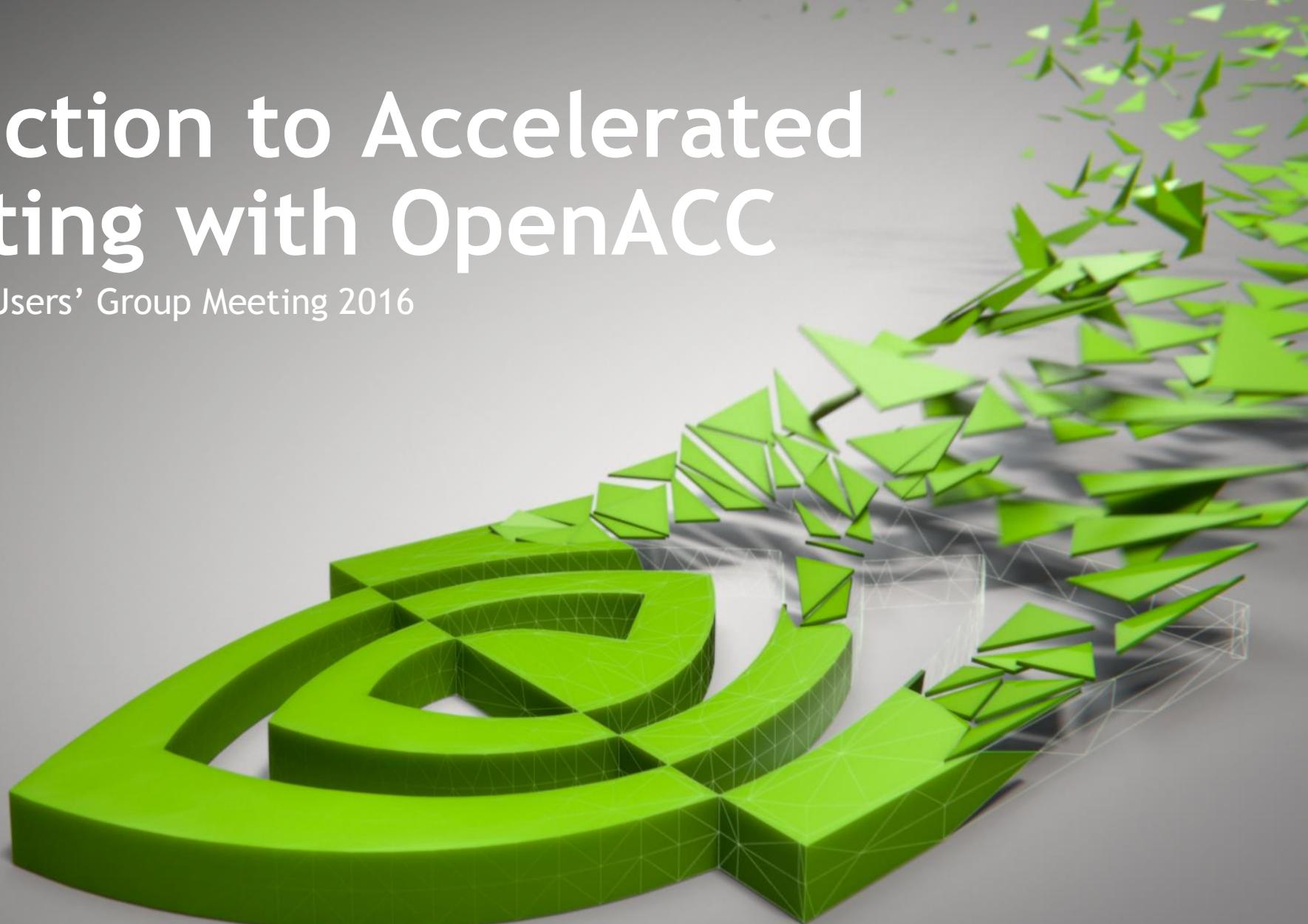


Introduction to Accelerated Computing with OpenACC

Jeff Larkin, OLCF Users' Group Meeting 2016



AGENDA

NVIDIA Introduction

Introduction to Accelerated Computing

Accelerated Computing with OpenACC

Next Steps

NVIDIA Introduction



GAMING



PRO VISUALIZATION



DATA CENTER



AUTO

THE WORLD LEADER IN VISUAL COMPUTING

Tesla Accelerates Discoveries

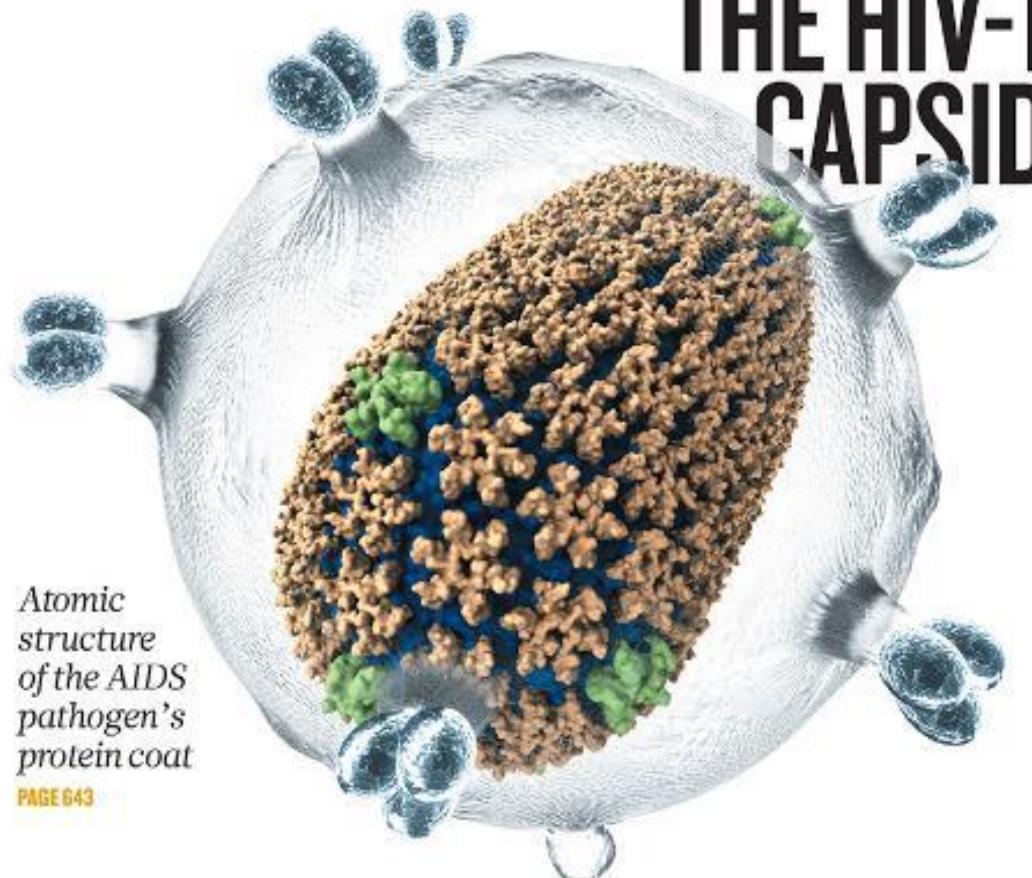
Using a supercomputer powered by the Tesla Platform with over 3,000 Tesla accelerators, University of Illinois scientists performed the first all-atom simulation of the HIV virus and discovered the chemical structure of its capsid — “the perfect target for fighting the infection.”

Without GPU, the supercomputer would need to be 5x larger for similar performance.

nature

THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE

THE HIV-1 CAPSID



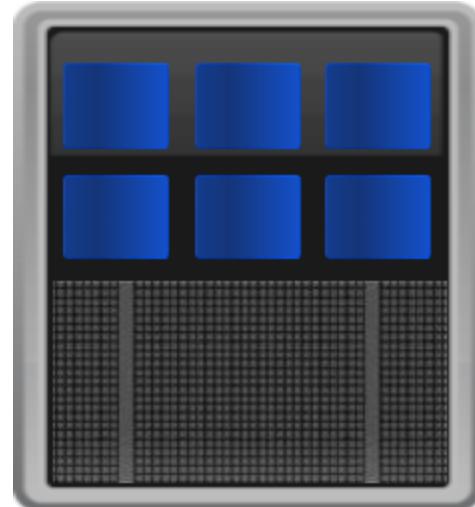
Introduction to Accelerated Computing

ACCELERATED COMPUTING

10X PERFORMANCE & 5X ENERGY EFFICIENCY FOR HPC

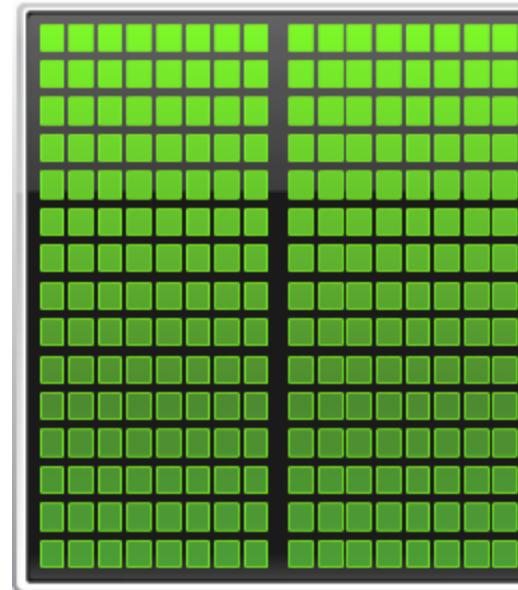
CPU

Optimized for
Serial Tasks



GPU Accelerator

Optimized for
Parallel Tasks



ACCELERATED COMPUTING

10X PERFORMANCE & 5X ENERGY EFFICIENCY FOR HPC



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

ACCELERATED COMPUTING

10X PERFORMANCE & 5X ENERGY EFFICIENCY FOR HPC

GPU Strengths

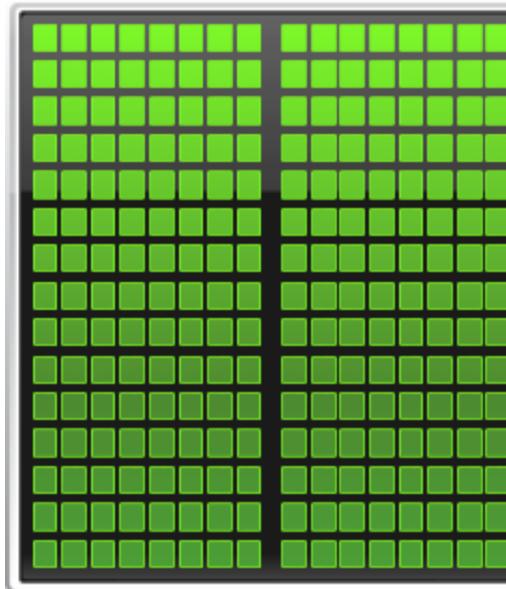
- High bandwidth main memory
- Latency tolerant via parallelism
- Significantly more compute resources
- High throughput
- High performance/watt

