# OpenACC Course

Lecture 1: Introduction to OpenACC

September 2015

**NVIDIA**

Course Objective:

Enable *you* to accelerate *your* applications with OpenACC.

# Course Syllabus

Oct 1: Introduction to OpenACC

Oct 6: Office Hours

Oct 15: Profiling and Parallelizing with the OpenACC Toolkit

Oct 20: Office Hours

Oct 29: Expressing Data Locality and Optimizations with OpenACC

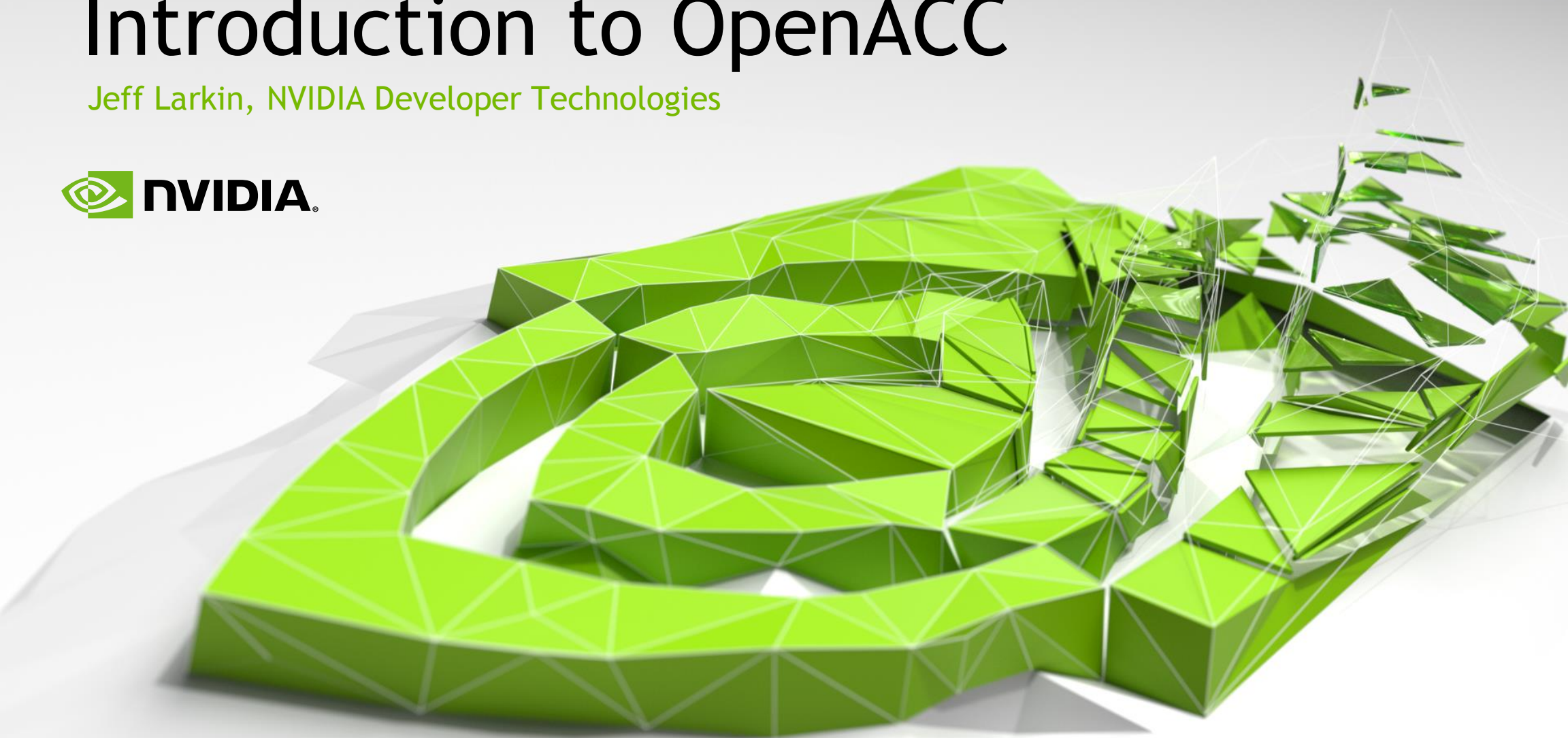Nov 3: Office Hours

Nov 12: Advanced OpenACC Techniques

Nov 24: Office Hours

**Recordings:**
https://developer.nvidia.com/openacc-course

# Introduction to OpenACC

Jeff Larkin, NVIDIA Developer Technologies

# Agenda

Why OpenACC?

