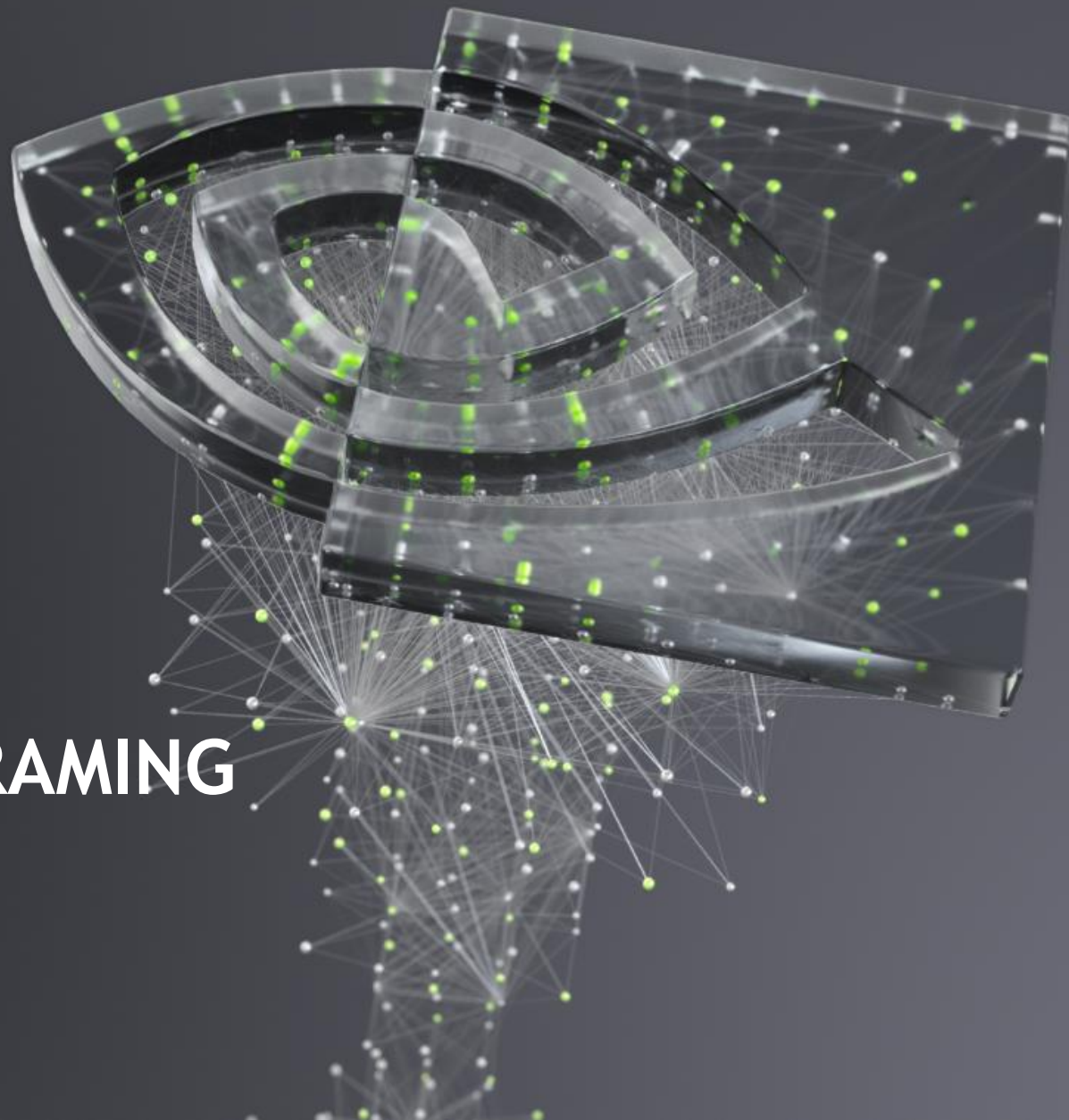# GPU ACCELERATION PRINCIPLE AND PROGRAMING - HPC SDK

Ryan Jeng | Tech Lead, DevTech APAC

# GPU PROGRAMMING IN 2020 AND BEYOND

## Math Libraries | Standard Languages | Directives | CUDA

```cpp
std::transform(par, x, x+n, y, y,
    [=](float x, float y){
        return y + a*x;
});
```

```fortran
do concurrent (i = 1:n)
    y(i) = y(i) + a*x(i)
enddo
```

**GPU Accelerated**
**C++ and Fortran**

```cpp
#pragma acc data copy(x,y)
{

...

std::transform(par, x, x+n, y, y,
    [=](float x, float y){
        return y + a*x;
});

...

}
```

**Incremental Performance**
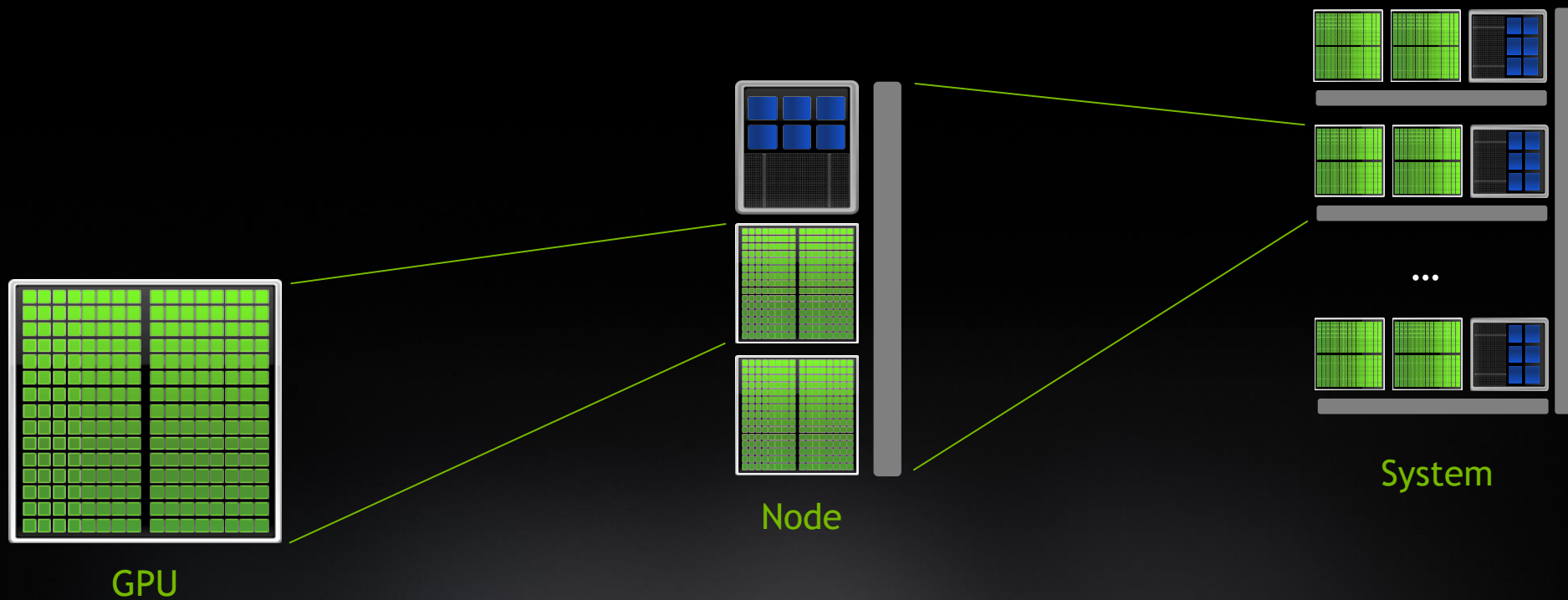**Optimization with Directives**

```cpp
__global__
void saxpy(int n, float a,
           float *x, float *y) {
    int i = blockIdx.x*blockDim.x +
        threadIdx.x;
    if (i < n) y[i] += a*x[i];
}

int main(void) {
    ...
    cudaMemcpy(d_x, x, ...);
    cudaMemcpy(d_y, y, ...);

    saxpy<<<(N+255)/256,256>>>(...);

    cudaMemcpy(y, d_y, ...);
```

