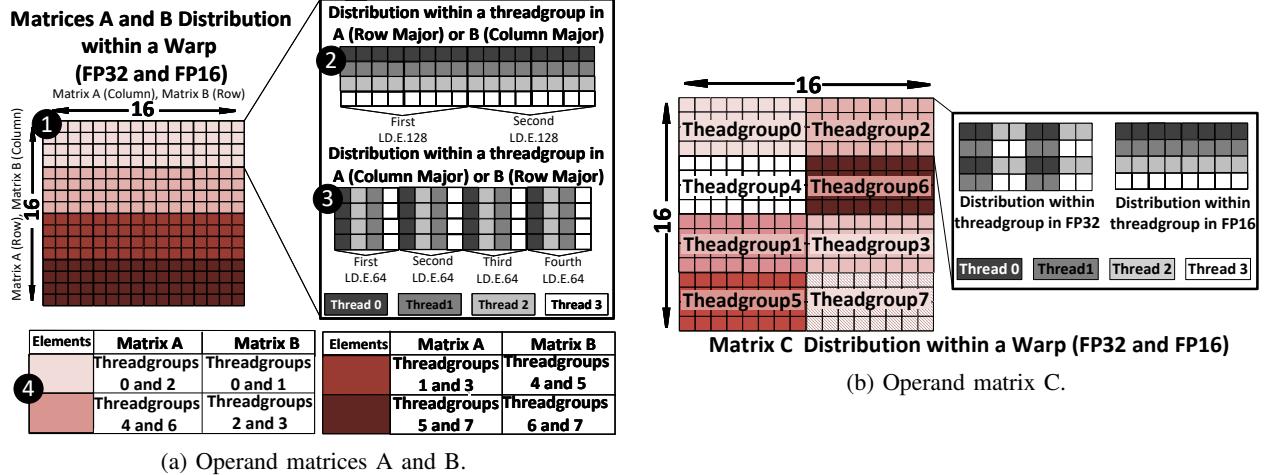


(b) Operand matrix C.

Figure 7: Distribution of operand matrix elements to threads for Tensor Cores in the Titan V (Volta).

support three new precision modes: 1-bit, 4-bit and 8-bit, along with three new tile sizes: $32 \times 8 \times 16$ and $8 \times 32 \times 16$ for 8 and 16-bit modes and $8 \times 8 \times 32$ for 4-bit mode. Support for 1-bit mode was only enabled very recently as of this writing and did not appear to work on our system. Thus, no analysis is provide for 1-bit mode in the rest of this paper. We found Turing has a simpler distribution of elements to threads than Volta. Specifically, each operand matrix element is loaded only once. Both tile size $32 \times 8 \times 16$ and $8 \times 32 \times 16$ employ the same distribution. For all modes and configurations, each row or column (depending on the mode and operand matrix) is loaded by a *threadgroup* and consecutive *threadgoups* load consecutive rows or columns.

### C. Machine ISA interface

This section summarizes what we learned about how Tensor Cores are accessed at the machine instruction set architecture level. This level is typically called SASS for NVIDIA GPUs. The analysis here is based upon examining SASS disassembly using NVIDIA's cuobjdump tool.

We found that wmma.load and wmma.store PTX instructions are implemented by being broken into a group of normal SASS load (LD.E.64, LD.E.128, LD.E.SYS) and store (ST.E.SYS) instructions. This suggests that Tensor Cores access operand matrix fragments directly from the normal GPU register file. In more detail, we found the wmma.load.c PTX instruction is broken into a group of LD.E.SYS instructions. For operand matrices A and B, depending on whether the operand matrix layout is row major or column major, wmma.load PTX instructions are broken into either four 64-bit loads (LD.E.64) or two 128-bit loads (LD.E.128), respectively.

Figure 9 illustrates the SASS code for Volta corresponding with a single wmma.mma PTX instruction. As can be seen in this figure, matrix-multiply accumulate operations are implemented via a new SASS instruction, HMMA. Each HMMA instruction has four operands and each operand uses a pair of registers. By comparing the registers used by the HMMA and the loads and stores, we have inferred that a pair of adjacent registers accessed by different memory operations are encoded in the HMMA instruction using a single register identifier. For example, R8 in the first HMMA instruction in Figure 9 appears from our analysis to represent the register pair <R8,R7>. The higher register identifier in the register pair is the one encoded in the instruction. For example, for the HMMA instruction on the first line of of Figure 9, the destination register R8 actually represents the pair <R8,R7>. Similarly, the remaining register identifiers actually represent three pairs of source operand registers (<R24,R23>, <R22,R21> and <R8, R7>). Each of the four pairs of registers corresponds to operand matrices A, through D.

Some registers are annotated with "reuse" in Figure 9. Gray [42] analyzed NVIDIA's SASS instruction set for the earlier Maxwell architecture where a similar annotation often appears. Based upon his analysis and related papers from NVIDIA on register file caching for GPUs [43], we believe the "reuse" notation indicates the associated operand is reused in the next step and therefore cached in the operand reuse cache to avoid a register fetch and possibly to reduce bank conflicts.

*1) Volta Tensor Cores:* Each wmma.mma PTX instruction is broken into a group of HMMA instructions.

Figure 9a illustrates the SASS code for mixed precision mode. In this mode, each PTX wmma.mma instruction is broken into 16 HMMA instructions. These are organized as four *sets* of four HMMA instructions. Each HMMA instruction is annotated with "STEP<n>" where <n> ranges from 0 to 3. Thus, each set comprises four steps. Figure 9b illustrates the SASS code for FP16 mode in which a single PTX wmma.mma instruction is broken into four sets consisting of only 2 steps.
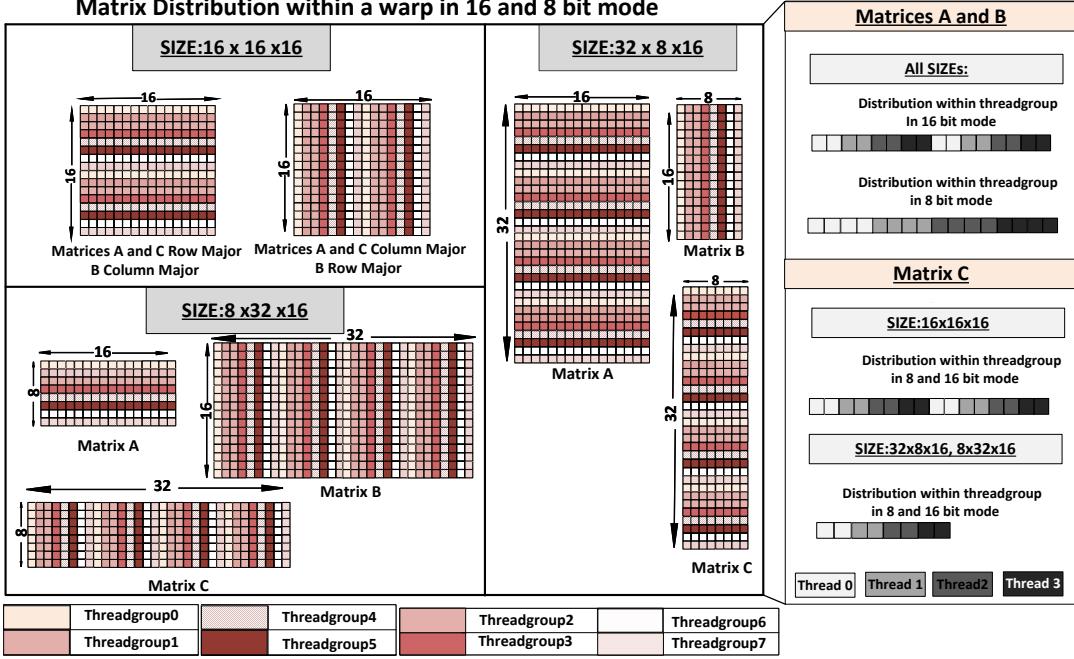
## Matrix Distribution within a warp in 16 and 8 bit mode



