# CUDA 9.2 UPDATE

OLCF User Group Call, 4/25/2018

# CUDA ECOSYSTEM 2018

**CUDA DOWNLOADS IN 2017**
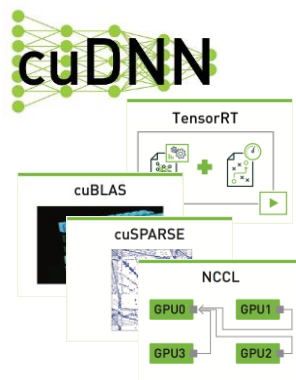**3,500,000**

**CUDA REGISTERED DEVELOPERS**
**800,000**

**GTC ATTENDEES**
**8,000+**

NVIDIA

# CUDA DEVELOPMENT ECOSYSTEM

## From Ease of Use to Specialized Performance
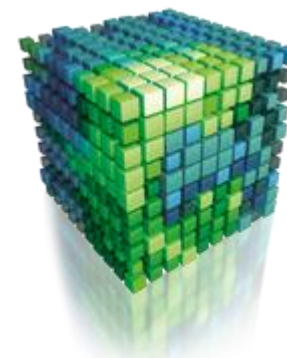


**Applications**

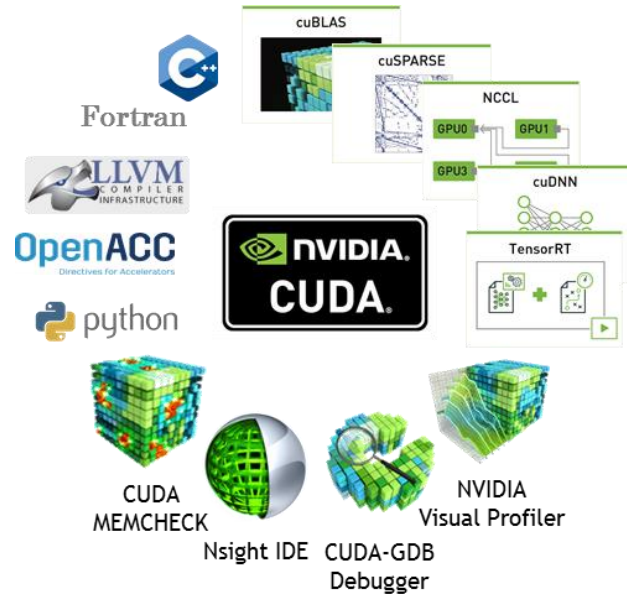**Frameworks**

**Libraries**

**Directives and Standard Languages**

**Specialized Languages**

# CUDA RELEASES
## Accelerating the Pace

**Four CUDA releases per year**

Faster release cadence for new features and improved stability for existing users

Upcoming limited decoupling of display driver and CUDA release for ease of deployment