Accelerated Computing Fundamentals

OpenACC Programming Cycle

Installing the OpenACC Toolkit

Accessing QwikLabs

Week 1 Homework

# Why OpenACC?

# 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

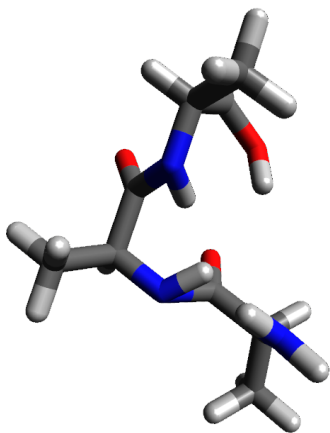RIKEN Japan
NICAM- Climate Modeling



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

**8000+**

Developers

using OpenACC

http://www.cray.com/sites/default/files/resources/OpenACC_213462.12_OpenACC_Cosmo_CS_FNL.pdf
http://www.hpcwire.com/off-the-wire/first-round-of-2015-hackathons-gets-underway
http://on-demand.gputechconf.com/gtc/2015/presentation/S5297-Hisashi-Yashiro.pdf
http://www.openacc.org/content/experiences-porting-molecular-dynamics-code-gpus-cray-xk7

# LS-DALTON

Large-scale application for calculating high-accuracy molecular energies



" *OpenACC makes GPU computing approachable for domain scientists. Initial OpenACC implementation required only minor effort, and more importantly, no modifications of our existing CPU implementation.* "
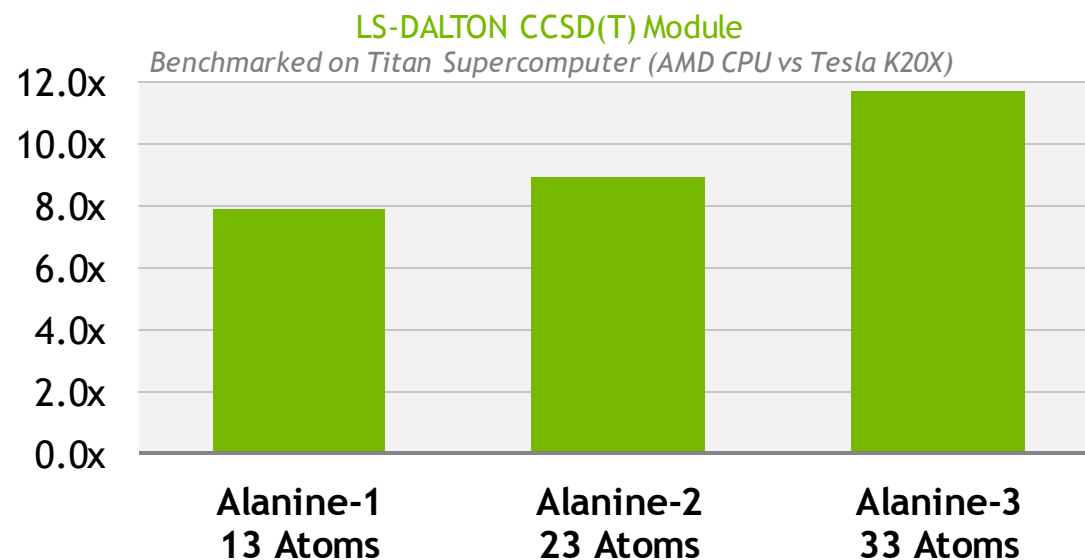
*Janus Juul Eriksen, PhD Fellow*
*qLEAP Center for Theoretical Chemistry, Aarhus University*

