# NVIDIA'S SOFTWARE ECOSYSTEM FOR HPC
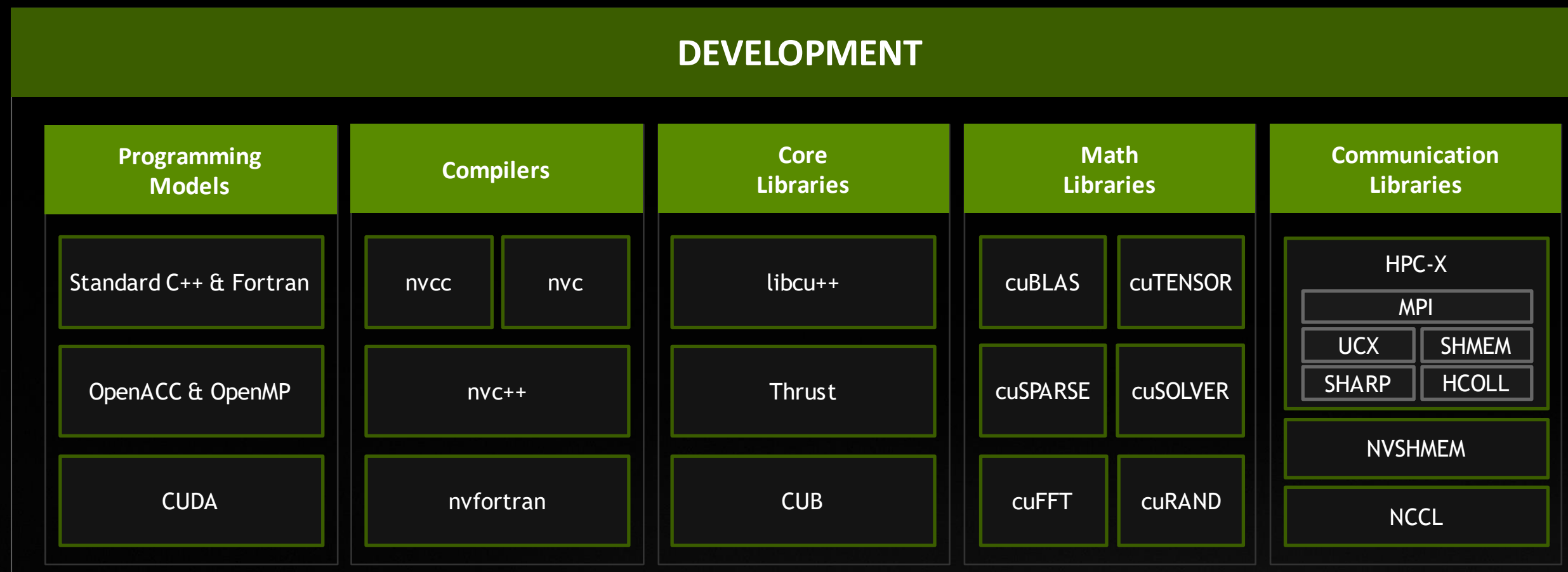
Max Katz

May 5, 2021
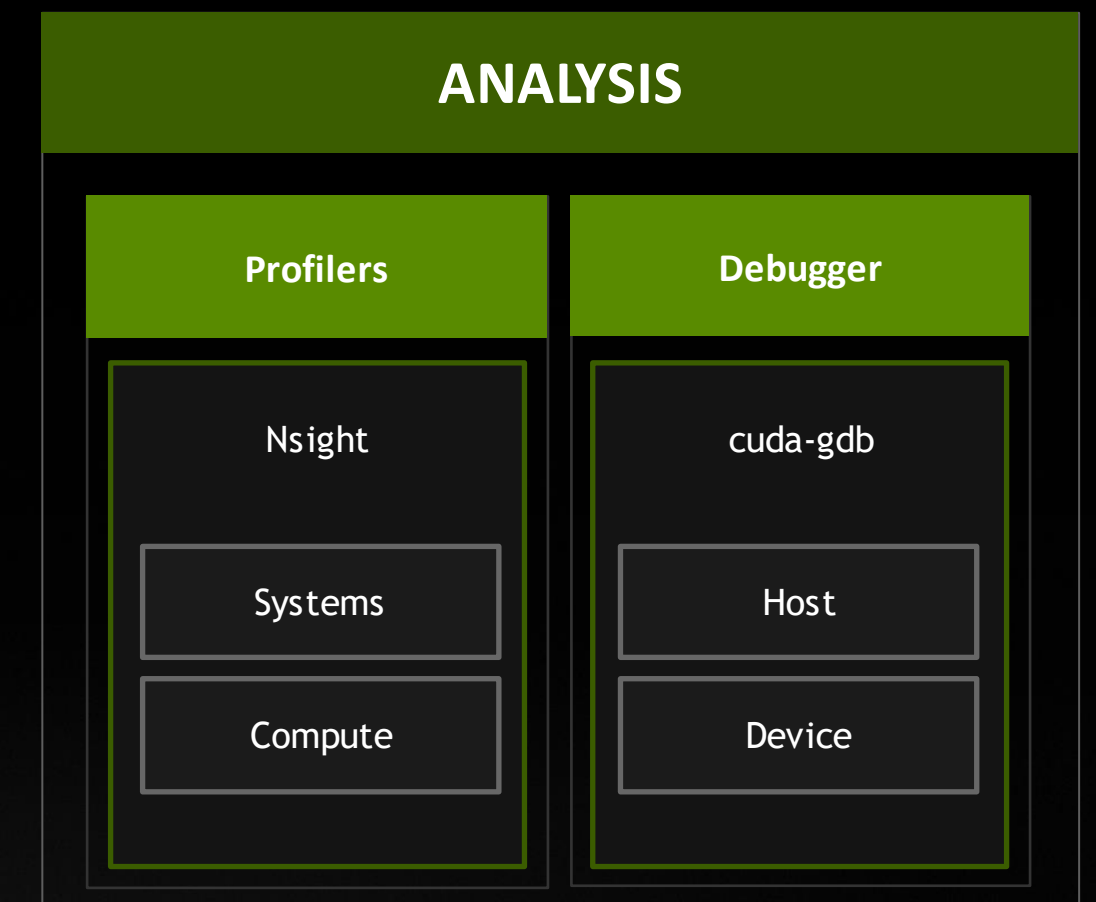
# NVIDIA HPC SDK

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

Develop for the NVIDIA Platform: GPU, CPU and Interconnect
Libraries | Accelerated C++ and Fortran | Directives | CUDA
7-8 Releases Per Year | Freely Available

# PROGRAMMING THE NVIDIA PLATFORM

## CPU, GPU and Network

### Accelerated Standard Languages

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


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


import legate.numpy as np
…
def saxpy(a, x, y):
    y[:] += a*x
```

### Incremental Portable Optimization

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

...

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

...

}
```

### Platform Specialization

```
