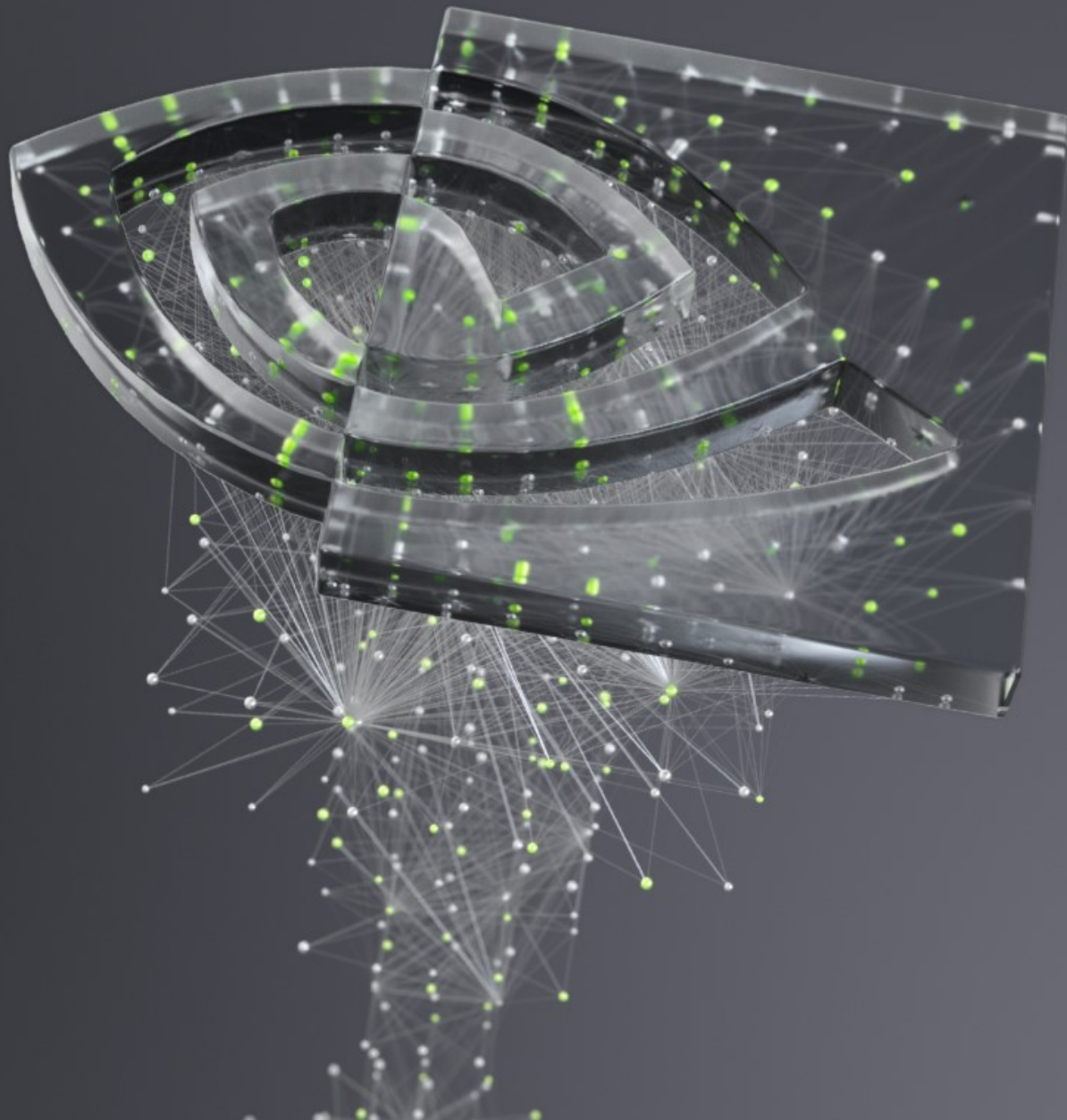# N-WAYS TO GPU COMPUTING

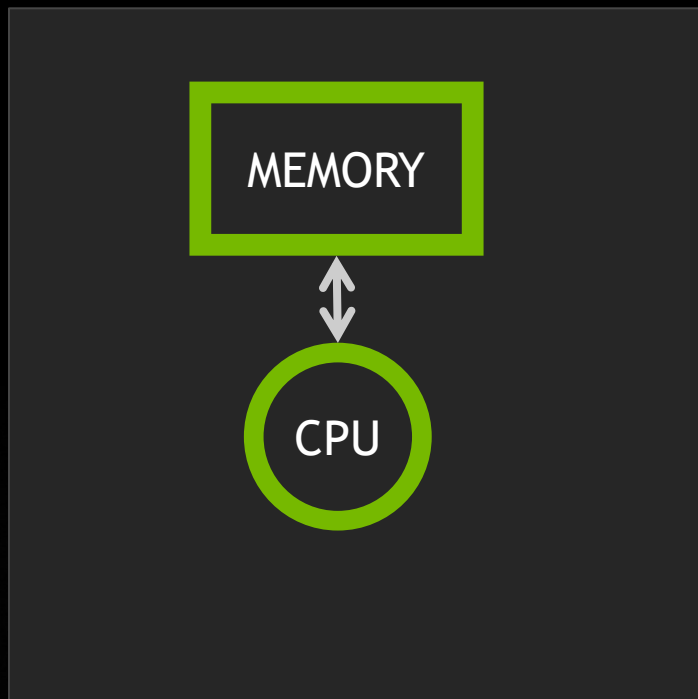# INTRODUCTION TO GPU COMPUTING

## What to expect?
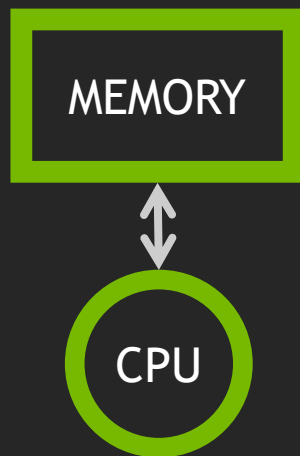
- Broad view on GPU Stack