## Minimal Effort

| Lines of Code Modified | # of Weeks Required | # of Codes to Maintain |
|---|---|---|
| **<100 Lines** | **1 Week** | **1 Source** |

## Big Performance

### LS-DALTON CCSD(T) Module
*Benchmarked on Titan Supercomputer (AMD CPU vs Tesla K20X)*



| | Alanine-1 13 Atoms | Alanine-2 23 Atoms | Alanine-3 33 Atoms |
|---|---|---|---|

# OpenACC Directives

Manage
Data
Movement

Initiate
Parallel
Execution

Optimize
Loop
Mappings

```
#pragma acc data copyin(a,b) copyout(c)
{
    ...
#pragma acc parallel
{
#pragma acc loop gang vector
    for (i = 0; i < n; ++i) {
        z[i] = x[i] + y[i];
        ...
    }
}
    ...
}
```

**OpenACC**
Directives for Accelerators

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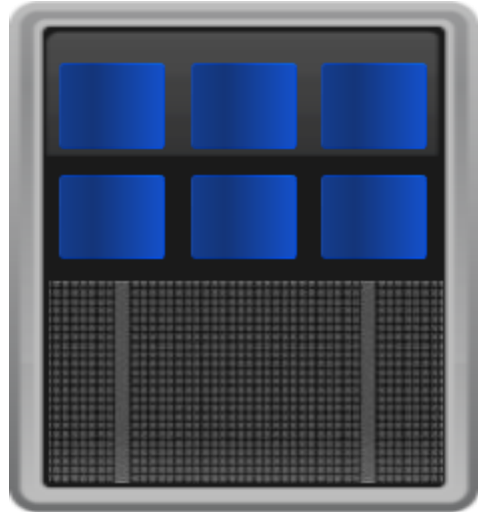- Interoperable
- Performance portable
- CPU, GPU, MIC

# Accelerated Computing Fundamentals

# Accelerated Computing
*10x Performance & 5x Energy Efficiency for HPC*
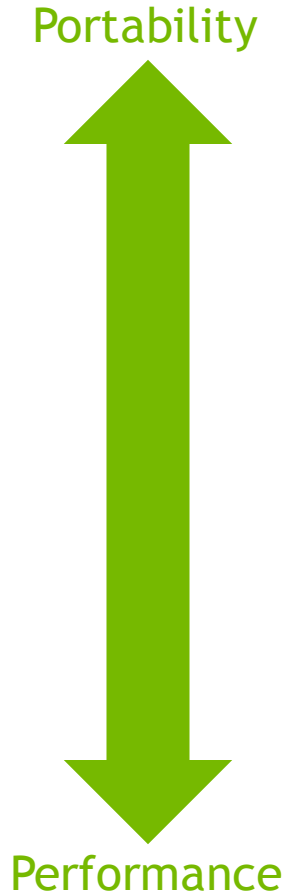
**CPU**
Optimized for
Serial Tasks

**GPU Accelerator**
Optimized for
Parallel Tasks

+

# What is Heterogeneous Programming?



GPU

Compute-Intensive Functions

A few % of Code
A large % of Time

Rest of Sequential
CPU Code

CPU

+

# Portability & Performance

Portability

Performance

## Accelerated Libraries

High performance with little or no code change

Limited by what libraries are available

## Compiler Directives

High Level: Based on existing languages; simple, familiar, portable

High Level: Performance may not be optimal

## Parallel Language Extensions

Greater flexibility and control for maximum performance

Often less portable and more time consuming to implement

NVIDIA.