__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, ...);
```

| Core | Math | Communication | Data Analytics | AI |
|------|------|---------------|----------------|-----|

Acceleration Libraries

# ACCELERATED STANDARDS

Parallel performance for wherever you code needs to run

```
std::transform(std::execution::par, x, x+n, y, y,
   [=] (auto xi, auto yi) { return y + a*xi; });
```

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

CPU

GPU

nvc++ –stdpar=multicore
nvfortran -stdpar=multicore

nvc++ –stdpar=gpu
nvfortran -stdpar=gpu

# ACCELERATED 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
...
!$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)
```
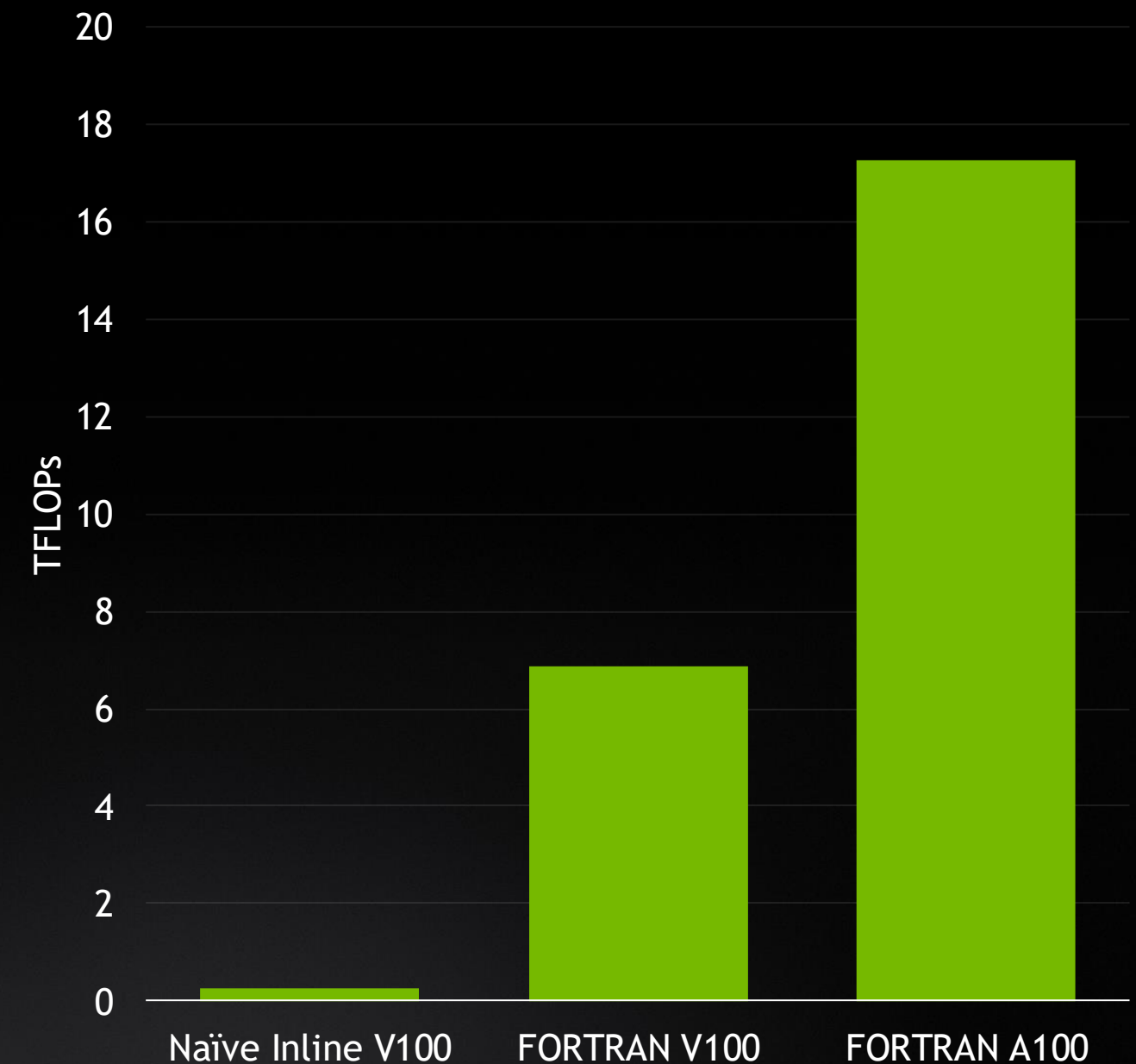
```fortran
real(8), dimension(ni,nk) :: a
real(8), dimension(nk,nj) :: b
real(8), dimension(ni,nj) :: c