Figure 8: Distribution of operand matrix elements to threads for tensor cores in the RTX 2080 (Turing).

```
                                                                    Cumulative
                                                                    Clock Cycles
       HMMA.884.F32.F32.STEP0 R8, R24.reuse.COL, R22.reuse.ROW, R8;    10
SET1   HMMA.884.F32.F32.STEP1 R10, R24.reuse.COL, R22.reuse.ROW, R10;  12
       HMMA.884.F32.F32.STEP2 R4, R24.reuse.COL, R22.reuse.ROW, R4;    14
       HMMA.884.F32.F32.STEP3 R6, R24.COL, R22.ROW, R6;                18
       HMMA.884.F32.F32.STEP0 R8, R20.reuse.COL, R18.reuse.ROW, R8;    20
SET2   HMMA.884.F32.F32.STEP1 R10, R20.reuse.COL, R18.reuse.ROW, R10;  22
       HMMA.884.F32.F32.STEP2 R4, R20.reuse.COL, R18.reuse.ROW, R4;    24
       HMMA.884.F32.F32.STEP3 R6, R20.COL, R18.ROW, R6;                28
       HMMA.884.F32.F32.STEP0 R8, R14.reuse.COL, R12.reuse.ROW, R8;    30
SET3   HMMA.884.F32.F32.STEP1 R10, R14.reuse.COL, R12.reuse.ROW, R10;  32
       HMMA.884.F32.F32.STEP2 R4, R14.reuse.COL, R12.reuse.ROW, R4;    34
       HMMA.884.F32.F32.STEP3 R6, R14.COL, R12.ROW, R6;                38
       HMMA.884.F32.F32.STEP0 R8, R16.reuse.COL, R2.reuse.ROW, R8;     40
SET4   HMMA.884.F32.F32.STEP1 R10, R16.reuse.COL, R2.reuse.ROW, R10;   42
       HMMA.884.F32.F32.STEP2 R4, R16.reuse.COL, R2.reuse.ROW, R4;     44
       HMMA.884.F32.F32.STEP3 R6, R16.COL, R2.ROW, R6;                 54
```

(a) Disassembled SASS instructions for Mixed precision mode

```
                                                              Cumulative
                                                              Clock Cycles
SET1   HMMA.884.F16.F16.STEP0 R4, R22.reuse.T, R12.reuse.T, R4;   12
       HMMA.884.F16.F16.STEP1 R6, R22.T, R12.T, R6;               21
SET2   HMMA.884.F16.F16.STEP0 R4, R16.reuse.T, R14.reuse.T, R4;   25
       HMMA.884.F16.F16.STEP1 R6, R16.T, R14.T, R6;               34
SET3   HMMA.884.F16.F16.STEP0 R4, R18.reuse.T, R8.reuse.T, R4;    38
       HMMA.884.F16.F16.STEP1 R6, R18.T, R8.T, R6;                47
SET4   HMMA.884.F16.F16.STEP0 R4, R2.reuse.T, R10.reuse.T, R4;    51
       HMMA.884.F16.F16.STEP1 R6, R2.T, R10.T, R6;                64
```

(b) Disassembled SASS instructions for FP16 mode

Figure 9: Disassembled SASS instructions corresponding to WMMA:MMA API

Figure 9 also shows the cumulative clock cycles for the Volta Tensor Cores. The latency of wmma.mma API in mixed precision mode is ten cycles lower than in FP16 mode.

*2) Turing Tensor Cores:* For Turing, each PTX wmma.mma instruction is broken into a group of four HMMA instructions for all modes except 4-bit where it is converted into

a single HMMA instruction. Table I shows the cumulative clock cycles for HMMA instructions on the Turing architecture. For $16 \times 16 \times 16$ tile size, the latency of wmma.mma in mixed precision mode on Turing, 99 cycles (Table I), is more than on Volta, 54 cycles (Figure 9a). The latency of mixed precision mode is more than FP16 mode. 8-bit mode is fastest. The latency of 4-bit mode is the highest, which may be because it is an experimental feature on the 2080 RTX.

| TILE SIZE (MxNxK) | PRECISION | Average Cumulative Clock Cycles | | | |
|---|---|---|---|---|---|
| | | SET 1 | SET 2 | SET 3 | SET 4 |
| 16x16x16 | 16Bit (FP32 Acc) | 42 | 56 | 78 | 99 |
| | 16Bit (FP16 Acc) | 44 | 52 | 60 | 74 |
| | 8Bit | 40 | 44 | 47 | 59 |
| 32x8x16 | 16Bit (FP32 Acc) | 48 | 60 | 81 | 104 |
| | 16Bit (FP16 Acc) | 44 | 52 | 60 | 74 |
| | 8Bit | 52 | 55 | 59 | 73 |
| 8x32x16 | 16Bit (FP32 Acc) | 42 | 56 | 77 | 99 |
| | 16Bit (FP16 Acc) | 42 | 50 | 58 | 72 |
| | 8Bit | 38 | 42 | 46 | 56 |
| 8x8x32 | 4Bit | 230 | - | - | - |

Table I: Average cycles to execute all HMMA instructions up to SET $n$ on Turing. "Acc" is accumulation mode.

### D. HMMA Instruction Analysis

This section explores HMMA execution in greater detail.