# Code for Portability & Performance

**Libraries**

- Implement as much as possible using portable libraries

**Directives**

- Use directives for rapid and portable development

**Languages**

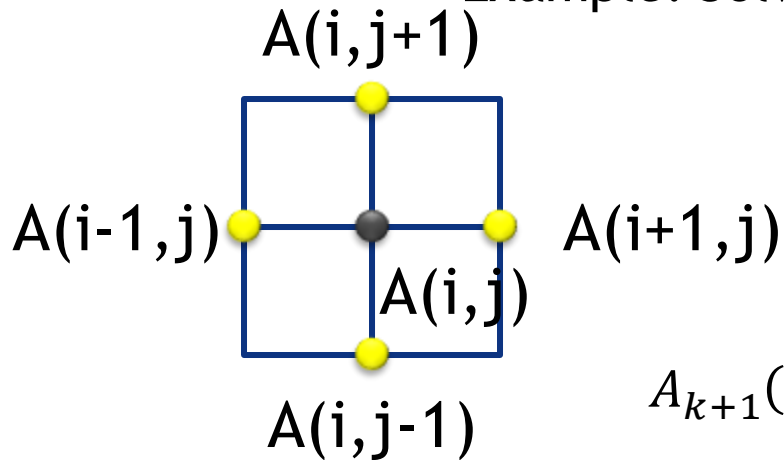- Use lower level languages for important kernels

# OpenACC Programming Cycle

```
                    ┌─────────────────┐
                    │     Identify    │
                    │    Available    │
                    │   Parallelism   │
                    └─────────────────┘

  ┌─────────────────┐                   ┌─────────────────┐
  │     Optimize    │                   │     Express     │