**Maximize GPU Performance with**
**CUDA C++/Fortran**

**GPU Accelerated Libraries**

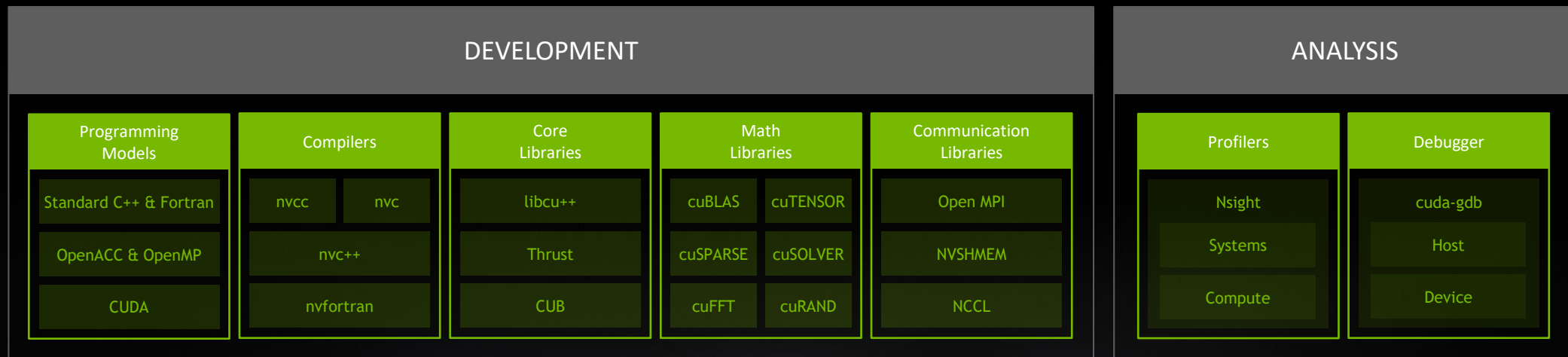# PROGRAMMING GPU-ACCELERATED HPC SYSTEMS

GPU | CPU | Interconnect

GPU

Node

System

# AVAILABLE NOW: THE NVIDIA HPC SDK

Available at developer.nvidia.com/hpc-sdk, on NGC, and in the Cloud

## NVIDIA HPC SDK

| DEVELOPMENT | | | | | ANALYSIS | |
|---|---|---|---|---|---|---|
| **Programming Models** | **Compilers** | **Core Libraries** | **Math Libraries** | **Communication Libraries** | **Profilers** | **Debugger** |
| Standard C++ & Fortran | nvcc · nvc | libcu++ | cuBLAS · cuTENSOR | Open MPI | Nsight | cuda-gdb |
| OpenACC & OpenMP | nvc++ | Thrust | cuSPARSE · cuSOLVER | NVSHMEM | Systems | Host |
| CUDA | nvfortran | CUB | cuFFT · cuRAND | NCCL | Compute | Device |

Develop for the NVIDIA HPC Platform: GPU, CPU and Interconnect

HPC Libraries | GPU Accelerated C++ and Fortran | Directives | CUDA

7-8 Releases Per Year | Freely Available

NVIDIA.

# HPC COMPILERS & STANDARD LANGUAGE PARALLELISM

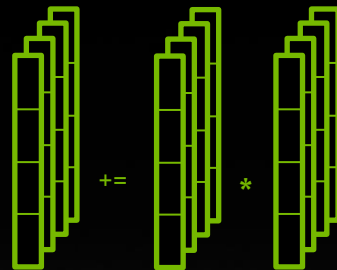# HPC COMPILERS

## NVC | NVC++ | NVFORTRAN
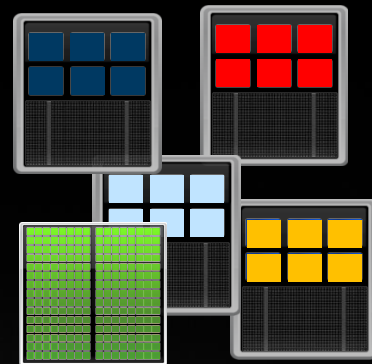
**Accelerated**
A100
Automatic

**Programmable**
Standard Languages
Directives
CUDA

**Multicore**
Directives
Vectorization

**Multi-Platform**
x86_64
Arm
OpenPOWER

# HPC PROGRAMMING IN ISO C++

ISO is the place for portable concurrency and parallelism

## C++17

**Parallel Algorithms**

➤  In NVC++

➤  Parallel and vector concurrency

**Forward Progress Guarantees**

➤ Extend the C++ execution model for accelerators

**Memory Model Clarifications**

➤ Extend the C++ memory model for accelerators

## C++20

**Scalable Synchronization Library**

➤  Express thread synchronization that is portable and scalable across CPUs and accelerators

➤  In libcu++:

  ➤  `std::atomic<T>`

  ➤  `std::barrier`

  ➤  `std::counting_semaphore`

  ➤  `std::atomic<T>::wait/notify_*`

  ➤  `std::atomic_ref<T>`

## C++23 and Beyond

**Executors**

➤  Simplify launching and managing parallel work across CPUs and accelerators

`std::mdspan/mdarray`

➤  HPC-oriented multi-dimensional array abstractions.