*1) Volta:* We examine the operation of each "set" of HMMA instructions in Figure 9. As shown in Figure 10a, irrespective of mode, when executing the HMMA instructions in a set, each *threadgroup* multiplies a $4 \times 4$ sub-tile of
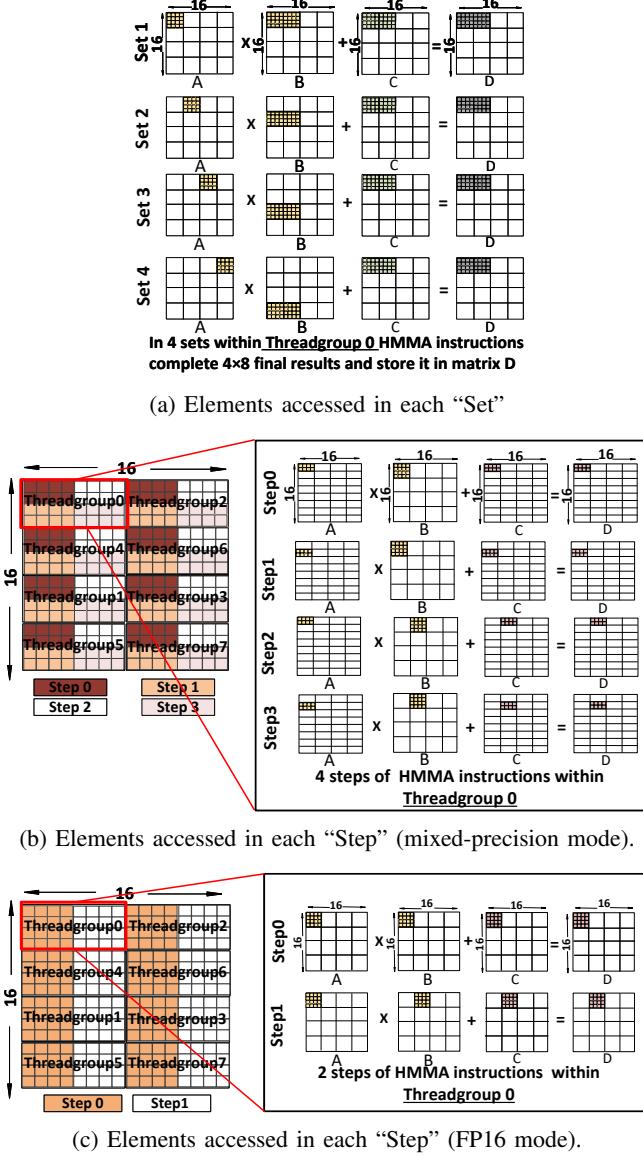
(a) Elements accessed in each "Set"



(b) Elements accessed in each "Step" (mixed-precision mode).



(c) Elements accessed in each "Step" (FP16 mode).

Figure 10: HMMA instruction analysis for Volta (Titan V).

operand matrix A with a $4 \times 8$ sub-tile of operand matrix B and accumulates the result with operand matrix C. For example, when *threadgroup* 0 executes the first set of HMMA instructions (Set 1) it multiplies the sub-tile consisting of the first four rows and columns of operand matrix A with the sub-tile consisting of the first four rows and first eight columns of operand matrix B. The result is accumulated with a $4 \times 8$ sub-tile of operand matrix C and stored in a $4 \times 8$ sub-tile of operand matrix D.

Figure 10b shows the detailed operation of each HMMA "step" within a "set" for threadgroup 0 for mixed-precision mode. Each "set" of HMMA instructions contains four "steps". We find in each step, a $2 \times 4$ sub-tile of operand matrix A

is multiplied with a $4 \times 4$ sub-tile of operand matrix B, accumulated with a $2 \times 4$ sub-tile of operand matrix C.

Similarly, Figure 10c shows the detailed operation of each HMMA "step" within a "set" for threadgroup 0 for FP16 mode. Each set of HMMA instructions contains two "steps". In each step, every threadgroup multiplies a $4 \times 4$ sub-tile of operand matrix A with a $4 \times 4$ sub-tile of operand matrix B and accumulates the result with matrix C.

*2) Turing:* Figure 11 illustrates the elements accessed by HMMA instructions on the Turing GPU architecture. The "step" annotation found on HMMA SASS instructions in Volta is not present in Turing. Given the latency results in Table I do not suggest increased parallelism one possibility is similar "steps" are sequenced by the microarchitecture using a state-machine. We make the following observations:

- The elements accessed for a particular mode are similar for different tile configurations.
- In FP16 and mixed precision mode the computation pattern is the product between two subtiles where one of the subtile is $8 \times 8$ and the other subtile is either $16 \times 8$ or $8 \times 16$. For example, for tile size $16 \times 16 \times 16$ or $32 \times 8 \times 16$, the computation in SET 1 is the product between the $16 \times 8$ subtile of matrix A with the $8 \times 8$ subtile of matrix B whereas for tile size $8 \times 32 \times 16$ the product is between the $8 \times 8$ subtile of matrix A with the $8 \times 16$ subtile of matrix B.
- For the 8-bit mode the computation pattern is the product between the $8 \times 16$ subtile of matrix A with $16 \times 8$ subtile of matrix B.
- In 4-bit mode, each wmma.mma PTX instruction is implemented with a single HMMA SASS instruction so we omit 4-bit mode in Figure 11.

### E. Discussion

In this section, we provide our analysis of the results presented above for Volta and infer a possible rationale for why execution is broken into "sets" and "steps".

Recall, each element of the input matrix is loaded by two different *threadgroups*. We wrote a microbenchmark to help determine how the fragments loaded by different threads are used by a HMMA instruction. For example, to determine how operand matrix elements loaded by thread 0 are used, we altered these values and observed how the result is affected. We found that *threadgroups* work in pairs to compute $8 \times 8$ subtiles of the result. We call each such pair of threadgroups an octet. There are four octets in a warp.

Table II shows the pair of *threadgroup* constituting each *octet*, which in general can be formulated as octet X = *threadgroup* X $\bigcup$ *threadgroup* X+4  where X lies in between 0 to 3. Table II also uses the notation [Row_Start : Row_End, Col_Start : Col_End ] to show the subtile of the operand matrix A and B accessed by the threads inside each octet. The elements loaded by the octet remain the same
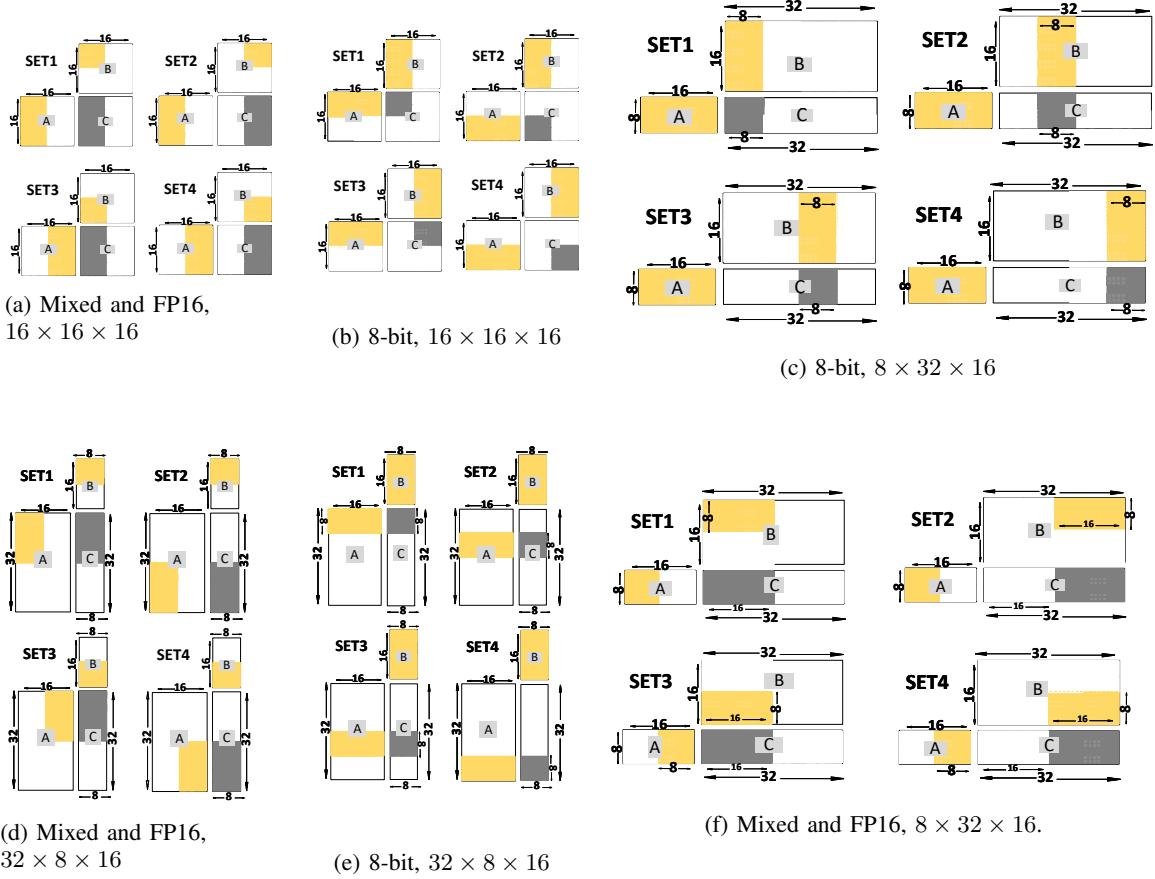
(a) Mixed and FP16, $16 \times 16 \times 16$