...

do nt = 1, ntimes
  d = c + matmul(a,b)
end do
```

**Inline FP64 matrix multiply**

**MATMUL FP64 matrix multiply**



TFLOPs

| | |
|---|---|
| Naïve Inline V100 | FORTRAN V100 | FORTRAN A100 |

# HPC PROGRAMMING IN ISO FORTRAN

## Examples of Patterns Accelerated in NVFORTRAN
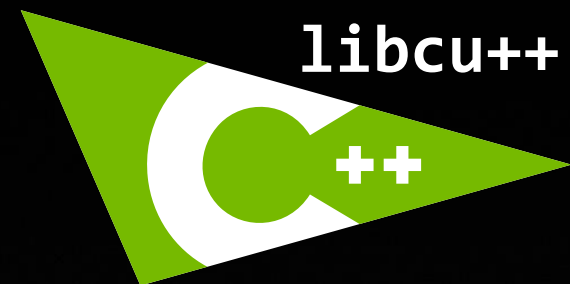
```
d = 2.5 * ceil(transpose(a)) + 3.0 * abs(transpose(b))
d = 2.5 * ceil(transpose(a)) + 3.0 * abs(b)
d = reshape(a,shape=[ni,nj,nk])
d = reshape(a,shape=[ni,nk,nj])
d = 2.5 * sqrt(reshape(a,shape=[ni,nk,nj],order=[1,3,2]))
d = alpha * conjg(reshape(a,shape=[ni,nk,nj],order=[1,3,2]))
d = reshape(a,shape=[ni,nk,nj],order=[1,3,2])
d = reshape(a,shape=[nk,ni,nj],order=[2,3,1])
d = reshape(a,shape=[ni*nj,nk])
d = reshape(a,shape=[nk,ni*nj],order=[2,1])
d = reshape(a,shape=[64,2,16,16,64],order=[5,2,3,4,1])
d = abs(reshape(a,shape=[64,2,16,16,64],order=[5,2,3,4,1]))
c = matmul(a,b)
c = matmul(transpose(a),b)
c = matmul(reshape(a,shape=[m,k],order=[2,1]),b)
c = matmul(a,transpose(b))
c = matmul(a,reshape(b,shape=[k,n],order=[2,1]))
```