GPU Weaknesses

- Relatively low memory capacity
- Low per-thread performance

GPU Accelerator

Optimized for
Parallel Tasks



Speed v. Throughput

Speed

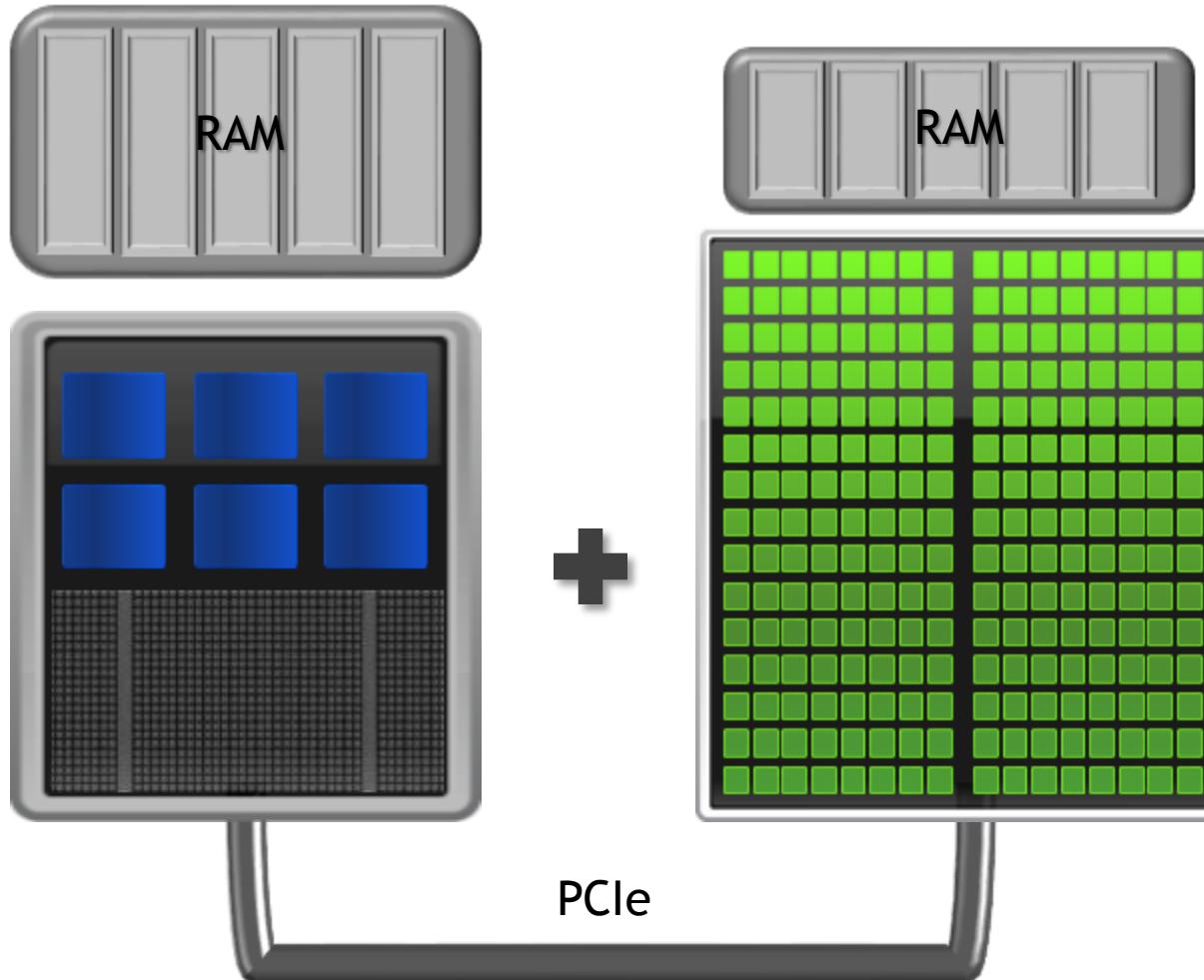


Throughput



Which is better depends on your needs...

Accelerator Nodes



CPU and GPU have distinct memories

- CPU generally larger and slower
- GPU generally smaller and faster

Execution begins on the CPU

- Data and computation are offloaded to the GPU

CPU and GPU communicate via PCIe

- Data must be copied between these memories over PCIe
- PCIe Bandwidth is much lower than either memories

Simplicity & Performance

Simplicity



Accelerated Libraries

Little or no code change for standard libraries; high performance

Limited by what libraries are available

Compiler Directives

High Level: Based on existing languages; simple and familiar

High Level: Performance may not be optimal

Parallel Language Extensions

Expose low-level details for maximum performance

Often more difficult to learn and more time consuming to implement

Performance

Code for Simplicity & Performance

Libraries

- Implement as much as possible using portable libraries.

Directives

- Use directives to rapidly accelerate your code.

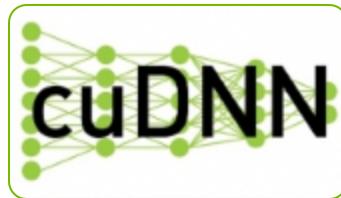
Languages

- Use lower level languages for important kernels.

GPU Accelerated Libraries

“Drop-in” Acceleration for Your Applications

Domain-specific
Deep Learning, GIS, EDA,
Bioinformatics, Fluids



Visual Processing
Image & Video



Linear Algebra
Dense, Sparse, Matrix



Math Algorithms
AMG, Templates, Solvers



OpenACC

Simple | Powerful | Portable

Fueling the Next Wave of
Scientific Discoveries in HPC

```
main()
{
    <serial code>
    #pragma acc kernels
    //automatically runs on GPU
    {
        <parallel code>
    }
}
```