(b) 8-bit, $16 \times 16 \times 16$

(c) 8-bit, $8 \times 32 \times 16$

(d) Mixed and FP16, $32 \times 8 \times 16$

(e) 8-bit, $32 \times 8 \times 16$

(f) Mixed and FP16, $8 \times 32 \times 16$.

Figure 11: HMMA instruction analysis for Turing (RTX2080).



(a) Elements of operand matrices accessed by each octet

(b) Outer product formulation during sets and steps in an octet

(c) Cycles to execute parallel HMMA operations versus number of warps per SM
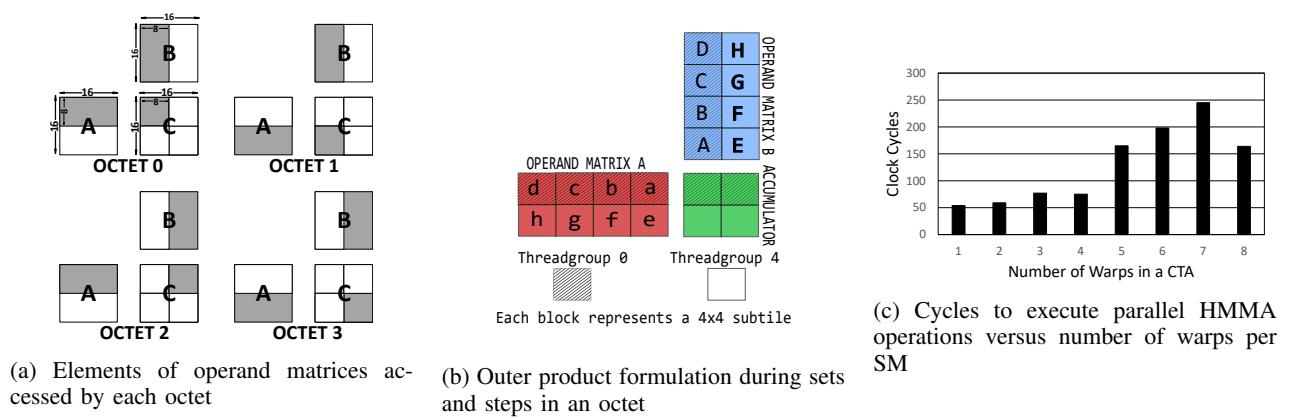
Figure 12

irrespective of the layout in which the operand matrices are stored.

Table II shows each element of the operand matrices A and B is loaded twice by threads in a different *threadgroup*. This enables each octet to work independently. Specifically, each octet reads an $8 \times 16$ subtile of operand matrix A, an $16 \times 8$ subtile of operand matrix B and an $8 \times 8$ subtile of operand matrix C as shown in Figure 12a.

To better understand the organization of threads into octets, we analyzed the calculation performed by octets in different "sets" and "steps". As shown in Figure 12b, in each set, every octet performs the outer product between input subtiles. For example, in Set 1 the outer product between input subtile $[a], [e]$ and $[A], [E]$ is completed to generate the partial result $[aA], [aE], [eA]$ and $[eE]$. Here each $[a], [e], [A], [E]$ represents a $4 \times 4$ subtile. To compute $[aE]$, *threadgroup* 0 needs operand matrix B subtile $[E]$ which is only loaded by *threadgroup* 4. Similarly, to compute $[eA]$, *threadgroup* 4 needs operand matrix B subtile $[A]$ which is only loaded by *threadgroup* 0. Thus, while *threadgroups* cannot, octets *can* work independently. Table III expands upon Figure 12b to tabulate all the outer product computations performed in different sets and steps.

| Octet | Threadgroup | Matrix A | Matrix B |
|---|---|---|---|
| 0 | 0 and 4 | [0:7,0:15] | [0:15,0:7] |
| 1 | 1 and 5 | [8:15,0:15] | [0:15,0:7] |
| 2 | 2 and 6 | [0:7,0:15] | [0:15,8:15] |
| 3 | 3 and 7 | [8:15,0:15] | [0:15,8:15] |

Table II: Octet composition and elements accessed

| SET | STEP | Threadgroup X | Threadgroup X+4 |
|---|---|---|---|
| 1 | 0 | $a[0:1] \times A$ | $e[0:1] \times A$ |
| | 1 | $a[2:3] \times A$ | $e[2:3] \times A$ |
| | 2 | $a[0:1] \times E$ | $e[0:1] \times E$ |
| | 3 | $a[2:3] \times E$ | $e[2:3] \times E$ |
| 2 | 0 | $b[0:1] \times B$ | $f[0:1] \times B$ |
| | 1 | $b[2:3] \times B$ | $f[2:3] \times B$ |
| | 2 | $b[0:1] \times F$ | $f[0:1] \times F$ |
| | 3 | $b[2:3] \times F$ | $f[2:3] \times F$ |
| 3 | 0 | $c[0:1] \times C$ | $g[0:1] \times C$ |
| | 1 | $c[2:3] \times C$ | $g[2:3] \times C$ |
| | 2 | $c[0:1] \times G$ | $g[0:1] \times G$ |
| | 3 | $c[2:3] \times G$ | $g[2:3] \times G$ |
| 4 | 0 | $d[0:1] \times D$ | $h[0:1] \times D$ |
| | 1 | $d[2:3] \times D$ | $h[2:3] \times D$ |
| | 2 | $d[0:1] \times H$ | $h[0:1] \times H$ |
| | 3 | $d[2:3] \times H$ | $h[2:3] \times H$ |

Table III: Octet computation details

## IV. A TENSOR CORE MICROARCHITECTURE

In this section we present a tensor core microarchitecture consistent with the observations made for Volta earlier in the paper.