**Monthly cuDNN & other library updates**

Rapid innovation in library performance and functionality

Library Meta Packages independent of toolkit for easy deployment

# CUDA 9.2 NEW FEATURES AT A GLANCE

## SYSTEM & PERFORMANCE

Unified Memory + ATS on IBM-POWER9

Launch Latency Optimizations

## DEVICE CODE IMPROVEMENTS

New WMMA sizes for Tensor Cores

Heterogeneous Half-Precision Datatypes

Volta Independent Thread Scheduling Control

## MATH LIBRARIES

New CUDA Library Meta-Packages

Volta Architecture-Optimized Algorithms
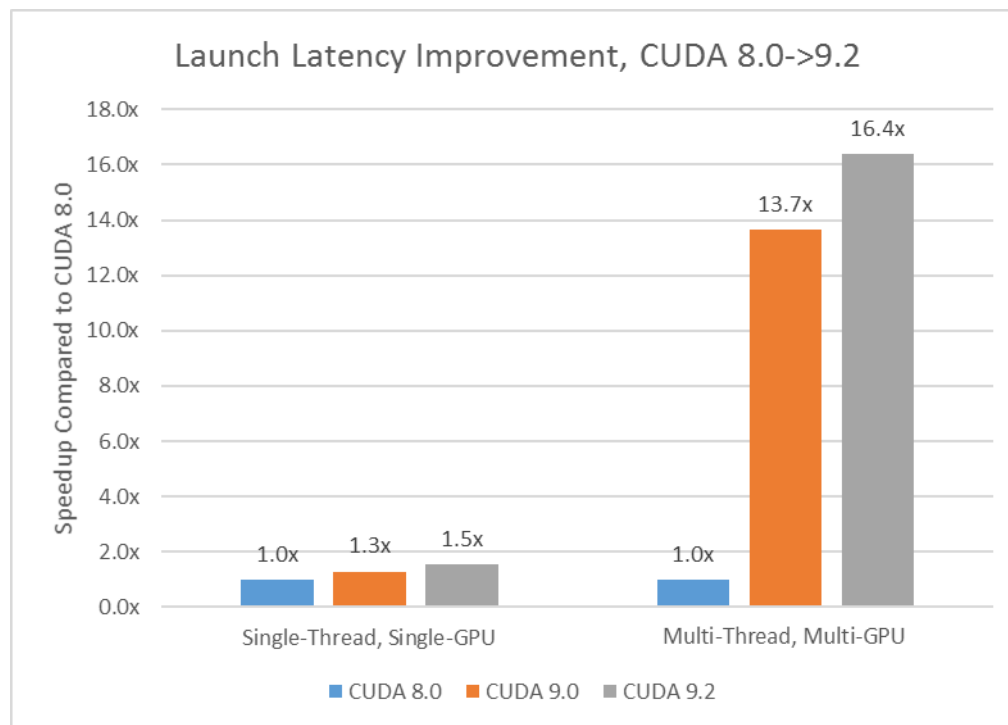
## TOOLS

Unified Nsight Product Family

Single-Pass Tracing & Profiling

# LAUNCH LATENCY IMPROVEMENTS

## Multi-GPU Launches & Kernels With Many Arguments: Now Much Faster

### Lower overhead for short kernels

- Significant factor for deep learning inference workloads

- Significant factor for small computational workloads (e.g. small FFT, small vector ops)



Launch Latency Improvement, CUDA 8.0->9.2

(Speedup Compared to CUDA 8.0)

Single-Thread, Single-GPU: 1.0x, 1.3x, 1.5x
Multi-Thread, Multi-GPU: 1.0x, 13.7x, 16.4x

CUDA 8.0 ■ CUDA 9.0 ■ CUDA 9.2

NVIDIA.