**Linear Algebra**

➤  C++ standard algorithms API to linear algebra

➤  Maps to vendor optimized BLAS libraries

**Extended Floating Point Types**

➤  First-class support for formats new and old: `std::float16_t/float64_t`

# HPC PROGRAMMING IN ISO C++

## C++ Parallel Algorithms

```
std::sort(std::execution::par, c.begin(), c.end());

std::unique(std::execution::par, c.begin(), c.end());
```

- ➤ Introduced in C++17
- ➤ Parallel and vector concurrency via execution policies

    `std::execution::par, std::execution::par_seq, std::execution::seq`

- ➤ Several new algorithms in C++17 including

    - ➤ `std::for_each_n(POLICY, first, size, func)`

- ➤ Insert `std::execution::par` as first parameter when calling algorithms
- ➤ NVC++ 20.5: automatic GPU acceleration of C++17 parallel algorithms

    - ➤ Leverages CUDA Unified Memory

NVIDIA.

# PARALLEL C++

```cpp
static inline
void CalcHydroConstraintForElems(Domain &domain, Index_t length,
                    Index_t *regElemlist, Real_t dvovmax, Real_t& dthydro)
{
#if _OPENMP
  const Index_t threads = omp_get_max_threads();
  Index_t hydro_elem_per_thread[threads];
  Real_t dthydro_per_thread[threads];
#else
  Index_t threads = 1;
  Index_t hydro_elem_per_thread[1];
  Real_t dthydro_per_thread[1];
#endif
#pragma omp parallel firstprivate(length, dvovmax)
  {
    Real_t dthydro_tmp = dthydro ;
    Index_t hydro_elem = -1 ;
#if _OPENMP
    Index_t thread_num = omp_get_thread_num();
#else
    Index_t thread_num = 0;
#endif
#pragma omp for
    for (Index_t i = 0 ; i < length ; ++i) {
      Index_t indx = regElemlist[i] ;

      if (domain.vdov(indx) != Real_t(0.)) {
        Real_t dtdvov = dvovmax / (FABS(domain.vdov(indx))+Real_t(1.e-20)) ;

        if ( dthydro_tmp > dtdvov ) {
          dthydro_tmp = dtdvov ;
          hydro_elem = indx ;
        }
      }
    }
    dthydro_per_thread[thread_num] = dthydro_tmp ;
    hydro_elem_per_thread[thread_num] = hydro_elem ;
  }
  for (Index_t i = 1; i < threads; ++i) {
    if(dthydro_per_thread[i] < dthydro_per_thread[0]) {
      dthydro_per_thread[0] = dthydro_per_thread[i];
      hydro_elem_per_thread[0] = hydro_elem_per_thread[i];
    }
  }
  if (hydro_elem_per_thread[0] != -1) {
    dthydro = dthydro_per_thread[0] ;
  }
  return ;
}
```

**C++ with OpenMP**

➢ Composable, compact and elegant

➢ Easy to read and maintain

➢ ISO Standard

➢ Portable – nvc++, g++, icpc, MSVC, …

```cpp
static inline void CalcHydroConstraintForElems(Domain &domain, Index_t length,
                                            Index_t *regElemlist,
                                            Real_t dvovmax,
                                            Real_t &dthydro)
{
  dthydro = std::transform_reduce(
    std::execution::par, counting_iterator(0), counting_iterator(length),
    dthydro, [](Real_t a, Real_t b) { return a < b ? a : b; },
    [=, &domain](Index_t i)
  {
    Index_t indx = regElemlist[i];
    if (domain.vdov(indx) == Real_t(0.0)) {
      return std::numeric_limits<Real_t>::max();
    } else {
      return dvovmax / (std::abs(domain.vdov(indx)) + Real_t(1.e-20));
    }
  });
}
```
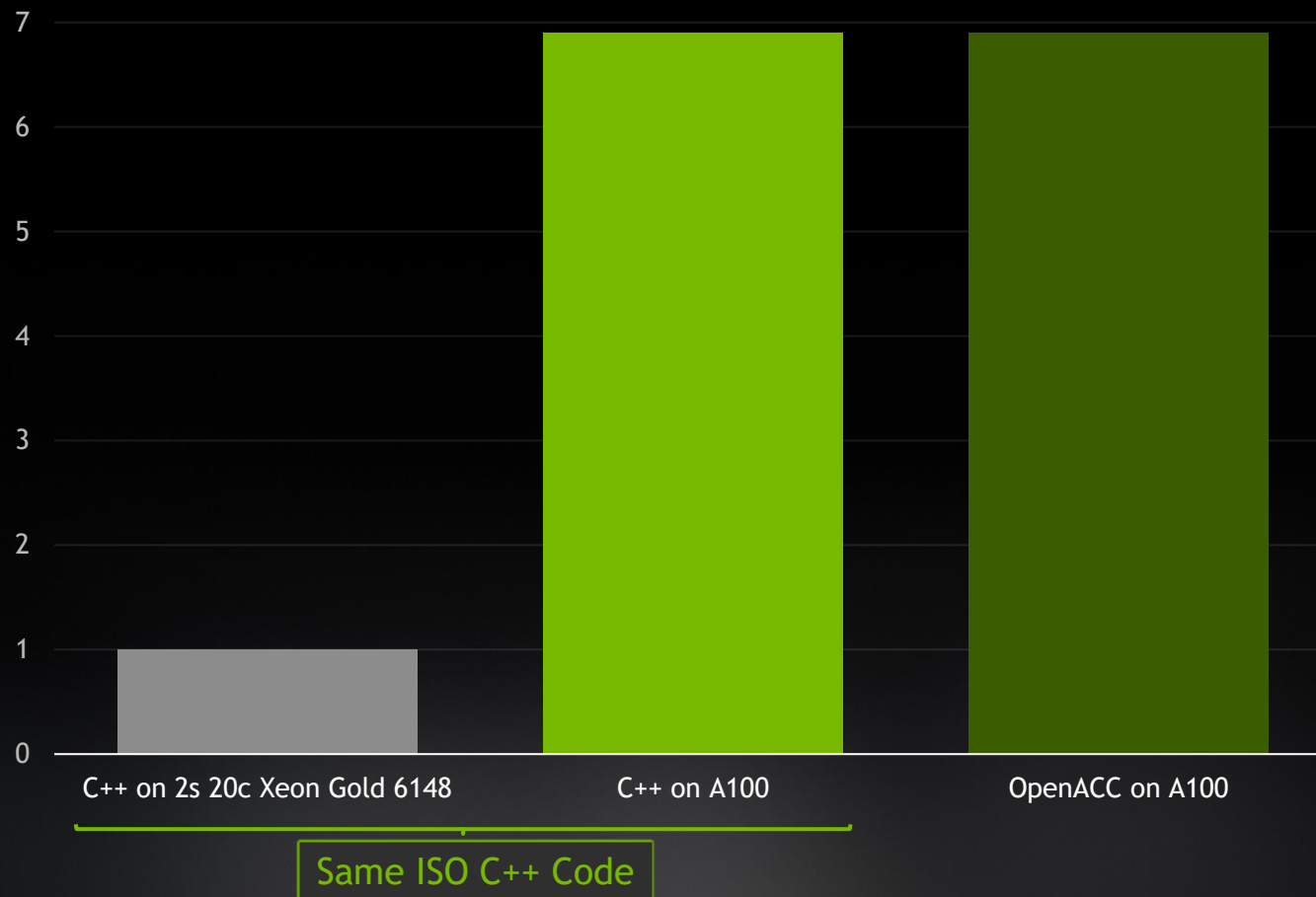
**Parallel C++17**