University of Illinois
PowerGrid- MRI Reconstruction



70x Speed-Up
2 Days of Effort

RIKEN Japan
NICAM- Climate Modeling



7-8x Speed-Up
5% of Code Modified

8000+
Developers
using OpenACC

Common Programming Models Across Multiple CPUs

Libraries

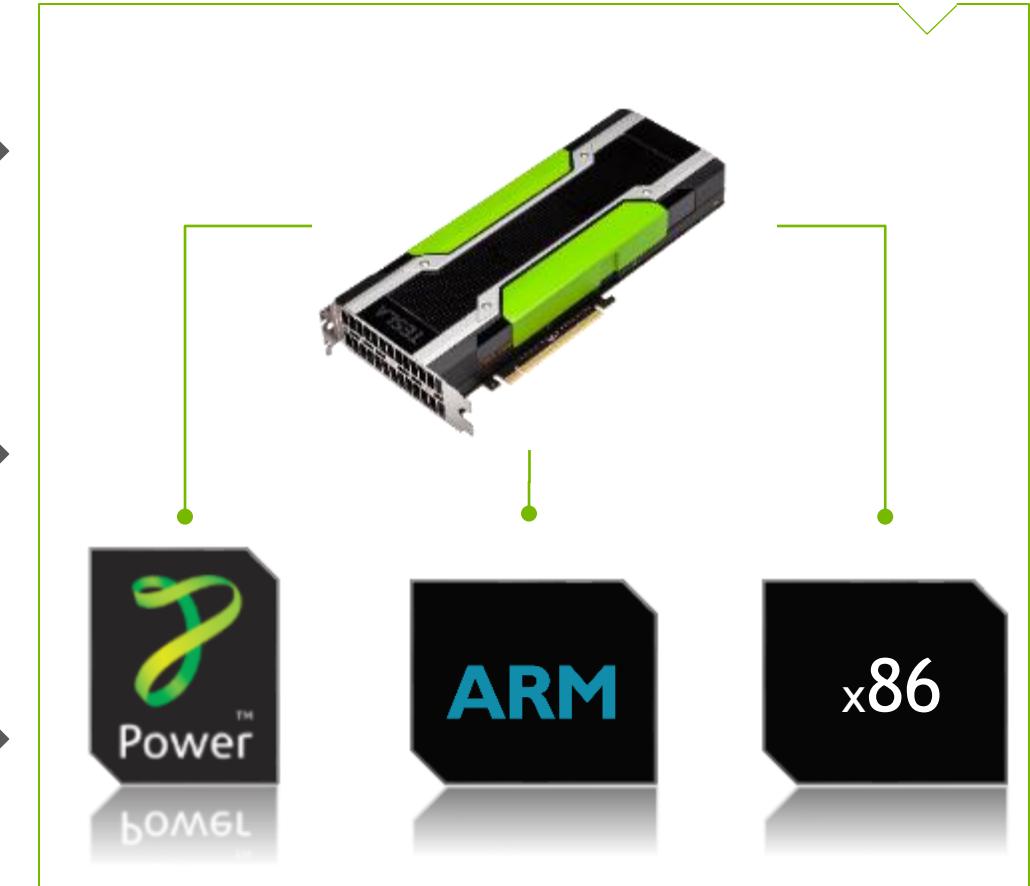


Compiler Directives

OpenACC



Programming Languages



GPU Developer Eco-System

Numerical Packages

MATLAB
Mathematica
LabView

Debuggers & Profilers

CUDA-GDB
NV Visual Profiler
NVIDIA Nsight
Visual Studio
Allinea
TotalView

Languages & Directives

C
C++
Fortran
Java
Python
OpenACC
OpenMP

Cluster Tools

GPUDirect RDMA
Datacenter
GPU Manager

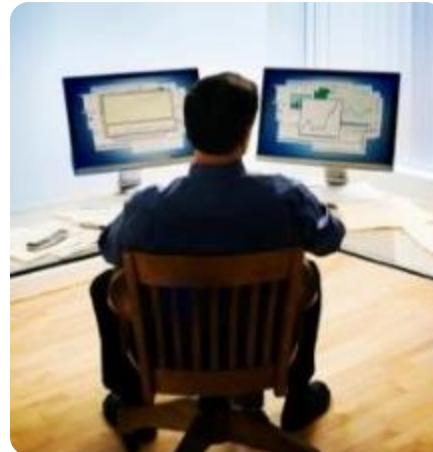
Libraries

FFT
BLAS
SPARSE
LAPACK
NPP
Video
Imaging

Consultants & Training



ANEO GPU Tech



OEM Solution Providers



CRAY

ASUS

SUPERMICRO

sgi

FUJITSU

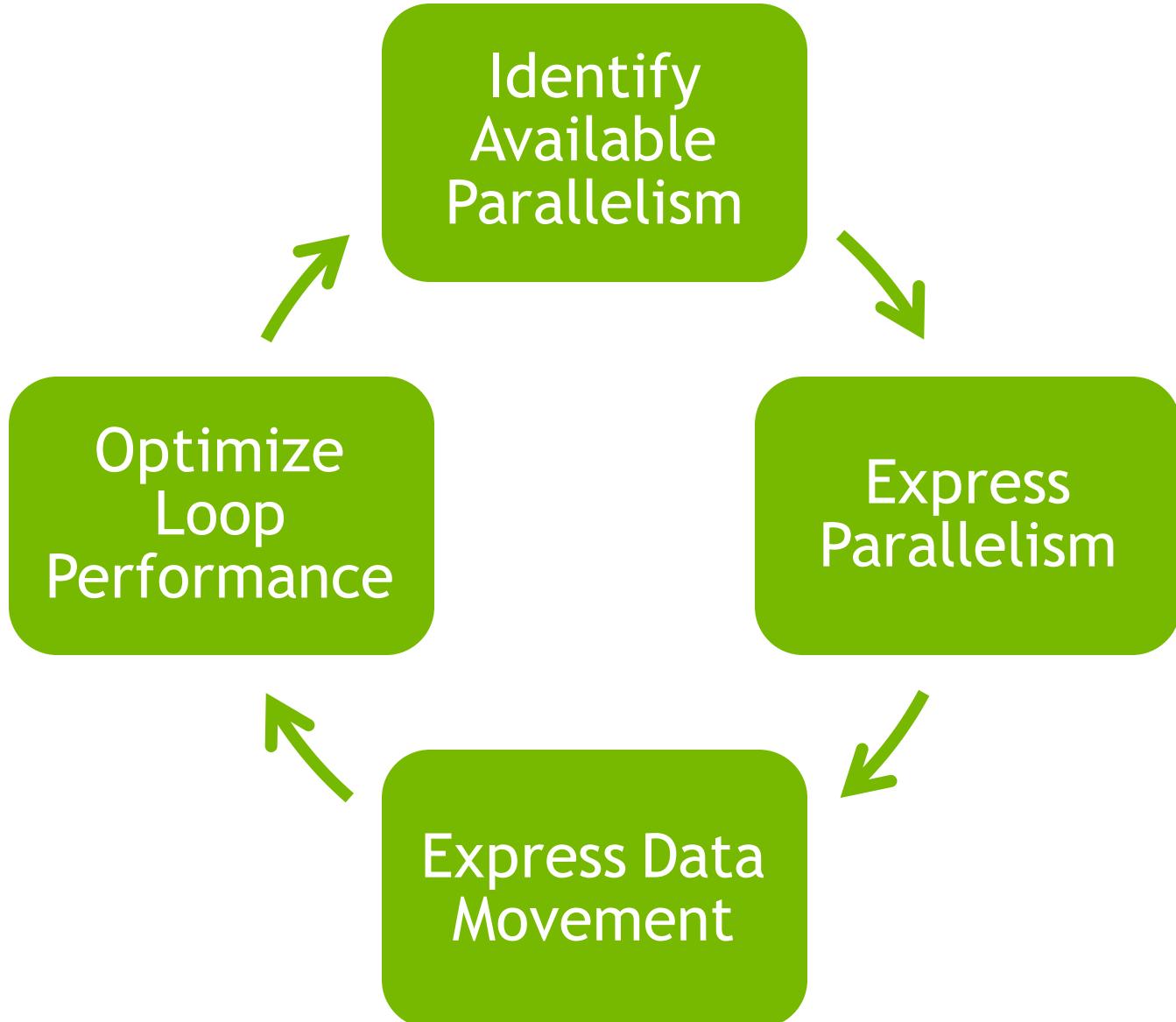
BULL



lenovo 联想

NEC

Accelerated Computing with OpenACC

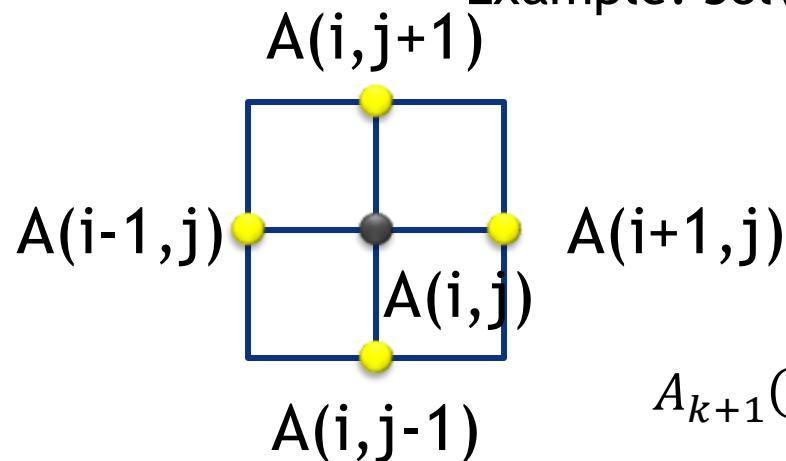


Example: Jacobi Iteration

Iteratively converges to correct value (e.g. Temperature), by computing new values at each point from the average of neighboring points.

Common, useful algorithm

Example: Solve Laplace equation in 2D: $\nabla^2 f(x, y) = 0$



$$A_{k+1}(i, j) = \frac{A_k(i - 1, j) + A_k(i + 1, j) + A_k(i, j - 1) + A_k(i, j + 1)}{4}$$

Jacobi Iteration: C Code

```
while ( err > tol && iter < iter_max ) {  
    err=0.0;  
  
    for( int j = 1; j < n-1; j++) {  
        for(int i = 1; i < m-1; i++) {  
  
            Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +  
                                  A[j-1][i] + A[j+1][i]);  
  
            err = max(err, abs(Anew[j][i] - A[j][i]));  
        }  
    }  
  
    for( int j = 1; j < n-1; j++) {  
        for( int i = 1; i < m-1; i++ ) {  
            A[j][i] = Anew[j][i];  
        }  
    }  
  
    iter++;  
}
```

