

INTRODUCTION TO OPENACC

Speaker

LECTURE 1 OUTLINE

Topics to be covered

- What is OpenACC and Why Should You Care?
- Profile-driven Development
- First Steps with OpenACC
- Week 1 Lab
- Where to Get Help

INTRODUCTION TO OPENACC

3 WAYS TO ACCELERATE APPLICATIONS

Applications

Libraries

Easy to use
Most Performance

Compiler
Directives

Easy to use
Portable code

OpenACC

Programming
Languages

Most Performance
Most Flexibility

OPENACC IS...

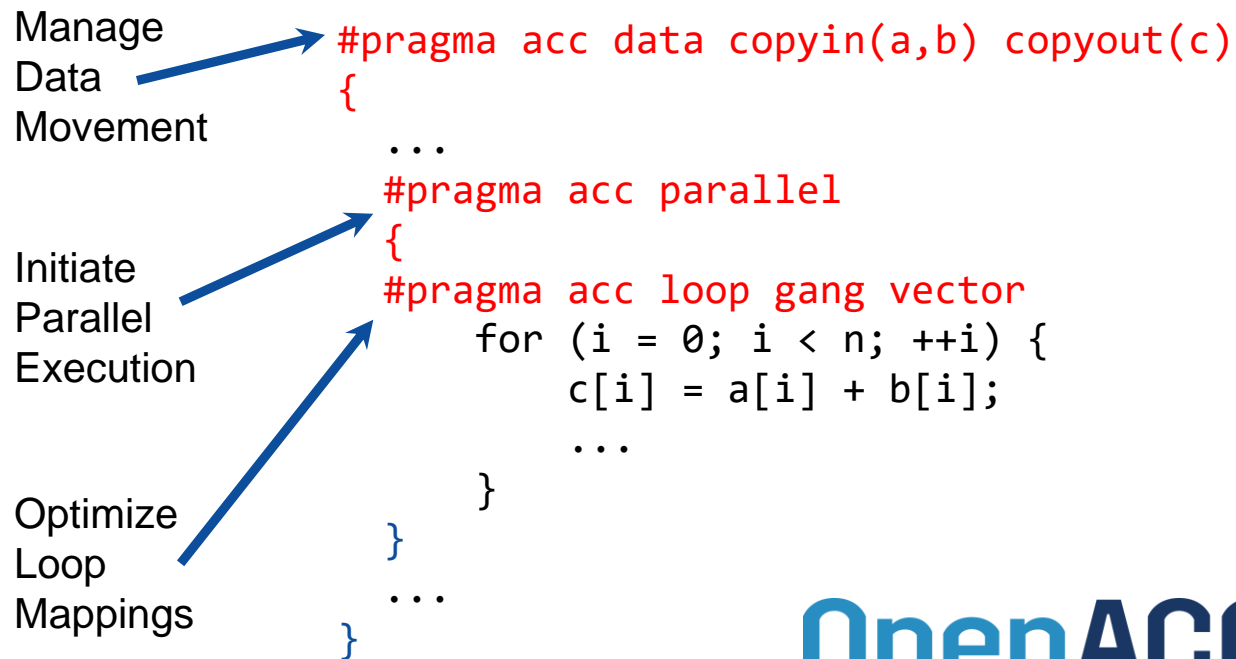
a directives-based **parallel programming model**
designed for **performance**
and **portability**.

Add Simple Compiler Directive

```
main()
{
    <serial code>
    #pragma acc kernels
    {
        <parallel code>
    }
}
```



OpenACC Directives



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

OpenACC
Directives for Accelerators

OPENACC

Incremental

- Maintain existing sequential code
- Add annotations to expose parallelism
- After verifying correctness, annotate more of the code

Single Source

- Rebuild the same code on multiple architectures
- Compiler determines how to parallelize for the desired machine
- Sequential code is maintained

Low Learning Curve

- OpenACC is meant to be easy to use, and easy to learn
- Programmer remains in familiar C, C++, or Fortran
- No reason to learn low-level details of the hardware.

OPENACC

Incremental

- Maintain existing sequential code
- Add annotations to expose parallelism
- After verifying correctness, annotate more of the code

Enhance Sequential Code

```
#pragma acc parallel loop
for( i = 0; i < N; i++ )
{
    < loop code >
}

#pragma acc parallel loop
for( i = 0; i < N; i++ )
{
    < loop code >
}
```

Begin with a working sequential code.

Parallelize it with OpenACC.

