Inside the NVIDIA Ampere Architecture

Ronny Krashinsky, Olivier Giroux GPU Architects GTC 2020



UNPRECEDENTED ACCELERATION AT EVERY SCALE





UNIFIED AI ACCELERATION



BERT Large Inference uses TRT 7.1 for T4/V100, with INT8/FP16 at batch size 256. Pre-production TRT for A100, uses batch size 94 and INT8 with sparsity

ACCELERATING HPC



All results are measured

Except BerkeleyGW, V100 used is single V100 SXM2. A100 used is single A100 SXM4 More apps detail: AMBER based on PME-Cellulose, GROMACS with STMV (h-bond), LAMMPS with Atomic Fluid LJ-2.5, NAMD with v3.0a1 STMV_NVE Chroma with szscl21_24_128, FUN3D with dpw, RTM with Isotropic Radius 4 1024^3, SPECFEM3D with Cartesian four material model BerkeleyGW based on Chi Sum and uses 8xV100 in DGX-1, vs 8xA100 in DGX A100

🕺 NVIDIA



2 NUDIA

A100 SM



Third-generation Tensor Core Faster and more efficient Comprehensive data types Sparsity acceleration

Asynchronous data movement and synchronization

Increased L1/SMEM capacity

1. New Tensor Core

2. Strong Scaling

3. Elastic GPU

4. Productivity



2. Strong Scaling

3. Elastic GPU

4. Productivity

	INPL	JT OPERANDS	AC	CUMULATOR	TOPS	X-factor vs. FFMA
V100	FP32	000000000000000000000000000000000000000	FP32	C	15.7	1x
	FP16		FP32		125	8x

V100 125 8x vs. TOPS FFMA FF16/FP32 Mixedprecision



	INPL	JT OPERANDS	AC	CUMULATOR	TOPS	X-factor vs. FFMA
V100	FP32	000000000000000000000000000000000000000	FP32	000000	15.7	1x
	FP16	000000000000000000000000000000000000000	FP32	000000	125	8x
A100	FP32		FP32	000000000000000000000000000000000000000	19.5	1x
	FP16		FP32	000000	312	16x

V100→A100 2.5x 2x TOPS TOPS/ FF16/FP32 SM Mixedprecision

	INPUT OPERANDS		AC	CUMULATOR	TOPS	X-factor vs. FFMA
V100	FP32		FP32		15.7	1x
	FP16		FP32		125	8 x
A100	FP32		FP32		19.5	1x
	TF32		FP32		156	8x
	FP16		FP32		312	16x
	BF16		FP32		312	16x

TF32 accelerates FP32 in/out data \rightarrow 10x vs. V100 FP32 BFloat16 (BF16) at same rate as FP16

	INPUT OPERANDS		ACCUMULATOR		TOPS	X-factor vs. FFMA		
V100	FP32		FP32		15.7	1x		
V 100	FP16		FP32		125	8x		
	FP32		FP32		19.5	1x		
	TF32		FP32		156	8x		
	FP16	00000000000	FP32		312	16x		
A100	BF16		FP32		312	16x		TODO
	FP16	1 /2	FP16		312	16x	2 2v	IOPS
	INT8		INT32		624	32x		track
	INT4	····· //2	INT32		1248	64x		operand
	BINARY	1/4	INT32		4992	256x	∂ 4X	width

Inference data types

	INPUT OPERANDS		ACCUMULATOR		TOPS	X-factor vs. FFMA	SPARSE TOPS	SPARSE X-factor vs. FFMA
V100	FP32		FP32		15.7	1x	-	-
V100	FP16	00000	FP32		125	8x	-	-
	FP32		FP32		19.5	1x	-	-
	TF32		FP32		156	8x	312	16x
	FP16		FP32		312	16x	624	32x
1100	BF16		FP32		312	16x	624	32x
ATUU	FP16	00000	FP16		312	16x	624	32x
	INT8		INT32	[624	32x	1248	64x
	INT4		INT32		1248	64x	2496	128x
	BINARY	0	INT32		4992	256x	-	-

With Sparsity another 2x, INT8/INT4 reach petaops

	INPL	JT OPERANDS	AC	CUMULATOR	TOPS	X-factor vs. FFMA	SPARSE TOPS	SPARSE X-factor vs. FFMA
V100	FP32		FP32		15.7	1x	-	-
V 100	FP16		FP32		125	8x	-	-
	FP32		FP32		19.5	1x	-	-
	TF32		FP32		156	8 x	312	16x
	FP16		FP32		312	16x	624	32x
	BF16		FP32		312	16x	624	32x
A100	FP16		FP16		312	16x	624	32x
	INT8		INT32		624	32x	1248	64x
	INT4		INT32		1248	64x	2496	128x
	BINARY ^{II}		INT32		4992	256x	<u>V100-</u>	<u>>A100</u>
	IEEE FF	264 (mmm)			19.5	1x	Z.3X f	
							TOL	

	INPL	JT OPERANDS	AC	CUMULATOR	TOPS	X-factor vs. FFMA	SPARSE TOPS	SPARSE X-factor vs. FFMA
V100	FP32	000000000000000000000000000000000000000	FP32		15.7	1x	-	-
V 100	FP16	000000000000000000000000000000000000000	FP32		125	8x	-	-
	FP32	000000000000000000000000000000000000000	FP32		19.5	1x	-	-
	TF32		FP32		156	8x	312	16x
	FP16	00000000000	FP32		312	16x	624	32x
	BF16		FP32		312	16x	624	32x
A100	FP16	00000000000	FP16		312	16x	624	32x
	INT8		INT32		624	32x	1248	64x
	INT4		INT32		1248	64x	2496	128x
	BINARY D		INT32		4992	256x	-	-
	IEEE F	P64 000000000000000000000000000000000000		19.5	1x	-	-	

INSIDE A100 TensorFloat-32 (TF32)



Range of FP32 with precision of FP16



FP32 input/output

FP32 storage and math for all activations, gradients, ... everything outside tensor cores

Out-of-the-box tensor core acceleration for DL

Easy step towards maximizing tensor core performance with mixed-precision (FP16, BF16)

Up to 4x speedup on linear solvers for HPC

 \rightarrow S22082: Mixed-Precision Training of Neural Networks, 5/20 2:45pm PDT

ightarrowS21681: How CUDA Math Libraries can help you unleash the power of the new NVIDIA A100 GPU (recording available) ightarrow

INSIDE A100 SPARSE TENSOR CORE



~No loss in inferencing accuracy

Evaluated across dozens of networks: vision, object detection, segmentation, natural language modeling, translation

→ S22085: Accelerating Sparsity in the NVIDIA Ampere Architecture, 5/20 1:30pm PDT

1. New Tensor Core

2. Strong Scaling

3. Elastic GPU

4. Productivity

DL STRONG SCALING

DL networks: Long chains of sequentiallydependent compute-intensive layers





runs in same time

Strong scaling





Fixed network runs ~2.5x faster

9 🥺 🥺 NIDIA

HOW TO KEEP TENSOR CORES FED?

Math bandwidth (MACs/clock/SM)





A100 STRONG SCALING INNOVATIONS





A100 TENSOR CORE 2x throughput vs. V100, >2x efficiency



V100 TC Instruction (1024 MACs, 8 cycles)





A100 TC Instruction



16x16x16 matrix multiply	FFMA	V100 TC	A100 TC	A100 vs. V100 (improvement)	A100 vs. FFMA (improvement)
Thread sharing	1	8	32	4x	32x
Hardware instructions	128	16	2	8x	64x
Register reads+writes (warp)	512	80	28	2.9x	18x
Cycles	256	32	16	2x	16x

Tensor Cores assume FP16 inputs with FP32 accumulator, V100 Tensor Core instruction uses 4 hardware instructions

22

A100 SM DATA MOVEMENT EFFICIENCY 3x SMEM/L1 bandwidth, 2x in-flight capacity



A100 L2 BANDWIDTH



Split L2 with hierarchical crossbar -2.3x increase in bandwidth over V100, lower latency



A100 DRAM BANDWIDTH



→ S21819: Optimizing Applications for NVIDIA Ampere GPU Architecture, 5/21 10:15am PDT

A100 COMPUTE DATA COMPRESSION





Math

RF

SMEM/L1

L2

DRAM

NVLINK

Up to 4x DRAM+L2 bandwidth and 2x L2 capacity for fine-grained unstructured sparsity



→ S21819: Optimizing Applications for NVIDIA Ampere GPU Architecture, 5/21 10:15am PDT

A100 NVLINK BANDWIDTH

Third Generation NVLink

Math

RF

SMEM/L1

L2

DRAM

NVLINK

50 Gbit/sec per signal pair 12 links, 25 GB/s in/out, 600 GB/s total 2x vs. V100





A100 ACCELERATES CUDA GRAPHS



Grid launches:

- CPU-to-GPU
- GPU grid-to-grid





32-node graphs of empty grids, DGX1-V, DGX-A100

One-shot CPU-to-GPU graph submission and graph reuse Microarchitecture improvements for grid-to-grid latencies

→S21760: CUDA New Features And Beyond, 5/19 10:15am PDT

A100 STRONG SCALING INNOVATIONS Delivering unprecedented levels of performance

A100 improvements over V100





1. New Tensor Core

2. Strong Scaling

3. Elastic GPU

4. Productivity

NVLINK: ONE BIG GPU

- InfiniBand/Ethernet: travels a long distance, consistency is the responsibility of software
- PCI Express: hardware consistency for I/O, not for programming language memory models
- NVLINK: hardware consistency for programming language memory models, like system bus