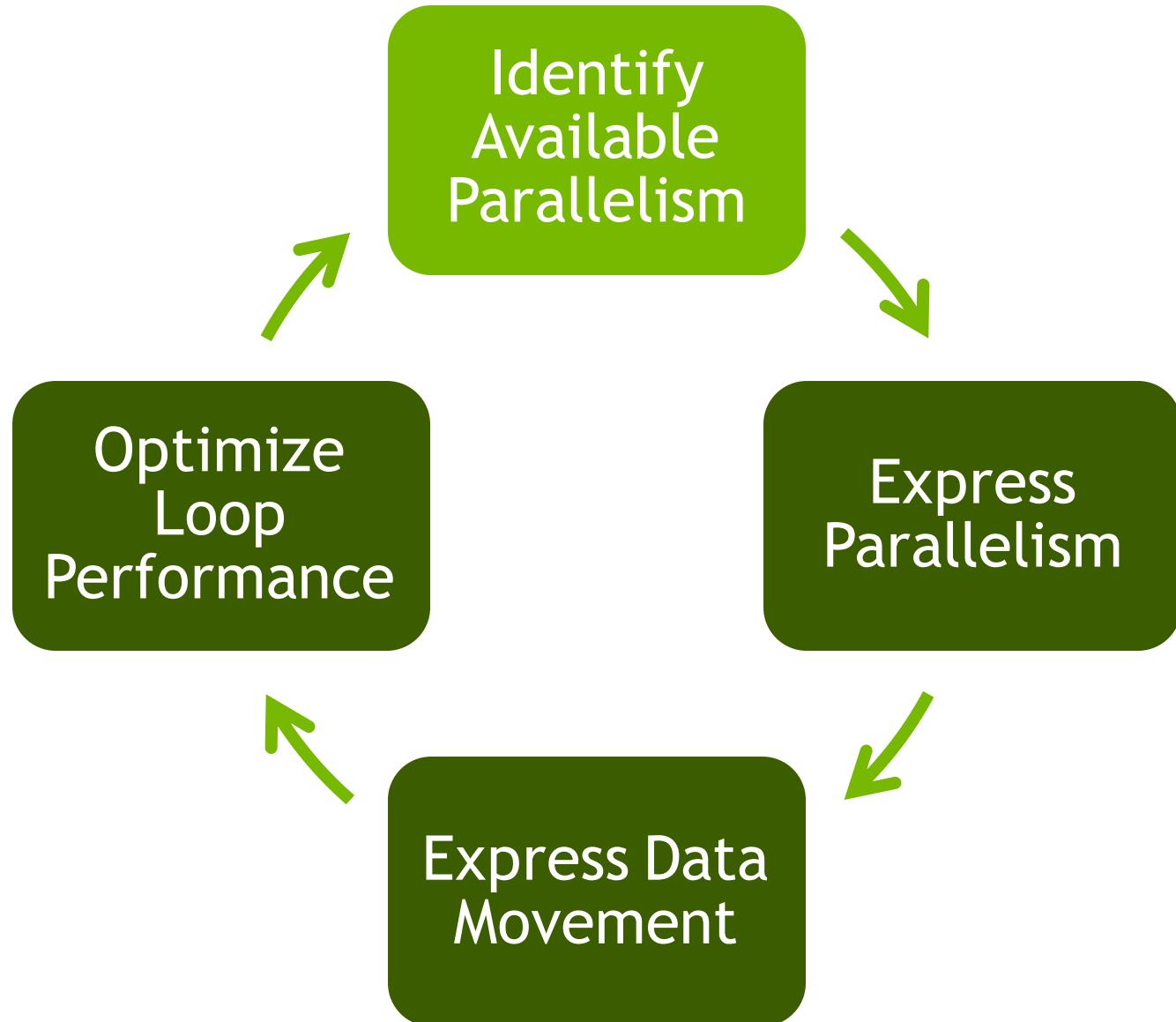
Iterate until converged

Iterate across matrix elements

Calculate new value from neighbors

Compute max error for convergence

Swap input/output arrays



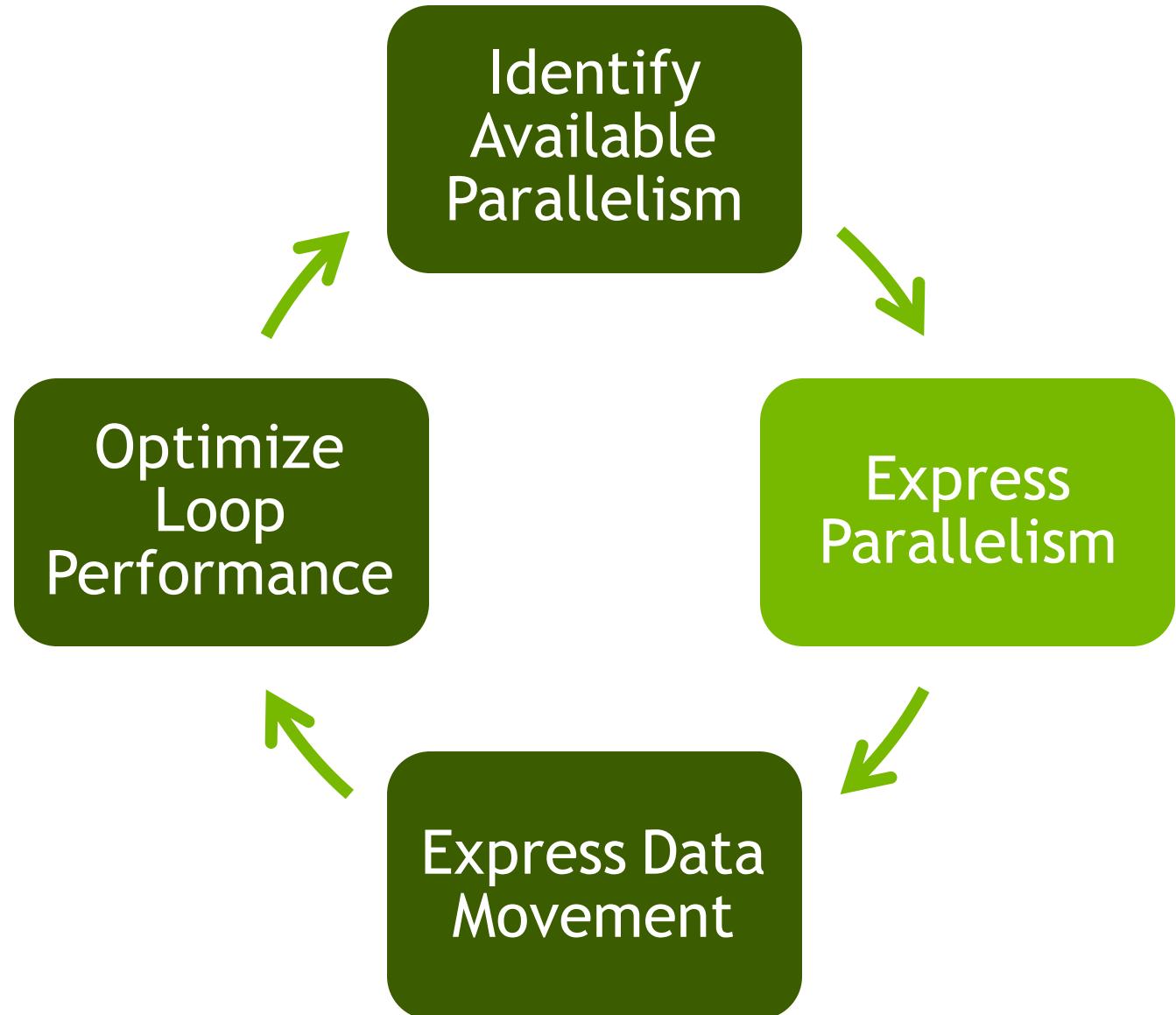
Identify Parallelism

```
while ( err > tol && iter < iter_max ) {  
    err=0.0;  
  
    for( int j = 1; j < n-1; j++) {  
        for(int i = 1; i < m-1; i++) {  
  
            Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +  
                                  A[j-1][i] + A[j+1][i]);  
  
            err = max(err, abs(Anew[j][i] - A[j][i]));  
        }  
    }  
  
    for( int j = 1; j < n-1; j++) {  
        for( int i = 1; i < m-1; i++ ) {  
            A[j][i] = Anew[j][i];  
        }  
    }  
  
    iter++;  
}
```

Data dependency
between iterations.

Independent loop
iterations

Independent loop
iterations



OpenACC Directives

```
Manage Data Movement → #pragma acc data copyin(a,b) copyout(c)
{ ... }

Initiate Parallel Execution → #pragma acc kernels
{
# pragma acc loop gang vector
for (i = 0; i < n; ++i) {
    z[i] = x[i] + y[i];
    ...
}
}

Optimize Loop Mappings → ...
}
```

- Incremental
- Single source
- Interoperable
- Performance portable
- CPU, GPU, MIC

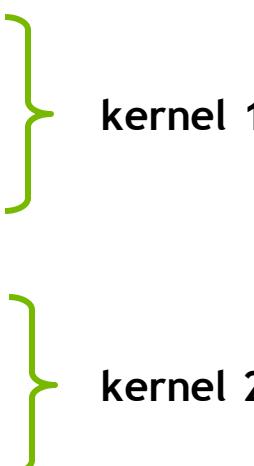
OpenACC
Directives for Accelerators

OpenACC kernels Directive

Identifies a region of code where I think the compiler can turn *loops* into *kernels*

```
#pragma acc kernels
{
for(int i=0; i<N; i++)
{
    x[i] = 1.0;
    y[i] = 2.0;
}

for(int i=0; i<N; i++)
{
    y[i] = a*x[i] + y[i];
}
```



The compiler identifies
2 parallel loops and
generates 2 kernels.

OpenACC kernels Directive (Fortran)

Identifies a region of code where I think the compiler can turn *loops* into *kernels*

```
!$acc kernels
do i=1,N
    x(i) = 1.0
    y(i) = 2.0
end do

y(:) = a*x(:) + y(:)

!$acc end kernels
```

}

kernel 1

}

kernel 2

The compiler identifies
2 parallel loops and
generates 2 kernels.

Loops vs. Kernels

```
for (int i = 0; i < 16384; i++)  
{  
    C[i] = A[i] + B[i];  
}
```

```
function loopBody(A, B, C, i)  
{  
    C[i] = A[i] + B[i];  
}
```

Loops vs. Kernels

```
for (int i = 0; i < 16384; i++)  
{  
    C[i] = A[i] + B[i];  
}
```

```
function loopBody(A, B, C, i)  
{  
    C[i] = A[i] + B[i];  
}
```

Calculate 0 -16383 in order.

Loops vs. Kernels

```
for (int i = 0; i < 16384; i++)  
{  
    C[i] = A[i] + B[i];  
}
```

Calculate 0 -16383 in order.

```
function loopBody(A, B, C, i)  
{  
    C[i] = A[i] + B[i];  
}
```

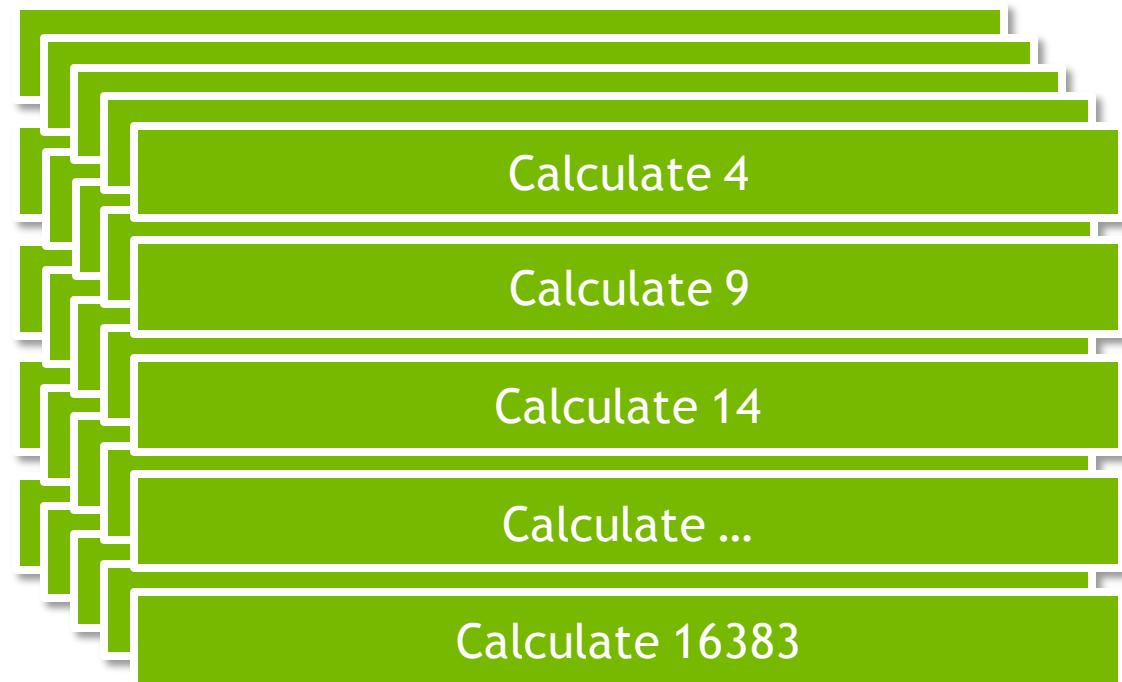
Calculate 0

Loops vs. Kernels

```
for (int i = 0; i < 16384; i++)  
{  
    C[i] = A[i] + B[i];  
}
```

Calculate 0 -16383 in order.

```
function loopBody(A, B, C, i)  
{  
    C[i] = A[i] + B[i];  
}
```



The Kernels Directive

Identifies a region of code where I think the compiler can turn *loops* into *kernels*

```
#pragma acc kernels
{
for (int i = 0; i < 16384;
     i++)
{
    C[i] = A[i] + B[i];
}
}
```

The Compiler will...