Recall each tensor core completes a $4 \times 4$ matrix-multiply and accumulate each cycle. To achieve this, each tensor core must be able to perform sixteen four-element dot-products (FEDPs) each cycle. As shown in Figure 9 and 10b, in steady state, a *threadgroup* takes two cycles to generate a $2 \times 4$ subtile of the output matrix. Thus, across all threads in a warp a HMMA instruction is executing 32 FEDP per cycle. Since each tensor core can only complete 16 FEDP per cycle it follows that full throughput requires two tensor cores per sub-core within an SM. To confirm this we wrote a microbenchmark that repeatedly executes HMMA operations, varies the number of warps per thread block and the number of thread blocks executing concurrently constant. As shown in Figure 12c, this microbenchmark shows that only four warps can concurrently execute on a single SM, but the Titan V SM has 8 tensor cores per SM. Thus, each warp appears to utilize two tensor cores.

Next, we consider register access bandwidth. The data in Figure 9a suggests the minimum initiation interval of an HMMA instruction is two cycles. There are three source operands and as noted in Section III-C for each source operand a pair of 32-bit registers is read. Taking all these factors into account the total register fetch bandwidth is $32 \times 2 \times 3 \times 32 = 6$kb every 2 cycles per warp. This bandwidth is sufficient for a warp to fetch the following every two cycles: eight $2 \times 4$ FP16 subtiles for operand A, eight $4 \times 4$ FP16 subtiles for operand B, and eight $2 \times 4$ FP32 subtiles or eight $4 \times 4$ FP16 subtiles for operand C. Given every warp accesses two tensor cores, the register bandwidth per tensor core is 1.5kb per warp per clock cycle.

NVIDIA states that in Volta INT and FP32 instructions can be co-issued [26]. On the other hand tensor core operations reportedly cannot be co-issued with integer and floating-point arithmetic instructions [44]. We believe the reason is that the tensor cores may be using the register file access ports associated with the INT and FP32 cores. There are 64 INT and 64 FP32 ALUs inside Titan V SM for a total of 128 ALUs. With eight tensor cores inside an SM sharing access to the register file each tensor core should be able to access $\frac{128}{8} \times 32 = 16 \times 32 = 512$ bits per source operand per cycle. Assuming three source operands per ALU (to support multiply-accumulate operations) this means each tensor core can access 1.5kb/cycle.

Figure 13 illustrates our proposed tensor core microarchitecture. Each warp utilizes two tensor cores. We assume two octets within a warp access each tensor core. Sixteen SIMD lanes are dedicated to each tensor core, eight to each octet, and four to each *threadgroup*. Each *threadgroup* lane fetches the operands into internal buffers. For operand matrix A and C, each *threadgroup* fetches the operands to its separate buffer whereas for operand matrix B both the *threadgroups* fetches to a shared buffer. The mode of operation and steps determine the *threadgroup* lane from which each operand is fetched. The buffers feed sixteen FP16 FEDP units. Inside
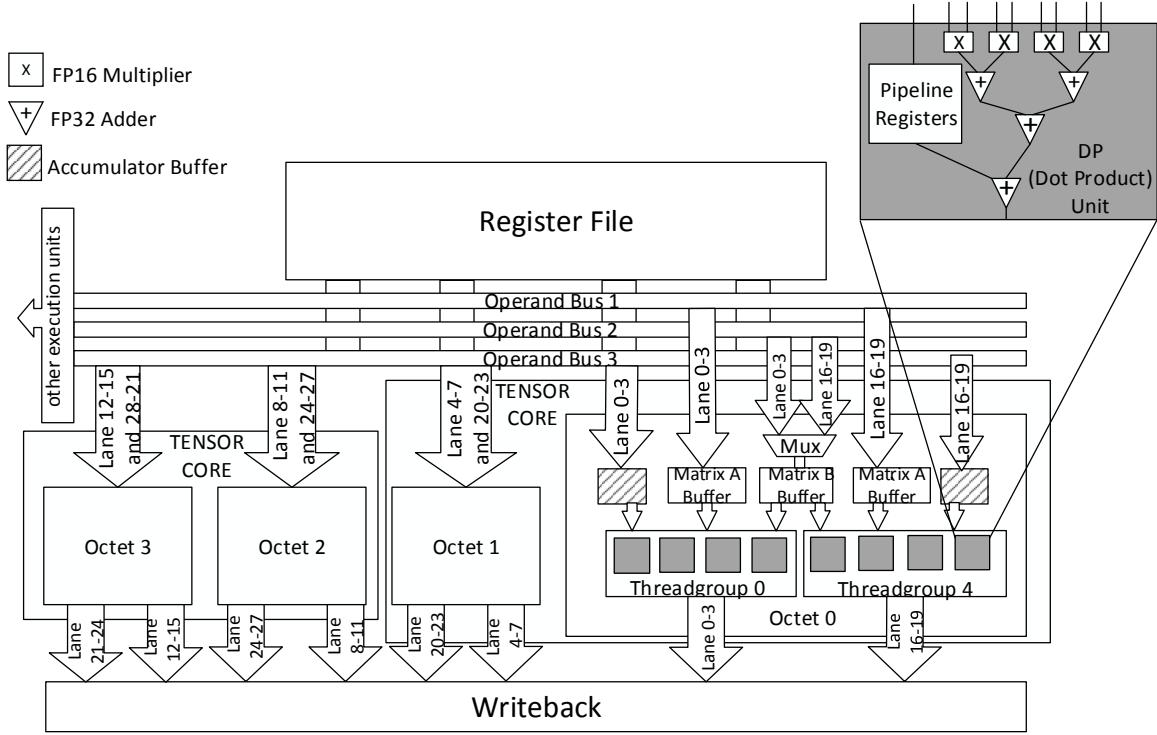
Figure 13: Proposed Tensor Core Microarchitecture

each FEDP unit, multiplication is performed in parallel in the first stage and accumulation occurs over three stages for a total of four pipeline stages. As each tensor core consists of sixteen FP16 FEDP units, it is capable of completing one $4 \times 4$ matrix multiplication each cycle.

## V. MODELING AND EVALUATION

### A. Modelling Tensor Cores

Our changes to model the tensor cores in Volta are available in the "dev" branch of GPGPU-Sim [24] on github[3]. We extended the current version of GPGPU-Sim to support 16-bit floating-point by using a half-precision C++ header-only library [45]. The library provides an efficient implementation of 16-bit floating-point conforming to the IEEE 754 half-precision format. It provides common arithmetic operations and type conversion. GPGPU-Sim currently only supports SASS execution for the G90 architecture; therefore, we only model tensor core operations at the PTX level. To do so, we added functional and timing models for the `wmma.load`, `wmma.mma` and `wmma.store` PTX instructions described in Section II-C.