```
c = matmul(transpose(a),transpose(b))
c = matmul(transpose(a),reshape(b,shape=[k,n],order=[2,1]))
d = spread(a,dim=3,ncopies=nk)
d = spread(a,dim=1,ncopies=ni)
d = spread(a,dim=2,ncopies=nx)
d = alpha * abs(spread(a,dim=2,ncopies=nx))
d = alpha * spread(a,dim=2,ncopies=nx)
d = abs(spread(a,dim=2,ncopies=nx))
d = transpose(a)
d = alpha * transpose(a)
d = alpha * ceil(transpose(a))
d = alpha * conjg(transpose(a))
c = c + matmul(a,b)
c = c - matmul(a,b)
c = c + alpha * matmul(a,b)
d = alpha * matmul(a,b) + c
d = alpha * matmul(a,b) + beta * c
```
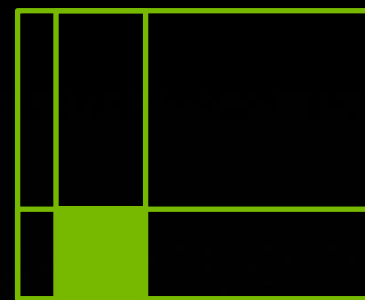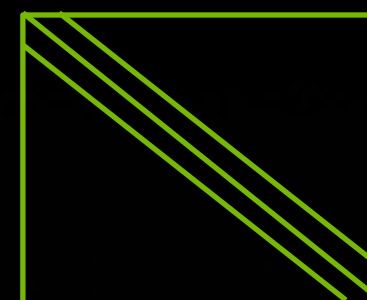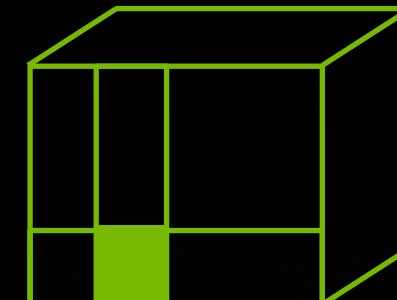
# NVIDIA PERFORMANCE LIBRARIES

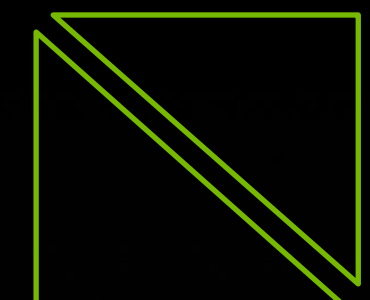## Core and Math

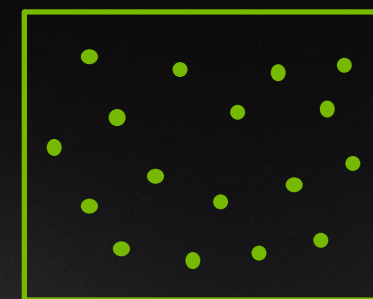| Core | Math |