- Fundamentals of GPU Architecture

- Good starting point
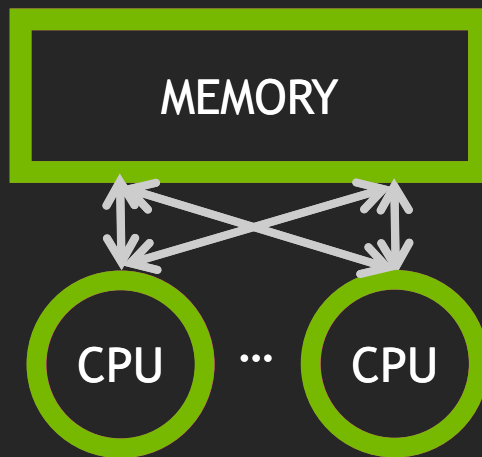
NVIDIA.

# HPC SYSTEM EVOLUTION



Sequential

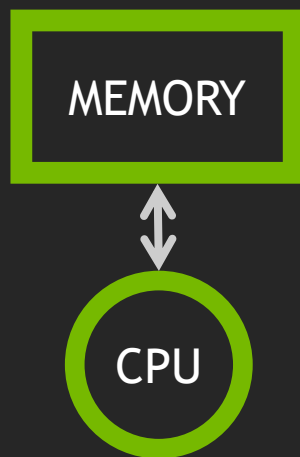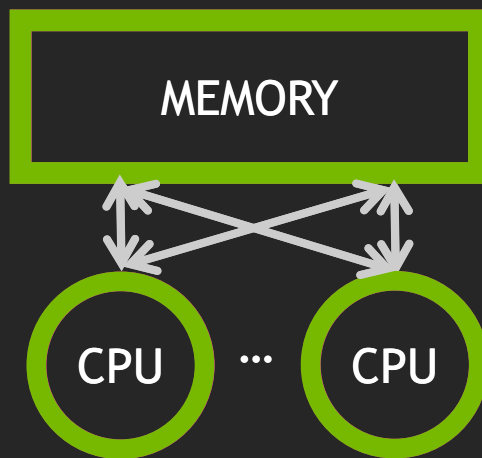3 NVIDIA.