HGX A100: 3RD GEN NVLINK

HGX A100 4-GPU: fully-connected system with 100GB/s all-to-all BW

HGX A100: 3RD GEN NVLINK & SWITCH

- HGX A100 4-GPU: fully-connected system with 100GB/s all-to-all BW
- New NVSwitch: 6B transistors in TSMC 7FF, 36 ports, 25GB/s each, per direction
- HGX A100 8-GPU: 6x NVSwitch in a fat tree topology, 2.4TB/s full-duplex bandwidth



DGX A100: PCIE4 CONTROL & I/O





35

CLOUD SMALL INSTANCE USAGE

- Small workloads can under-utilize GPU cloud instances, provisioned at whole GPU level
- CSPs can't use MPS for GPU space-sharing, because it doesn't provide enough isolation


NEW: MULTI-INSTANCE GPU (MIG)

- Up to 7 instances total, dynamically reconfigurable
- Compute instances: compute/fault isolation, but share/compete for memory
- **GPU instances:** separate and isolated paths through the entire memory system



ELASTIC GPU COMPUTING

- Each A100 is 1 to 7 GPUs
- Each DGX A100 is 1 to 56 GPUs
- Each GPU can serve a different user, with full memory isolation and QoS



UNIFIED AI ACCELERATION

→S21975: Inside NVIDIA's Multi-Instance GPU Feature (recording available)

 \rightarrow S21884: Under the Hood of the new DGX A100 System Architecture (recording available soon)

 \rightarrow S21702: Introducing NVIDIA DGX A100: The Universal AI System for Enterprise, 5/20 9:00am PDT

1. New Tensor Core

2. Strong Scaling

3. Elastic GPU

4. Productivity

COMPUTE CAPABILITY

Programming Model Development at NVIDIA



GPU PROGRAMMING IN 2020 AND BEYOND Math Libraries | Standard Languages | Directives | CUDA

global

INVIDIA



GPU Accelerated Math Libraries

ightarrowS21766 Inside the NVIDIA HPC SDK: the Compilers, Libraries and Tools for Accelerated Computing, 5/19 1:30PM PST 41

PROGRAMMING MODEL WANTED Software pipelining to hide latency is hard.

```
_device___ void exhibit_A1()
memcpy(/* ... */); //< blocks here</pre>
/* more work */
compute(); //< needed here</pre>
/* more work */
                Data
```

```
_device___ void exhibit_B1()
```

```
compute_head();
__syncthreads(); //< blocks here
/* more work */
```

```
compute_tail(); //< needed here
/* more work */</pre>
```

Compute

PROGRAMMING MODEL WANTED Software pipelining to hide latency is hard.



```
__device__ void exhibit_B2()
{
    compute_head();
    __syncthreads(); //< blocks here
    /* compute_head();
    __syncthreads(); */
    compute_tail(); */< needed here
    /* compute_tail(); */</pre>
```

Compute



CO-DESIGNED: A100 & C++20 barrier Key to asynchronous programming in compute_80

#include <cuda/barrier> // ISO C++20 conforming extension
using barrier = cuda::barrier<cuda::thread_scope_block>;

```
class barrier { // synopsis
  //...
  void arrive_and_wait();
  arrival_token arrive(ptrdiff_t = 1); Nonblocking
  void wait(arrival_token &&) const;
  //...
};
```

ASYNCHRONOUS COPY + BARRIER

Capability	PTX ISA	CUDA C++ API
Asynchronous barrier	<pre>mbarrier.{<basis functions="">}</basis></pre>	cuda::barrier<>
Asynchronous copy	cp.async.ca + cp.async.mbarrier.arrive	<pre>cuda::memcpy_async()</pre>
+Cache-bypass	cp.async.cg	
+Zero-fill ragged edge	cp.async.* … wr-size, rd-size;	CUDA 11 preview library in experimental:: namespace
+User-level tracking	cp.async.mbarrier.arrive.noinc	
+Single-threaded mode	<pre>cp.async.{commit_group, wait_group}</pre>	

ASYNCHRONOUS PROGRAMMING MODEL

#include <cuda/barrier> // ISO C++20 conforming extension
using barrier = cuda::barrier<cuda::thread_scope_block>;



MULTI-BUFFERING PIPELINES IN C++

```
#include <cuda/barrier> // ISO C++20 conforming extension
using barrier = cuda::barrier<cuda::thread scope block>;
global void exhibit C(/* ... */) {
  __shared__ barrier b[2];
                                                                                  Data
  // ^^initialization omitted
  barrier::arrival token t[2];
   da::memcpy_async(/* ... */, b[v]);
 t[0] = b[0].arrive();
 for(int step = 0, next = 1; step < steps; ++step, ++next)</pre>
   if(next < steps) {</pre>
     b[next & 1].wait(t[next & 1]);
     cude::mem.py_async(/* ... */, b[next & 1]);
     t[nekt & 1] = b[next & 1].arrive();
                                                                               Compute
   b[step & 1].uait(t[step & 1]);
   compute();
   t[step & 1] = b[step & 1].arrive();
```

MULTI-BUFFERING PIPELINES IN C++



→S21760: CUDA New Features And Beyond, 5/19 10:15am PDT

OUR PRODUCTIVITY GAINS FROM A100



50



UNPRECEDENTED ACCELERATION AT EVERY SCALE





Whitepaper: NVIDIA A100 Tensor Core GPU Architecture www.nvidia.com/nvidia-ampere-architecture-whitepaper

CS 295: Modern Systems GPU Computing Introduction



Sang-Woo Jun Spring 2019



Graphic Processing – Some History

1990s: Real-time 3D rendering for video games were becoming common
 o Doom, Quake, Descent, ... (Nostalgia!)

□ 3D graphics processing is immensely computation-intensive





Texture mapping

Shading

Warren Moore, "Textures and Samplers in Metal," Metal by Example, 2014 Gray Olsen, "CSE 470 Assignment 3 Part 2 - Gourad/Phong Shading," grayolsen.com, 2018

Graphic Processing – Some History

□ Before 3D accelerators (GPUs) were common

CPUs had to do all graphics computation, while maintaining framerate!

 \circ $\,$ Many tricks were played $\,$



Doom (1993) : "Affine texture mapping"

- Linearly maps textures to screen location, disregarding depth
- Doom levels did not have slanted walls or ramps, to hide this

Graphic Processing – Some History

Before 3D accelerators (GPUs) were common

CPUs had to do all graphics computation, while maintaining framerate!

Many tricks were played



Quake III arena (1999) : "Fast inverse square root" magic!

```
float Q_rsqrt( float number )
{
    const float x2 = number * 0.5F;
    const float threehalfs = 1.5F;
    union {
      float f;
      uint32_t i;
    } conv = {number}; // member 'f' set to value of 'number'.
      conv.i = 0x5f3759df - ( conv.i >> 1 );
      conv.f *= ( threehalfs - ( x2 * conv.f * conv.f ) );
    return conv.f;
}
```

Introduction of 3D Accelerator Cards

- Much of 3D processing is short algorithms repeated on a lot of data
 pixels, polygons, textures, ...
- Dedicated accelerators with simple, massively parallel computation





A Diamond Monster 3D, using the Voodoo chipset (1997) (Konstantin Lanzet, Wikipedia)



Peak Performance vs. CPU



System Architecture Snapshot With a GPU (2019)



High-Performance Graphics Memory

- Modern GPUs even employing 3D-stacked memory via silicon interposer
 - \circ $\,$ Very wide bus, very high bandwidth
 - $\circ~$ e.g., HBM2 in Volta



Graphics Card Hub, "GDDR5 vs GDDR5X vs HBM vs HBM2 vs GDDR6 Memory Comparison," 2019

Massively Parallel Architecture For Massively Parallel Workloads!

- □ NVIDIA CUDA (Compute Uniform Device Architecture) 2007
 - \circ A way to run custom programs on the massively parallel architecture!
- OpenCL specification released 2008
- Both platforms expose synchronous execution of a massive number of threads
 GPU Threads



CUDA Execution Abstraction

Block: Multi-dimensional array of threads

- o 1D, 2D, or 3D
- Threads in a block can synchronize among themselves
- $\circ~$ Threads in a block can access shared memory
- CUDA (Thread, Block) ~= OpenCL (Work item, Work group)
- Grid: Multi-dimensional array of blocks
 - \circ 1D or 2D
 - $\circ~$ Blocks in a grid can run in parallel, or sequentially
- □ Kernel execution issued in grid units
- Limited recursion (depth limit of 24 as of now)

Simple CUDA Example





Simple CUDA Example



More Complex Example: Picture Blurring

□ Slides from NVIDIA/UIUC Accelerated Computing Teaching Kit

Another end-to-end example <u>https://devblogs.nvidia.com/even-easier-introduction-cuda/</u>

Great! Now we know how to use GPUs – Bye?

Matrix Multiplication Performance Engineering



Coleman et. al., "Efficient CUDA," 2017



Volta Execution Architecture

- General Gen
 - Specialization to make use of chip space...?
- Not much on-chip memory per thread
 - \circ 96 KB Shared memory
 - 1024 Registers per FP32 core
- □ Hard limit on compute management
 - 32 blocks AND 2048 threads AND 1024 threads/block
 - e.g., 2 blocks with 1024 threads, or 4 blocks with 512 threads
 - Enough registers/shared memory for all threads must be available (all context is resident during execution)





Resource Balancing Details

- □ How many threads in a block?
- □ Too small: 4x4 window == 16 threads
 - \circ 128 blocks to fill 2048 thread/SM
 - SM only supports 32 blocks -> only 512 threads used
 - SM has only 64 cores... does it matter? Sometimes!
- □ Too large: 32x48 window == 1536 threads
 - Threads do not fit in a block!
- □ Too large: 1024 threads using more than 64 registers
- Limitations vary across platforms (Fermi, Pascal, Volta, ...)