Rerun the code to verify correctness and performance

OPENACC

Supported Platforms

POWER

Sunway

x86 CPU

AMD GPU

NVIDIA GPU

PEZY-SC

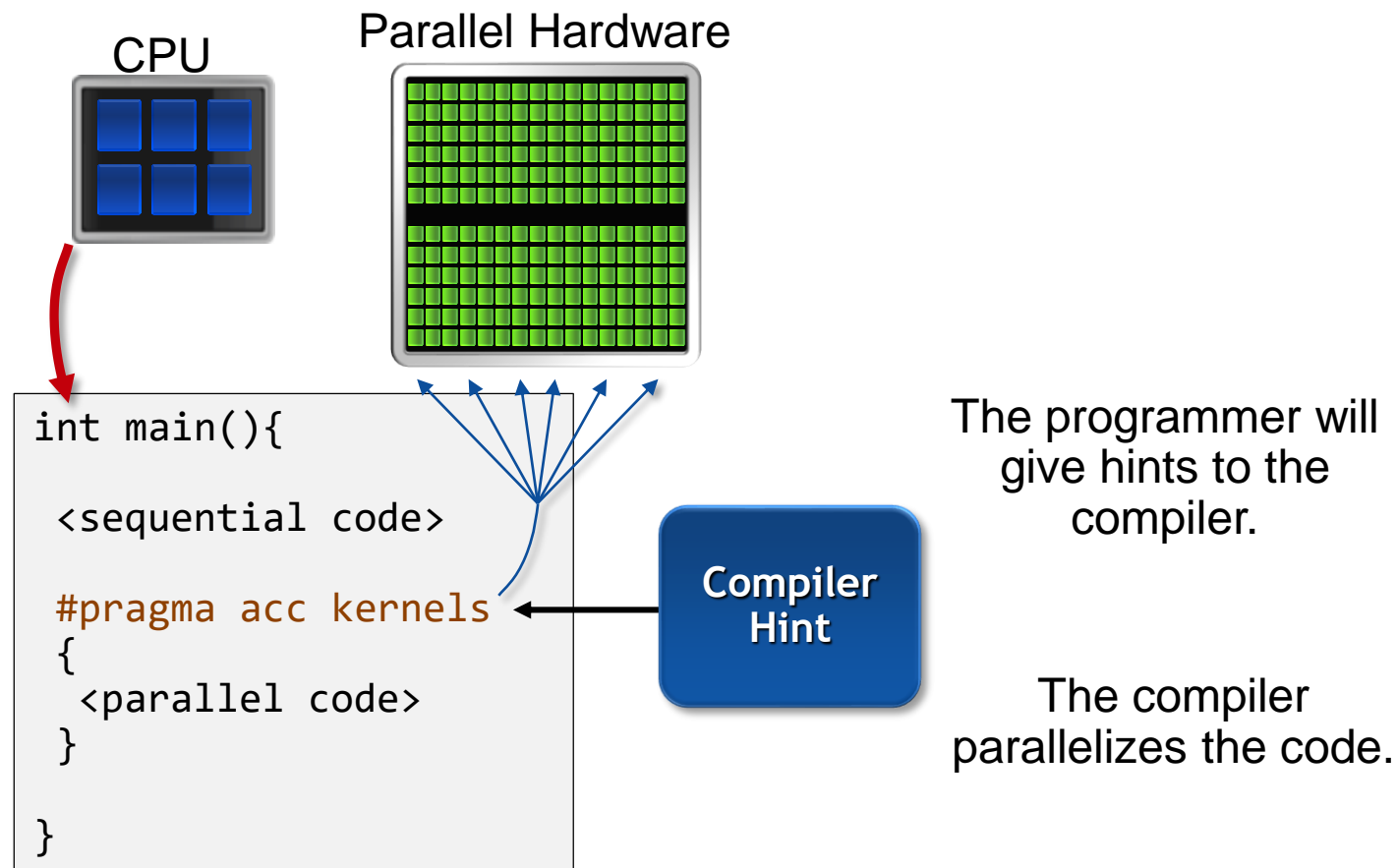
Single Source

- Rebuild the same code on multiple architectures
- Compiler determines how to parallelize for the desired machine
- Sequential code is maintained

The compiler can **ignore** your OpenACC code additions, so the same code can be used for **parallel** or **sequential** execution.

```
int main(){  
  
...  
  
#pragma acc parallel loop  
for(int i = 0; i < N; i++)  
    < loop code >  
}
```

OPENACC



Low Learning Curve

- OpenACC is meant to be easy to use, and easy to learn
- Programmer remains in familiar C, C++, or Fortran
- No reason to learn low-level details of the hardware.

OPENACC

Incremental

- Maintain existing sequential code
- Add annotations to expose parallelism
- After verifying correctness, annotate more of the code

Single Source

- Rebuild the same code on multiple architectures
- Compiler determines how to parallelize for the desired machine
- Sequential code is maintained

Low Learning Curve

- OpenACC is meant to be easy to use, and easy to learn
- Programmer remains in familiar C, C++, or Fortran
- No reason to learn low-level details of the hardware.