# VOLTA NANOSLEEP TIMER
## For Polling & Synchronization Operations
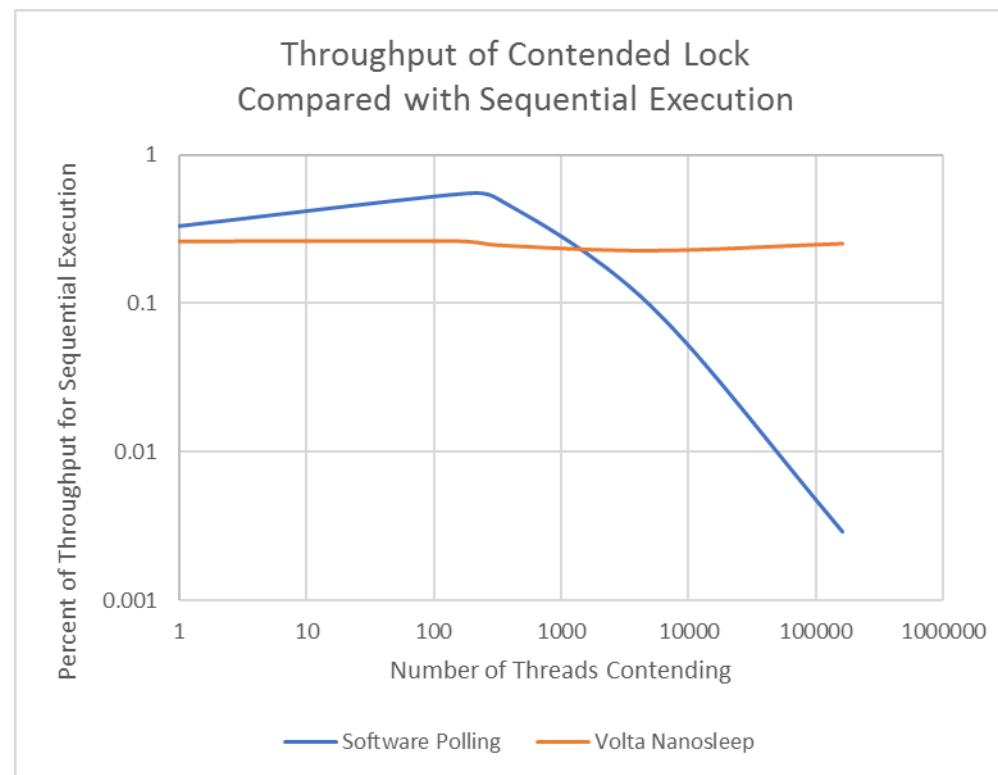
## New *nanosleep()* instruction

```
__device__ void __nanosleep(unsigned int ns);
```

Sleeps a thread for an amount of time

Sleeping thread yields execution to other active threads

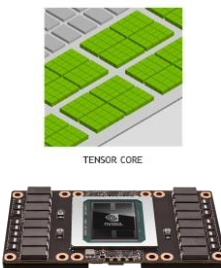Integrated into hardware thread scheduler

Ideal for timed-backoff polling

Throughput of Contended Lock
Compared with Sequential Execution



⬢ nvidia.

# MATH LIBRARIES: WHAT'S NEW

## VOLTA PLATFORM SUPPORT

Volta architecture optimized GEMMs, & GEMM extensions for Volta Tensor Cores (cuBLAS)

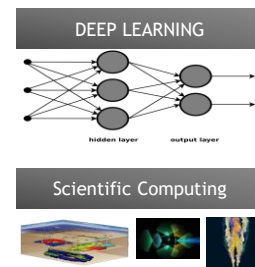Out-of-box performance on Volta (all libraries)

## PERFORMANCE

GEMM optimizations for RNNs (cuBLAS)

Faster image processing (NPP)

Prime factor FFT performance (cuFFT)
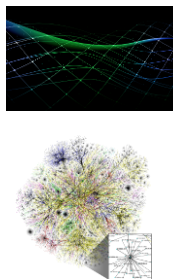
SpMV performance (cuSPARSE)

## NEW ALGORITHMS

Mixed-precision Batched GEMM for attention models (cuBLAS)

Image Augmentation and batched image processing routines (NPP)

Batched pentadiagonal solver (cuSPARSE)

## MEMORY & FOOTPRINT OPTIMIZATION

Large FFT sizes on multi-GPU systems (cuFFT)

Modular functional blocks with small footprint (NPP)