Warp Scheduling Unit

- □ Threads in a block are executed in 32-thread "warp" unit
 - Not part of language specs, just architecture specifics
 - A warp is SIMD Same PC, same instructions executed on every core
- □ What happens when there is a conditional statement?
 - Prefix operations, or control divergence
 - More on this later!
- □ Warps have been 32-threads so far, but may change in the future
Memory Architecture Caveats

□ Shared memory peculiarities

- Small amount (e.g., 96 KB/SM for Volta) shared across all threads
- $\circ~$ Organized into banks to distribute access
- $\circ~$ Bank conflicts can drastically lower performance
- □ Relatively slow global memory
 - Blocking, caching becomes important (again)
 - If not for performance, for power consumption...



8-way bank conflict 1/8 memory bandwidth

Are Mobile DNN Accelerators Accelerating DNNs?

Qingqing Cao*, Alexandru E. Irimiea‡, Mohamed Abdelfattah◊, Aruna Balasubramanian*, Nicholas D. Lane†◊



5th International Workshop on Embedded and Mobile Deep Learning (EMDL 2021)

Deep Learning Applications Exploding on Mobile









Challenges in Running **DNNs** on Mobile Processors

Mobile processors often run on battery-powered devices



Mobile processors often have limited processing power



Are Mobile DNN Accelerators the Solutions?



Vision processing unit (VPU) Tensor processing unitIntelligence processing unit(EdgeTPU)(IPU)

Many customized DNN accelerators, xPUs: designed for better energy efficiency and greater DNN processing power

Background: Architecture of Mobile CPUs

CPU

Control	ALU	ALU
Control	ALU	ALU
С	ache	

• Single instruction multiple data (SIMD) capabilities to parallelize compute intensive operations

 Designed for more general tasks, complex control logic and lower compute density

Background: Architecture of Mobile GPUs

GPU



- Stronger SIMD capabilities (many GPU shader cores/more ALUs), high compute density
- Mobile GPUs often share memory with CPUs, memory accesses are energyconsuming

Background: Architecture of DSPs



DSP

• Enhance SIMD by introducing vector execution unit in addition to ALUs

 Integer based operations, hardwareassisted multithreading, more energy efficient

Background: Differences between CPU, GPU and DSP

CPU: general purpose, SIMD instructions for parallel processing

GPU: originally designed for graphics processing, massive parallelism (more SIMD processing units), power hungry

DSP: originally for photo/video/audio processing, integer-only vector operations, energy efficient

Goal of Mobile DNN Accelerator Empirical Study

Understanding the DNN performance on mobile accelerators in comparison to conventional processors like CPU, GPU, and DSP architectures

Mobile DNN Accelerators: The architecture of NCS VPU

VPU of Neural Compute Stick (NCS)

Interfaces (SPI, US	B, Ethernet, etc)
Low Power Me	mory Fabric
Optimised Configurable	Vision HW Engines
<u>.</u>	RISC-RT Schedule
12 Shave Processors (128-bit vector VLIW)	POSIX RTOS

- 128-bit vector processing cores (SHAVE processors), SIMD support controlled via VLIW instructions
- NCS1: 12 SHAVE cores, 2MB on-chip memory (SRAM)
- NCS2: 16 SHAVE cores, 2.5 MB SRAM, additional neural compute engine dedicated for DNN workloads with higher efficiency

Why NCS Matters? Relationships to DNN Accelerators

CPUs and GPUs (SIMD/SIMT)



Memory access is the bottleneck for CPUs and GPUs



Example: AlexNet has **724M** MACs , but **2896M** DRAM accesses required!

Why NCS Matters? Relationships to DNN Accelerators



Partial sum accumulation does NOT have to access DRAM

Accelerators (Dataflow Processing)



Credit: Eyeriss [ISCA 2016]

DNN Accelerators Should Run Faster and Be More Energy-Efficient

Device	Maximum On-Chip GOPs/s	Ext. Memory (kB)	Maximum BW (GB/s)	Power (W)
CPU	17.6/core	1536		
GPU	519	1024	27.8	5
DSP	128	512		
NCS1	1000	2048	7	1
NCS2	4000	2560	11.9	1
			1	

Study Setup

Hardware:

- Android Smartphone (OnePlus 3),
- NCS1 and NCS2,
- Jetson TX2 board (w/ CUDA GPU)

Software:

- OpenVINO SDK for NCS
- TensorFlow Lite for CPU
- Snapdragon Neural Processing Engine(SNPE) for GPU and DSP

Power measurements:

- Monsoon power monitor for OnePlus 3
- MakerHawk UM34C for NCS on PC
- On-board power module for TX2



Study Methodology

DNN Workloads



- → 4 Popular CNNs: SqueezeNet, MobileNetV2, ResNet50, InceptionV3
- → 4 NLP Transformer (BERT) models: BERT-tiny, BERT-mini, BERT-small, BERTmedium

Measurement Setup



- Measure inference latency and energy and average across 10 runs
- Scale the voltage and current measurements to make them comparable

GPU and DSP Are Faster than NCS1 for CNN Workloads



NCS1 is slower than GPU and DSP, likely because the SHAVE cores are not optimized for CNNs.

Due to dedicated neural compute engine, NCS2 consistently beats NCS1 as expected, and is often faster than GPUs.

DSP is still the fastest accelerator for all studied CNN workloads.

CNN: GPU and DSP Are More Energy-Efficient than NCS



The NCS1/2 consume more energy than the GPU, but much less than the CPU in all cases.

DSP consumes least energy for all studied CNN workloads.

Roofline Analysis of CNNs



NCS devices are memory bound (except SqueezeNet)

- Existing models cannot fit in memory (even with better memory architecture)
- Parallelization does not help because the workload is memory bound

Can changing models help: NCS2 "Vectorization" Effects



Reason: CNN filter size exceeds the vector instruction parallelism in the NCS2

The effects happen much less on other devices such as CPU/GPU (hence our focus on NCS2)

NCS-Aware CNN Optimizations

We decreased the output filter depths in InceptionV3 to the closest power of two and then retrained the model

		Latency	Power	Perf/Watt
Original	NCS2	37 ms	633 mW	43 fps/W
	DSP	36 ms	587 mW	47 fps/W
Pruned	NCS2	20 ms	592 mW	84 fps/W
	DSP	23 ms	478 mW	91 fps/W

NCS2 becomes 15% faster than DSP after pruning, even though it was 3% before pruning.

NCS Devices Are Slower than CPU/GPU for NLP Models



NCS1 and NCS2 are 5 ~ 10x slower than the CPU and are 20 ~ 60 slower than the GPU

Hypothesis: NCS accelerators are not optimized for NLP models like BERT;

BERT models use attention mechanism which has different properties compared to CNNs



Performance Characterization of DNNs on commodity mobile DNN accelerators -- NCS

Development of Preliminary Accelerator Specific Neural Network Inference Optimizations

Major Takeaway: Mobile DNN Accelerators are not ready yet to accelerator DNNs, and are just like conventional mobile processors, only as good as the software and algorithms that drive them.



NVIDIA CUDA ARCHITECTURE

João Paulo Navarro - Solutions Architect jpnavarro@nvidia.com

• NVIDIA history

- Why GPU and accelerated computing
- GPU architecture
- GPU data-center architecture



FORCES SHAPING COMPUTING



BEYOND MOORE'S LAW

FORCES SHAPING COMPUTING



BEYOND MOORE'S LAW - 1000X EVERY 10 YEARS

ACCELERATED COMPUTING

FORCES SHAPING COMPUTING



BEYOND MOORE'S LAW - 1000X EVERY 10 YEARS

ACCELERATED COMPUTING

COMPUTERS WRITING SOFTWARE

25 YEARS OF ACCELERATED COMPUTING



25 YEARS OF ACCELERATED COMPUTING



CHALLENGES: ACCELERATING BIG AND SMALL

AI Advances Demand Exponentially Higher Compute



3000X Higher Compute Required to Train Largest Models Since Volta

AI Applications Demand Distributed Pervasive Acceleration



TODAY'S HYPERCONVERGED DATA CENTER

Impossible to Optimally Design Server Mix for Unpredictable Demand



REIMAGINING THE GPU

Three Breakthroughs to Fuel the Next Era of Modern Accelerated Data Centers





A GIANT LEAP IN PERFORMANCE



UNIFIED AI TRAINING AND INFERENCE ACCELERATION

1-50

SCALABILITY FOR THE ELASTIC DATACENTER

HIGH-PERFORMANCE COMPUTING WITH NVIDIA

RISE OF GPU COMPUTING



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp
HOW GPU ACCELERATION WORKS



BEYOND MOORE'S LAW

Progress Of Stack In 6 Years

2013

 cuBLAS: 5.0

 cuFFT: 5.0

 cuRAND: 5.0

 cuSPARSE: 5.0

 NPP: 5.0

 Thrust: 1.5.3

 CUDA: 5.0

 Resource Mgr: r304

 Base OS: CentOS 6.2



Accelerated Server With Fermi



Measured performance of Amber, CHROMA, GTC, LAMMPS, MILC, NAMD, Quantum Espresso, SPECFEM3D 2019





Accelerated Server with Volta

NVIDIA DATA CENTER PLATFORM

Single Platform Drives Utilization and Productivity



MOST ADOPTED PLATFORM FOR ACCELERATING HPC



600+ Applications Accelerated



125 Systems on Top 500



World's #1 Summit: 149 PF World's #2 Sierra: 95 PF Europe's #1 Piz Daint: 21 PF Japan's #1 ABCI: 20 PF Industrial #1 Total Pangea 3: 18 PF

ALL TOP 15 APPLICATIONS ACCELERATED

NEW HIGHS IN TOP 500 LIST

NVIDIA POWERS WORLD'S FASTEST SUPERCOMPUTER

Summit Becomes First System To Scale The 100 Petaflops Milestone





27,648 Volta Tensor Core GPUs

NVIDIA POWERS FASTEST SUPERCOMPUTERS IN US, EUROPE, JAPAN, INDUSTRY

17 of World's 20 Most Energy-efficient Supercomputers



ORNL Summit World's Fastest 27,648 GPUs| 122 PF LLNL Sierra US 2nd Fastest 17,280 GPUs| 72 PF

ABCI Japan's Fastest 4,352 GPUs | 20 PF Piz Daint Europe's Fastest 5,320 GPUs | 20 PF ENI HPC4 Fastest Industrial 3,200 GPUs| 12 PF

WEAK NODES

Lots of Nodes Interconnected with Vast Network Overhead

STRONG NODES

Few Lightning-Fast Nodes with Performance of Hundreds of Weak Nodes



GPU ARCHITECTURE

SINGLE AND DOUBLE PRECISION THROUGHPUT

Floating-Point Operations per Second for the CPU and GPU



SINGLE AND DOUBLE PRECISION THROUGHPUT

Memory Bandwidth for the CPU and GPU



Why does GPU have memory?