# HPC SYSTEM EVOLUTION



Sequential

Multithreaded
P-Thread/OpenMP

# HPC SYSTEM EVOLUTION



Sequential

Multithreaded
P-Thread/OpenMP

Distributed
MPI

# GPU ARCHITECTURE CONTINUES TO DELIVER PERFORMANCE

APPLICATIONS

ALGORITHMS

SYSTEMS

CUDA

CPU — GPU

ARCHITECTURE

GPU-Computing perf
1.5X per year

1.1X per year

1.5X per year

Single-threaded perf

# ACCELERATED COMPUTING PILLARS



X-FACTOR SPEED UP

GPU

DPU

CPU

FULL STACK

CUDA

DATA-CENTER SCALE

Developer Productivity

DEVELOPERS ++

INSTALLED BASE ++

PERFORMANCE ++

CUDA
EVERYWHERE

NVIDIA.

# HIERARCHY OF SCALES



**Multi-System Rack**
Unlimited Scale

**Multi-GPU System**
8 GPUs

**Multi-SM GPU**
108 Multiprocessors

**Multi-Core SM**
2048 threads

# CUDA PLATFORM: TARGETS EACH LEVEL OF THE HIERARCHY

## The CUDA Platform Advances State Of The Art From Data Center To The GPU



### System Scope

FABRIC MANAGEMENT

DATA CENTER OPERATIONS

DEPLOYMENT

MONITORING

COMPATIBILITY

SECURITY

### Node Scope

GPU-DIRECT

NVLINK

LIBRARIES

UNIFIED MEMORY

ARM

MIG

### Program Scope

CUDA C++

OPENACC

STANDARD LANGUAGES

SYNCHRONIZATION

PRECISION

SCOPE OF THIS SESSION

# ACCELERATED PLATFORM

**BUSINESS APPLICATIONS**

| Customer Engagement | Patient Diagnostics | Fraud Detection | Quality Assurance | Industrial Automation | Precision Marketing | Molecular Simulations | ++ |

**SOFTWARE HUB**

HELM

*Certified Containers*

*Pre-trained Models*

*SDKs*

**APPLICATION FRAMEWORKS**

| SMART CITY *Metropolis* | CONVERSATIONAL AI *Jarvis* | AUTONOMOUS VEHICLES *Drive* | RECOMMENDATION SYSTEMS *Merlin* | HEALTHCARE *Clara* | ++ ... |

**DEVELOPER TOOLKITS**

| ML & DATA ANALYTICS RAPIDS  dmlc XGBoost | AI TRAINING & INFERENCE TensorFlow  TensorRT  PYTORCH  mxnet | HIGH PERFORMANCE COMPUTING *NVIDIA HPC SDK* | RENDERING & VISUALIZATION *IndeX*  *OptiX* |

**ACCELERATION LIBRARIES**

| COMPUTE *CUDA-X* | NETWORKING, STORAGE & SECURITY *DOCA*  *MAGNUM IO* |

**OPERATIONS**

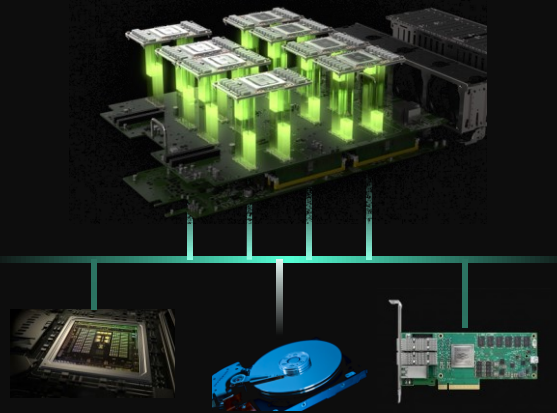TRITON INFERENCE SERVER     FLEET COMMAND

NVIDIA GPU Operator

Red Hat

vmware

**SERVERS & CLOUD**

DGX HGX
Purpose Built

EGX
Mainstream & Edge

CSP Instances

**VALIDATED SOLUTIONS**

**MONITORING**

DCGM  Grafana

Prometheus  UFM

**HARDWARE TECHNOLOGIES**

GPU     NVSwitch     BlueField DPU     SMART NIC     Mellanox Switch

NVIDIA

# HOW GPU ACCELERATION WORKS



**GPU**

5% of Code

Rest of Sequential
CPU Code

**CPU**

+

NVIDIA.

# ACCELERATED COMPUTING

**CPU**
Optimized for
Serial Tasks

**GPU Accelerator**
Optimized for
Parallel Tasks

NVIDIA.

# SILICON BUDGET

- The three components of any processor

| Less | **ALU** | More |
| More | **Control** | Less |
| More | **Cache** | Less |

NVIDIA.