# CUDA TENSOR CORE PROGRAMMING

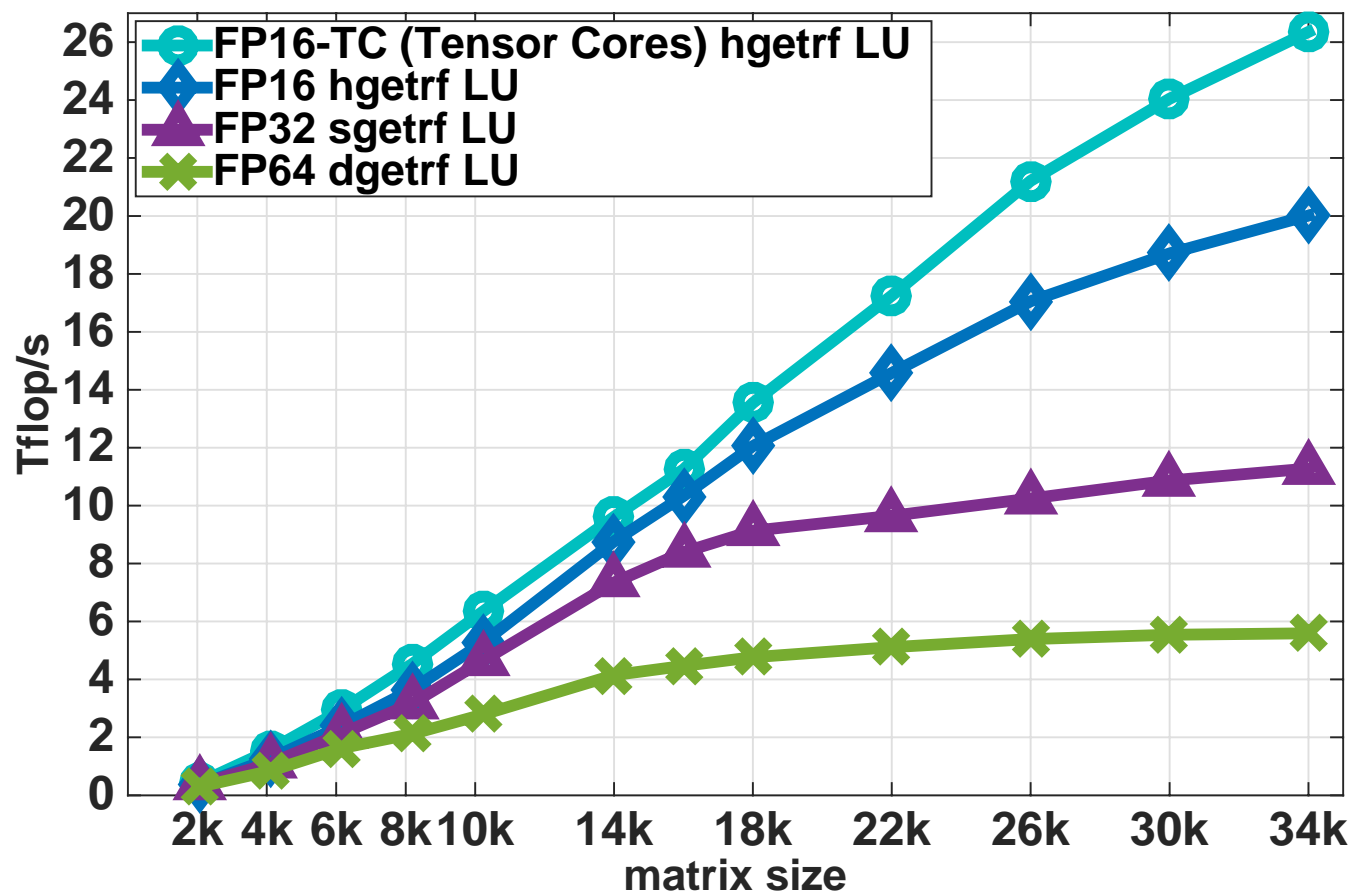## 16x16x16 Warp Matrix Multiply and Accumulate (WMMA)

```
wmma::mma_sync(Dmat, Amat, Bmat, Cmat);
```

$$D = \begin{pmatrix} & \\ & \end{pmatrix} \begin{pmatrix} & \\ & \end{pmatrix} + \begin{pmatrix} & \\ & \end{pmatrix}$$