CPU VS GPU

The GPU Devotes More Transistors to Data Processing



20-Series Architecture (Fermi)



512 Scalar Processor (SP) cores execute parallel thread instructions



16 Streaming Multiprocessors (SMs) each contains [•] 32 scalar processors [•] 32 fp32 / int32 ops / clock, [•] 16 fp64 ops / clock

4 Special Function Units (SFUs)

Shared register file (128KB)

48 KB / 16 KB Shared memory
16KB / 48 KB L1 data cache



		Instr	ache							
	Schedu	ler		Scheduler						
	Dispate	ch		Dispatch						
		1								
Core	Core	Core	Core	L/S L/S	SFU					
Core	Core	Core	Core							
Core	Core	Core	Core		SFL					
Core	Core	Core Core	Core	L/S L/S L/S	SFU					
Core	Core	Core	Core		SFU					
Core	Core L1	Core	Core	LIS	Iure					
	Share	d Memor	у	Unit	Text					

Kepler cc 3.5 SM (GK110)



"SMX" (enhanced SM) 192 SP units ("cores") 64 DP units LD/ST units, 64K registers 4 warp schedulers Each warp scheduler is dualissue capable K20: 13 SMX's, 5GB K20X: 14 SMX's, 6GB K40: 15 SMX's, 12GB

SMX	SMX																		
	Warp Scheduler Warp Scheduler Warp Scheduler Warp Scheduler																		
Di	spate	h act	Dispat	tch	D	Dispatch Dispatch					Dispatch Dispatch					Dispatch Dispatch			
	+		+		+ +					+ +									
							Regi	ster	File (6	65,53	6 x 3	2-bit))						
Core	Core	Com	DP Unit	Core	Core	Core	DP Unit	LD/ST	SEU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LDVST	SEU
-			-																
Core	Соге	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	OP Unit	Core	Core	Core	OP Unit	LD/ST	SFU
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LDIST	SFU
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LDIST	SFU
Gore	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Gore	Core	Core	DP Unit	LDIST	SFU
Core	Core	Core	EP Link	Core	Core	Core	DP Linit	LDST	SELL	Core	Core	Com	DP Linit	Core	Core	Corp	DP Linit	LDIST	SEL
	Cons	Cont		Core	Core					COTE	CORE	COIC		COID	CUIE	COID			
Core	Соге	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	OP Unit	Core	Core	Core	DP Unit	LD/ST	SFU
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU
Core	Core	Core	BP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	OP Unit	LDIST	SFU
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	OP Unit	Core	Core	Core	DP Unit	LDIST	SFU
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	OP Unit	Core	Core	Core	DP Unit	LD/ST	SFU
Core	Core	Gore	DP Unit	Core	Core	Core	DP Unit	LD/ST	SEU	Core	Core	Core	OP Unit	Core	Core	Core	DP Unit	LDIST	SEU
									OFU									10000	0.00
Core	Core	Core	OP Unit	Core	Core	Core	DP Unit	LD/ST	SPU	Core	Core	Core	UP Unit	Core	Core	Core	DP Unit	LDIST	SFU
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU
							CAND	Inter	conne	ct Nel	twork	0							
							04 KB	Snar	ad Me	mor	y / L1	Gac	ne						
		_			_	-	48 K	B Re	ad-Or	ny D	ata C	ache	1000				-		
-	Tex		Tex			Tex		Te>	{		Tex		Tex			Tex		Tex	
	Tex		Tex			Tex		Te	4		Tex		Tex			Tex		Tex	

Maxwell/Pascal cc5.2, cc6.1 SM



"SMM" (enhanced SM) 128 SP units ("cores") 4 DP units LD/ST units cc 6.1: INT8 4 warp schedulers Each warp scheduler is dualissue capable M40: 24 SMM's, 12/24GB P40: 30 SM's, 24GB P4: 20 SM's, 8GB

	Verte	x Felch			roly Morph Tesse	ittingine 3.1 elletor	0	v	lewport T	ransform			
			Attrib	ute Setup			Stream O	hitput					
				_	Instruction	an Cache	_	_	_				
-	- 1	nstructi	on Butte	e.		Instruction Buffer							
p	kpatch GN	1		Reparton Ch		spatist Give	t and	2	Inpatch Ur	it .			
	Registe	er File (1	16,384 x	32-bit)		Regist	er File (1	16,384 x	32-bit)				
Core	Cáre	Core	Core	LOIST	SFU	Gore	Core	Core	Core	LDIST	50		
Core	Core	Core	Core	LDIST	sru	Com	Core	Core	Corn	LDST	56		
Core	Gure	Gure	Com	LDIST	SFU	Gore	Core	Core	Core	LOIST	SE		
Core	Core	Core	Core	LOIST	BFU	Core	Core	Core	Corp	LOIST	SF		
Core	Cere	Core	Core	LEVET	sru	Core	Con	Core	Core	LDIST	se		
Core	Core	Gere	Core	LOST	aru.	Gore	Gore	Gore	Core	LOIST	56		
Gore	Core	Gare	Core	LDIST	BFU	Core	Core	Care	Core	LEAST	SF		
Core	Core	Gire	Gire	LOIST	sru	Core	Gere	Gura	Core	LOIST	SF		
_													
	Tex		1	Tex	+ COURS I	L1 Game	Тех			Тяк			
k	Tex	netructé	an Buffr	Tex			Tex	nstructi	on Bulle	Tex.			
	Tex I	nutrucii Wary B	an Duffr	Tex If			Tex	nstructi Warp Se	on Buffe	Tex of			
D	Tex Inpetch Un	notrució Warp 19 a	an Duffe Declara	Tex It Source D It It			Tex Inputch Dr Peoplet	natructi Warp Se 4	on Buffe cheduler	Tex Mapeter S	vet		
C	Tex Inpetch Un Regist	netructi Ware F er File (an Duffr 16,384 x	Tex R Search D			Tex Inset: the Regist	nstructi Worp Se 4 er File (on Duffe hersder 16,384 x	Tex Maputan S 32-bit)	vet		
Cere	Tex Inputth Un Regist Corv	netructi Warp I a er File (Cure	an Duffe Predent 16,384 x Core	Tex F Sector D 32-bit) LDIST	an a	Core	Tex Regist	nstruct) Worp Si It Eore	on Duff Instaler 16,384 x Core	Tex Counter U 32-bit) LDIST	ert Her		
Core Core	Tex Inpatch Unit Regist Core	netructi Worp 1 er File (Core Core	an Dulle 2600er 16,384 x Core Córe	Tex Reaction D 32-bit() LDIST	SPU SPU	Curs Curs	Tex Instant Dr Rogist Care Care	Nern 3 Warp 3 L Cere Cere	on Bulfe heifuler 16,384 x Care Core	Tex F Sustant U 32-bit) LDST	nt SF		
Care Core Core	Tex Regist Core Core Core	nstructi Warp So ar File (Core Core	an Duffe December 16,384 x Core Core Core	T+x F Search U 32-bit) LDST LDST	sru sru sru	Curre Curre Corre	Tex Nearth Dir Rogist Care Care Care	nstruct) Word S H Core Core Core	on Duff Definition 16,384 x Core Core	Tex Second to 32-bit) Lost Lost Lost	rit SF SF		
Care Core Core Core	Tex Regist Core Core Core	er File (Core Core	an Bulli Andular 16,384 x Core Core Core	Tex	SPU SPU SPU SPU	Core Core Core Core	Tex Regist Care Care Care	nstructi Warp 3 I Core Core Core	an Duff heifuiar 16,384 x Care Core Core	Tex 32-bit) LDIST LDIST LDIST	rt Sf Sf		
Core Core Core Core Core	Tex Regist Core Core Core Core	netructi Warp I er File (Core Core Core Core	an Duffe Decement 16, 354 x Core Core Core Core	Tex R Source U 32-bit) LDST LDST LDST LDST	SPU SPU SPU SPU SPU	Curs Curs Cors Cors Cors	Tex Regist Core Core Core Core	outnict) Worg 3 4 Core Core Core Core Core	on Built Induin 16,384 × Core Core Core Core	Tex Sector D 32-bit) LDIST LDIST LDIST	SF SF SF		
Core Core Core Core Core Core	Tex Regist Carv Carv Care Care Care Care	er File (Core Core Core Core Core	an Duffe Heldunt 16,354 x Care Care Care Care Care Care	Tex Terret D	SPU SPU SPU SPU SPU SPU	Core Core Core Core Core	Tex Regist Core Core Core Core Core	retruction Worm of Core Core Core Core Core	on Bulli Definition 16,384 × Core Core Core Core Core	Tex 32-bit) LDST LDST LDST LDST LDST	rt SF SF SF		
Corre Corre Corre Corre Corre Corre	Tex Regist Core Core Core Core Core	Core Core Core Core Core Core Core	Core Core Core Core Core Core Core	Tex Control D Control D Control Con	SPU SPU SPU SPU SPU SPU SPU	Curs Curs Curs Curs Curs Curs Curs Curs	Tex Regist Care Care Care Care Care Care Care Care	er File (Core Core Core Core Core Core	en Dudfo Trestater Core Core Core Core Core Core	Tes 32-bit) Lost Lost Lost Lost Lost Lost	sr sr sr sr sr		
Corre Corre Corre Corre Corre Corre Corre	Tex Registri Uri Cars Cars Cars Cars Cars Cars Cars Cars	Anterucia Vicey Ba en File (Corre Corre Corre Corre Corre Corre Corre	an Buff Notes To 354 x Core Core Core Core Core Core	Tex Sector 32-bili LDest LDest LDEST LDEST LDEST	8FU 8FU 8FU 8FU 8FU 8FU 8FU 8FU 8FU	Curs Curs Cars Cars Cars Cars Cars Cars Cars Ca	Tex Regist Care Care Care Care Care Care Care Care	er File (Gore Gore Gore Core Core Core Core Core	an Buth Declar Core Core Core Core Core Core Core Cor	Tex a Support Source Loss Loss Loss Loss Loss Loss Loss Loss Loss Loss	sr sr sr sr sr sr		
Cors Cors Cors Cors Cors Cors Cors	Tex Regist Cars Cars Cars Cars Cars Cars Cars Cars	Antruche Vorge 5 a corre Corre Corre Corre Corre Corre Corre Corre	an Budh Stadam Care Care Care Care Care Care Care Care	Tex r Sector D S2-bil) LDST LDST LDST LDST LDST	SPU SPU SPU SPU SPU SPU SPU SPU SPU	Core Core Core Core Core Core Core Core	Tex Pediation Regist Care Care Care Care Care Care Care Care	nstruct) Word 33 er File (Care Care Care Care Care Care Care Care	in Buffe Defaue Core Core Core Core Core Core Core Cor	Tes s s s s s s s s s s s s s	SF SF SF SF SF SF		
Corre Corre Corre Corre Corre Corre Corre	Tex Regist Care Care Care Care Care Care Care Care	Virage Da Virage	an Builli Malant 16,384 x Care Care Care Care Care Care Care Care	Tex ar Seventin 10 Seventin	SPU SPU SPU SPU SPU SPU SPU SPU SPU	Curra Curra Curra Curra Curra Curra Curra Curra Curra Curra Curra Curra Curra	Tex Regist Care Care Care Care Care Care Care Care	netruction Verage Sale Corres Corres Corres Corres Corres Corres Corres	an Dudie Charlier Core Core Core Core Core Core Core Co	Tex Sector 1 32-bit LOST LOST LOST LOST LOST LOST LOST LOST	sr sr sr sr sr sr		