1. Analyze the code to determine if it contains parallelism
2. Identify data that needs to be moved to/from the GPU memory
3. Generate kernels
4. Run on the GPU

Parallelize with OpenACC kernels

```
while ( err > tol && iter < iter_max ) {
    err=0.0;

#pragma acc kernels
{
    for( int j = 1; j < n-1; j++) {
        for(int i = 1; i < m-1; i++) {

            Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +
                                  A[j-1][i] + A[j+1][i]);

            err = max(err, abs(Anew[j][i] - A[j][i]));
        }
    }

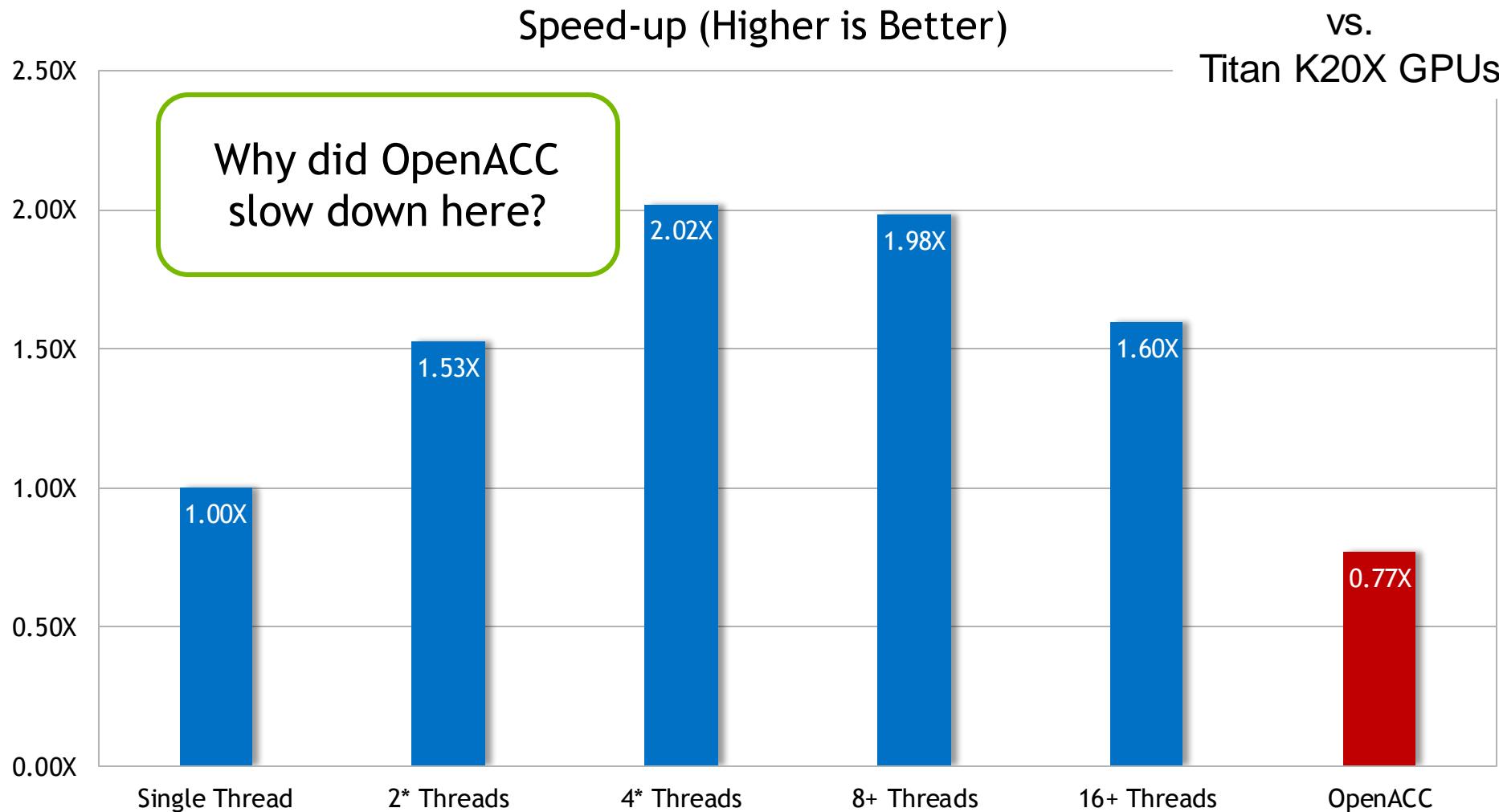
    for( int j = 1; j < n-1; j++) {
        for( int i = 1; i < m-1; i++ ) {
            A[j][i] = Anew[j][i];
        }
    }
}
iter++;
}
```

Look for parallelism
within this region.