# LULESH PERFORMANCE

## Speedup – Higher is Better



Chart with vertical axis from 0 to 7:

- C++ on 2s 20c Xeon Gold 6148: ≈1
- C++ on A100: ≈6.9
- OpenACC on A100: ≈6.9

Same ISO C++ Code (covering C++ on 2s 20c Xeon Gold 6148 and C++ on A100)

# HPC PROGRAMMING IN ISO FORTRAN

ISO is the place for portable concurrency and parallelism

## Fortran 2018

**Array Syntax and Intrinsics**
- NVFORTRAN 20.5
- Accelerated matmul, reshape, spread, etc

**DO CONCURRENT**
- NVFORTRAN 20.x
- Auto-offload & multi-core

**Co-Arrays**
- Coming Soon
- Accelerated co-array images

## Fortran 202x

**DO CONCURRENT Reductions**
- REDUCE subclause added
- Support for +, *, MIN, MAX, IAND, IOR, IEOR.
- Support for .AND., .OR., .EQV., .NEQV on LOGICAL values
- Atomics

# HPC PROGRAMMING IN ISO FORTRAN

## NVFORTRAN Accelerates Fortran Intrinsics with cuTENSOR Backend

```fortran
real(8), dimension(ni,nk) :: a
real(8), dimension(nk,nj) :: b
real(8), dimension(ni,nj) :: c, d

...

!$acc enter data copyin(a,b,c) create(d)

do nt = 1, ntimes
    !$acc kernels
    do j = 1, nj
        do i = 1, ni
            d(i,j) = c(i,j)
            do k = 1, nk
                d(i,j) = d(i,j) + a(i,k) * b(k,j)
            end do
        end do
    end do
    !$acc end kernels
end do

!$acc exit data copyout(d)
```

**Inline FP64 matrix multiply**

```fortran
real(8), dimension(ni,nk) :: a
real(8), dimension(nk,nj) :: b
real(8), dimension(ni,nj) :: c, d

...

!$acc enter data copyin(a,b,c) create(d)

...

!$acc host_data use_device(a,b,c,d)
do nt = 1, ntimes
    d = c + matmul(a,b)
end do
!$acc end host_data

...

!$acc exit data copyout(d)
```

**MATMUL FP64 matrix multiply**



Bar chart, y-axis TFLOPs (0 to 20), x-axis categories: Naïve Inline V100 (≈0.3), FORTRAN V100 (≈6.8), FORTRAN A100 (≈17.3)

HPC LIBRARIES

# A100 FEATURES IN MATH LIBRARIES

## Automatic Acceleration of Critical Routines in HPC and AI

**cuBLAS**

BF16, TF32 and FP64 Tensor Cores

**cuSPARSE**

Increased memory BW, Shared Memory and L2

**cuTENSOR**

BF16, TF32 and FP64 Tensor Cores

**cuSOLVER**

BF16, TF32 and FP64 Tensor Cores

**cuFFT**

Increased memory BW, Shared Memory and L2

**CUDA Math API**

BF16 Support

# A100 TENSOR CORES IN LIBRARIES

## cuBLAS

- Automatic Tensor Core acceleration
- Removed matrix size restrictions for Tensor Core acceleration

**DGEMM on A100**

- Up to 19 TFLOPs, 2.4x V100

### FP64 Matrix Multiply: A100 vs V100



Legend: A100 FP64 Tensor Core (DMMA) — V100 FP64

Y-axis: TFLOPS (0, 4, 8, 12, 16, 20)
X-axis: Matrix Size (m=n=k) (0, 1024, 2048, 3072, 4096, 5120, 6144, 7168, 8192)

# A100 TENSOR CORES IN LIBRARIES

## cuBLAS

NVIDIA V100 FP32        NVIDIA A100 Tensor Core TF32        NVIDIA V100 FP64        NVIDIA A100 Tensor Core FP64

**10X**
THROUGHPUT

**2.5X**
THROUGHPUT

TF32 MMA Dimensions: m,n,k = 16x8x8        DMMA Dimensions: m,n,k = 8x8x4

# 10X        2.5X

# cuSPARSELt

## Extension Library with Sparse Matmul

- High-performance library for general matrix-matrix operations in which at least one operand is a sparse matrix
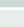- Ampere Sparse MMA tensor core support
- Mixed-precision support
- Matrix pruning and compression functionalities
- Auto-tuning functionality

Sparse Tensor Core

select

mux    mux

dot-product

Input activations

Dense trained weights

Fine-grained structured pruning (2:4 non-zero)

Fine-tuning weights

zero

Compress

Non-zero data    Non-zero indices

Output activations

# A100 TENSOR CORE

| | INPUT OPERANDS | | ACCUMULATOR | | TOPS | X-factor vs. FFMA | SPARSE TOPS | SPARSE X-factor vs. FFMA |
|---|---|---|---|---|---|---|---|---|
| **V100** | FP32 | | FP32 | | 15.7 | 1x | - | - |
| | FP16 | | FP32 | | 125 | 8x | - | - |
| **A100** | FP32 | | FP32 | | 19.5 | 1x | - | - |
| | TF32 | | FP32 | | 156 | 8x | 312 | 16x |
| | FP16 | | FP32 | | 312 | 16x | 624 | 32x |
| | BF16 | | FP32 | | 312 | 16x | 624 | 32x |
| | FP16 | | FP16 | | 312 | 16x | 624 | 32x |
| | INT8 | | INT32 | | 624 | 32x | 1248 | 64x |
| | INT4 | | INT32 | | 1248 | 64x | 2496 | 128x |
| | BINARY | | INT32 | | 4992 | 256x | - | - |
| | IEEE FP64 | | | | 19.5 | 1x | - | - |