|---|---|

**Core**
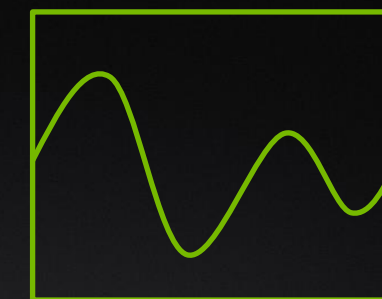
libcu++

Thrust

CUB

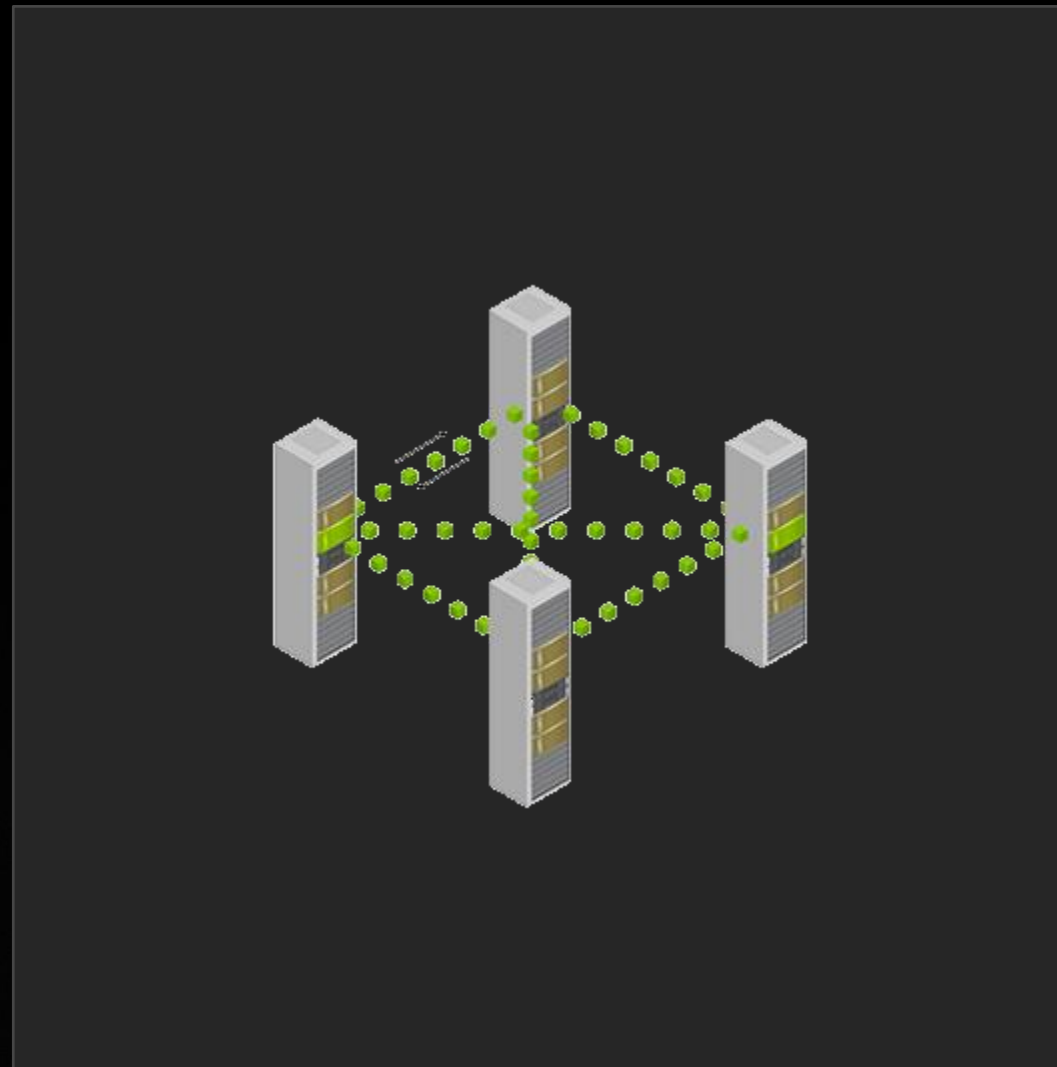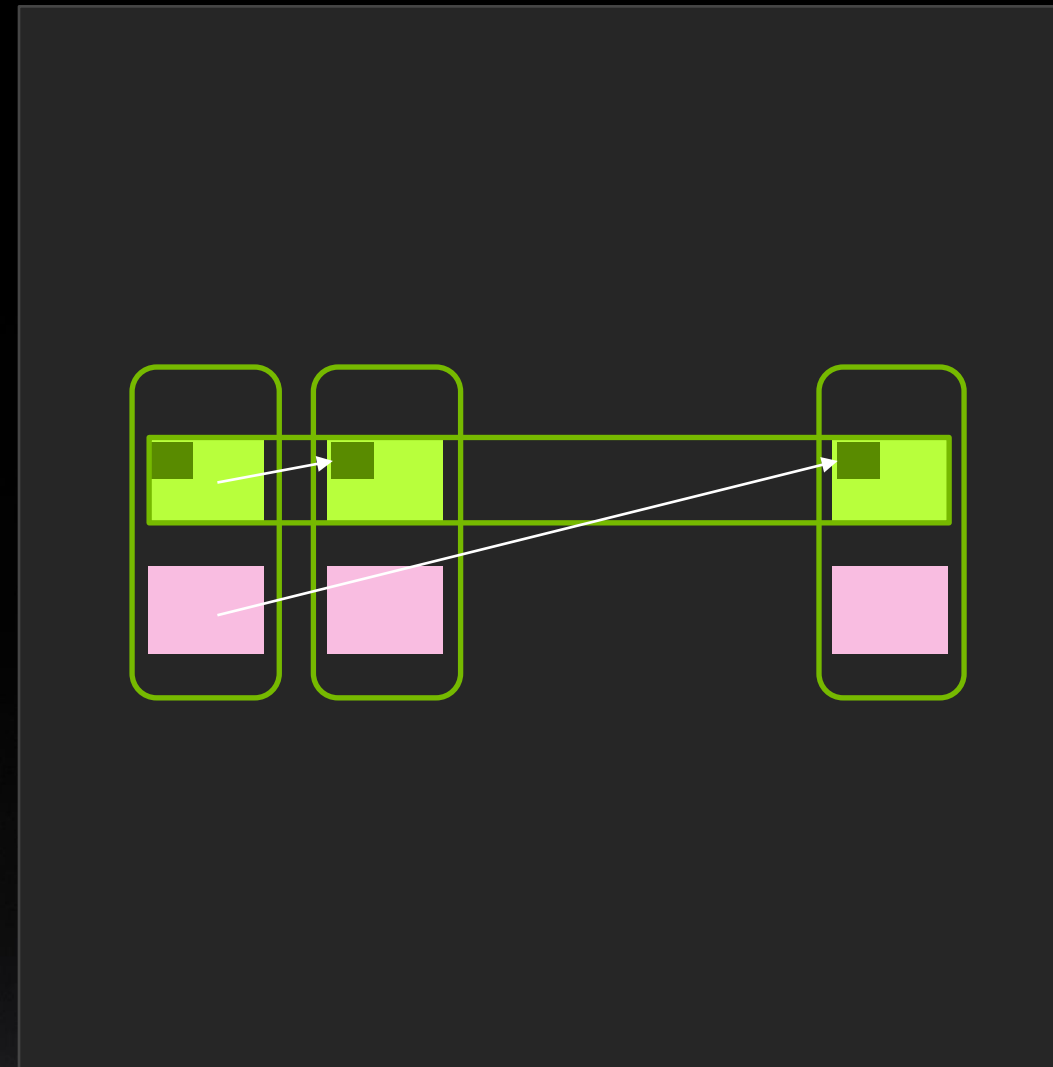**Math**

cuBLAS

cuSPARSE

cuTENSOR

cuSOLVER

cuRAND

cuFFT

CUDA Math API

# NVIDIA COMMUNICATION LIBRARIES



**HPC-X**
Optimized whole-system communications

**NVSHMEM**
Low-latency PGAS programming

**NCCL**
Multi-node collectives for accelerators

# INTRODUCING LEGATE
## Accelerated and Distributed

A framework for programming large numbers of GPUs as if they were a single processor

Pass data between Legate libraries without worrying about distribution or synchronization requirements

Legate NumPy and Pandas aim to transparently scale existing Numpy and Pandas workloads

Legate Numpy and Legate Pandas available now and opensource!