  │       Loop      │                   │   Parallelism   │
  │   Performance   │                   │                 │
  └─────────────────┘                   └─────────────────┘

                    ┌─────────────────┐
                    │   Express Data  │
                    │    Movement     │
                    └─────────────────┘
```

# 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(i,j+1)

A(i-1,j)

A(i+1,j)

A(i,j)

A(i,j-1)

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

# Jacobi Iteration: C Code

```c
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++;
}
```
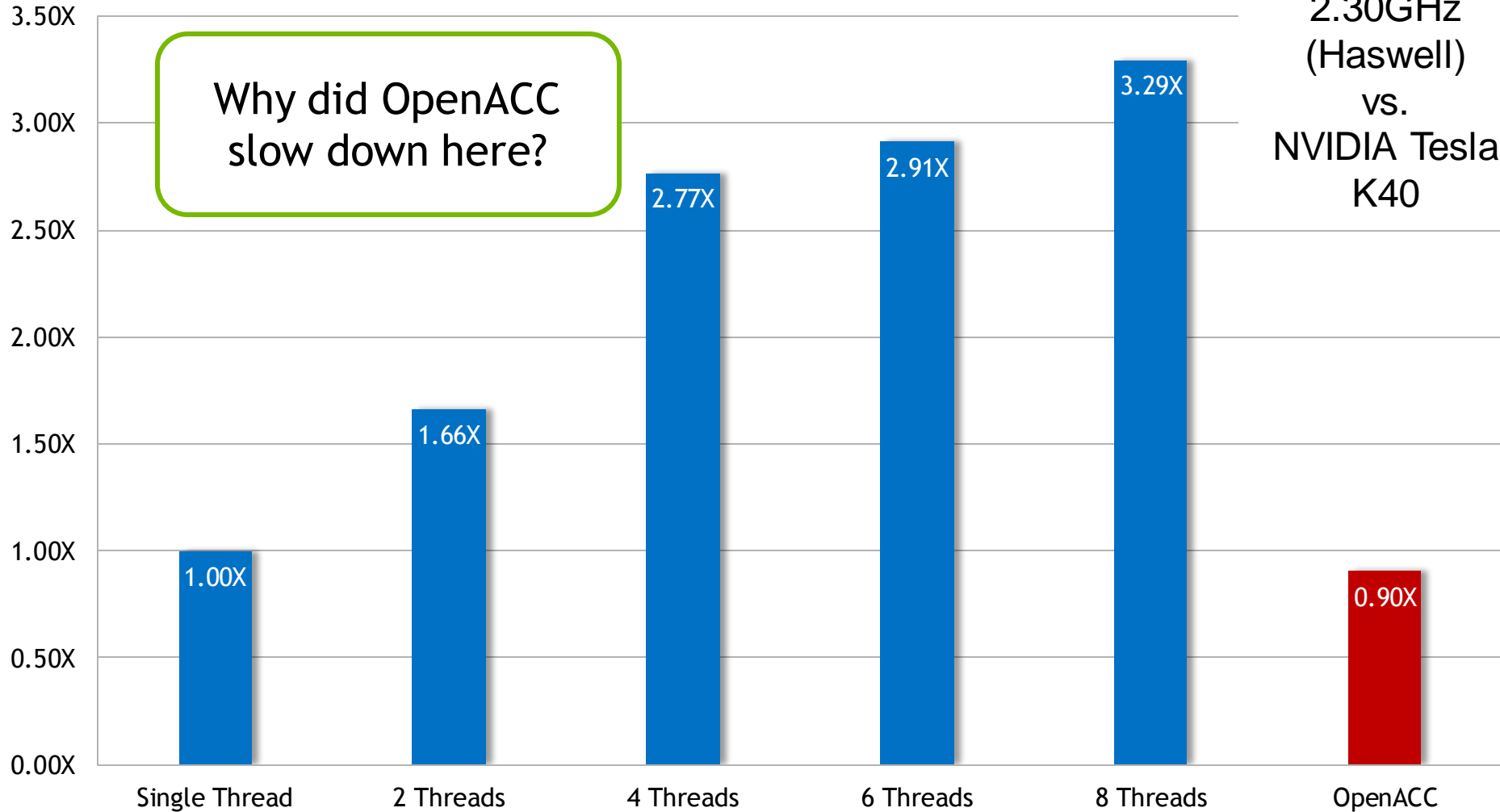
- Iterate until converged