FP16 or FP32   FP16     FP16    FP16 or FP32

$$D = AB + C$$

# LINEAR ALGEBRA + TENSOR CORES



**Double Precision LU Decomposition**

- Compute initial solution in FP16

- Iteratively refine to FP64

Achieved FP64 Tflops: **26**

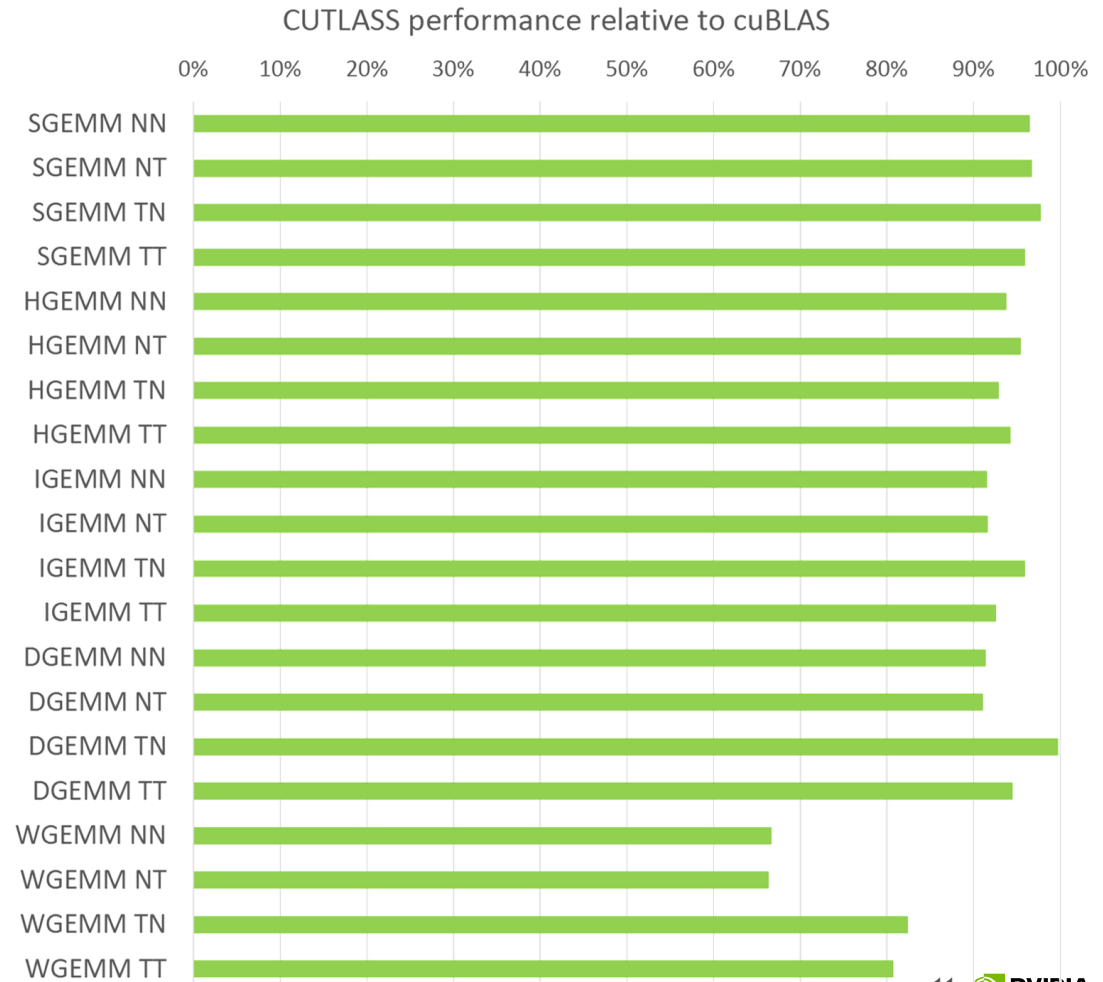Device FP64 Tflops: **7.5**

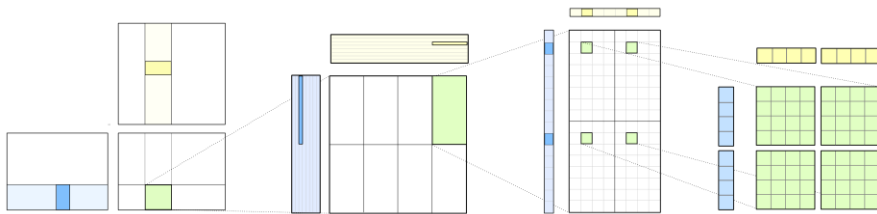# CUTLASS

Template library for linear algebra operations in CUDA C++

>90% CUBLAS performance

Open Source (3-clause BSD License)
https://github.com/NVIDIA/cutlass
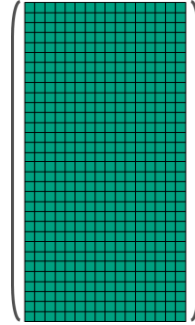
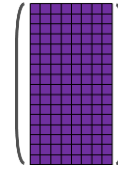CUTLASS performance relative to cuBLAS
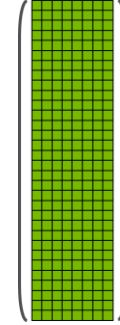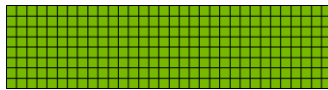


NVIDIA.

# NEW WMMA MATRIX SIZES

## WMMA 32x8x16

D
32x8
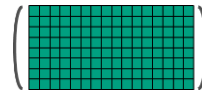
=

A
32x16

B
16x8

+

C
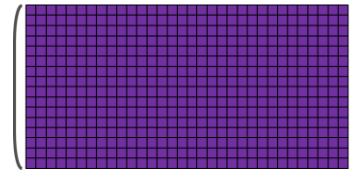32x8

## WMMA 8x32x16

D
8x32

=

A
8x16

B
16x32

+

C
8x32

NVIDIA.

# TOOLS UPDATES FOR CUDA 9.2

## NVPROF

New Metrics: Tensor Cores, L2, Memory Instructions Per Load/Store

PCIe Topology Display

Single-Pass Tracing & Profiling

## CUPTI

New Activity Kind: PCIE

New Attribute: Profiling Scope (Device-Level, Context-Level)

Exposes New Metrics

## VISUAL PROFILER

Summary View for Memory Hierarchy

Improved Handling of Segments for UVM Data on the Timeline

## DEBUGGER

Lightweight Coredump Files

User-Induced Coredumps

Coredump Support on Volta-MPS

NVIDIA.