Our functional model of the `wmma.load` and `wmma.store` PTX instructions support all possible layout combinations for operand matrix A, B and C. Our functional model follows

the operand matrix element to thread mapping shown in Figure 7. We have verified the timing model generates the exact same number of coalesced memory transactions generated by the Titan V GPU for these operations.

Our functional model of the `wmma.mma` instruction supports all 32 possible configurations supported on the Titan V GPU. A timing model for the tensor core functional unit is added to the GPU pipeline. We interface our tensor core timing model to the operand collector unit modeled in GPGPU-Sim. Each `wmma.mma` instruction is issued to the tensor core unit after all of its source operands are ready in the operand collector. We updated the scoreboard to check for RAW and WAW hazard associated with `wmma.mma` instructions.

We validate our tensor core model by comparing against an NVIDIA Tesla Titan V with CUDA Capability 7.0, hosted by an Intel Core i7-4771 3.50GHz based workstation with Ubuntu 16.04.4 LTS, CUDA Toolkit Version 9.0, NVIDIA 410.48 GPU driver, and gcc 4.9.4. Figure 14a compares the cycles required to execute a WMMA based matrix-multiply and accumulate kernel on the Titan V GPU and GPGPU-Sim as matrix size varies. We find GPGPU-Sim tracks real hardware very accurately with a standard deviation of less than 5%. This is despite the fact our model is implemented at the PTX level.

---

(a) WMMA-based GEMM kernel cycle count as matrix size varies.

(b) Instructions per cycle (IPC) correlation of CUTLASS GEMM kernel on GPGPU-Sim vs Titan V.

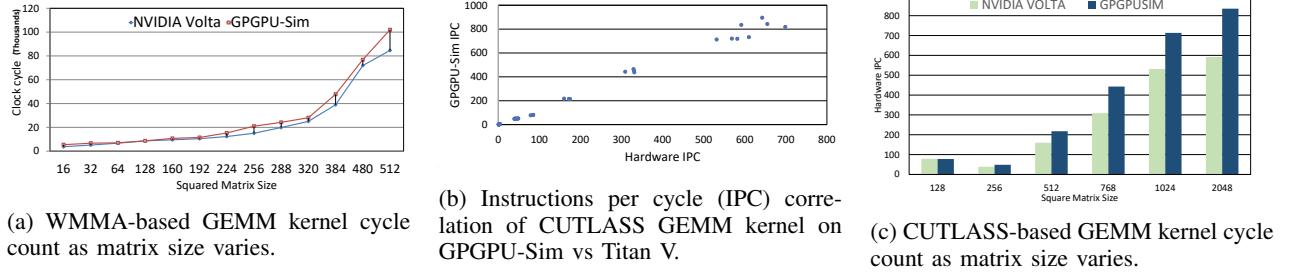(c) CUTLASS-based GEMM kernel cycle count as matrix size varies.

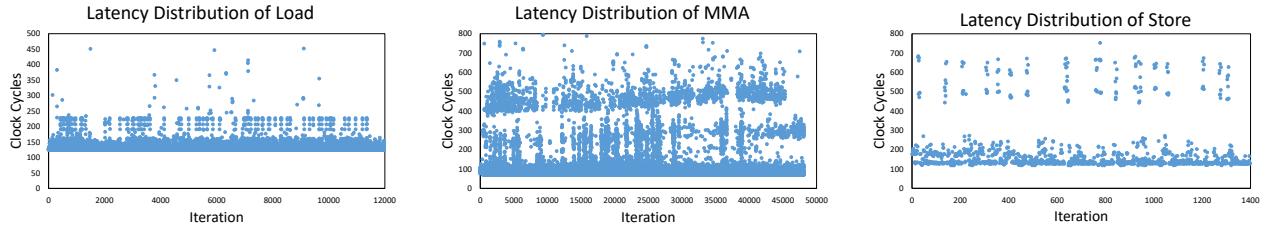Figure 14: Comparison of simulated and actual performance



Figure 15: Distribution of `wmma.load`, `wmma.mma` and `wmma.store` latency for matrix size $1024 \times 1024$ GEMM using shared memory
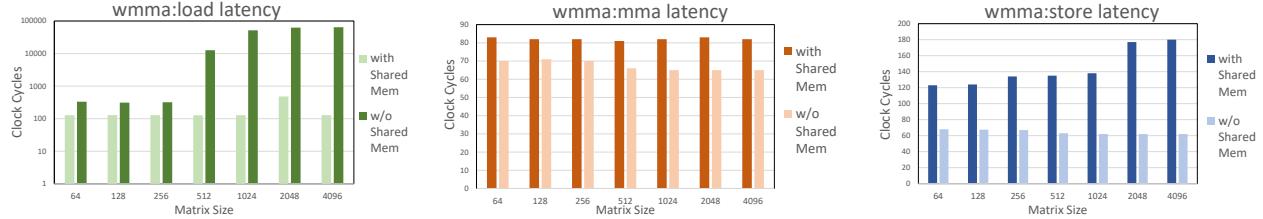


Figure 16: Variation of latency of `wmma.load`, `wmma.mma` and `wmma.store` with matrix size
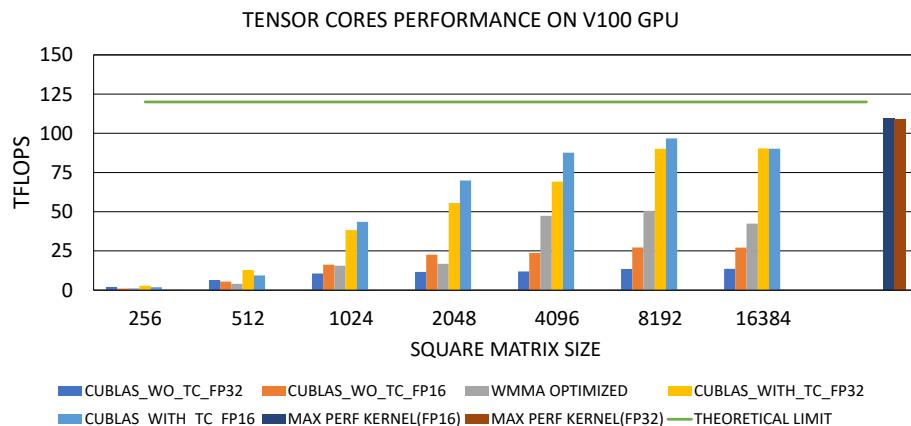


Figure 17: Tensor Cores Performance

### B. CUTLASS

CUTLASS is an open-source CUDA C++ template library for efficient linear algebra in C++. It provides basic building block for implementing high-performance fused matrix-multiply kernels for deep learning.

We modified GPGPU-Sim to enable it to run CUTLASS including adding missing API calls and PTX instruction definitions. NVIDIA developed a unit-test suite for CUTLASS library consisting of around 680 test cases. We verified these test cases run with our modifications to GPGPU-Sim[4]. Figure 14b shows a comparison of Instructions Per Cycle (IPC) measured on GPGPU-Sim versus a real NVIDIA Titan V GPU for a tensor core enabled kernel developed using CUTLASS. This data shows an IPC correlation of 99.60%. Figure 14c shows GPGPU-Sim tends to have higher performance versus hardware as matrix size increases.