Pascal/Volta cc6.0/7.0

64 SP units ("cores") 32 DP units LD/ST units FP16 @ 2x SP rate cc7.0: TensorCore 4 warp schedulers Each warp scheduler is dualissue capable P100: 50 SM's, 16GB V100: 80 SM's, 16/32GB

М														
_	_	_	_	_		L1 Instru	ction	Cache	_	_	_	_		
		L0 I	nstruct	ion C	ache	_				L0 Ir	istruc	tion C	ache	
	War	p Scl	heduler	(32 th	hread/clk)	1			War	p Sch	edule Linit	r (32 t	hread/clk)	
Register File (16,384 x 32-bit)									Reg	ister	File (16,384	4 x 32-bit)	8
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
FP64	INT	INT	FP32	FP32	TENCOR	TENCOR		FP64	INT	INT	FP32	FP32	TENSO	TENCOR
FP64	INT	INT	FP32	FP32	CORE	CORE		FP64	INT	INT	FP32	FP32	CORE	CORE
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
LD/ LD/	LD/	LD/	LD/ ST	LD/	LD/ LD/	SFU			LD/	LD/	LD/	LD/ ST		SFU
51 51	al	91	91	aı	31 31		╎└╴	31 31	al	а	91	31	al al	
		L0 I	nstructi	ion C	ache					LO Ir	nstruc	tion C	ache	
	War	p Sci	heduler	(32 ti	hread/cik)		Warp Scheduler (32 thread/clk)							
	Di	spate	n Unit (32 m	read/cik)				U	spatci		(32 m	read/cik)	
	Reg	ister	File (1	6,384	4 x 32-bit)				Reg	ister	File (16,384	4 x 32-bit)	
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
FP64	INT	INT	FP32	FP32	TENSOR	TENSOR		FP64	INT	INT	FP32	FP32	TENSOR	TENSOR
FP64	INT	INT	FP32	FP32	CORE	CORE		FP64	INT	INT	FP32	FP32	CORE	CORE
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
FP64	INT	INT	FP32	FP32				FP64	INT	INT	FP32	FP32		
LD/ LD/ ST ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ LD/ ST ST	SFU		LD/ LD/ ST ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ LD/ ST ST	SFU
					128KE	L1 Data Ca	che /	Shared M	emory	,				
	-			-	Terr		1	à	Toy		-	-		

INTRODUCING TESLA V100



The Fastest and Most Productive GPU for Deep Learning and HPC



TESLA V100

21B transistors 815 mm²

80 SM 5120 CUDA Cores 640 Tensor Cores

16 GB HBM2 900 GB/s HBM2 300 GB/s NVLink



*full GV100 chip contains 84 SMs

VOLTA GV100 SM

	GV100
FP32 units	64
FP64 units	32
INT32 units	64
Tensor Cores	8
Register File	256 KB
Unified L1/Shared memory	128 KB
Active Threads	2048