DIRECTIVE-BASED HPC PROGRAMMING

Who's Using OpenACC

3 OF TOP 5 HPC APPS



5 OF 13 CAAR CODES



2 OF LAST 9 FINALISTS



450 DOMAIN EXPERTS



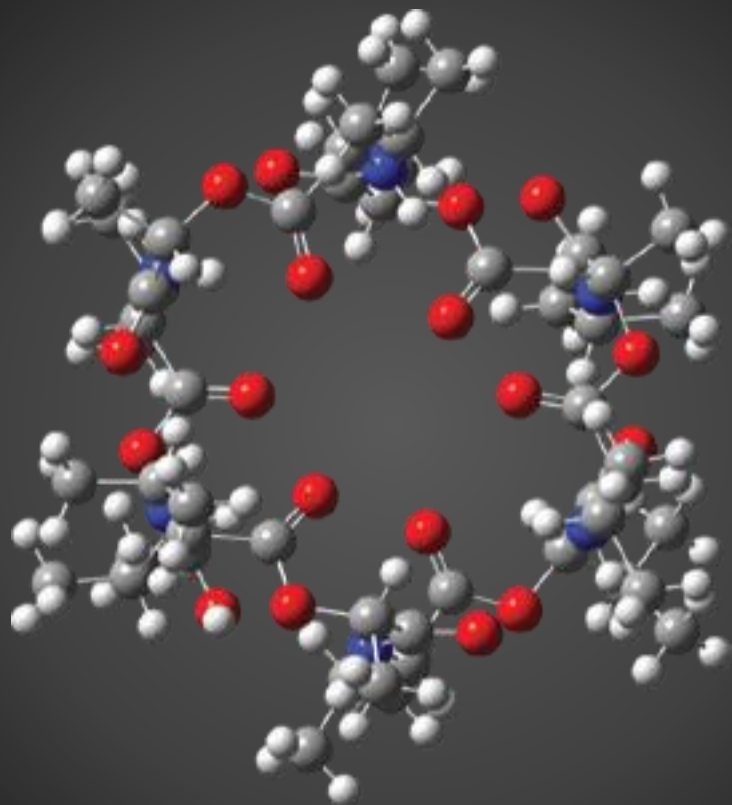
ACCELERATED APPS



100,000 DOWNLOADS



GAUSSIAN 16



Mike Frisch, Ph.D.
President and
CEO
Gaussian, Inc.



Using OpenACC allowed us to continue development of our fundamental algorithms and software capabilities simultaneously with the GPU-related work. In the end, we could use the same code base for SMP, cluster/network and GPU parallelism. PGI's compilers were essential to the success of our efforts.