# CPU IS A LATENCY REDUCING ARCHITECTURE

**CPU**
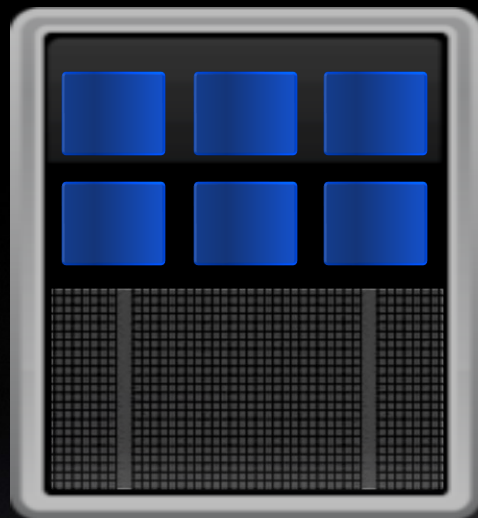Optimized for
Serial Tasks



## CPU Strengths

- Very large main memory
- Very fast clock speeds
- Latency optimized via large caches
- Small number of threads can run very quickly

## CPU Weaknesses

- Relatively low memory bandwidth
- Cache misses very costly
- Low performance/watt

NVIDIA.

# GPU IS ALL ABOUT HIDING LATENCY

## GPU Strengths

- High bandwidth main memory
- Significantly more compute resources
- Latency tolerant via parallelism
- High throughput
- High performance/watt

## GPU Weaknesses

- Relatively low memory capacity
- Low per-thread performance

## GPU Accelerator
### Optimized for Parallel Tasks

# LOW LATENCY VS HIGH THROUGHPUT

- CPU architecture must minimize latency within each thread

- GPU architecture hides latency with computation (data-parallelism, to 30k threads!)

**GPU – High Throughput Processor**

T₄ T₃ T₂ T₁

**CPU core – Low Latency Processor**

T₁ T₂ T₃ T₄

**Computation Thread**

Tₙ — Processing

Waiting for data

Ready to be processed

NVIDIA.

# SPEED V. THROUGHPUT

## Speed

## Throughput



Which is better depends on your needs...

# HUGE BREADTH OF PLATFORMS, SYSTEMS, LANGUAGES

Data Analytics

AI Training

Scientific Computing

Video Analytics

Cloud Gaming

Machine Learning

Graphics

5G Networks

Genomics

AI Inference

# NVIDIA HPC SDK

Download at developer.nvidia.com/hpc-sdk

## NVIDIA HPC SDK

### DEVELOPMENT

| Compilers | Math Libraries | Communication Libraries | Programming Models |
|---|---|---|---|
| nvcc    nvc | cuBLAS    cuTENSOR | Open MPI | Standard C++ & Fortran |
| nvc++ | cuSPARSE    cuSOLVER | NVSHMEM | OpenACC & OpenMP |
| nvfortran | cuFFT    cuRAND | NCCL | CUDA |

### ANALYSIS

| Profilers | Debugger |
|---|---|
| Nsight | cuda-gdb |
| Systems | Host |
| Compute | Device |

### DEPLOYMENT

| Container |
|---|
| HPC Container Maker / NVIDIA Container Runtime |

Develop for the NVIDIA HPC Platform: GPU, CPU and Interconnect
HPC Libraries | GPU Accelerated C++ and Fortran | Directives | CUDA

NVIDIA.

# N-WAYS TO GPU PROGRAMMING

## 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
```

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

...

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

...

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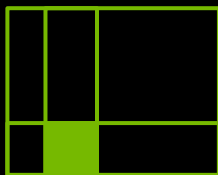
```cpp
__global__
void saxpy(int n, float a,
           float *x, float *y) {
    int i = blockIdx.x*blockDim.x +
            threadIdx.x;
    if (i < n) y[i] += a*x[i];
}