M							L1 Instru	ctio	on Cache								
	10/-	L0 I	nstruc	tion C	ache						L0 Ir	nstruc	tion C	ache			
	vva Di	rp Scr ispatc	h Unit	(32 th	nread/ read/c	cik) ik)				Di	rp Sch spatcl	h Unit	(32 th	nread read/o	/cik) :lk)		
Register File (16,384 x 32-bit)									Register File (16,384 x 32-bit)								
FP64	INT	INT	FP32	FP32					FP64	INT	INT	FP32	FP32	-			
FP64	INT	INT	FP32	FP32					FP64	INT	INT	FP32	FP32				
FP64	INT	INT	FP32	FP32					FP64	INT	INT	FP32	FP32	Ħ			
FP64	INT	INT	FP32	FP32	TEN	SOR	TENSOR		FP64	INT	INT	FP32	FP32	TEN	SOR	TENSOR	
FP64	INT	INT	FP32	FP32	со	RE	CORE		FP64	INT	INT	FP32	FP32	cc	RE	CORE	
FP64	INT	INT	FP32	FP32					FP64	INT	INT	FP32	FP32				
FP64	INT	INT	FP32	FP32					FP64	INT	INT	FP32	FP32				
FP64	INT	INT	FP32	FP32					FP64	INT	INT	FP32	FP32	H			
LD/ LD ST ST	/ LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	SFU		LD/ LD ST S1	/ LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	SFU	
								┛┖									
		101	astruc	tion C	ache			٦٢			1.0.1	astruc	tion C	ache			
	Wa	L0 lı rp Scł	nstruc 1edule	tion C r (32 t	ache hread/	clk)				Wai	L0 Ir rp Sch	nstruc Iedule	tion C r (32 t	ache hread	/clk)		
	Wa Di	L0 II rp Sch spatc	nstruc nedule h Unit	tion C r (32 t (32 th	ache hread/ read/c	cik) ik)				Wa Di	L0 Ir rp Sch spatcl	nstruc Iedule h Unit	tion C r (32 t (32 th	ache hread read/o	/clk) :lk)		
	Wa Di Reç	L0 li rp Scl ispatc gister	nstruc hedule h Unit File ('	tion C r (32 t (32 th 16,384	ache hread/ read/c 4 x 32·	clk) lk) -bit)				War Di Reg	L0 Ir rp Sch spatcl jister	nstruc Iedule h Unit File ('	tion C r (32 t (32 th 16,384	ache hread read/c 4 x 32	/clk) :lk) !-bit)		
FP64	Wa Di Reç INT	L0 II rp Sch ispatc gister INT	nstruc nedule h Unit File (' FP32	tion C r (32 t (32 th 16,384 FP32	ache hread/ read/c 4 x 32	clk) lk) -bit)			FP64	War Di Reg	L0 Ir rp Sch spatcl jister INT	nstruc Iedule h Unit File (1 FP32	tion C r (32 t (32 th 16,384 FP32	ache hread read/d 4 x 32	/clk) :lk) !-bit)		
FP64 FP64	Wa Di Reg INT	L0 II rp Sch ispatc gister INT INT	nstruc hedule h Unit File (' FP32 FP32	tion C r (32 t (32 th 16,384 FP32 FP32	ache hread/ read/c 4 x 32	clk) Ik) -bit)			FP64	War Di Reg INT	L0 Ir rp Sch spatcl jister INT INT	nstruc nedule h Unit File (* FP32 FP32	tion C r (32 t (32 th 16,384 FP32 FP32	ache hread read/o 4 x 32	/clk) :lk) !-bit)		
FP64 FP64 FP64	Wa Di Reg INT INT	LO II rp Sch spatc gister INT INT INT	nstruc hedule h Unit File (' FP32 FP32 FP32	tion C r (32 th (32 th 16,384 FP32 FP32 FP32	ache hread/ read/c 4 x 32	clk) lk) -bit)			FP64 FP64 FP64	Wan Di Reg INT INT INT	L0 Ir rp Sch spatcl jister INT INT	nstruc nedule h Unit File (* FP32 FP32 FP32	tion C r (32 th (32 th 16,384 FP32 FP32 FP32	ache hread/ read/ 4 x 32	/clk) :lk) !-bit)		
FP64 FP64 FP64 FP64	Wa Di Reg INT INT INT	LO II rp Sch spatc jister INT INT INT INT	nstruc nedule h Unit File (* FP32 FP32 FP32 FP32	tion C r (32 t (32 th 16,384 FP32 FP32 FP32 FP32	ache hread/ read/c 4 x 32 4 x 32	clk) lk) -bit) SOR	TENSOR		FP64 FP64 FP64 FP64	Wai Di Reg INT INT INT INT	L0 Ir rp Sch spatcl jister INT INT INT	nstruc nedule h Unit File (* FP32 FP32 FP32	tion C r (32 ti (32 th 16,384 FP32 FP32 FP32 FP32	ache hread read/d 4 x 32	/clk) :lk) !-bit) SOR	TENSOR	
FP64 FP64 FP64 FP64 FP64	Wa Di Reç INT INT INT INT INT	L0 II rp Sch spatc gister INT INT INT INT	nstruc nedule h Unit File (* FP32 FP32 FP32 FP32 FP32	tion C r (32 t (32 th 16,384 FP32 FP32 FP32 FP32 FP32	ache hread/ read/c 4 x 32 4 x 32 TEN CO	clk) ik) -bit) SOR RE	TENSOR		FP64 FP64 FP64 FP64 FP64	Wat Di Reg INT INT INT INT INT	L0 Ir rp Sch spatcl jister INT INT INT INT	nstruc nedule h Unit File (* FP32 FP32 FP32 FP32 FP32	tion C r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32	ache hread read/o 4 x 32	/clk) :lk) ?-bit) SOR PRE	TENSOR	
FP64 FP64 FP64 FP64 FP64 FP64	Wa Di Reg INT INT INT INT INT INT	LOIN rp Sch spatci sispatci ister INT INT INT INT INT	File (FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	tion C r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32	ache hread/c 4 x 32- TENS CO	clk) lk) -bit) SOR RE	TENSOR		FP64 FP64 FP64 FP64 FP64 FP64	Wan Di Reg INT INT INT INT INT INT	L0 Irr p Schur spatcl sister INT INT INT INT INT	nstruc nedule h Unit File (' FP32 FP32 FP32 FP32 FP32	tion C r (32 t (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32	ache hread/(4 x 32 TEN CC	/clk) :lk) !-bit) SOR PRE	TENSOR	
FP64 FP64 FP64 FP64 FP64 FP64 FP64	Wa Di Reg INT INT INT INT INT INT	L0 In rp Scl spatc ister INT INT INT INT INT INT INT INT	nstruc nedule h Unit FIIe (' FP32 FP32 FP32 FP32 FP32 FP32	tion C r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	ache hread/c 4 x 32: TEN: CO	cik) ik) -bit) SOR RE	TENSOR		FP64 FP64 FP64 FP64 FP64 FP64 FP64	Wan Di Reg INT INT INT INT INT INT INT	LO Irr p Sch spatcl jister INT INT INT INT INT INT	restruction of the second seco	tion C r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	ache hread/ 4 x 32	/clk) slk) t-bit) SOR RE	TENSOR	
FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	Wa Di Reg INT INT INT INT INT INT INT INT	L0 II rp Sch ispatc jister INT INT INT INT INT INT INT	netruc nedule h Unit File (' FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	tion C r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	ache hread/c 4 x 32: TEN: CO	clk) lk) -bit) SOR RE	TENSOR		FP64 FP64 FP64 FP64 FP64 FP64 FP64	Waa Di Reg INT INT INT INT INT INT INT	LO Ir rp Sch spatcl jister INT INT INT INT INT INT INT	restruc edule h Unit File (' FP32 FP32 FP32 FP32 FP32 FP32 FP32	tion C r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	ache hread/ 4 x 32 TEN CC	/clk) blk) t-bit) SOR RE	TENSOR	
FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	Wa DI Reg INT INT INT INT INT INT INT INT INT INT	LOII rp Sch spatc spatc ister INT INT INT INT INT INT INT INT	nstruc nedule h Unit File (' FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	tion C r (32 t (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	ache hread/ 4 x 32: TEN: CO	cik) ik) -bit) SOR RE	TENSOR CORE		FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	Waand Di Di Regg INT INT INT INT INT INT INT (LD/ ST	LO Ir rp Sch spatcl sister INT INT INT INT INT INT INT INT INT	redule edule h Unit File (* FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	tion C r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	ache hread/ 4 x 32	/clk) llk) l-bit) SOR RE	TENSOR CORE	
FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	Wa Di Reç INT INT INT INT INT INT INT INT INT	L0 II rp Sch spatc: INT INT INT INT INT INT INT INT INT INT	nstruc nedule h Unit File (' FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	tion C r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	ache hread/c 4 x 32: TEN: CO	clk) lk) -bit) SOR RE	TENSOR CORE SFU	che	FP64 FP64 FP64 FP64 FP64 FP64 FP64 FP64	Wala Di Reg INT INT INT INT INT INT INT INT INT INT	LO In rp Schel spatcl ister INT INT INT INT INT INT INT INT	redule h Unit File (' FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	tion C r (32 th (32 th 16,384 FP32 FP32 FP32 FP32 FP32 FP32 FP32 FP32	ache hread 4 x 32 TEN CC	/clk) llk) L-bit) SOR RE LD/ ST	TENSOR CORE	

32 📀 nvidia

VOLTA GV100 SM REDESIGNED FOR PRODUCTIVITY

Twice the schedulers

Simplified Issue Logic

Large, fast L1 cache

Improved SIMT model

Tensor acceleration

=

The easiest SM to program yet

SM						I d In comu	tion Co									
	_		- c tru c	tion C	coho			cne				tion C	eebe	_		
	Wa	n Sch	istruc iedule	r (32 fl	ache bread/clk)	_	Warp Scheduler (32 thread/clk)									
	Di	spatcl	h Unit	(32 th	read/clk)		Warp Scheduler (32 thread/clk) Dispatch Unit (32 thread/clk)									
	Register File (16,384 x 32-bit)															
FP64	INT	INT	FP32	FP32			F	P64	INT	INT	FP32	FP32	+			
FP64	INT	INT	FP32	FP32			F	P64	INT	INT	FP32	FP32				
FP64	INT	INT	FP32	FP32			F	P64	INT	INT	FP32	FP32				
FP64	INT	INT	FP32	FP32	TENSOR	TENSOR	F	P64	INT	INT	FP32	FP32	TENSOR		TENSOR	
FP64	INT	INT	FP32	FP32	CORE	CORE	F	P64	INT	INT	FP32	FP32	co	RE	CORE	
FP64	INT	INT	FP32	FP32			F	P64	INT	INT	FP32	FP32	\pm			
FP64	INT	INT	FP32	FP32			F	P64	INT	INT	FP32	FP32				
FP64	INT	INT	FP32	FP32			F	P64	INT	INT	FP32	FP32				
LD/ LD/ ST ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ LD/ ST ST	SFU	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	LD/ ST	SFU	
		L0 Ir	nstruc	tion C	ache					L0 Ir	nstruc	tion C	ache			
	Wa	r <mark>p Sc</mark> h	nedule	r (32 t	hread/clk)				War	p Sch	edule	r (32 t	hread	'clk)		
	Di	spatcl	h Unit	(32 th	read/clk)				Dis	spatcl	n Unit	(32 th	read/c	lk)		
	Reg	ister	File (16,384	4 x 32-bit)		Register File (16,384 x 32-bit)									
FP64	INT	INT	FP32	FP32			F	P64	INT	INT	FP32	FP32	+			
FP64	INT	INT	FP32	FP32			F	P64	INT	INT	FP32	FP32	\pm			
FP64	INT	INT	FP32	FP32			F	P64	INT	INT	FP32	FP32				
FP64	INT	INT	FP32	FP32	TENSOR	TENSOR	F	P64	INT	INT	FP32	FP32	TEN	SOR	TENSOR	
EDGA	INT	INT	FP32	FP32	CORE	CORE	F	P64	INT	INT	FP32	FP32	co	RE	CORE	
FF04							F		INT	INT	5000	5000				
FP64	INT	INT	FP32	FP32				⁵ 64			FP32	FP32				
FP64 FP64	INT INT	INT INT	FP32 FP32	FP32 FP32			F	964 964	INT	INT	FP32	FP32 FP32				
FP64 FP64 FP64	INT INT INT	INT INT INT	FP32 FP32 FP32	FP32 FP32 FP32			FI	264 264 264	INT	INT INT	FP32 FP32 FP32	FP32 FP32 FP32				
FP64 FP64 FP64 LD/ LD/ ST LD/	INT INT INT LD/ ST	INT INT INT LD/ ST	FP32 FP32 FP32 LD/ ST	FP32 FP32 FP32 LD/ ST	LD/ ST ST	SFU	FI FI LD/ ST	264 264 LD/ ST	INT INT LD/ ST	INT INT LD/ ST	FP32 FP32 FP32 LD/ ST	FP32 FP32 FP32 LD/ ST	LD/ ST	LD/ ST	SFU	
FP64 FP64 FP64 LD/ LD/ ST ST	INT INT INT LD/ ST	INT INT INT LD/ ST	FP32 FP32 FP32 LD/ ST	FP32 FP32 FP32 LD/ ST	LD/ ST ST 128K	SFU B L1 Data Cae	FI LD/ ST che / Sha	264 264 264 264 LD/ ST	INT INT LD/ ST	INT INT LD/ ST	FP32 FP32 FP32 LD/ ST	FP32 FP32 FP32 LD/ ST	LD/ ST	LD/ ST	SFU	

33 📀 nvidia

GPU PERFORMANCE COMPARISON

	P100	V100	Ratio
Training acceleration	10 TOPS	120 TOPS	12x
Inference acceleration	21 TFLOPS	120 TOPS	6x
FP64/FP32	5/10 TFLOPS	7.5/15 TFLOPS	1.5x
HBM2 Bandwidth	720 GB/s	900 GB/s	1.2 x
NVLink Bandwidth	160 GB/s	300 GB/s	1 .9 x
L2 Cache	4 MB	6 MB	1.5x
L1 Caches	1.3 MB	10 MB	7.7x

MEMORY HIERARCHY

NEW HBM2 MEMORY ARCHITECTURE

1.5x Delivered Bandwidth





VOLTA MEMORY SUBSYSTEM

Tesla V100



80 Streaming Multiprocessors 256KB register file (20 MB)

Unified Shared Mem / L1 Cache 128KB, variable split (10MB Total, 14 TB/s), Volta caches L1 writes

6 MB L2 Cache, L2 is write back

16/32 GB HBM2 (900 GB/s)

L1, L2 CACHES Why do GPU have caches?