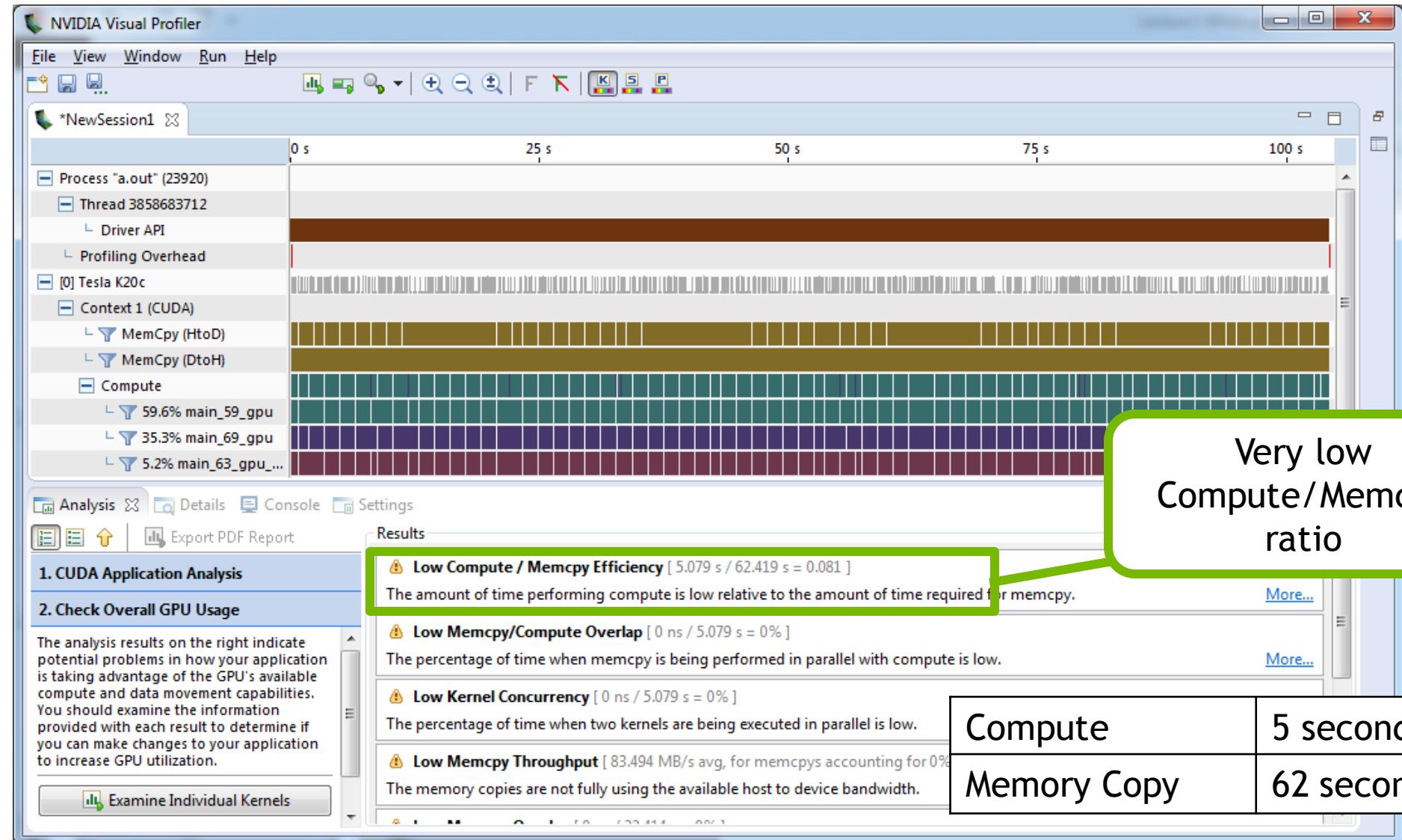
Building the code

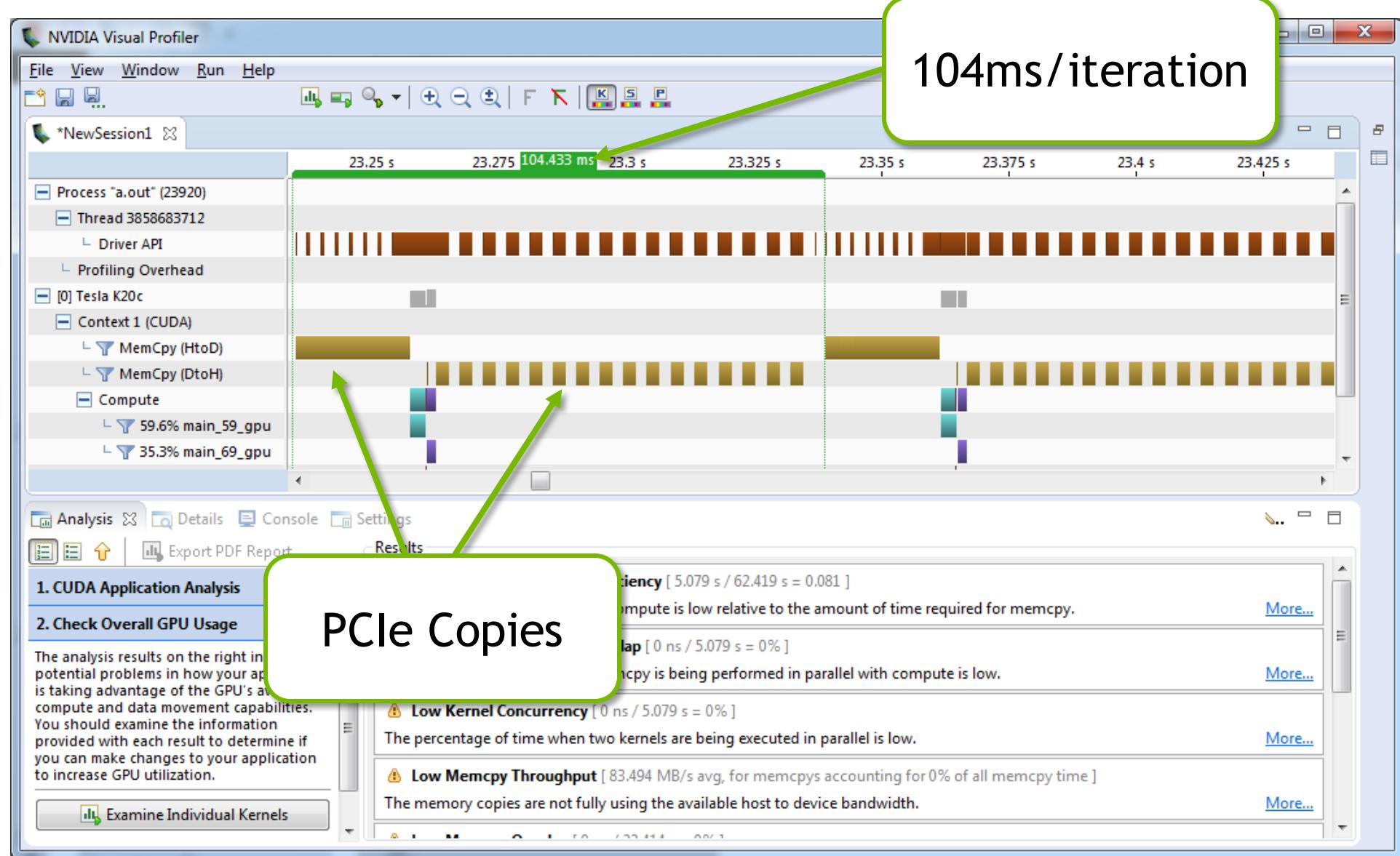
```
$ pgcc -fast -ta=tesla -Minfo=all laplace2d.c
main:
  40, Loop not fused: function call before adjacent loop
      Generated vector sse code for the loop
  51, Loop not vectorized/parallelized: potential early exits
  55, Generating copyout(Anew[1:4094][1:4094])
      Generating copyin(A[:, :])
      Generating copyout(A[1:4094][1:4094])
      Generating Tesla code
  57, Loop is parallelizable
  59, Loop is parallelizable
      Accelerator kernel generated
      57, #pragma acc loop gang /* blockIdx.y */
      59, #pragma acc loop gang, vector(128) /* blockIdx.x threadIdx.x */
      63, Max reduction generated for error
  67, Loop is parallelizable
  69, Loop is parallelizable
      Accelerator kernel generated
      67, #pragma acc loop gang /* blockIdx.y */
      69, #pragma acc loop gang, vector(128) /* blockIdx.x threadIdx.x */
```

Titan IL-16 CPUs
vs.
Titan K20X GPUs



* Compute Module Mode, + Default mode





Excessive Data Transfers

```
while ( err > tol && iter < iter_max )
{
    err=0.0;
```

A, Anew resident on host

These copies happen every iteration of the outer while loop!

• •

```
#pragma acc kernels
```

A, Anew resident on accelerator

```

for( int j = 1; j < n-1; j++) {
    for(int i = 1; i < m-1; i++) {
        Anew[j][i] = 0.25 * (A[j][i+1] +
                               A[j][i-1] + A[j-1][i] +
                               A[j+1][i]);
    }
    err = max(err, abs(Anew[j][i] -
                        A[j][i]));
}

```

}

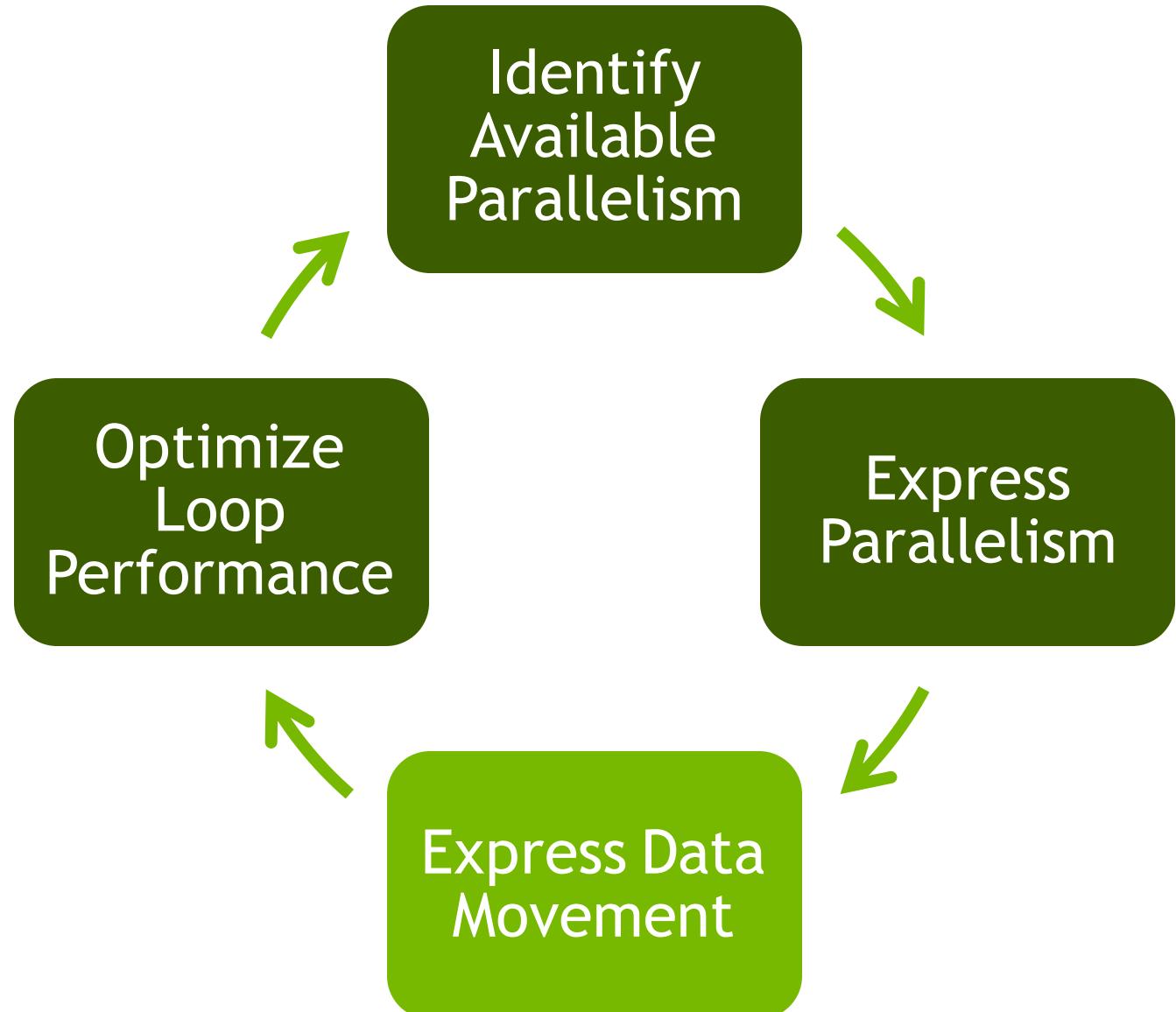
A, Anew resident on accelerator

Identifying Data Locality

```
while ( err > tol && iter < iter_max ) {  
    err=0.0;  
  
#pragma acc kernels  
{  
    for( int j = 1; j < n-1; j++) {  
        for(int i = 1; i < m-1; i++) {  
  
            Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +  
                                  A[j-1][i] + A[j+1][i]);  
  
            err = max(err, abs(Anew[j][i] - A[j][i]));  
        }  
    }  
  
    for( int j = 1; j < n-1; j++) {  
        for( int i = 1; i < m-1; i++ ) {  
            A[j][i] = Anew[j][i];  
        }  
    }  
}  
  
iter++;  
}
```

Does the CPU need the data between these loop nests?

Does the CPU need the data between iterations of the convergence loop?