int main(void) {
    cudaMallocManaged(&x, ...);
    cudaMallocManaged(&y, ...);
    ...
    saxpy<<<(N+255)/256,256>>>(...,x, y)
    cudaDeviceSynchronize();
    ...
}
```

**GPU Accelerated
C++ and Fortran**

**Incremental Performance
Optimization with Directives**

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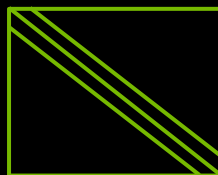CUDA C++/Fortran**

**GPU Accelerated Math Libraries**

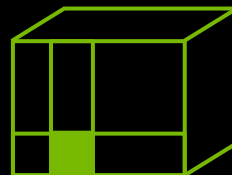NVIDIA.

# GPU ACCELERATED MATH LIBRARIES

## cuBLAS
BF16, TF32 and FP64 Tensor Cores

## cuSPARSE
Increased memory BW, Shared Memory & L2

## cuTENSOR
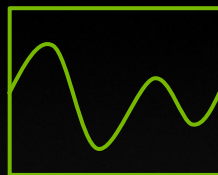BF16, TF32 and FP64 Tensor Cores

## cuSOLVER
BF16, TF32 and FP64 Tensor Cores

## nvJPEG
Hardware Decoder

## cuFFT
BF16, TF32 and FP64 Tensor Cores

## CUDA Math API
Increased memory BW, Shared Memory & L2

## CUTLASS
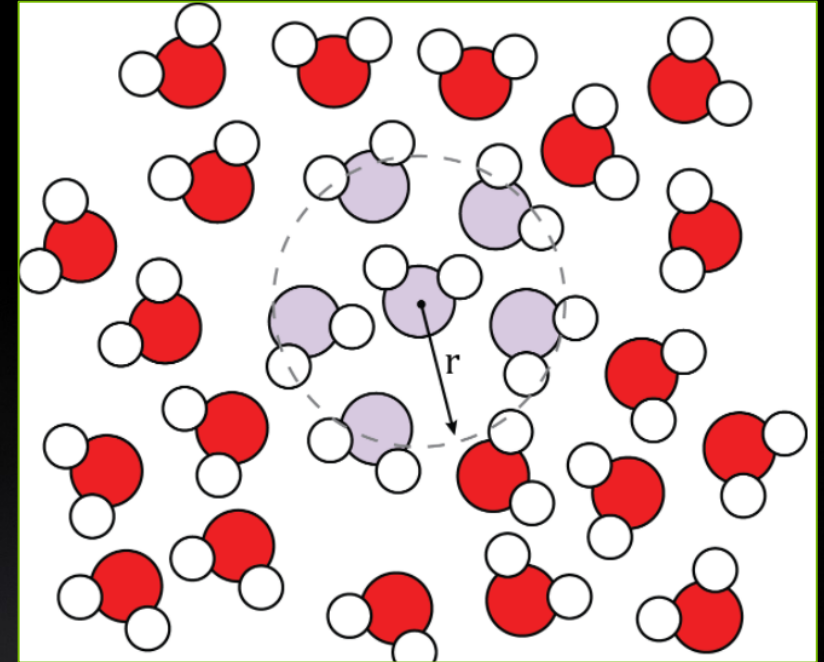BF16 & TF32 Support

NVIDIA.

# APPLICATION
## Molecular Simulation

### RDF

The radial distribution function (RDF) denoted in equations by g(r) defines the probability of finding a particle at a distance r from another tagged particle.

NVIDIA.

# RDF
## Pseudo Code - C

```c
for (int frame=0;frame<nconf;frame++){

    for(int id1=0;id1<numatm;id1++)
    {
        for(int id2=0;id2<numatm;id2++)
        {
            dx=d_x[]-d_x[];
            dy=d_y[]-d_y[];
            dz=d_z[]-d_z[];
            r=sqrtf(dx*dx+dy*dy+dz*dz);

            if (r<cut) {
                ig2=(int)(r/del);
                d_g2[ig2] = d_g2[ig2] +1 ;
            }
        }
    }
}
```

▹ Across Frames

▹ Find Distance

▹ Reduction

NVIDIA

# RDF
## Pseudo Code - Fortran

```fortran
do iconf=1,nframes
   if (mod(iconf,1).eq.0) print*,iconf

   do i=1,natoms
      do j=1,natoms
           dx=x(iconf,i)-x(iconf,j)
           dy=y(iconf,i)-y(iconf,j)
           dz=z(iconf,i)-z(iconf,j)

           r=dsqrt(dx**2+dy**2+dz**2)
           if(r<cut)then
                g(ind)=g(ind)+1.0d0
           endif
      enddo
   enddo
enddo
```

▷ **Across Frames**

▷ **Find Distance**

▷ **Reduction**

NVIDIA.

THANK YOU

BACKUP