```python
import legate.numpy as np

def cg_solve(A, b, tol=1e-10):
    x = np.zeros(A.shape[1])
    r = b - A.dot(x)
    p = r
    rsold = r.dot(r)
    for i in xrange(b.shape[0]):
        Ap = A.dot(p)
        alpha = rsold / (p.dot(Ap))
        x = x + alpha * p
        r = r - alpha * Ap
        rsnew = r.dot(r)
        if np.sqrt(rsnew) < tol:
            break
        beta = rsnew / rsold
        p = r + beta * p
        rsold = rsnew
    return x
```

# LEGATE NUMPY

## Results from "CFD Python"
https://github.com/barbagroup/CFDPython

```python
import legate.numpy as np

for _ in range(iter):
    un = u.copy()

    vn = v.copy()
    b = build_up_b(rho, dt, dx, dy, u, v)
    p = pressure_poisson_periodic(b, nit, p, dx, dy)

    u[1:-1, 1:-1] = (
        un[1:-1, 1:-1]
        - un[1:-1, 1:-1] * dt / dx * (un[1:-1, 1:-1] - un[1:-1, 0:-2])
        - vn[1:-1, 1:-1] * dt / dy * (un[1:-1, 1:-1] - un[0:-2, 1:-1])
        - dt / (2 * rho * dx) * (p[1:-1, 2:] - p[1:-1, 0:-2])
        + nu
        * (
            dt
            / dx ** 2
            * (un[1:-1, 2:] - 2 * un[1:-1, 1:-1] + un[1:-1, 0:-2])
            + dt
            / dy ** 2
            * (un[2:, 1:-1] - 2 * un[1:-1, 1:-1] + un[0:-2, 1:-1])
        )
        + F * dt
    )
```

## Distributed NumPy Performance
### (weak scaling)

Extracted from "CFD Python" course at https://github.com/barbagroup/CFDPython
Barba, Lorena A., and Forsyth, Gilbert F. (2018). CFD Python: the 12 steps to Navier-Stokes equations.
*Journal of Open Source Education*, **1**(9), 21, https://doi.org/10.21105/jose.00021

# LEGATE PANDAS

## Pandas join micro-benchmark – 300M rows/GPU

```python
import legate.numpy as np
import legate.pandas as pd

size = num_rows_per_gpu * num_gpus

key_l = np.arange(size)
payload_l = np.random.randn(size) * 100.0
lhs = pd.DataFrame({"key": key_l, "payload": payload_l})

key_r = key_l // 3 * 3    # selectivity: 0.33
payload_r = np.random.randn(size) * 100.0
rhs = pd.DataFrame({"key": key_r, "payload": payload_r})