# MATH LIBRARY DEVICE EXTENSIONS

## Introducing cuFFTDx: Device Extension

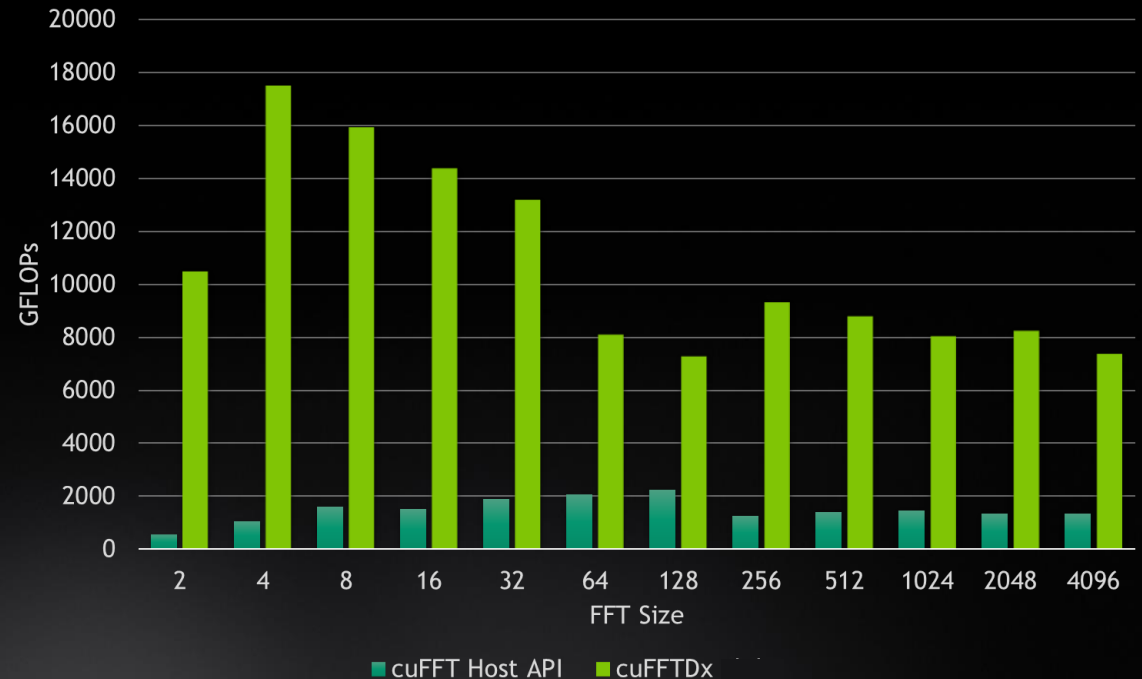**Available in Math Library EA Program**

Device callable library

Retain and reuse on-chip data

Inline FFTs in user kernels
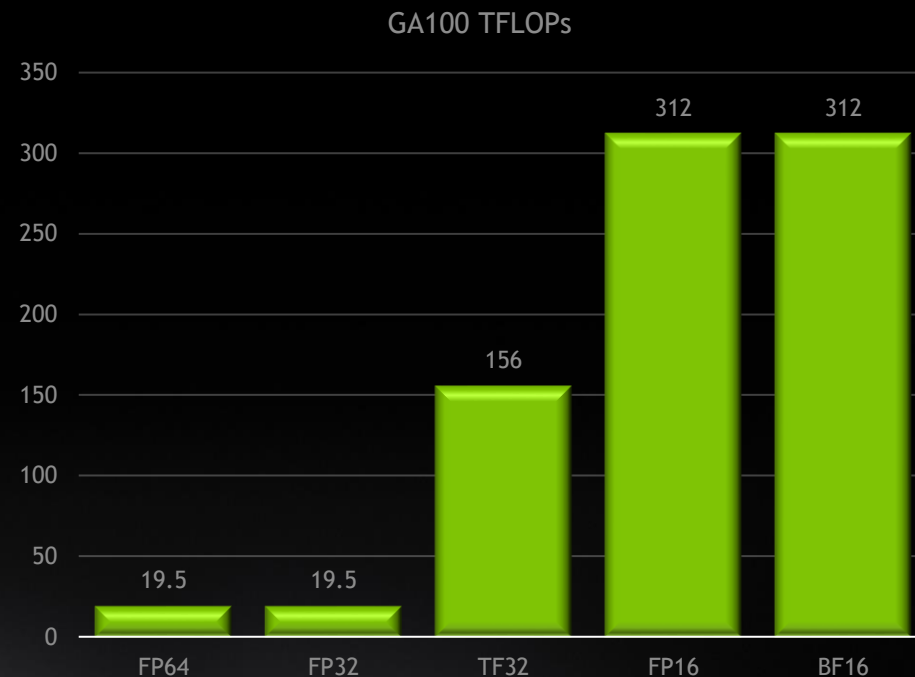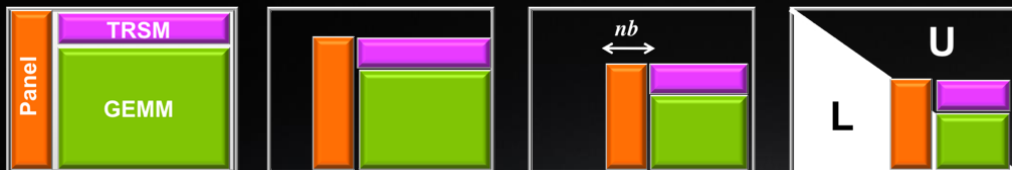
Combine multiple FFT operations

### cuFFTDx Device API V100 Performance
Small-size FFTs



Bar chart with y-axis labeled "GFLOPs" (0 to 20000) and x-axis labeled "FFT Size" (2, 4, 8, 16, 32, 64, 128, 256, 512, 1024, 2048, 4096). Legend: cuFFT Host API, cuFFTDx.