- 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**

Identify Available Parallelism → Express Parallelism → Express Data Movement → Optimize Loop Performance → (back to Identify Available Parallelism)

# OpenACC kernels Directive

The kernels directive identifies a region that may contain *loops* that the compiler can turn into parallel *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];
}
}
```

kernel 1

kernel 2

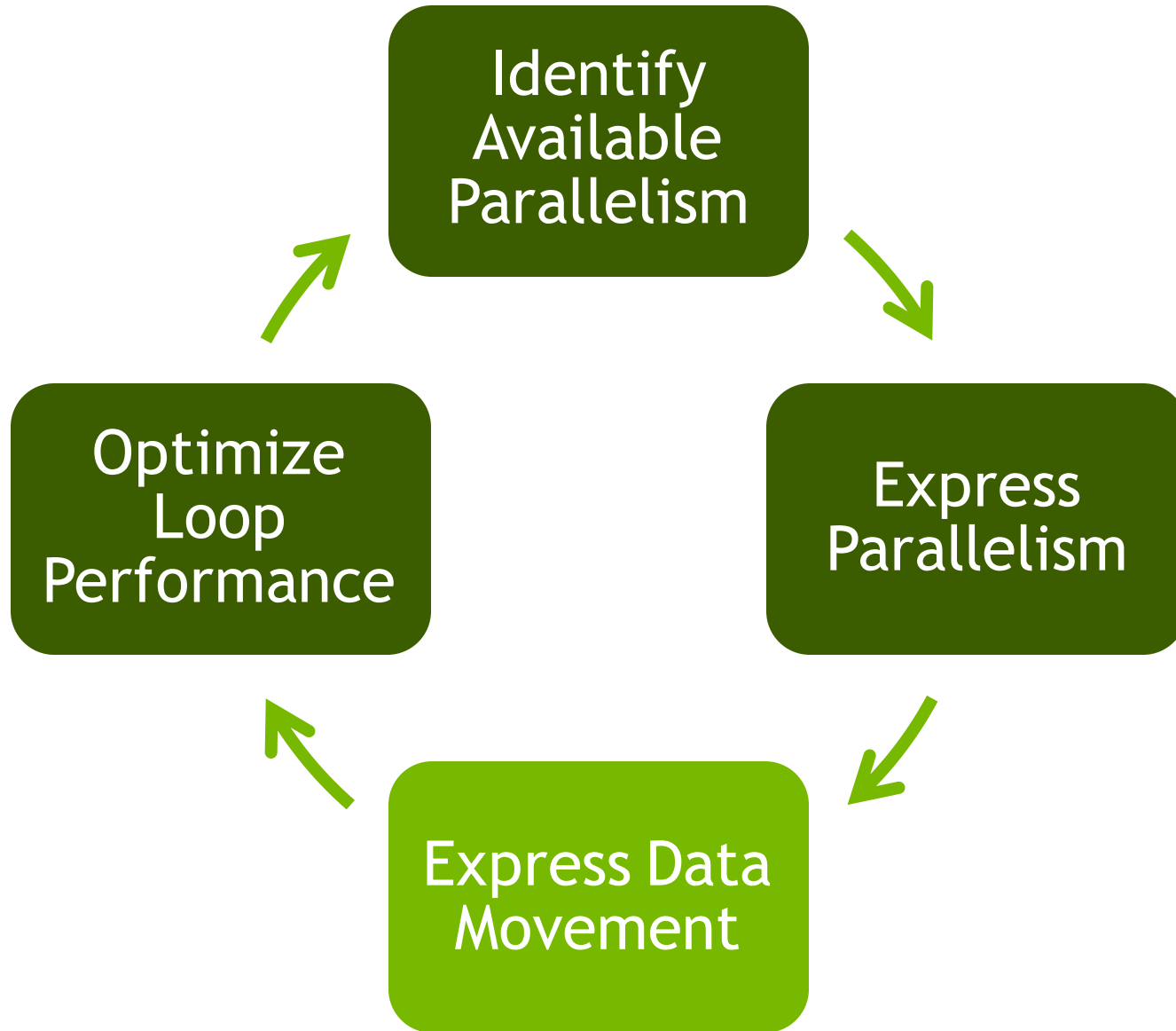The compiler identifies 2 parallel loops and generates 2 kernels.

# Parallelize with OpenACC kernels

```c
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 */
```

# Speed-up (Higher is Better)

Intel Xeon E5-2698 v3 @ 2.30GHz (Haswell) vs. NVIDIA Tesla K40

Why did OpenACC slow down here?

| Single Thread | 2 Threads | 4 Threads | 6 Threads | 8 Threads | OpenACC |
|---|---|---|---|---|---|
| 1.00X | 1.66X | 2.77X | 2.91X | 3.29X | 0.90X |

# Excessive Data Transfers

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

```
#pragma acc kernels
```

| A, Anew resident on host | → | 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]);
  }
}
...
```