Caches on GPUs can help with:

"Smoothing" irregular, unaligned access patterns Caching common data accessed by many threads

Faster register spills, local memory

Can help in codes that don't use shared memory

SHARED MEMORY

Scratch-pad memory on each SM

User-managed cache, hardware does not evict data

Data written to SMEM stays there until this the code overwrites the data or threadblock finishes execution

Useful for: Storing frequently-accessed data, to reduce DRAM accesses Communication among threads of a threadblock

Performance benefits compared to DRAM: 20-40x lower latency ~15x higher bandwidth

UNIFIED SHARED MEM / L1 CACHE

Variable split



How to specify the L1 / Smem split:

cudaFuncSetAttribute (MyKernel, cudaFuncAttributePreferredSharedMemoryCarveout, carveout);

The driver usually does a pretty good job at choosing the right split.

To overcome 48 KB per threadblock limitation call: cudaFuncSetAttribute (MyKernel, cudaFuncAttributeMaxDynamicSharedMemorySize, maxsize);

RECAP: PASCAL L1 AND SHARED MEMORY



UNIFYING KEY TECHNOLOGIES



VOLTA L1 AND SHARED MEMORY

Volta Streaming L15 :

Low cache hit latency 4x more bandwidth 5x more capacity

Volta Shared Memory :

Unified storage with L1 Configurable up to 96KB



"SCALABILITY OF CPU AND GPU SOLUTIONS OF THE PRIME ELLIPTIC CURVE DISCRETE LOGARITHM PROBLEM"

Jairo Panetta (ITA), Paulo Souza (ITA), Luiz Laranjeira (UnB), Carlos Teixeira Jr (UnB)



VOLTA WARP EXECUTION MODEL

PASCAL WARP EXECUTION MODEL

if (threadIdx.x < 4) {
 A;
 B;
} else {
 X;
 Y;
}</pre>





PASCAL WARP EXECUTION MODEL



WARP IMPLEMENTATION



32 thread warp

Volta



VOLTA WARP EXECUTION MODEL


VOLTA WARP EXECUTION MODEL



Software synchronization also supported, e.g. locks for doubly-linked list!

VOLTA TENSOR CORE



TENSOR CORE

Mixed Precision Matrix Math 4x4 matrices

B_{0,2} C_{0,3} B_{0,0} **B**_{0,1} **B**_{0.3} A_{0,1} A_{0,2} A_{0,3} C_{0,2} $C_{0,1}$ C_{0,0} A_{0,0} B_{1,1} B_{1,3} C1,1 **C**_{1,0} C_{1,2} A_{1,2} A_{1,3} **B**_{1,0} **B**_{1,2} C_{1,3} **A**_{1,0} **A**_{1,1} C_{2,0} B_{2,1} C_{2,1} B_{2,2} B_{2,3} C_{2,3} A_{2,1} A_{2,2} C_{2,2} **B**_{2,0} A_{2.3} A_{2,0} B_{3,0} B_{3,1} B_{3,2} A_{3,2} **B**_{3.3} A_{3,0} $C_{2,1}$ A_{3.3} $C_{3.0}$ C_{2} A_{3.1} 2 FP16 or FP32 **FP16 FP16** FP16 or FP32

D = AB + C



TENSOR SYNCHRONIZATION

FULL WARP 16X16 MATRIX MATH



VOLTA TENSOR OPERATION



Also supports FP16 accumulator mode for inferencing

50

A GIANT LEAP FOR DEEP LEARNING



Inside the NVIDIA Ampere Architecture



UNPRECEDENTED ACCELERATION AT EVERY SCALE



UNIFIED AI ACCELERATION



BERT Large Inference uses TRT 7.1 for T4/V100, with INT8/FP16 at batch size 256. Pre-production TRT for A100, uses batch size 94 and INT8 with sparsity

ACCELERATING HPC



All results are measured

Except BerkeleyGW, V100 used is single V100 SXM2. A100 used is single A100 SXM4 More apps detail: AMBER based on PME-Cellulose, GROMACS with STMV (h-bond), LAMMPS with Atomic Fluid LJ-2.5, NAMD with v3.0a1 STMV_NVE Chroma with szscl21_24_128, FUN3D with dpw, RTM with Isotropic Radius 4 1024^3, SPECFEM3D with Cartesian four material model BerkeleyGW based on Chi Sum and uses 8xV100 in DGX-1, vs 8xA100 in DGX A100

💿 nvidia.



A100 SM

L1 Instruction Cache							
L0 Instruction Cache							
Warp Scheduler (32 th	read/clk)	Warp Sche	eduler (32 thread/clk)				
Dispatch Unit (32 thr	ead/clk)	Dispatch Unit (32 thread/clk)					
Register File (16,384 x 32-bit)		Register File (16,384 x 32-bit)					
INT32 INT32 FP32 FP32 FP64		INT32 INT32 FP32 FP32	FP64				
INT32 INT32 FP32 FP32 FP64		INT32 INT32 FP32 FP32	FP64				
INT32 INT32 FP32 FP32 FP64	3 rd Gen.	INT32 INT32 FP32 FP32	^{FP64} 3 rd Gen.				
INT32 INT32 FP32 FP32 FP64	TENSOR	INT32 INT32 FP32 FP32	FP64 TENSOR				
INT32 INT32 FP32 FP32 FP64	CORE	INT32 INT32 FP32 FP32 INT32 INT32 FP32 FP32	FP64 CORE				
INT32 INT32 FP32 FP32 FP64		INT32 INT32 FP32 FP32	FP64				
INT32 INT32 FP32 FP32 FP64		INT32 INT32 FP32 FP32	FP64				
LD/ LD/ LD/ LD/ LD/ LD/ ST ST ST ST ST ST ST	LD/ LD/ SFU	LD/ LD/ LD/ LD/ ST ST ST ST	LD/ LD/ LD/ LD/ ST ST ST				
L 0 Instruction Cache							
Warp Scheduler (32 th	read/clk)	Warp Scheduler (32 thread/clk)					
Dispatch Unit (32 thr	ead/clk)	Dispatch Unit (32 thread/clk)					
Register File (16,384 x 32-bit) Register File (16,384 x 32-bit)							
INT32 INT32 FP32 FP32 FP64		INT32 INT32 FP32 FP32	FP64				
INT32 INT32 FP32 FP32 FP64		INT32 INT32 FP32 FP32	FP64				
INT32 INT32 FP32 FP32 FP64	3 rd Gen.	INT32 INT32 FP32 FP32	^{FP64} 3 rd Gen.				
INT32 INT32 FP32 FP32 FP64	TENCOD	INT32 INT32 FP32 FP32	TENCOD				
INT32 INT32 FP32 FP32 FP64	TENSOR	INT32 INT32 FP32 FP32	FP64				
INT32 INT32 FP32 FP32 FP64 INT32 INT32 FP32 FP32 FP64	CORE	INT32 INT32 FP32 FP32 INT32 INT32 FP32 FP32	FP64 FP64 CORE				
INT32 INT32 FP32 FP64 INT32 INT32 FP32 FP32 FP64 INT32 INT32 FP32 FP32 FP64	CORE	INT32 INT32 FP32 FP32 INT32 INT32 FP32 FP32 INT32 INT32 FP32 FP32	FP64 FP64 FP64				
INT32 INT32 FP32 FP32 FP64		INT32 INT32 FP32 FP32 INT32 INT32 FP32 FP32 INT32 INT32 FP32 FP32 INT32 INT32 FP32 FP32	FP64 FP64 FP64 FP64				
INT32 INT32 FP32 FP32 FP32 INT32 INT32 FP32 FP32 FP64 INT32 INT32 INT32 FP32 FP32 INT32 INT32 INT32 FP32 FP32 INT32 INT32 INT32 INT32 INT32 INT32 INT32	CORE	INT32 INT32 FP32 FP32 INT32 INT32 INT32 FP32 INT32 INT32 INT32 FP32 INT32 INT32 INT32 INT32 INT32 INT32	FP64 FP64 FP64 CORE FP64 FP64 FP64 FP64 Lby Lby St St				
INT32 INT32 FP32 FP32 FP32 FP64 INT32 INT32 FP32 FP32 FP64 LD/ LD/ LD/ ST ST ST LD/ LD/ LD/ ST ST ST ST LD/ LD/ LD/ ST ST ST ST	I Data Cac	INT32 INT32 FP32 FP32 INT32 INT32 INT32 FP32 INT32 INT32 INT32 FP32 INT32 INT32 INT32 FP32 INT32 INT32 INT32 INT32 INT32 INT32	FP64 FP64 FP64 FP64 LDY LDY LDY SFU Memory				

Third-generation Tensor Core Faster and more efficient Comprehensive data types Sparsity acceleration