# TENSOR CORE ACCELERATED LINEAR SOLVERS

## Mixed Precision Dense Linear Solvers

- Common HPC Solvers dominated by matrix multiplication
  - LU, QR

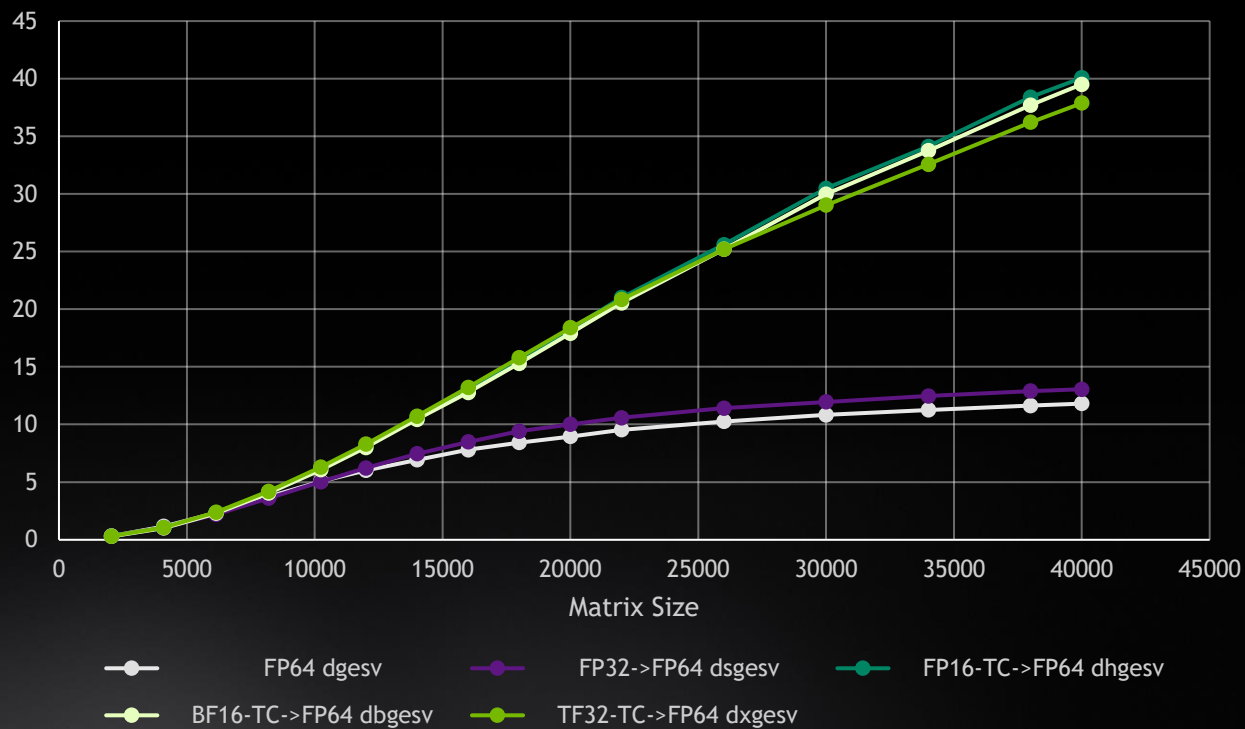- Can we accelerate with FP16 Tensor Core and retain FP64 accuracy? Yes!

**GA100 TFLOPs**

| Precision | TFLOPs |
|-----------|--------|
| FP64 | 19.5 |
| FP32 | 19.5 |
| TF32 | 156 |
| FP16 | 312 |
| BF16 | 312 |

# TENSOR CORE ACCELERATED LINEAR SOLVERS

## Mixed Precision Dense Linear Solvers

- ➤ LU & QR Solvers available in cuSOLVER

- ➤ FP64 input, FP64 output, black-box mixed precision acceleration

- ➤ Real and Complex, single and multi-RHS

### ZGETRF Performance



Legend:
- FP64 dgesv
- FP32->FP64 dsgesv
- FP16-TC->FP64 dhgesv
- BF16-TC->FP64 dbgesv
- TF32-TC->FP64 dxgesv

X-axis: Matrix Size

# MULTI GPU SUPPORT IN LIBRARIES

## Linear Algebra and FFT

### cuFFT
- Single Process Multi-GPU FFT
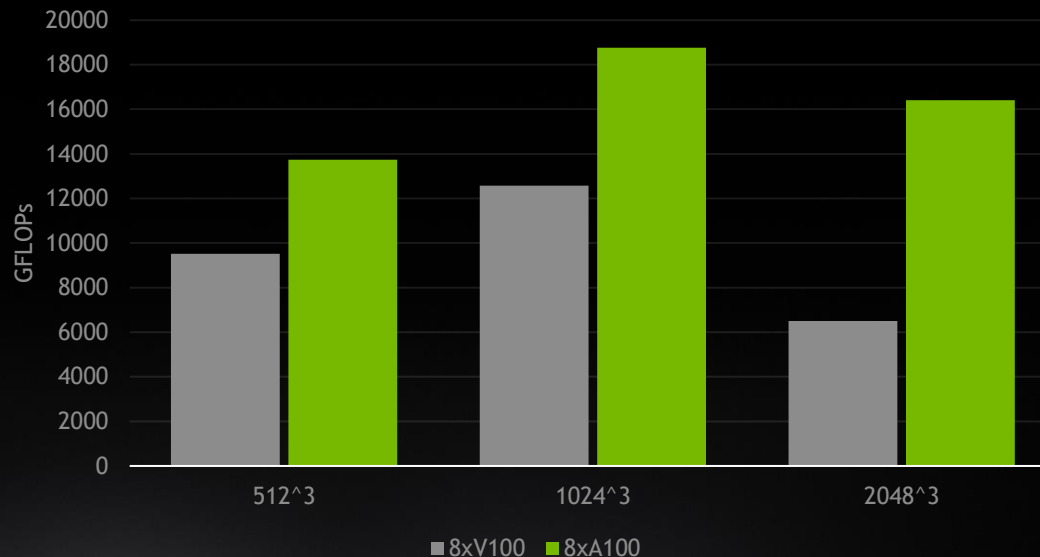- **Multi Node Multi-GPU FFT Coming Soon**

### cuSOLVER
- Single Process Multi-GPU Eigensolver
- Single Process Multi-GPU LU
- Single Process Multi-GPU Cholesky
- **Multi Node Multi-GPU LU Coming Soon**