### C. Profiling Tensor Cores

In this section we measure the performance gain obtained when employing tensor cores measured on a real NVIDIA Titan V GPU.

NVIDIA's documentation suggests that tensor cores can provide peak theoretical performance of 125 TFLOPs. The maximum performance we obtained for a GEMM kernel was around 96 TFLOPs. This performance was observed for $8192 \times 8192$ matrix using FP16 mode. To measure the maximum sustainable tensor core throughput we developed a kernel with repeated `wmma.mma` operations (computational intensity on the order of $10^8$). The performance obtained was 109.6 TFLOPs in FP16 mode and 108.7 TFLOPs in mixed-precision mode.

Figure 15 shows the results of profiling the latency of `wmma.load`, `wmma.mma` and `wmma.store` instructions during several iterations of a WMMA kernel. This kernel uses shared memory and performs matrix-multiply accumulate operations on a $1024 \times 1024$ matrix. All three graphs show occasional high latencies. These may result from some combination of warp scheduling policies and high memory traffic. We find the minimum latency of `wmma.load`, `wmma.store` and `wmma.mma` instructions is 125, 120 and 70 clock cycles respectively.

In Figure 16 we plot the median latency to analyze how `wmma.load`, `wmma.mma` and `wmma.store` latency varies with the matrix size for WMMA kernels. The `wmma.load` latency is plotted with a logarithmic axis. Using shared-memory reduces median `wmma.load` latency by more than $100\times$ when operating on a larger matrix.

Figure 17 shows the performance achieved by the tensor cores in different scenarios: In this figure we compare performance of a GEMM kernel implemented with various APIs (CUBLAS, WMMA) with (WITH) or without (WO) tensor cores (TC) using mixed-precision (FP32) or FP16

---

[4]https://github.com/gpgpu-sim/cutlass-gpgpu-sim

mode. In this graph "MAX PERF KERNEL" is our kernel designed to stress tensor core performance in FP16 or mixed-precision (FP32) mode. THEORETICAL LIMIT is the peak performance of 125 TFLOPs. The WMMA GEMM includes optimizations like using shared memory and proper memory layout. The performance gain obtained using the cuBLAS GEMM kernel is more than the WMMA GEMM implementation (both the kernels using tensor cores). cuBLAS is a highly optimized library which has optimizations to avoid shared memory bank conflicts and employs software pipelining. We find tensor cores provide a performance boost of about $3 - 6\times$ times that of SGEMM (Single Precision GEMM) kernel and about $3\times$ that of HGEMM (Half Precision GEMM).

### VI. RELATED WORK

This section briefly discusses related work. Wong et al. [46] performed a thorough analysis of the NVIDIA GT200 using an extensive set of microbenchmarks. They explored architectural details of the processing cores and the memory hierarchies. The describe previously undisclosed details of barrier synchronization and the memory hierarchy including TLB organization in GPUs. Jia et al. [40] explored tensor cores in detail. They decoded sets and steps for Volta tensor cores in mixed-precision mode. In contrast, we comprehensively investigated both modes of operation. We found that sets and steps behave differently in FP16 mode than in mixed precision mode. We uncovered the organization of theadgroups into octets. We determined the mapping of operand matrix elements to threads for the tensor cores in the Turing architecture and found they behaves differently from the Volta tensor cores. We also provide a methodology for uncovering the information presented (including describing our microbenchmarks). Markidis et al. [47] studied the impact of precision loss and programmability aspect of Tensor Cores for HPC application. Khairy, et al. [48] studied the memory system of modern GPUs including Volta and discovered many important design decisions in the memory system. They modeled it in GPGPU-Sim and achieve a very high correlation on a wide range of GPGPU workloads.

### VII. CONCLUSION

In this paper we investigated the design of the tensor core machine learning accelerators integrated into recent GPUs from NVIDIA. We performed a detailed characterization and analysis of the tensor cores implemented in NVIDIA's Volta and Turing architectures. This analysis guided the development of a detailed architectural model. We implemented a model for the Volta tensor cores in GPGPU-Sim and found its performance agreed well with hardware, obtaining a 99.6% IPC correlation versus a Titan V GPU. As part of our efforts we also enabled CUTLASS, NVIDIA's open-source CUDA C++ template library supporting tensor cores, on GPGPU-Sim. We believe that combined the above work

will serve as a promising starting point for further micro-architectural investigation of machine learning workloads.

REFERENCES

[1] M. M. Najafabadi, F. Villanustre, T. M. Khoshgoftaar, N. Seliya, R. Wald, and E. Muharemagic, "Deep learning applications and challenges in big data analytics," *Journal of Big Data*, vol. 2, p. 1, Feb 2015.

[2] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE international conference on*, pp. 6645–6649, IEEE, 2013.

[3] A. Bordes, X. Glorot, J. Weston, and Y. Bengio, "Joint learning of words and meaning representations for open-text semantic parsing," in *Artificial Intelligence and Statistics*, pp. 127–135, 2012.

[4] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems*, pp. 91–99, 2015.

[5] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.

[6] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3156–3164, 2015.

[7] K. Kavukcuoglu, P. Sermanet, Y.-L. Boureau, K. Gregor, M. Mathieu, and Y. L. Cun, "Learning convolutional feature hierarchies for visual recognition," in *Advances in Neural Information Processing Systems*, pp. 1090–1098, 2010.

[8] N. P. Jouppi, C. Young, N. Patil, D. Patterson, G. Agrawal, R. Bajwa, S. Bates, S. Bhatia, N. Boden, A. Borchers, *et al.*, "In-datacenter performance analysis of a tensor processing unit," in *Computer Architecture (ISCA), 2017 ACM/IEEE 44th Annual International Symposium on*, pp. 1–12, IEEE, 2017.

[9] T. Chen, Z. Du, N. Sun, J. Wang, C. Wu, Y. Chen, and O. Temam, "Diannao: A small-footprint high-throughput accelerator for ubiquitous machine-learning," *ACM Sigplan Notices*, vol. 49, no. 4, pp. 269–284, 2014.

[10] Y.-H. Chen, J. Emer, and V. Sze, "Eyeriss: A spatial architecture for energy-efficient dataflow for convolutional neural networks," in *ACM SIGARCH Computer Architecture News*, vol. 44, pp. 367–379, IEEE Press, 2016.

[11] S. Chakradhar, M. Sankaradas, V. Jakkula, and S. Cadambi, "A dynamically configurable coprocessor for convolutional neural networks," *ACM SIGARCH Computer Architecture News*, vol. 38, no. 3, pp. 247–257, 2010.

[12] M. M. Khan, D. R. Lester, L. A. Plana, A. Rast, X. Jin, E. Painkras, and S. B. Furber, "Spinnaker: mapping neural networks onto a massively-parallel chip multiprocessor," in *Neural Networks, 2008. IJCNN 2008.(IEEE World Congress on Computational Intelligence). IEEE International Joint Conference on*, pp. 2849–2856, IEEE, 2008.