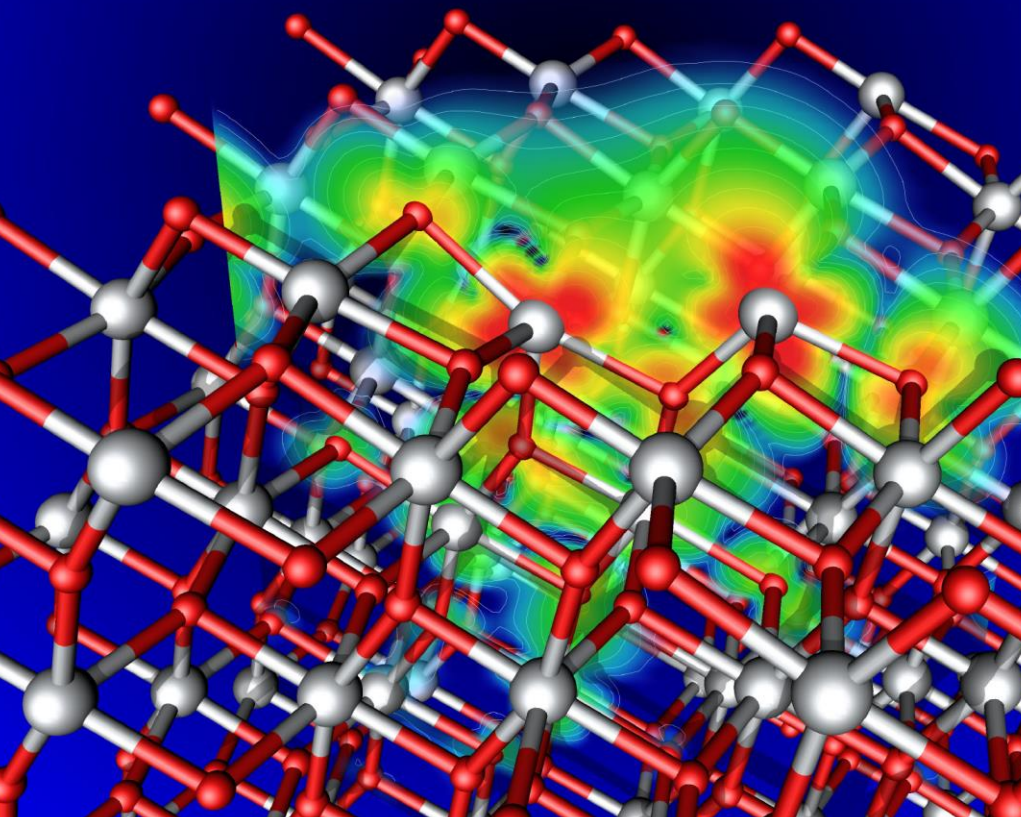
VASP



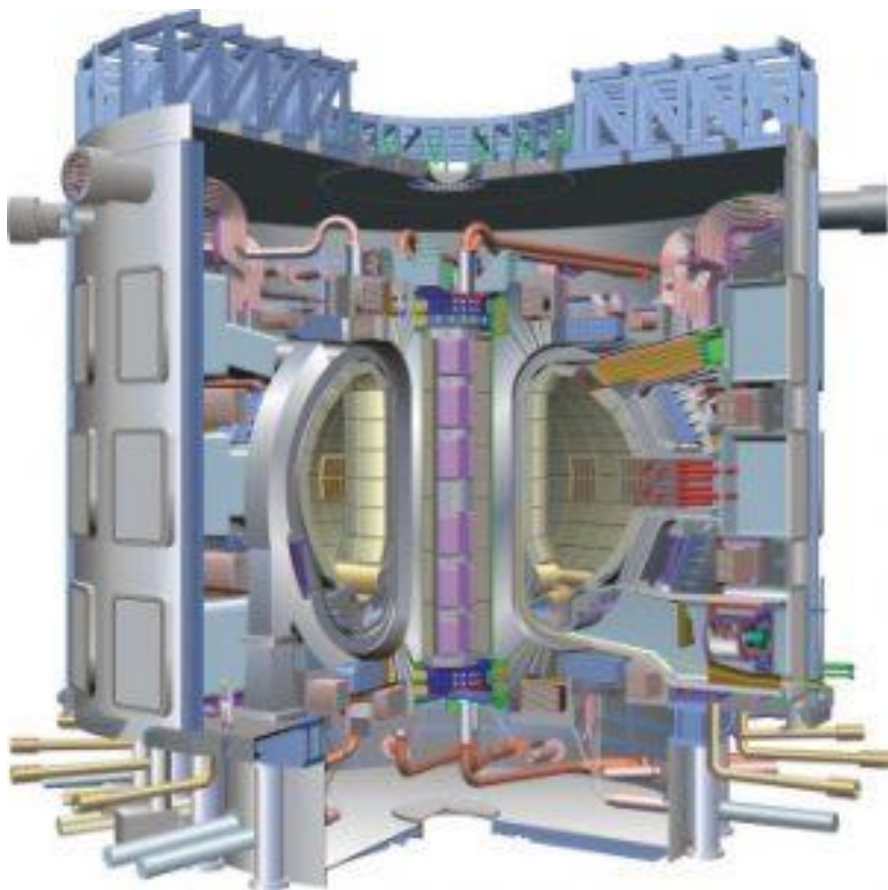
Prof. Georg Kresse
Computational Materials Physics
University of Vienna



For VASP, OpenACC is *the* way forward for GPU acceleration. Performance is similar and in some cases better than CUDA C, and OpenACC dramatically decreases GPU development and maintenance efforts. We're excited to collaborate with NVIDIA and PGI as an early adopter of CUDA Unified Memory.



GTC



Zhihong Lin
Professor and Principal Investigator
UC Irvine



Using OpenACC our scientists were able to achieve the acceleration needed for integrated fusion simulation with a minimum investment of time and effort in learning to program GPUs.



MAS



Ronald M. Caplan
Computational Scientist
Predictive Science Inc.



Adding OpenACC into MAS has given us the ability to migrate medium-sized simulations from a multi-node CPU cluster to a single multi-GPU server. The implementation yielded a portable single-source code for both CPU and GPU runs. Future work will add OpenACC to the remaining model features, enabling GPU-accelerated realistic solar storm modeling.



GAUSSIAN 16



Mike Friesch, Ph.D.
President and
CEO
Gaussian, Inc.

“ Using OpenACC allowed us to continue development of our fundamental algorithms and software capabilities simultaneously with the GPU-related work. In the end, we could use the same code base for SMP, cluster/ network and GPU parallelism. PGI's compilers were essential to the success of our efforts. ”

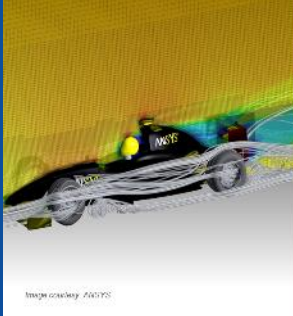


Image courtesy: ANSYS

ANSYS FLUENT



Sunil Sathya
Lead Software Developer
ANSYS Fluent

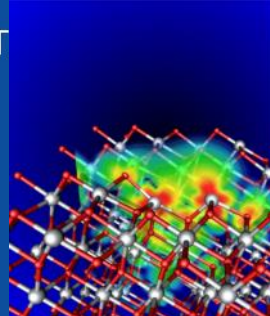
“ We've effectively used OpenACC for heterogeneous computing in ANSYS Fluent with impressive performance. We're now applying this work to more of our models and new platforms. ”

VASP



Prof. Georg Kresse
Computational Materials Physics
University of Vienna

“ For VASP, OpenACC is the way forward for GPU acceleration. Performance is similar and in some cases better than CUDA C, and OpenACC dramatically decreases GPU development and maintenance efforts. We're excited to collaborate with NVIDIA and PGI as an early adopter of CUDA Unified Memory. ”

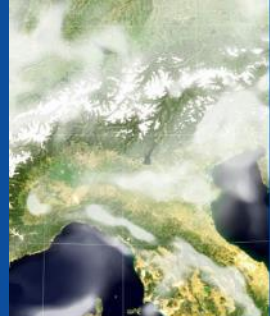


COSMO



Dr. Oliver Fuhrer
Senior Scientist
Matscorder

“ OpenACC made it practical to develop for GPU-based hardware while retaining a single source for almost all the COSMO physics code. ”



E3SM



Mark A. Taylor
Multiphysics Applications
Sandia

“ The CAAR project provided us with early access to Summit hardware and access to PGI compiler experts. Both of these were critical to our success. PGI's OpenACC support remains the best available and is competitive with much more intrusive programming model approaches. ”

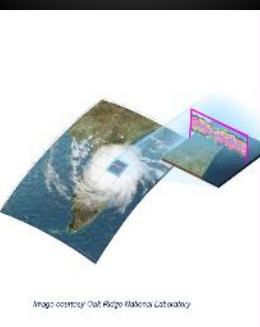


Image courtesy: Oak Ridge National Laboratory

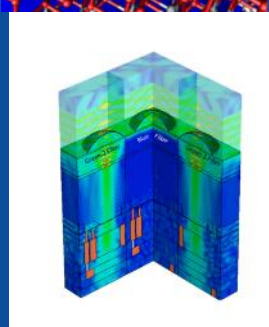


NUMECA FINE/Open



David Guizwiler
Lead Software Developer
NUMECA

“ Porting our unstructured C++ CFD solver FINE/Open to GPUs using OpenACC would have been impossible two or three years ago, but OpenACC has developed enough that we're now getting some really good results. ”



SYNOPSIS



Dr. Lutz Schneider
Senior R&D Engineer
Synopsis Inc.