Data regions

The **data** directive defines a region of code in which GPU arrays remain on the GPU and are shared among all kernels in that region.

```
#pragma acc data
{
#pragma acc kernels
...
#pragma acc kernels
...
}
```



Arrays used within the data region will remain on the GPU until the end of the data region.

Data Clauses

`copy (list)`

Allocates memory on GPU and copies data from host to GPU when entering region and copies data to the host when exiting region.

`copyin (list)`

Allocates memory on GPU and copies data from host to GPU when entering region.

`copyout (list)`

Allocates memory on GPU and copies data to the host when exiting region.

`create (list)`

Allocates memory on GPU but does not copy.

`present (list)`

Data is already present on GPU from another containing data region.

`deviceptr(list)`

The variable is a device pointer (e.g. CUDA) and can be used directly on the device.

Array Shaping

Compiler sometimes cannot determine size of arrays

Must specify explicitly using data clauses and array “shape”

C/C++

```
#pragma acc data copyin(a[0:nelem]) copyout(b[s/4:3*s/4])
```

Fortran

```
!$acc data copyin(a(1:end)) copyout(b(s/4:3*s/4))
```

Note: data clauses can be used on **data**, **parallel**, or **kernels**

Express Data Locality

```
#pragma acc data copy(A) create(Anew)
while ( err > tol && iter < iter_max ) {
    err=0.0;
#pragma acc kernels
{
    for( int j = 1; j < n-1; j++) {
        for(int i = 1; i < m-1; i++) {

            Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +
                                  A[j-1][i] + A[j+1][i]);

            err = max(err, abs(Anew[j][i] - A[j][i]));
        }
    }

    for( int j = 1; j < n-1; j++) {
        for( int i = 1; i < m-1; i++ ) {
            A[j][i] = Anew[j][i];
        }
    }
}
iter++;
}
```

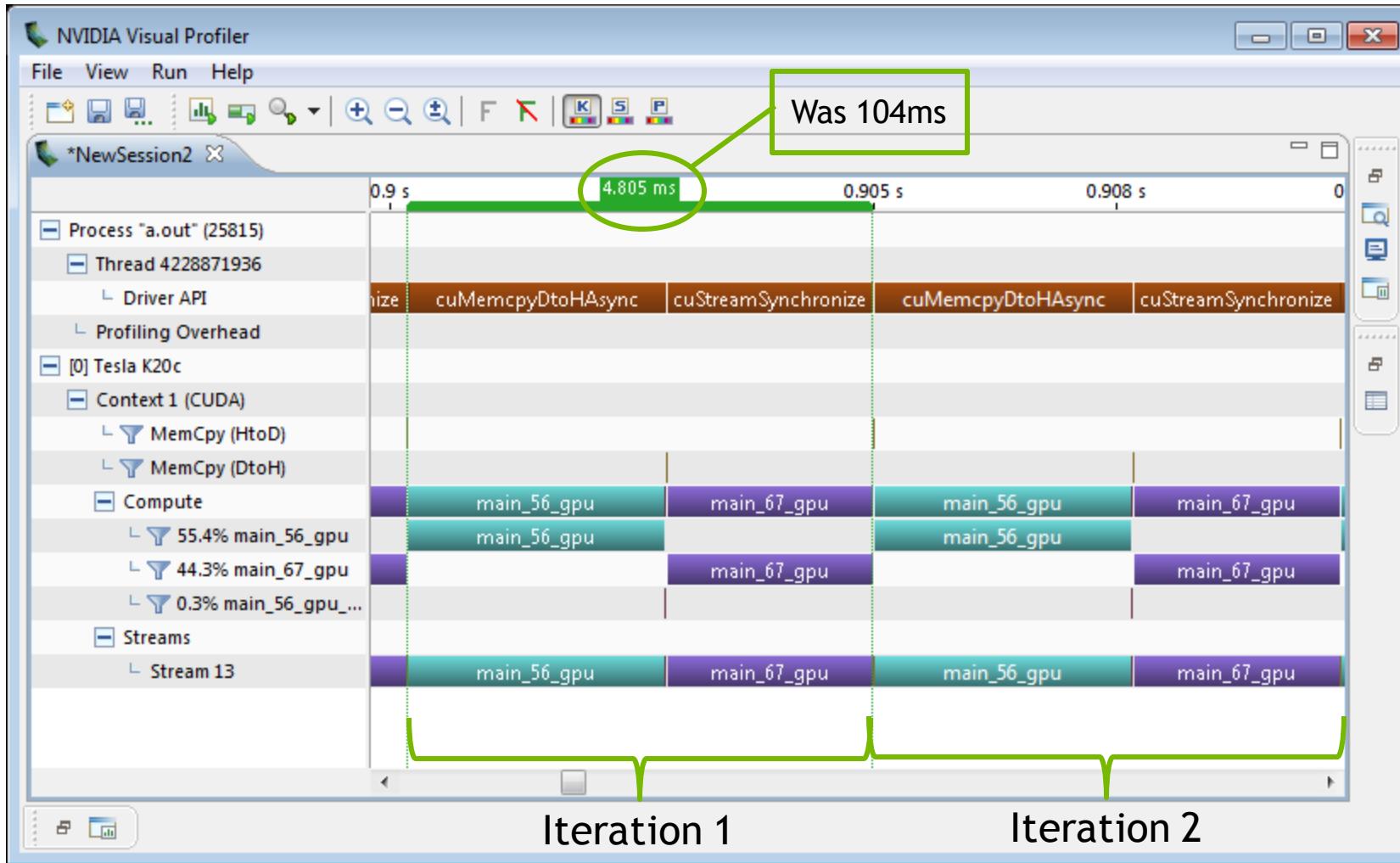
Copy A to/from the accelerator only when needed.

Create Anew as a device temporary.

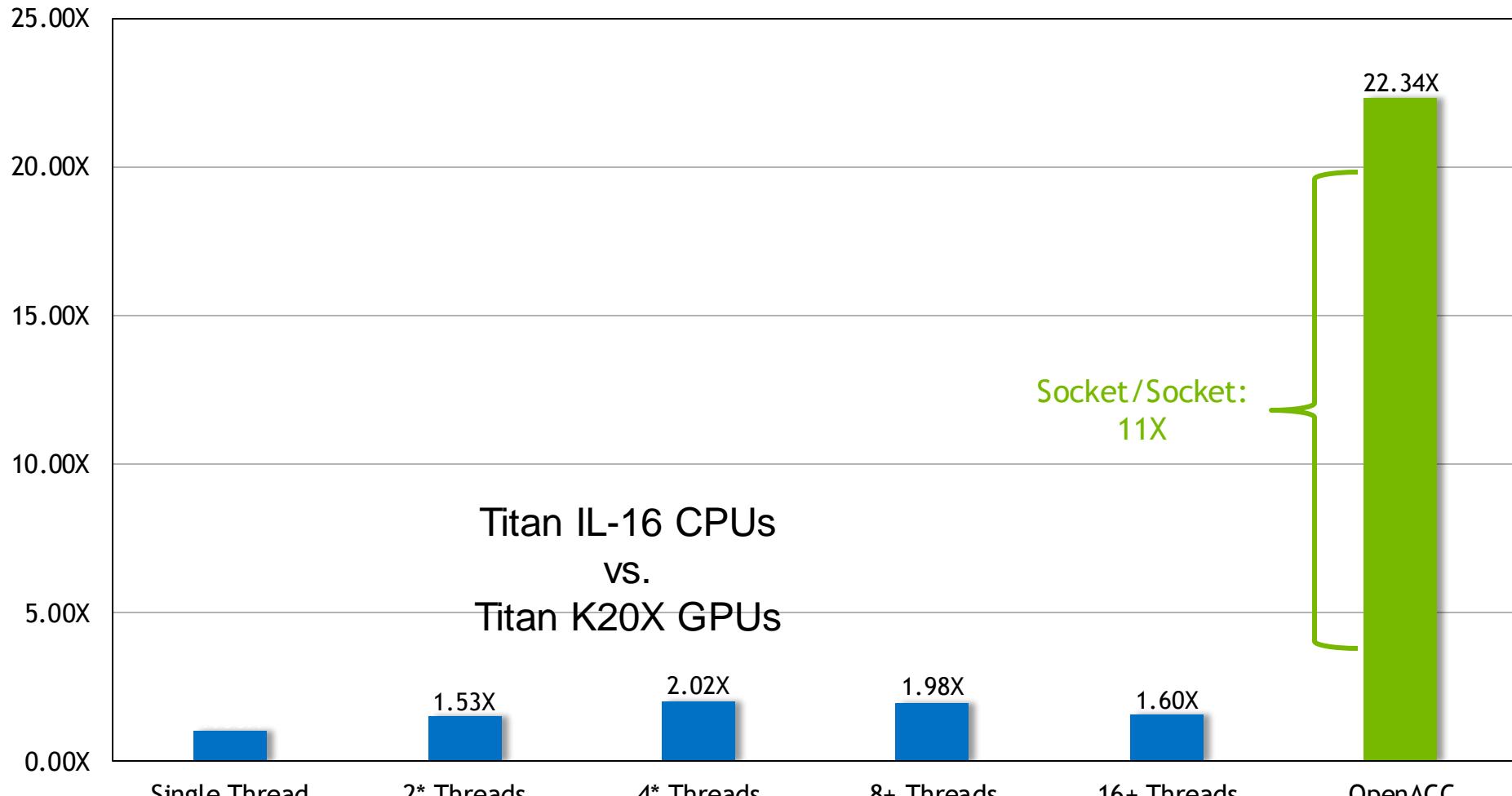
Rebuilding the code

```
$ pgcc -fast -acc -ta=tesla -Minfo=all laplace2d.c
main:
  40, Loop not fused: function call before adjacent loop
      Generated vector sse code for the loop
  51, Generating copy(A[:, :])
      Generating create(Anew[:, :])
      Loop not vectorized/parallelized: potential early exits
  56, Accelerator kernel generated
      56, Max reduction generated for error
      57, #pragma acc loop gang /* blockIdx.x */
      59, #pragma acc loop vector(256) /* threadIdx.x */
  56, Generating Tesla code
  59, Loop is parallelizable
  67, Accelerator kernel generated
      68, #pragma acc loop gang /* blockIdx.x */
      70, #pragma acc loop vector(256) /* threadIdx.x */
  67, Generating Tesla code
  70, Loop is parallelizable
```