### cuBLAS
- Improved Single Process Multi-GPU GEMM

Multi GPU cuFFT Performance, 8xV100 vs 8xA100

GFLOPs

20000
18000
16000
14000
12000
10000
8000
6000
4000
2000
0

512^3    1024^3    2048^3

■ 8xV100    ■ 8xA100

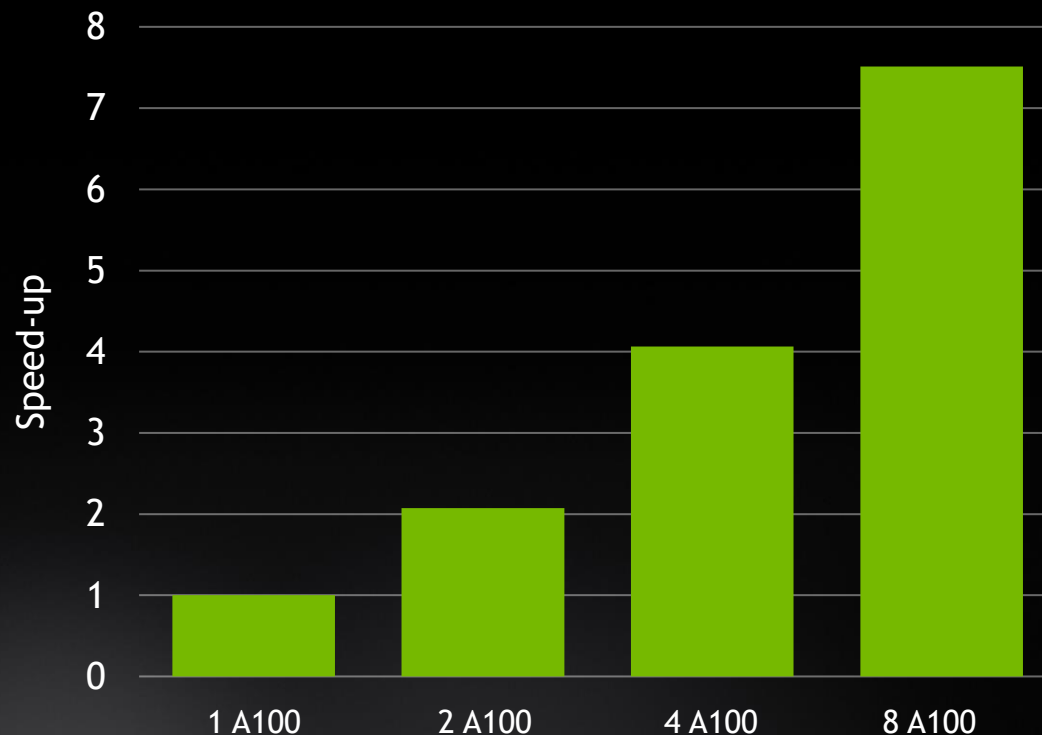# MULTI GPU WITH THE NVIDIA HPC SDK

Cloverleaf Hydrodynamics Mini-App

## Full Integration provided by HPC SDK
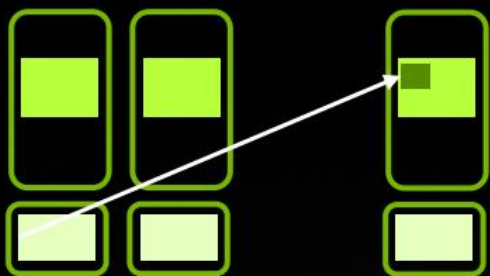
➤ Fortran + OpenACC + Open MPI

## Strong Scaling - Cloverleaf BM128

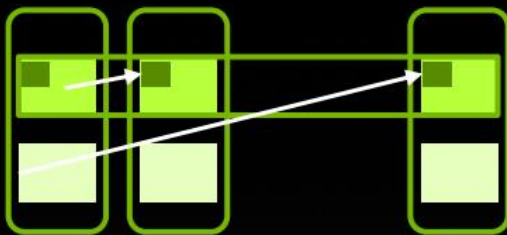➤ Perfect scaling to 4 A100 GPUs
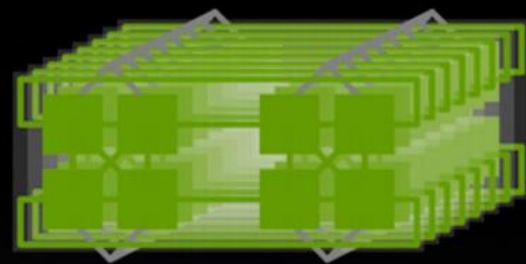➤ 7.5X speed-up on 8 A100 GPUs