out = lhs.merge(rhs, on="key")
```
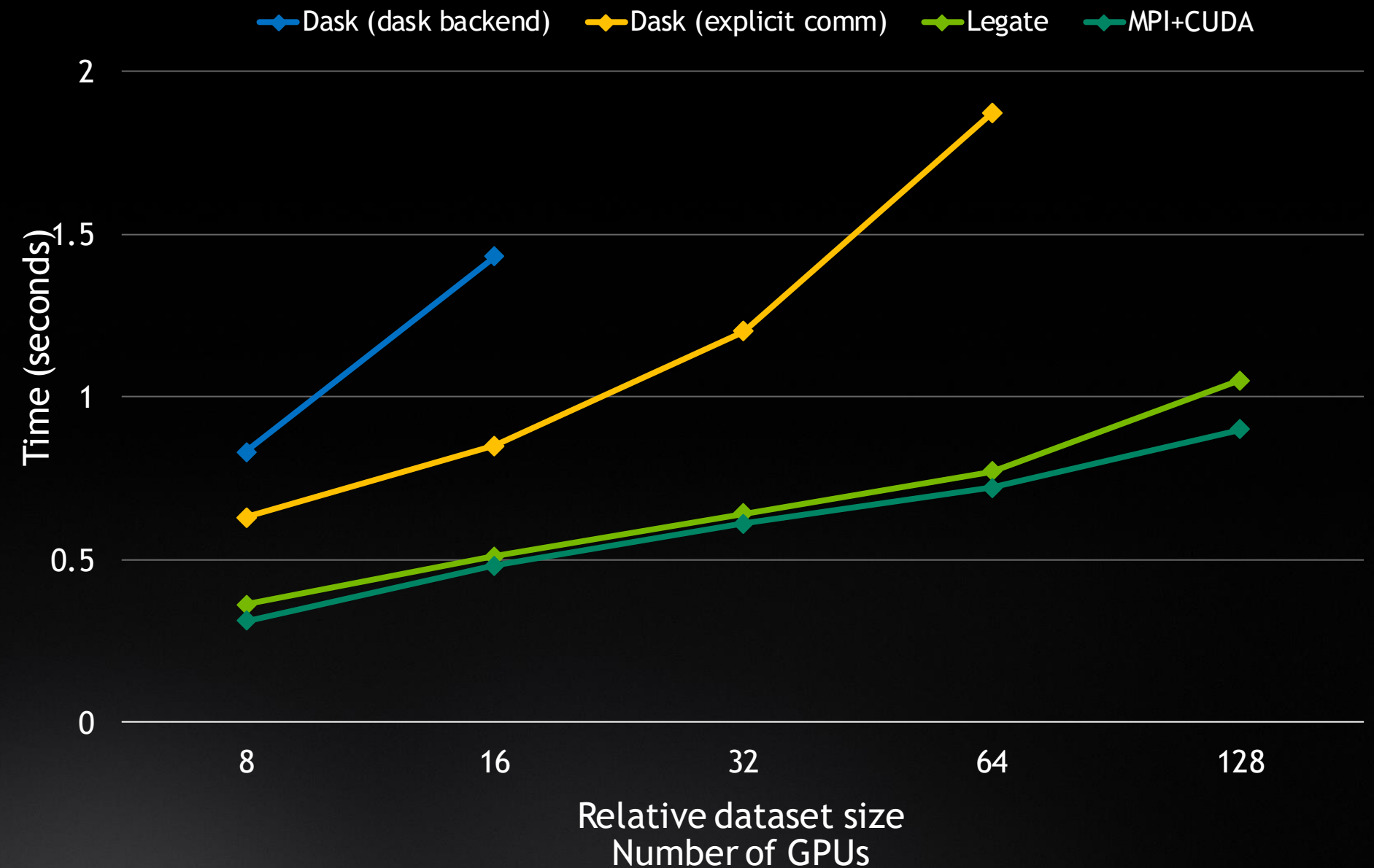
### Distributed Pandas Performance (weak scaling)

# NSIGHT PRODUCT FAMILY

Nsight Systems - Analyze application performance (system-wide)

Nsight Compute - Optimize CUDA kernel

Other developer tools:

Compute Sanitizer: memory debugging similar to valgrind for GPUs

cuda-gdb: CUDA-aware extension of gdb

Start here

**Nsight Systems**
Comprehensive workload-level performance

Re-check overall performance

Re-check overall performance

Dive into top CUDA kernels by using metrics/counter collection

Dive into graphics frames

**Nsight Compute**
Detailed CUDA kernel performance

**Nsight Graphics**
Detailed frame/render performance

NVIDIA Nsight Systems 2018.0.0

File   View   Help

Select device for profiling...        More info....

aggcpy-yeast-CUDA.qdrep

Timeline View        1x        1 warning, 5 messages

3s +637ms    +637.2ms    +637.4ms    +637.6ms    +637.8ms    +638ms    +638.2ms    +638.4ms    +638.6ms    +638.8ms    3s 638.962ms    +639.2ms    +639.4ms

System
MPD_TIMEST...  MCLKR Barrier [194.213 µs]   MPD_TIMESTEP [524.355 µs]   MCLKR Barri...   MPD_TIMESTEP [475.814 µs]   MCLKR Barrier [155....   MPD_TIMESTEP [466.524 µs]   MCLKR Barrier [167.96...   MPD_TIMESTEP [446.305 µs]

NVTX
MPD_S...   MPD_...   MPD_...   MPD_...

CUDA API
cudaSt...  MGPU_x_k...  MGPU_y_kernel  MGPU_z_ke...  MGPU_pre...  cuda...  cuda...  MGPU_...  MGPU_y_ke...  MGPU_z_ke...  MGPU_pre...  cud...  cuda...  MGPU_y_kernel  MGPU_z_ke...  MGPU_pre...  cudaM...  cuda...  MGP...  MGPU_prec...  cuda

[1999] lm

System
pt...   pthread_...   pthr...   threa...

NVTX
MPD_TIMEST...  MCLKR Barrier [193.247 µs]   MPD_TIMESTEP [452.747 µs]   MCLKR Barrier [174.799...]   MPD_TIMESTEP [526.627 µs]   MCLKR Barrie...   MPD_TIMESTEP [540.467 µs]   ...846 µs]

Waiting
Duration: 36.206 µs

Call stack at 3.639s:
libpthread-2.23.so:_pthread_mutex_lock
libcuda.so.384.125:0x7fd5683568d6
libcuda.so.384.125:0x7fd5683568825
libcuda.so.384.125:0x7fd5688268825
libcuda.so.384.125:0x7fd568269958
libcuda.so.384.125:cuMemcpyAsync
libcudart.so.9.0.176:0x7fd5749664fd
libcudart.so.9.0.176:0x7fd574943573
libcudart.so.9.0.176:cudaMemcpyAsync
lm!lm::rdme::MGPUMpdRdmeSolver::run_next_timestep(...)
lm!lm::rdme::MGPUMpdRdmeSolver::run_thread(...)
libpthread-2.23.so:start_thread
libc-2.23.so:__clone

MPD_S...   MPD_...

CUDA API
cud...  cudaSt...  MGPU_y_kernel  MGPU_z_ker...  MGPU_pre...  cud...  cud...  MG...  MGPU_y_ker...  MGPU_z_ke...  MGPU_prec...  cudaMemcpyAs...  cuda...  GPU...  MGPU_y_kernel  MGPU_z_ke...  MGPU_pre...  cud...  Memc...

[2002] lm

System
pthrea...   pt...   ...   p...  pt...   pthrea...   pth...

NVTX
MPD_TIMESTEP [506.5...   MCLKR Barrier [...   MPD_TIMESTEP [573.056 µs]   MCLK...   MPD_TIMESTEP [555.258 µs]   MCLKR B...   MPD_TIMESTEP [578.180 µs]   ...203 µs]

MPD_SYN...   MPD_SYN...   MPD_SYN...

CUDA API
cudaMem...  cudaStre...  MGPU_x...  MGPU_y_kernel  MGPU_z_ke...  MGPU_prec...  cudaMem...  cudaStre...  MGP...  MGPU_y_ke...  MGPU_z_ke...  MGPU_pre...  cudaMemcpy...  cudaStre...  MG...  MGPU_y_kernel  MGPU_z_ke...  MGPU_pre...  cudaMe

[1998] lm

System
pthrea...   ...   pt...   pt...   pt...

NVTX
MPD_TIMESTEP [511.14...   MCLKR Barrier [...   MPD_TIMESTEP [489.153 µs]   MCLKR Barrier [13...   MPD_TIMESTEP [466.427 µs]   MCLKR Barrier [164.4...   MPD_TIMESTEP [583.928 µs]   MCLK...   MPD_TIMESTEP [466.209 µs]

MPD_SY...   MPD_SY...   MPD_...   MPD_SYN...

CUDA API
cudaMemc...  cudaStre...  MG...  MGPU_y_kernel  MGPU_z_ke...  MGPU_pre...  cu...  cudaStr...  MGPU_y_ker...  MGPU_z_ke...  MGPU_pre...  cu...  cudaM...  cud...  MGPU_x...  MGPU_y_kernel  MGPU_z_ke...  MGPU_pre...  cudaMemcp...  cudaStre...  MGPU_prec...  cudaMemc

[1987] lm

11 threads hidden...

CUDA (Tesla V100-SXM2-32GB)
Stream 70
Memory
Kernels
MGPU...  MGPU_p...   MGPU_x_kernel_unpack   MGPU_y_k...  MGPU_z_k...  MGPU_p...   MGPU_x_kernel_unpack   MGPU_y_k...  MGPU_z_k...  MGPU_pr...   MGPU_x_kernel_unpack   MGPU_y_k...  MGPU_z_k...  MGPU_p...   MGPU_x_kernel_unpack   M...

1 stream(s) hidden...

CUDA (Tesla V100-SXM2-32GB)
CUDA (Tesla V100-SXM2-32GB)
CUDA (Tesla V100-SXM2-32GB)
CUDA (Tesla V100-SXM2-32GB)
CUDA (Tesla V100-SXM2-32GB)
CUDA (Tesla V100-SXM2-32GB)
CUDA (Tesla V100-SXM2-32GB)

Bottom-Up View        Process [1979] lm (10 of 21 threads)

Filter...    99.01% of data is shown due to applied filters.        Search...

| Symbol Name | Self, % | Module Name |
| --- | --- | --- |
| lm::rdme::MGPUMpdRdmeSolver::run_thread(int) | 38.59 | /opt/lm/bin/lm |
| 0x7fd568442726 | 8.26 | /usr/lib/x86_64-linux-gnu/libcuda.so.384.125 |
| 0x7fd568373448 | 2.97 | /usr/lib/x86_64-linux-gnu/libcuda.so.384.125 |
| 0x7fd568373436 | 2.25 | /usr/lib/x86_64-linux-gnu/libcuda.so.384.125 |

# NVTX
## NVIDIA Tools Extension for Profiling

NVTX can be used to manually instrument applications, for example in C:

```
nvtxRangePush("region_name");

nvtxRangePop();
```

These names are then shown on the nsys timeline (can also be used with ncu)

Headers provided with CUDA toolkit for C/C++; can also be used with Fortran, generic Python (e.g. provided by CuPy), TensorFlow, and PyTorch

# COLLECT A PROFILE WITH NSIGHT SYSTEMS

```
$ nsys profile -o report --stats=true ./myapp.exe
```

Generated file: `report.qdrep;` open for viewing in the Nsight Systems UI

Can be done inside a container if the container has `nsys`:

```
$ mpirun -n 4 singularity run --nv -B $CONTAINER_IMG nsys profile python train.py
```

# NSIGHT COMPUTE



CUDA kernel profiler

Targeted metric sections for various performance aspects

Customizable data collection and presentation (tables, charts, ...)

UI and Command Line

Python-based rules for guided analysis (or post-processing)

# NSIGHT COMPUTE



Source/PTX/SASS analysis and correlation

Source metrics per instruction and aggregated (e.g. PC sampling data)

Metric heatmap

# KERNEL PROFILES WITH NSIGHT COMPUTE

```
$ ncu –k mykernel –o report ./myapp.exe
```

Generated file: `report.ncu-rep;` open for viewing in the Nsight Compute UI

(Without the -k option, Nsight Compute will profile everything and take a long time!)