“ Using OpenACC, we've GPU-accelerated the Synopsis TCAD Sentaurus Device EMW simulator to speed up optical simulations of image sensors. GPUs are key to improving simulation throughput in the design of advanced image sensors. ”



Image courtesy: NCAR

MPAS-A



Richard Loft
Director, Technology Development
NCAR

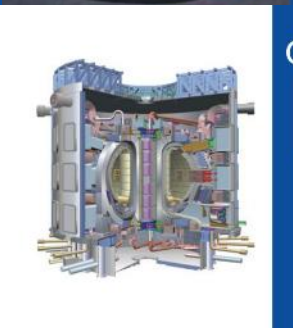
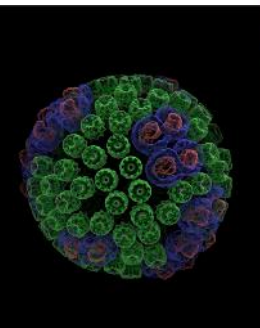
“ Our team has been evaluating OpenACC as a pathway to performance portability for the Model for Prediction (MPAS) atmospheric model. Using this approach on the MPAS dynamical core, we have achieved performance on a single P100 GPU equivalent to 2-7 dual socketed Intel Xeon nodes on our new Cheyenne supercomputer. ”

VMD



John Stone
Senior Research Programmer
Buckham Institute
University of Illinois

“ Due to Amdahl's law, we need to port more parts of our code to the GPU if we're going to speed it up. But the sheer number of routines poses a challenge. OpenACC directives give us a low-cost approach to getting at least some speed-up out of these second-tier routines. In many cases it's completely sufficient because with the current algorithms, GPU performance is bandwidth-bound. ”



GTC



Zhihong Lin
Professor and Principal Investigator
UC Irvine

“ Using OpenACC our scientists were able to achieve the acceleration needed for integrated fusion simulation with a minimum investment of time and effort in learning to program GPUs. ”

OpenACC
More Science. Less Programming

GAMERA



Takuma Yamaguchi, Kohji Fujita, Shojiro Ichimaru, Masaki Hori, Lutz Schneider
The University of Tokyo

“ With OpenACC and a compute node based on NVIDIA's Tesla P100 GPU, we achieved more than a 14X speed up over a K Computer node running our earthquake disaster simulation code. ”



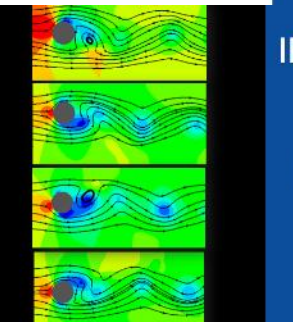
Map courtesy: University of Tokyo

SANJEEVINI



Abhilash Jayaraj
Project Scientist
Indian Institute of Technology
New Delhi

“ In an academic environment maintenance and speeding up existing codes is a tedious task. OpenACC provides a great platform for computational scientists to accomplish both tasks without involving a lot of efforts or manpower in speeding up the entire computational task. ”

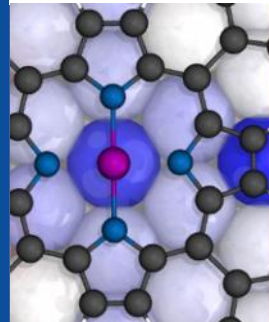


IBM-CFD



Sorenath Roy
Assistant Professor
Mechanical Engineering Department
Indian Institute of Technology Kharagpur

“ OpenACC can prove to be a handy tool for computational engineers and researchers to obtain fast solution of non-linear dynamics problem in immersed boundary incompressible CFD. We have obtained order of magnitude reduction in compiling time by porting several components of our legacy codes to GPU. Especially the routines involving search algorithm and matrix solvers have been well-accelerated to improve the overall scalability of the code. ”

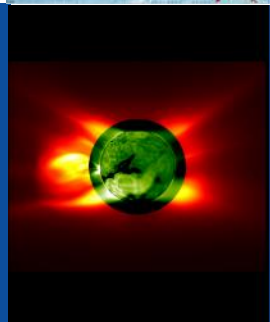


PWscf (Quantum ESPRESSO)



Filippo Sgalla
Senior Contributor
Quantum ESPRESSO group

“ CUDA Fortran gives us the full performance potential of the CUDA programming model and NVIDIA GPUs. While leveraging the potential of explicit data movement, ISCUP_KERNELS directives give us productivity and source code maintainability. It's the best of both worlds. ”



MAS



Ronald M. Caplan
Computational Scientist
Predictive Science Inc.