# HIERARCHICAL MEMORY STATISTICS

# NSIGHT PRODUCT FAMILY

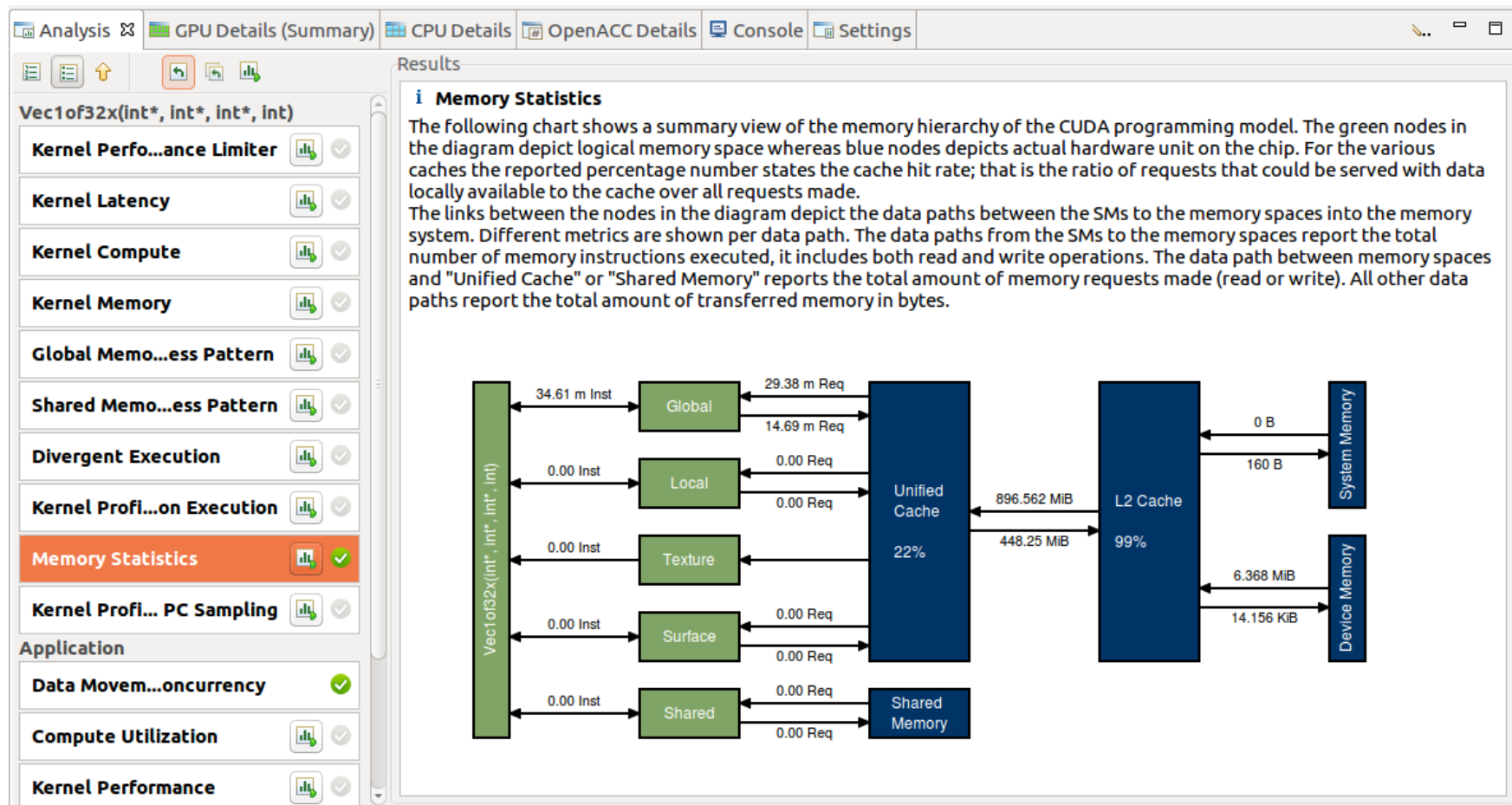## Standalone Performance Tools

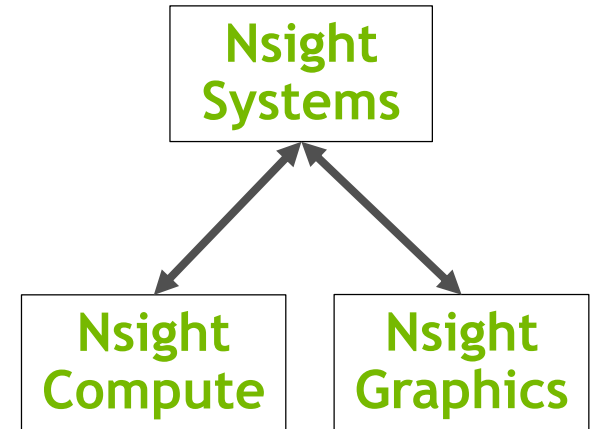**Nsight Systems** - System-wide application algorithm tuning