[13] C. Farabet, C. Poulet, J. Y. Han, and Y. LeCun, "Cnp: An fpga-based processor for convolutional networks," in *Field Programmable Logic and Applications, 2009. FPL 2009. International Conference on*, pp. 32–37, IEEE, 2009.

[14] D. Kim, J. Kung, S. Chai, S. Yalamanchili, and S. Mukhopadhyay, "Neurocube: A programmable digital neuromorphic architecture with high-density 3d memory," in *Computer Architecture (ISCA), 2016 ACM/IEEE 43rd Annual International Symposium on*, pp. 380–392, IEEE, 2016.

[15] 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*, pp. 161–170, ACM, 2015.

[16] D. Moloney, "Embedded deep neural networks: The cost of everything and the value of nothing ," in *Hot Chips 28 Symposium (HCS), 2016 IEEE*, pp. 1–20, IEEE, 2016.

[17] 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 *The 49th Annual IEEE/ACM International Symposium on Microarchitecture*, p. 20, IEEE Press, 2016.

[18] A. Tilley, "AI Chip Boom: This Stealthy AI Hardware Startup Is Worth Almost A Billion." https://www.forbes.com/sites/aarontilley/2017/08/31/ai-chip-cerebras-systems-investment/, Sep 2017.

[19] MLPerf., "MLPerf v0.5 Results." https://mlperf.org/results/, Dec 2018.

[20] NVIDIA Corporation, "NVIDIA Turing Architecture Whitepaper." https://www.nvidia.com/content/dam/en-zz/Solutions/design-visualization/technologies/turing-architecture/NVIDIA-Turing-Architecture-Whitepaper.pdf, June 2017.

[21] NVIDIA Corporation, "NVIDIA T4 Tensor Core GPUs for Accelerating Inference." https://www.nvidia.com/en-us/data-center/tesla-t4/, Dec 2018.

[22] NVIDIA Corporation, "Tensor Cores in NVIDIA Volta Architecture." https://www.nvidia.com/en-us/data-center/tensorcore/, Sep 2018.

[23] NVIDIA Corporation, "Gordon Bell Award." https://blogs.nvidia.com/blog/2018/09/17/nvidia-volta-tensor-core-gpus-gordon-bell-finalists/, Sep 2018.

[24] A. Bakhoda, G. L. Yuan, W. W. Fung, H. Wong, and T. M. Aamodt, "Analyzing cuda workloads using a detailed gpu simulator," in *Performance Analysis of Systems and Software, 2009. ISPASS 2009. IEEE International Symposium on*, pp. 163–174, IEEE, 2009.

[25] J. Lew, D. Shah, S. Pati, S. Cattell, M. Zhang, A. Sandhupatla, C. Ng, N. Goli, M. D. Sinclair, T. G. Rogers, and T. M. Aamodt, "Analyzing machine learning workloads using a detailed GPU simulator," *CoRR*, vol. abs/1811.08933, 2018.

[26] NVIDIA Corporation, "NVIDIA TESLA V100 GPU ARCHITECTURE." http://images.nvidia.com/content/volta-architecture/pdf/volta-architecture-whitepaper.pdf, June 2017.

[27] J. Choquette, O. Giroux, and D. Foley, "Volta: Performance and programmability," *IEEE Micro*, vol. 38, no. 2, pp. 42–52, 2018.

[28] NVIDIA Corporation, "CUDA C Programming Guide (CUDA 9.0)." https://docs.nvidia.com/cuda/archive/9.0/cuda-c-programming-guide/, Sep 2017.

[29] NVIDIA Corporation, "Programming Tensor Cores in CUDA 9." https://devblogs.nvidia.com/programming-tensor-cores-cuda-9/, Oct 2017.

[30] NVIDIA Corporation, "Parallel Thread Execution ISA Version 6.0." https://docs.nvidia.com/cuda/archive/9.0/parallel-thread-execution/index.html, Sep 2017.

[31] M. Wolfe, "More iteration space tiling," in *Proceedings of the 1989 ACM/IEEE Conference on Supercomputing*, Supercomputing '89, pp. 655–664, 1989.

[32] NVIDIA Corporation, "cuBLAS Developer Guide." https://docs.nvidia.com/cuda/cublas/index.html, Aug 2008.

[33] NVIDIA Corporation, "cuDNN Developer Guide." https://docs.nvidia.com/deeplearning/sdk/cudnn-developer-guide/index.html, Aug 2014.

[34] S. Chetlur, C. Woolley, P. Vandermersch, J. Cohen, J. Tran, B. Catanzaro, and E. Shelhamer, "cudnn: Efficient primitives for deep learning," *CoRR*, vol. abs/1410.0759, 2014.

[35] NVIDIA Corporation, "CUTLASS: Fast Linear Algebra in CUDA C++." https://devblogs.nvidia.com/cutlass-linear-algebra-cuda/, Dec 2017.

[36] NVIDIA Corporation, "CUTLASS: CUDA Templates for Linear Algebra Subroutines." https://github.com/NVIDIA/cutlass, June 2018.

[37] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, *et al.*, "Tensorflow: a system for large-scale machine learning.," in *OSDI*, vol. 16, pp. 265–283, 2016.

[38] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," in *NIPS-W*, 2017.

[39] NVIDIA Corporation, "CUDA C Programming Guide (CUDA 10)." https://docs.nvidia.com/cuda/archive/10.0/cuda-c-programming-guide/index.html, Sep 2018.

[40] Z. Jia, M. Maggioni, B. Staiger, and D. P. Scarpazza, "Dissecting the NVIDIA volta GPU architecture via microbenchmarking," *CoRR*, vol. abs/1804.06826, 2018.

[41] pancake, "Radare2 - A command line framework for reverse engineering binaries." https://rada.re/r/down.html, Feb 2006.

[42] Scott Gray, "SGEMM Implementation." https://github.com/NervanaSystems/maxas/wiki/SGEMM, Apr 2017.

[43] M. Gebhart, S. W. Keckler, and W. J. Dally, "A compile-time managed multi-level register file hierarchy," in *Microarchitecture (MICRO), 2011 44th Annual IEEE/ACM International Symposium on*, pp. 465–476, IEEE, 2011.

[44] NVIDIA Corporation, "Inside Volta: The World's Most Advanced Data Center GPU." https://devblogs.nvidia.com/inside-volta/, May 2017.

[45] C. Rau, "Half-precision floating point library." http://half.sourceforge.net/, Aug 2017.

[46] H. Wong, M.-M. Papadopoulou, M. Sadooghi-Alvandi, and A. Moshovos, "Demystifying gpu microarchitecture through microbenchmarking," in *Performance Analysis of Systems & Software (ISPASS), 2010 IEEE International Symposium on*, pp. 235–246, IEEE, 2010.

[47] S. Markidis, S. W. Der Chien, E. Laure, I. B. Peng, and J. S. Vetter, "Nvidia Tensor Core Programmability, Performance & Precision," in *2018 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pp. 522–531, IEEE, 2018.

[48] M. Khairy, A. Jain, T. M. Aamodt, and T. G. Rogers, "Exploring modern GPU memory system design challenges through accurate modeling," *CoRR*, vol. abs/1810.07269, 2018.