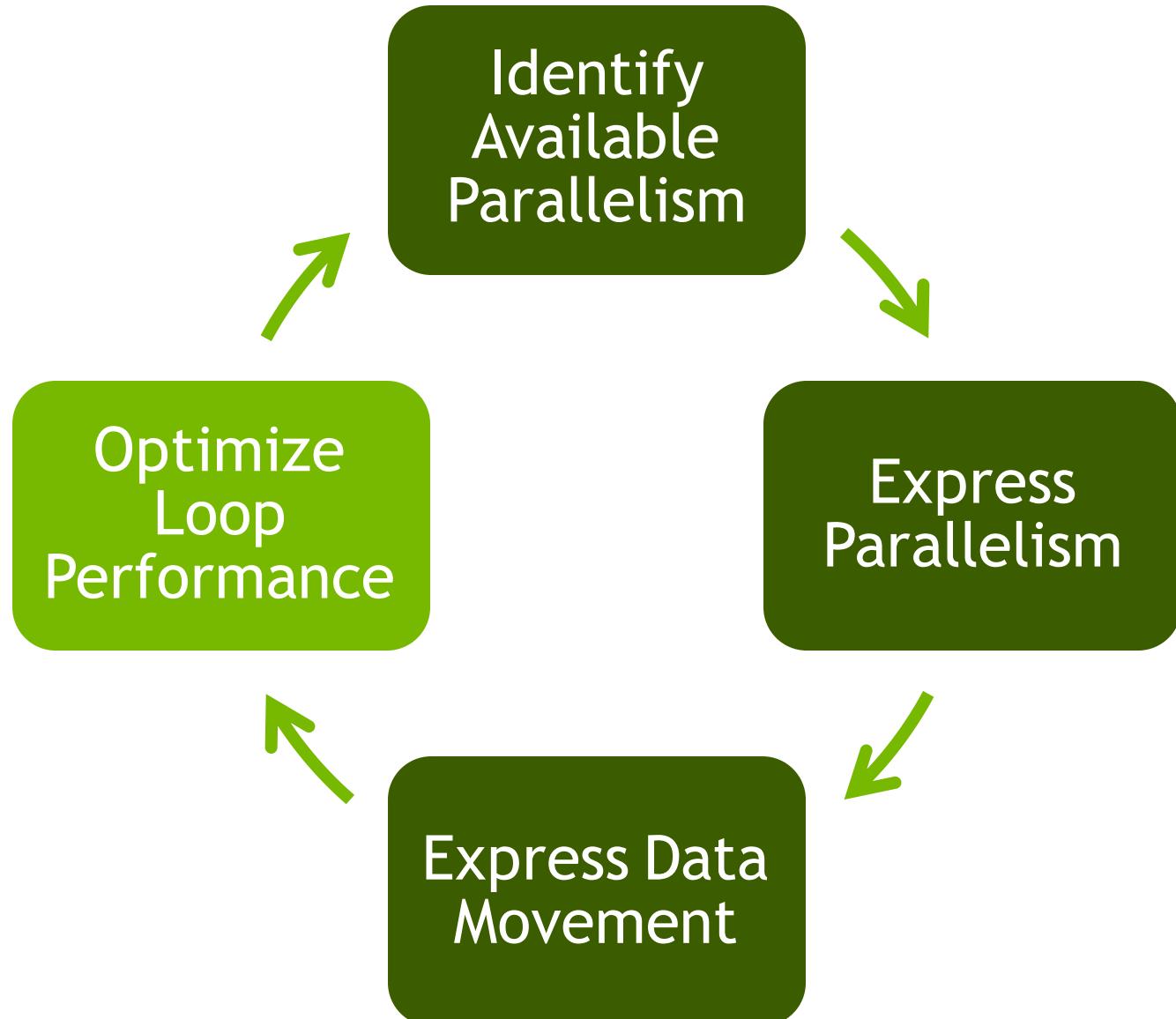
Visual Profiler: Data Region



Speed-Up (Higher is Better)



* Compute Module Mode, + Default mode



The loop Directive

The **loop** directive gives the compiler additional information about the *next* loop in the source code through several clauses.

- **independent** - all iterations of the loop are independent
- **collapse (N)** - turn the next N loops into one, flattened loop
- **tile(N[,M,...])** - break the next 1 or more loops into *tiles* based on the provided dimensions.

These clauses and more will be discussed in greater detail in a later class.

Optimize Loop Performance

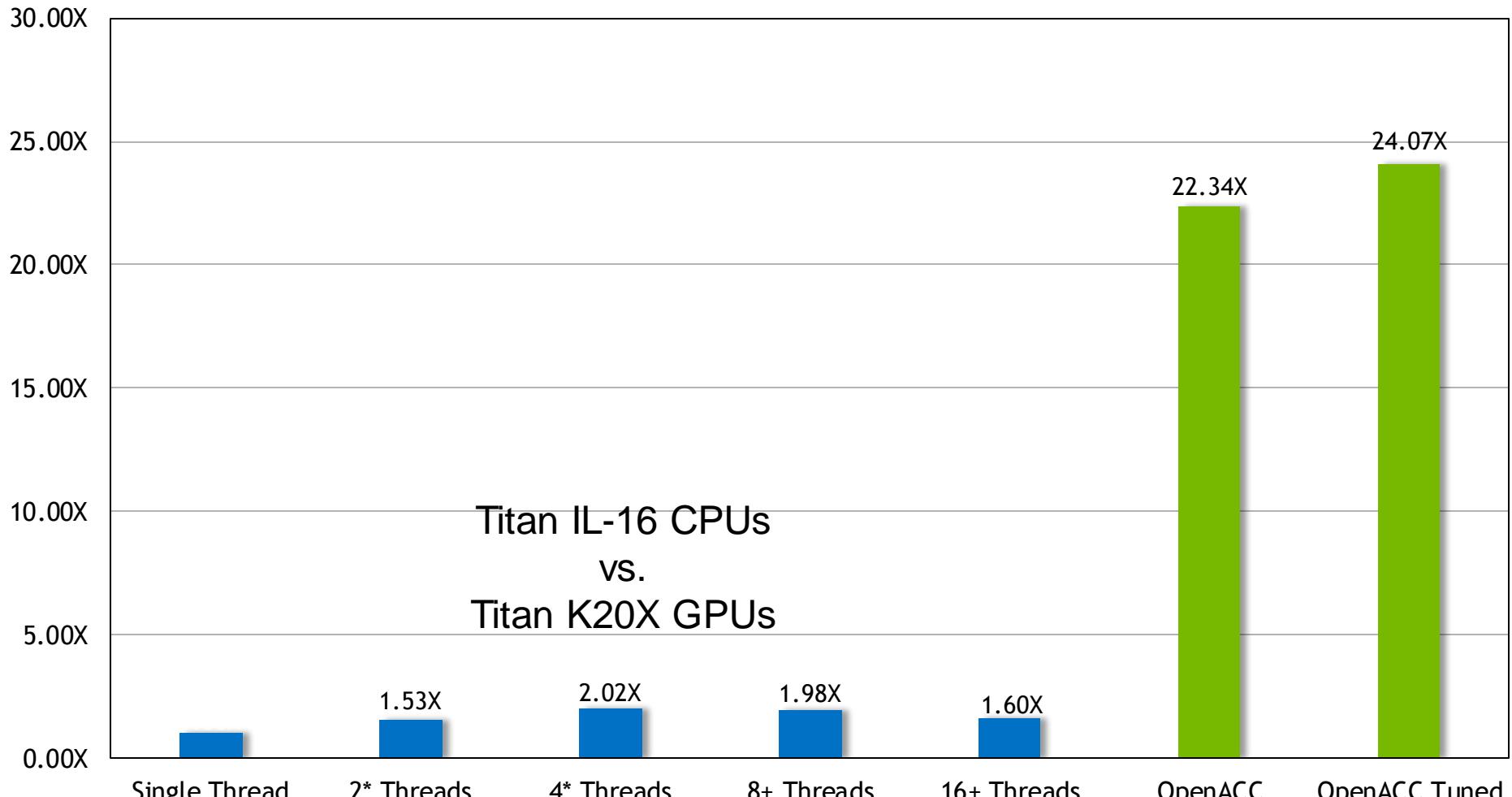
```
#pragma acc data copy(A) create(Anew)
while ( err > tol && iter < iter_max ) {
    err=0.0;
#pragma acc kernels
{
#pragma acc loop device_type(nvidia) tile(32, 4)
    for( int j = 1; j < n-1; j++) {
        for(int i = 1; i < m-1; i++) {

            Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +
                                  A[j-1][i] + A[j+1][i]);

            err = max(err, abs(Anew[j][i] - A[j][i]));
        }
    }
#pragma acc loop device_type(nvidia) tile(32, 4)
    for( int j = 1; j < n-1; j++) {
        for( int i = 1; i < m-1; i++ ) {
            A[j][i] = Anew[j][i];
        }
    }
}
iter++;
}
```

“Tile” the next two loops
into 32x4 blocks, but
only on NVIDIA GPUs.

Speed-Up (Higher is Better)



* Compute Module Mode, + Default mode

Next Steps

Online Training

<https://nvidia.qwiklab.com/>

NVIDIA offers a variety of hands-on labs online.

- OpenACC, Libraries, CUDA, Deep Learning

You can obtain free credits to take these labs by

- Visiting <https://developer.nvidia.com/qwiklabs-signup> and
- Signing up with the promo code OPENACC

NVIDIA Online OpenACC Courses

2015 course videos and slides are available online, the homework labs are in qwiklabs.

<https://developer.nvidia.com/openacc-courses>

Advanced OpenACC Course began May 19th.

<https://developer.nvidia.com/openacc-advanced-course>

2016 OpenACC course begins in September