“ Adding OpenACC into MAS has given us the ability to migrate medium-sized simulations from a multi node CPU cluster to a single multi-GPU server. The implementation yielded a portable single-source code for both CPU and GPU runs. Future work will add OpenACC to the remaining model features, enabling GPU accelerated realistic solar storm modeling. ”

OPENACC SYNTAX

OPENACC SYNTAX

Syntax for using OpenACC directives in code

C/C++

```
#pragma acc directive clauses  
<code>
```

Fortran

```
!$acc directive clauses  
<code>
```

- A ***pragma*** in C/C++ gives instructions to the compiler on how to compile the code. Compilers that do not understand a particular pragma can freely ignore it.
- A ***directive*** in Fortran is a specially formatted comment that likewise instructions the compiler in its compilation of the code and can be freely ignored.
- “***acc***” informs the compiler that what will come is an OpenACC directive
- ***Directives*** are commands in OpenACC for altering our code.
- ***Clauses*** are specifiers or additions to directives.

EXAMPLE CODE

LAPLACE HEAT TRANSFER

Introduction to lab code - visual

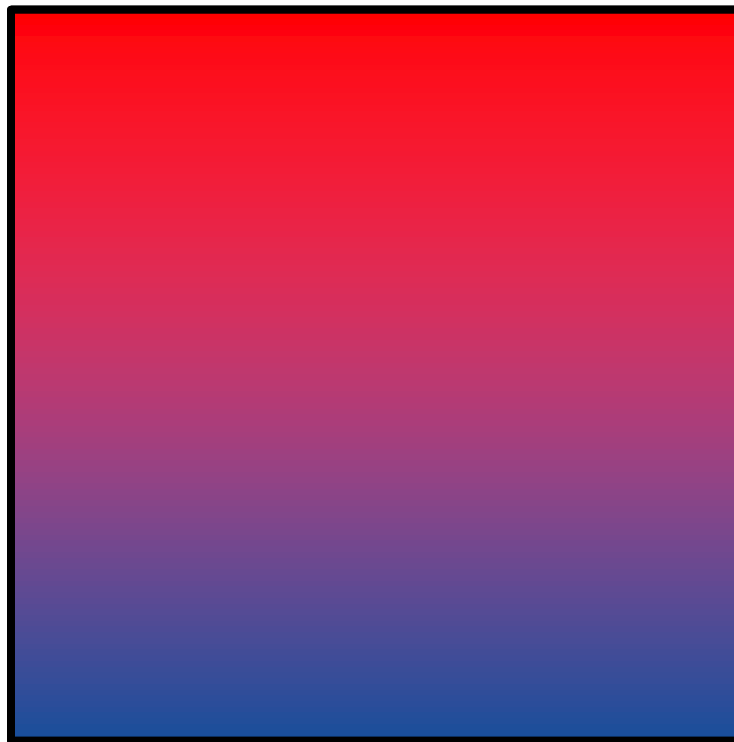
We will observe a simple simulation of heat distributing across a metal plate.

We will apply a consistent heat to the top of the plate.

Then, we will simulate the heat distributing across the plate.