**Nsight Compute** - Debug/optimize specific CUDA kernel

**Nsight Graphics** - Debug/optimize specific graphics shader

## IDE Plugins

**Nsight Visual Studio/Eclipse Edition** – editor, debugger, some perf analysis

## Workflow

```
        Nsight
        Systems
       ↙      ↘
  Nsight      Nsight
  Compute     Graphics
```

# UNIFIED MEMORY WITH ATS ON IBM POWER9

## IBM Power9 CPUs With NVLink Interconnect

Managed Memory With ATS

V100 1

V100 0

POWER9

NVLink Interconnect

### ALLOCATION

Automatic access to all system memory: malloc, stack, file system

### ACCESS

All data accessible concurrently from any processor, anytime

Atomic operations resolved directly over NVLink

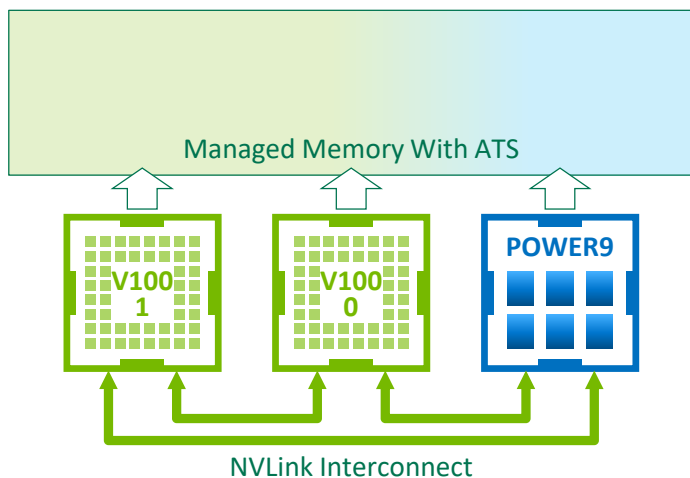# UNIFIED MEMORY WITH ATS ON IBM POWER9

## IBM Power9 CPUs With NVLink Interconnect

Managed Memory With ATS

V100 1

V100 0

POWER9

NVLink Interconnect

## ATS & POWER9 FEATURES

ATS allows GPUDirect RDMA to unified memory

Managed memory is cache-coherent between CPU & GPU

CPU has direct access to GPU memory without need for migration

NVIDIA.

# WHAT YOU CAN DO WITH UNIFIED MEMORY

## Works everywhere today

```
int *data;
cudaMallocManaged(&data, sizeof(int) * n);
kernel<<< grid, block >>>(data);
```

## Works on POWER9 + CUDA 9.2

```
int *data = (int*)malloc(sizeof(int) * n);
kernel<<< grid, block >>>(data);
```

```
int data[1024];
kernel<<< grid, block >>>(data);
```
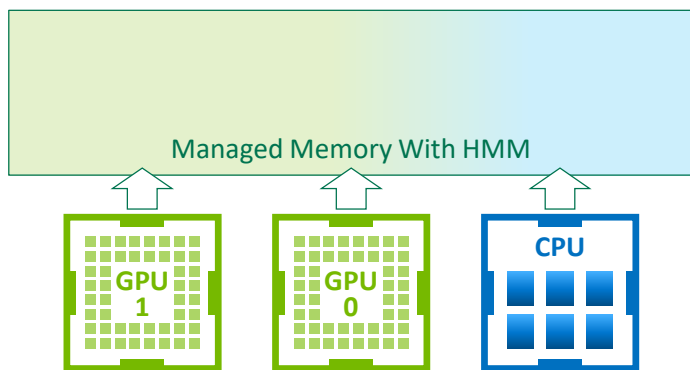
```
int *data = (int*)alloca(sizeof(int) * n);
kernel<<< grid, block >>>(data);
```

```
extern int *data;
kernel<<< grid, block >>>(data);
```

NVIDIA.

# BEYOND

# HETEROGENEOUS MEMORY ON x86-LINUX

## Feature Parity With POWER9 + ATS

Managed Memory With HMM

GPU 1
GPU 0
CPU

## ALLOCATION

Automatic access to <u>all</u> system memory: malloc, stack, file system

## ACCESS

All data accessible concurrently from any processor, anytime

Concurrent atomic operations permitted, resolved via page fault