# COMMUNICATION LIBRARIES

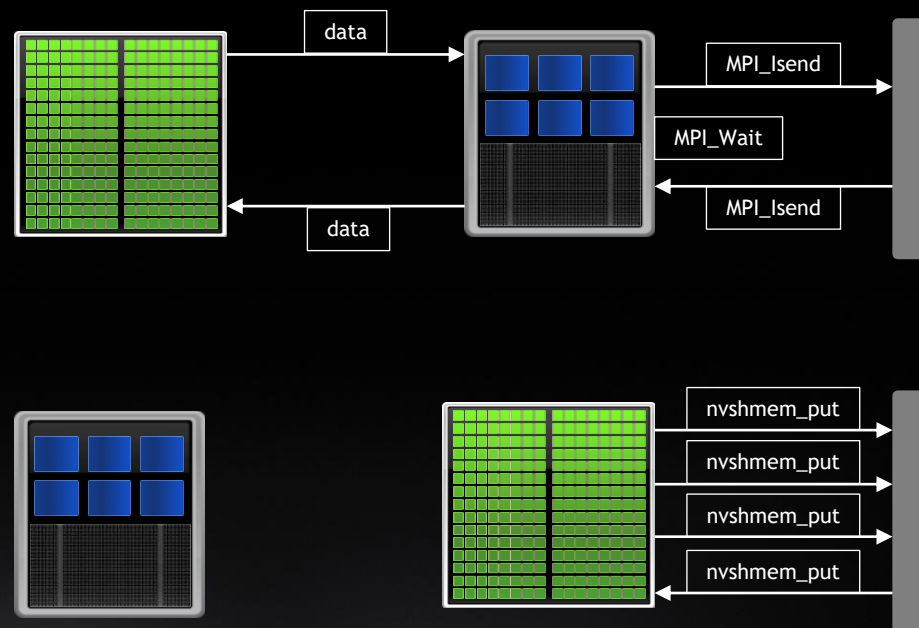Single GPU, Multi GPU, and Multi Node



**Open MPI + UCX**

**NVSHMEM**

**NCCL**

# INTRODUCING NVSHMEM

## GPU Optimized OpenSHMEM

- ➢ Initiate from CPU or GPU
- ➢ Initiate from within CUDA kernel
- ➢ Issue onto a CUDA stream
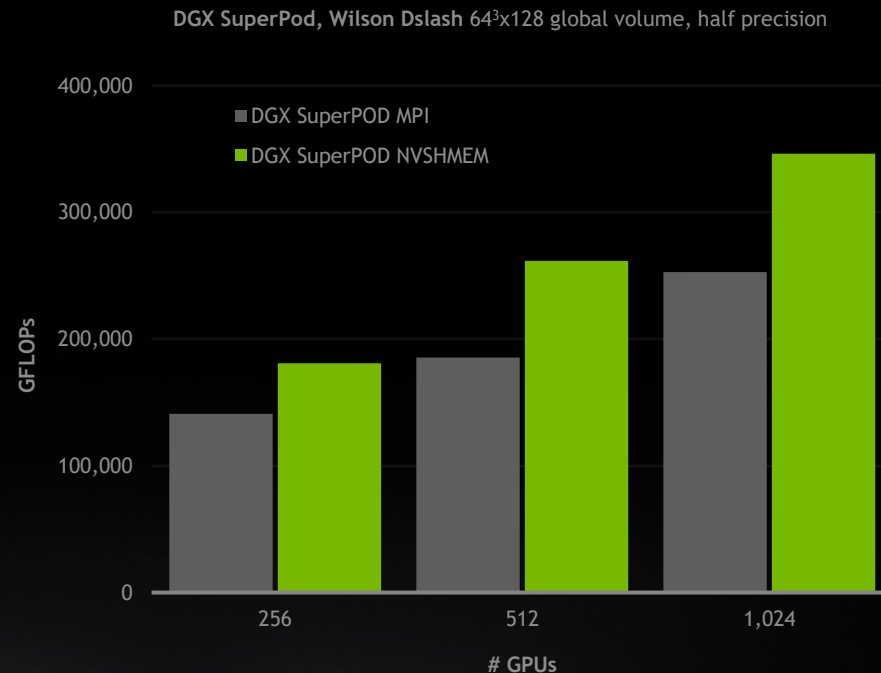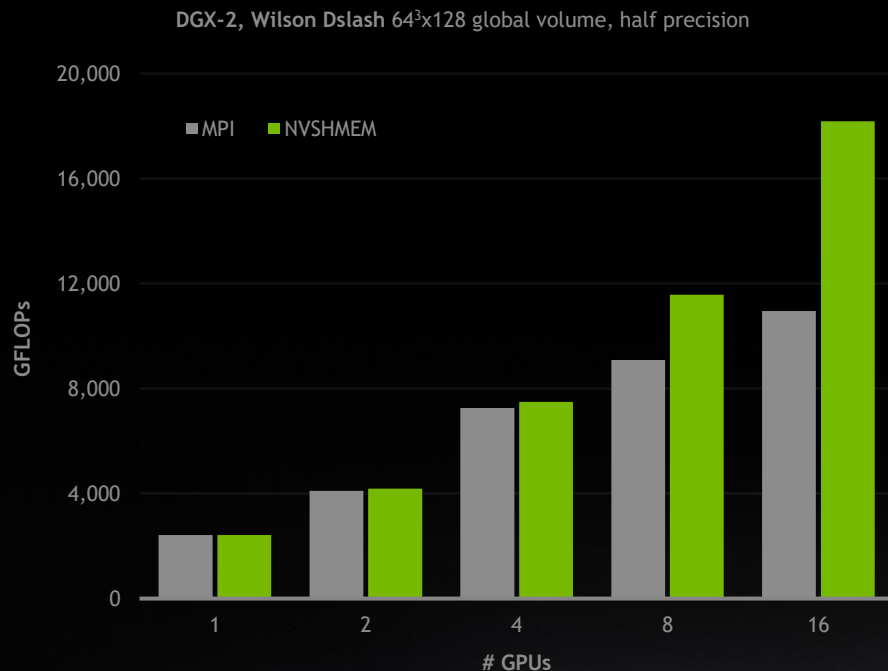- ➢ Interoperable with MPI & OpenSHMEM

Pre-release Impact

- ➢ LBANN, Kokkos/CGSolve, QUDA

data

MPI_Isend

MPI_Wait

MPI_Isend

data

nvshmem_put

nvshmem_put

nvshmem_put

nvshmem_put

# INTRODUCING NVSHMEM

## Impact in HPC Applications

**DGX-2, Wilson Dslash** $64^3$x128 global volume, half precision

**DGX SuperPod, Wilson Dslash** $64^3$x128 global volume, half precision

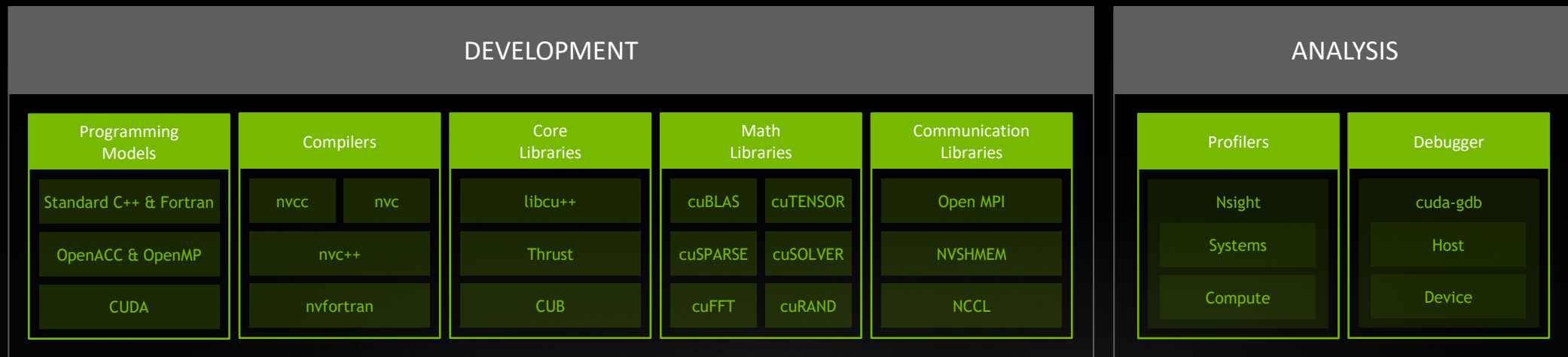**QUDA**: Quantum Chromodynamics on CUDA

➤ Up to 1.7X Single Node Speedup

➤ Up to 1.4X Multi Node Speedup

NVIDIA.

# AVAILABLE NOW: THE NVIDIA HPC SDK

Available at developer.nvidia.com/hpc-sdk, on NGC, and in the Cloud

## NVIDIA HPC SDK

| DEVELOPMENT | | | | | ANALYSIS | |
|---|---|---|---|---|---|---|
| **Programming Models** | **Compilers** | **Core Libraries** | **Math Libraries** | **Communication Libraries** | **Profilers** | **Debugger** |
| Standard C++ & Fortran | nvcc / nvc | libcu++ | cuBLAS / cuTENSOR | Open MPI | Nsight | cuda-gdb |
| OpenACC & OpenMP | nvc++ | Thrust | cuSPARSE / cuSOLVER | NVSHMEM | Systems | Host |
| CUDA | nvfortran | CUB | cuFFT / cuRAND | NCCL | Compute | Device |

Develop for the NVIDIA HPC Platform: GPU, CPU and Interconnect

HPC Libraries | GPU Accelerated C++ and Fortran | Directives | CUDA

7-8 Releases Per Year | Freely Available

# 加入 NVIDIA 開發者計畫

developer.nvidia.com/developer-program

請掃描左邊 QR code 加入 NVIDIA 開發者計畫，活動結束前完成註冊，即可換取精美小禮物。

請向 NVIDIA 工作人員出示主旨為「Your application for the program NVIDIA Developer Program is approv