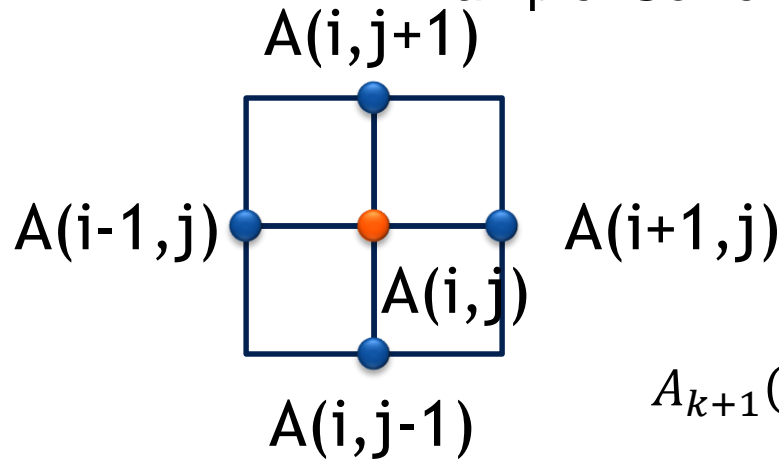
Very Hot

Room Temp



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

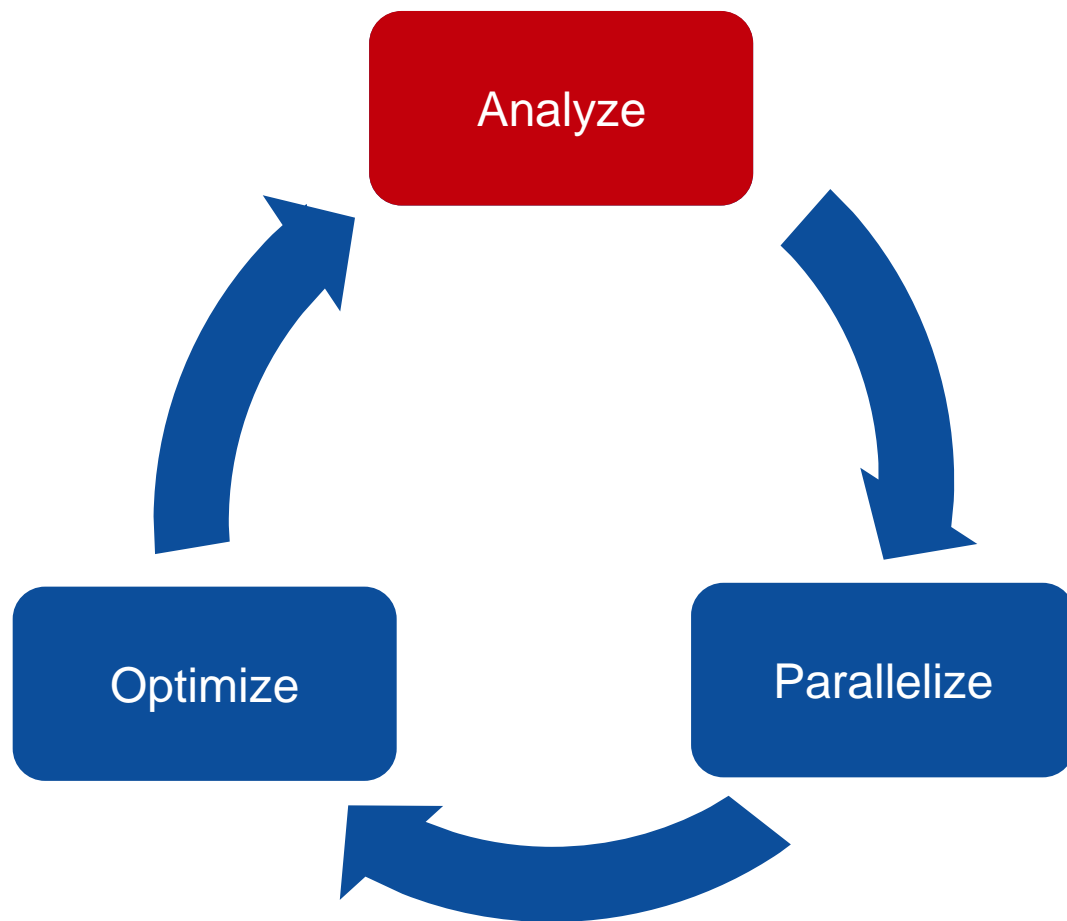


Swap input/output arrays

PROFILE-DRIVEN DEVELOPMENT

OPENACC DEVELOPMENT CYCLE

- **Analyze** your code to determine most likely places needing parallelization or optimization.
- **Parallelize** your code by starting with the most time consuming parts and check for correctness.
- **Optimize** your code to improve observed speed-up from parallelization.



PROFILING SEQUENTIAL CODE

Profile Your Code

Obtain detailed information about how the code ran.

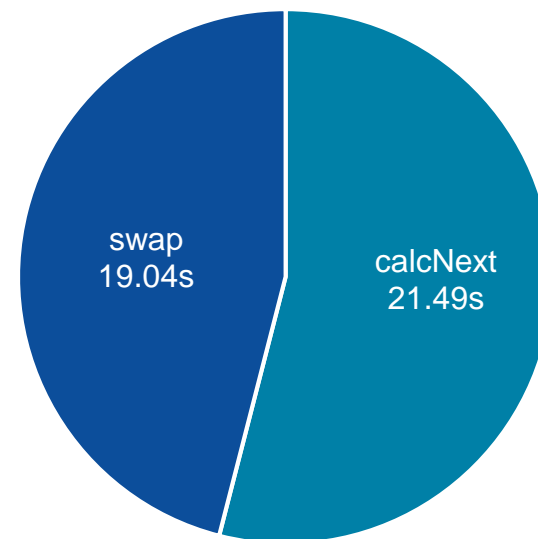
This can include information such as:

- Total runtime
- Runtime of individual routines
- Hardware counters

Identify the portions of code that took the longest to run. We want to focus on these “hotspots” when parallelizing.