Asynchronous data movement and synchronization

Increased L1/SMEM capacity

INTRODUCING NVIDIA A100 Greatest Generational Leap - 20X Volta

	Peak	Vs Volta
FP32 TRAINING	312 TFLOPS	20X
INT8 INFERENCE	1,248 TOPS	20X
FP64 HPC	19.5 TFLOPS	2.5X
MULTI INSTANCE GPU		7X GPUs



54B XTOR | 826mm2 | TSMC 7N | 40GB Samsung HBM2 | 600 GB/s NVLink

NVIDIA A100 DETAILED SPECS

	Peak Performance			
Transistor Count	54 billion			
Die Size	826 mm ²			
FP64 CUDA Cores	3,456			
FP32 CUDA Cores	6,912			
Tensor Cores	432			
Streaming Multiprocessors	108			
FP64	9.7 teraFLOPS			
FP64 Tensor Core	19.5 teraFLOPS			
FP32	19.5 teraFLOPS			
TF32 Tensor Core	156 teraFLOPS 312 teraFLOPS*			
BFLOAT16 Tensor Core	312 teraFLOPS 624 teraFLOPS*			
FP16 Tensor Core	312 teraFLOPS 624 teraFLOPS*			
INT8 Tensor Core	624 TOPS 1,248 TOPS*			
INT4 Tensor Core	1,248 TOPS 2,496 TOPS*			
GPU Memory	40 GB			
Interconnect	NVLink 600 GB/s PCIe Gen4 64 GB/s			
Multi-Instance GPUs	Various Instance sizes with up to 7MIGs @5GB			
Form Factor	4/8/16 SXM GPUs in HGX A1 <u>00</u>			
Max Power	400W (SXM)			



1. New Tensor Core

2. Strong Scaling

3. Elastic GPU

4. Productivity

1. New Tensor Core

2. Strong Scaling

3. Elastic GPU

4. Productivity

NEW TF32 TENSOR CORES



INSIDE A100 TensorFloat-32 (TF32)



Range of FP32 with precision of FP16



FP32 input/output

FP32 storage and math for all activations, gradients, ... everything outside tensor cores

Out-of-the-box tensor core acceleration for DL

Easy step towards maximizing tensor core performance with mixed-precision (FP16, BF16)

Up to 4x speedup on linear solvers for HPC

6622082: Mixed-Precision Training of Neural Networks, 5/202:45pmPDT

🐼 21681: How CUDA Math Libraries can help you unleash the power of the new NVIDIA A100 GPU (recording available) 16 🧧 💷 🚺

INSIDE A100 SPARSE TENSOR CORE



~No loss in inferencing accuracy

Evaluated across dozens of networks: vision, object detection, segmentation, natural language modeling, translation

• S22085: Accelerating Sparsity in the NVIDIAAmpere Architecture, 5/20 1:30pm PDT

1. New Tensor Core

2. Strong Scaling

3. Elastic GPU

4. Productivity

DL STRONG SCALING

DL networks: Long chains of sequentiallydependent compute-intensive layers





~2.5x larger network runs in same time







Fixed network runs ~2.5x faster

HOW TO KEEP TENSOR CORES FED?

Math bandwidth (MACs/clock/SM)





A100 STRONG SCALING INNOVATIONS



A100 TENSOR CORE 2x throughput vs. V100, >2x efficiency



Registers

V100 TC Instruction (1024 MACs, 8 cycles)

8-Thread 8-Thread 8-Thread 8-Thread

FFMA

(32 MACs, 2 cycles)

32 Threads (Warp)





A100 TC Instruction



16x16x16 matrix multiply	FFMA	V100 TC	A100 TC	A100 vs. V100 (improvement)	A100 vs. FFMA (improvement)
Thread sharing	1	8	32	4x	32x
Hardware instructions	128	16	2	8x	64x
Register reads+writes (warp)	512	80	28	2.9 x	18x
Cycles	256	32	16	2x	16x

Tensor Cores assume FP16 inputs with FP32 accumulator, V100 Tensor Core instruction uses 4 hardware instructions

💿 NVIDIA

A100 SM DATA MOVEMENT EFFICIENCY 3x SMEM/L1 bandwidth, 2x in-flight capacity



A100 L2 BANDWIDTH



Split L2 with hierarchical crossbar -2.3x increase in bandwidth over V100, lower latency



A100 DRAM BANDWIDTH



A100 COMPUTE DATA COMPRESSION

Activation sparsity due to ReLU



Math

RF

SMEM/L1

2

DRAM

NVLINK

Up to 4x DRAM+L2 bandwidth and 2x L2 capacity for fine-grained unstructured sparsity



S21819: Optimizing Applications for NVIDIA Ampere GPU Architecture, 5/21 10:15am FDT

A100 NVLINK BANDWIDTH

Third Generation NVLink

Math

RF

SMEM/L1

DRAM

NVLINK

50 Gbit/sec per signal pair 12 links, 25 GB/s in/out, 600 GB/s total 2x vs. V100



6621884: Under the Hood of the new DGX A100 System Architecture (recording available soon)

A100 STRONG SCALING INNOVATIONS Delivering unprecedented levels of performance

A100 improvements over V100





1. New Tensor Core

2. Strong Scaling

3. Elastic GPU

4. Productivity

NVLINK: ONE BIG GPU

- InfiniBand/Ethernet: travels a long distance, consistency is the responsibility of software
- PCI Express: hardware consistency for I/O, not for programming language memory models
- NVLINK: hardware consistency for programming language memory models, like system bus



HGX A100: 3RD GEN NVLINK

HGX A100 4-GPU: fully-connected system with 100GB/s all-to-all BW

HGX A100: 3RD GEN NVLINK & SWITCH

- HGX A100 4-GPU: fully-connected system with 100GB/s all-to-all BW
- New NVSwitch: 6B transistors in TSMC 7FF, 36 ports, 25GB/s each, per direction
- HGX A100 8-GPU: 6x NVSwitch in a fat tree topology, 2.4TB/s full-duplex bandwidth



DGX A100: PCIE4 CONTROL & I/O


CLOUD SMALL INSTANCE USAGE

- Small workloads can under-utilize GPU cloud instances, provisioned at whole GPU level
- CSPs can't use MPS for GPU space-sharing, because it doesn't provide enough isolation



NEW: MULTI-INSTANCE GPU (MIG)

- Up to 7 instances total, dynamically reconfigurable
- **Compute instances:** compute/fault isolation, but share/compete for memory
- **GPU instances:** separate and isolated paths through the entire memory system



ELASTIC GPU COMPUTING

- Each A100 is 1 to 7 GPUs
- Each DGX A100 is 1 to 56 GPUs
- Each GPU can serve a different user, with full memory isolation and QoS



UNIFIED AI ACCELERATION

- 6621975: Inside NVIDIA's Multi-Instance GPUFeature (recording available)
- 6621884: Under the Hood of the new DGX A100 System Architecture (recording available soon)
- 6621702: Introducing NVIDIA DGX A100: The Universal AI System for Enterprise, 5/20 9:00am PDT

7X HIGHER INFERENCE THROUGHPUT WITH MIG



BERT Large Inference | T4: TRT 7.1, Precision = INT8, Batch Size =256, V100: TRT 7.1, Precision = FP16, Batch Size =256 | A100 with 7 MIG instances of 1g.5gb : Pre-production TRT, Batch Size =94, Precision = INT8 with Sparsity

9X MORE PERFORMANCE IN 4 YEARS

Beyond Moore's Law With Full Stack Innovation



Geometric Mean of application speedups vs. P100 : Benchmark Application: Amber [PME-Cellulose_NVE], Chroma [szscl21_24_128], GROMACS [ADH Dodec], MILC [Apex Medium], NAMD [stmv_nve_cuda], PyTorch (BERT Large Fine Tuner], Quantum Espresso [AUSURF112-jR]; Random Forest FP32 [make_blobs (160000 x 64 : 10)], TensorFlow [ResNet-50], VASP 6 [Si Huge], IGPU node: with dual-socket CPUs with 4x P100, V100, or A100 GPUs.

1. New Tensor Core

2. Strong Scaling

3. Elastic GPU

4. Productivity



https://www.youtube.com/watch?v=TJcKYUTaBtg&t=3s

NVIDIA

One Platform. All challenges.



ONE ARCHITECTURE – CUDA



KEY ANNOUNCEMENT ASSETS

JHH Keynote On Demand

Tuesday, May 19th

Inside the NVIDIA Ampere Architecture 09:00 AM - 10:00 AM

CUDA New Features And Beyond 10:15 AM - 11:15 AM

<u>CUDA on NVIDIA Ampere GPU</u> <u>Architecture: Taking Your</u> <u>Algorithms to the Next Level of</u> <u>Performance</u> **11:30 AM - 12:30 PM**

Inside the NVIDIA HPC SDK: the Compilers, Libraries and Tools for Accelerated Computing 1:30 PM - 2:30 PM

Wednesday, May 20th

Introducing NVIDIA DGX A100: The Universal AI System for Enterprise 9:00 AM - 10:00 AM

Accelerating Deep Learning Inference With Sparse Tensor Cores of Ampere GPU Architecture 1:30 PM - 2:30 PM

Mixed-Precision Training of Neural Networks 2:45 PM - 3:45 PM

Thursday, May 21st

Tensor Core Performance on NVIDIA GPUs: The Ultimate Guide 9:00 AM - 10:00 AM

Optimizing Applications for NVIDIA Ampere GPU Architecture 10:15 AM - 11:15 AM

Developing CUDA kernels to push Tensor Cores to the Absolute Limit on NVIDIA A100 11:30 AM - 12:30 PM

High-Performance Next-Generation Deep-Learning Clusters 1:30 PM - 2:30