C
o
p
y

C
o
p
y

**These copies happen every iteration of the outer while loop!**

| A, Anew resident on host | ← | 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++;
}
```
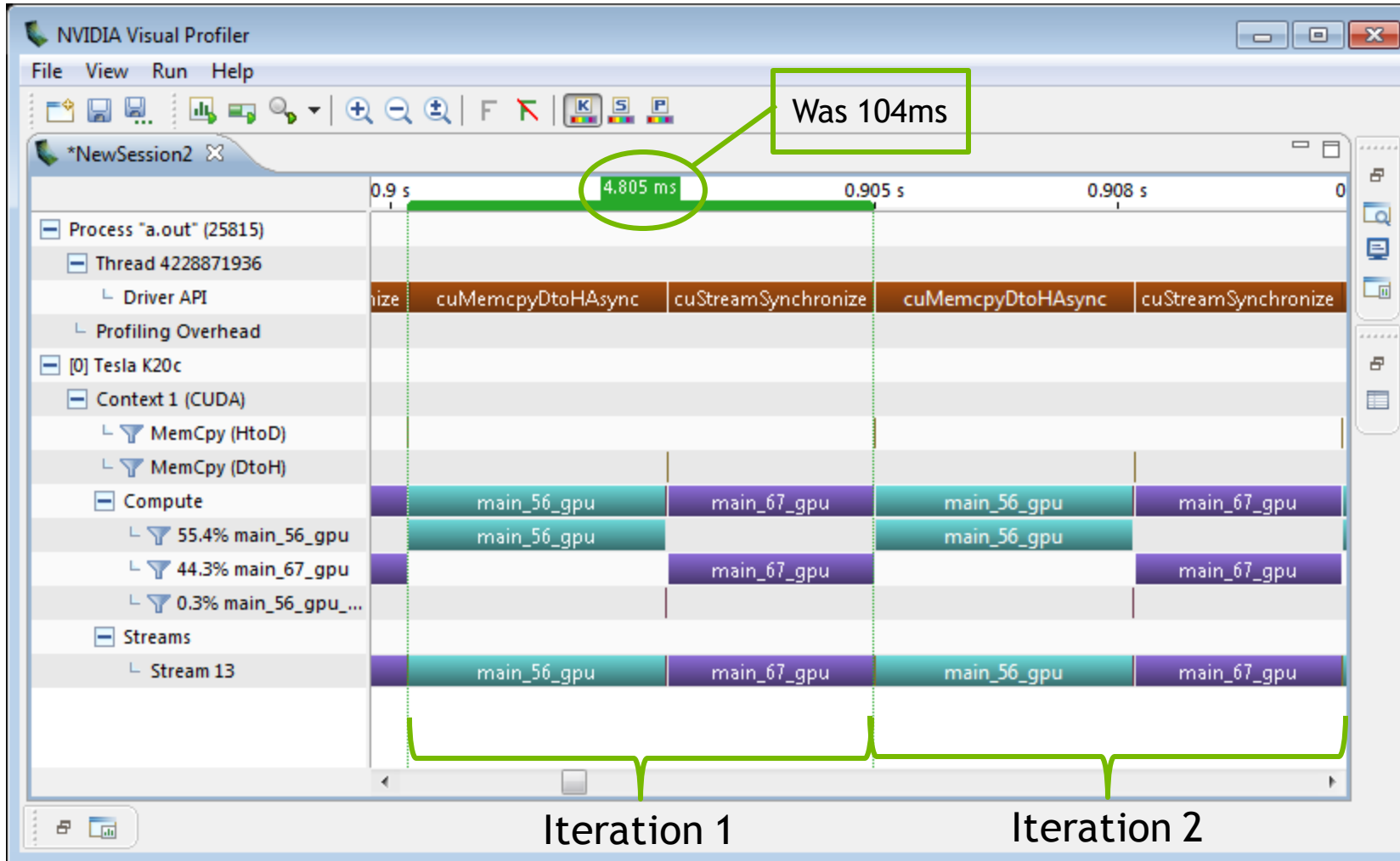
Does the CPU need the data between these loop nests?

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
...
}
```

Data Region

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)



Intel Xeon E5-2698 v3 @ 2.30GHz (Haswell)
vs.
NVIDIA Tesla K40

Socket/Socket:
5.2X

| | Single Thread | 2 Threads | 4 Threads | 6 Threads | 8 Threads | OpenACC |
|---|---|---|---|---|---|---|
| | | 1.90X | 3.20X | 3.74X | 3.83X | 19.89X |

25.00X
20.00X
15.00X
10.00X
5.00X
0.00X

Identify Available Parallelism

Express Parallelism

Express Data Movement

Optimize Loop Performance

# 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.

NVIDIA.

# 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)



Intel Xeon E5-2698 v3 @ 2.30GHz (Haswell)
vs.
NVIDIA Tesla K40

| | Single Thread | 2 Threads | 4 Threads | 6 Threads | 8 Threads | OpenACC | OpenACC Tuned |
|---|---|---|---|---|---|---|---|
| | | 1.90X | 3.20X | 3.74X | 3.83X | 19.89X | 21.22X |

# The OpenACC Toolkit

# Introducing the New OpenACC Toolkit
## Free Toolkit Offers Simple & Powerful Path to Accelerated Computing

OpenACC
Directives for Accelerators

http://developer.nvidia.com/openacc

**PGI Compiler**
Free OpenACC compiler for academia

**NVProf Profiler**
Easily find where to add compiler directives

**GPU Wizard**
Identify which GPU libraries can jumpstart code

**Code Samples**
Learn from examples of real-world algorithms

**Documentation**
Quick start guide, Best practices, Forums

# Download the OpenACC Toolkit

▷ Go to
https://developer.nvidia.com/openacc



NVIDIA.

# Download the OpenACC Toolkit

▸ Go to
https://developer.nvidia.com/openacc

▸ Register for the toolkit

  ▸ If you are an academic developer, be sure
    to click the check box at the bottom.

# Download the OpenACC Toolkit

▸ Go to
https://developer.nvidia.com/openacc

▸ Register for the toolkit

  ▸ If you are an academic developer, be sure to click the check box at the bottom.

▸ You will receive an email from NVIDIA

  ▸ Be sure to read the Quick Start Guide

Your OpenACC Toolkit License - Message (HTML)    ?  ⊞  —  □  ✕

FILE    MESSAGE

Thu 8/6/2015 10:30 AM

NVIDIA <no-reply@nvidia.eu>

**Your OpenACC Toolkit License**

To    Jeff Larkin

Retention Policy   Inbox (45 days)                    Expires   Never

⊘ NVIDIA.

Hello Jeff,

Welcome to the NVIDIA OpenACC Toolkit, which enables you to do more science with less programming.

Your OpenACC Toolkit contains various tools and best practice guides to help you get started with your accelerated computing journey.

1. Please refer the OpenACC Toolkit Quick Start Guide for an overview of the tools and documentation included in the toolkit.

2. Only the PGI Accelerator Compiler included in the toolkit requires a license key. Your

NVIDIA Your OpenACC Toolkit License

# Windows/Mac Developers

- The OpenACC Toolkit is only available on Linux, however…

- The PGI compiler is available on Mac and Windows from
  http://www.pgroup.com/support/trial.htm

  - You should still register for the OpenACC Toolkit to get the 90 day license.

- The CUDA Toolkit contains the libraries and profiling tools that will be used in this course.

  - https://developer.nvidia.com/cuda-zone

- The OpenACC Programming Guide is available from http://bit.ly/openacc-guide

  - Obtaining all examples and guides from the toolkit will still require downloading the full OpenACC toolkit.

# Using QwikLabs

# Getting access

Go to [nvidia.qwiklab.com](nvidia.qwiklab.com), log-in or create an account



Sign In or Create a New Account

# Homework

# Complete "2X in 4 Steps" Qwiklab

▸ C: http://bit.ly/nvoacclab1c

▸ F90: http://bit.ly/nvoacclab1f

▸ This lab is browser-based and should take you roughly 1 hour.

# Install the OpenACC Toolkit (Optional)

▸ Go to developer.nvidia.com/openacc

▸ Register for the OpenACC Toolkit

▸ Install on your personal machine. (Linux Only)

# Where to find help

- OpenACC Course Recordings - https://developer.nvidia.com/openacc-course

- OpenACC on StackOverflow - http://stackoverflow.com/questions/tagged/openacc

- OpenACC Toolkit - http://developer.nvidia.com/openacc

Additional Resources:

- Parallel Forall Blog - http://devblogs.nvidia.com/parallelforall/

- GPU Technology Conference - http://www.gputechconf.com/

- OpenACC Website - http://openacc.org/

NVIDIA.

# Course Syllabus

Oct 1:   Introduction to OpenACC

Oct 6:   Office Hours

Oct 15: Profiling and Parallelizing with the OpenACC Toolkit

Oct 20: Office Hours

Oct 29: Expressing Data Locality and Optimizations with OpenACC

Nov 3:   Office Hours

Nov 12: Advanced OpenACC Techniques

Nov 24: Office Hours

**Recordings:**
https://developer.nvidia.com/openacc-course