Lab Code: Laplace Heat Transfer

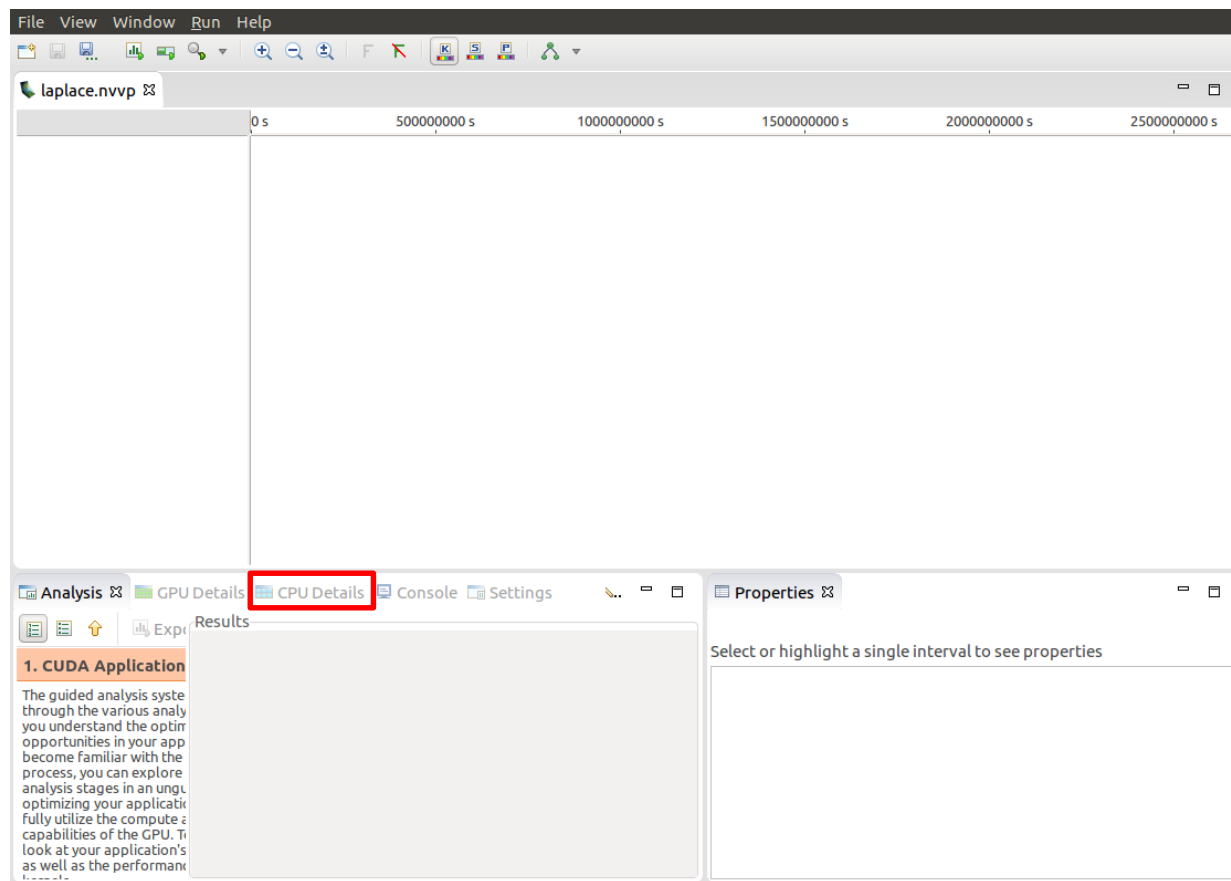
Total Runtime: 39.43 seconds



PROFILING SEQUENTIAL CODE

First sight when using PGPROF

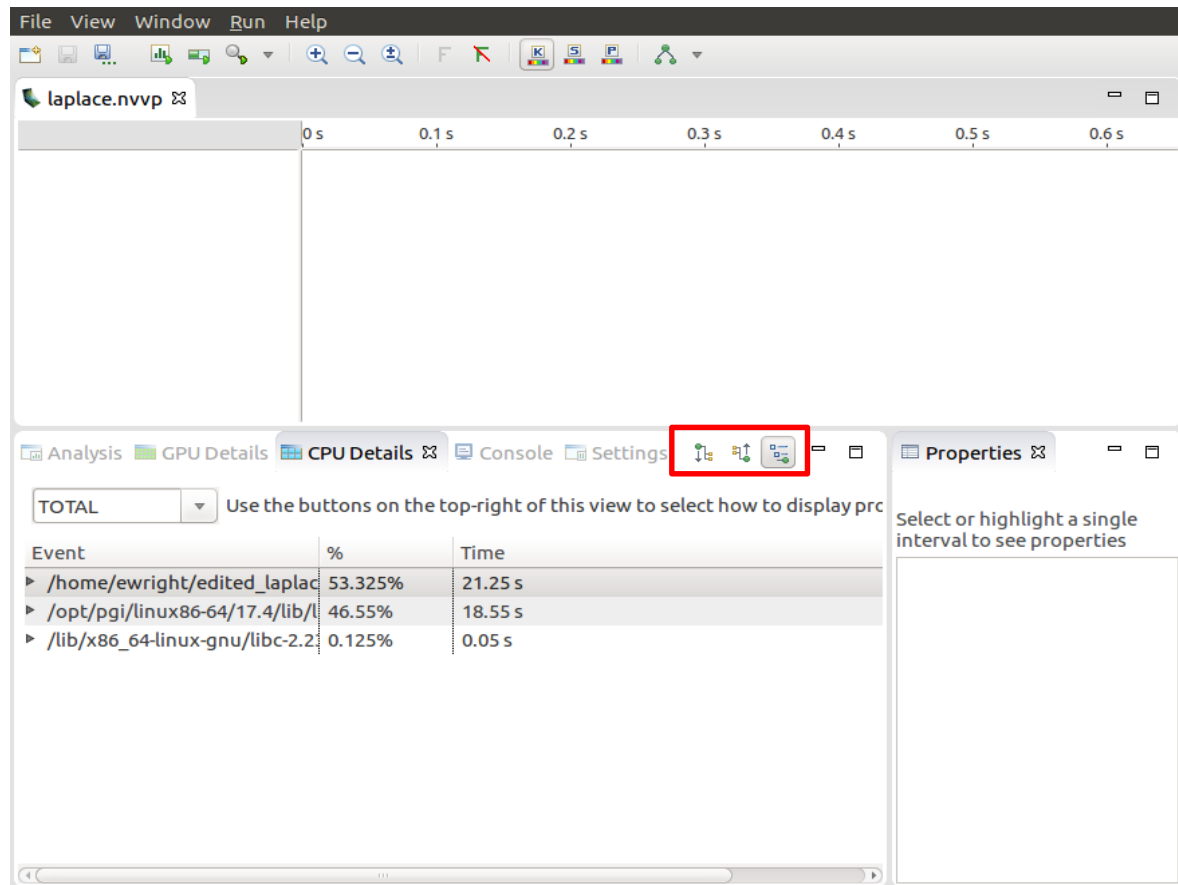
- Profiling a simple, sequential code
- Our sequential program will run on the CPU
- To view information about how our code ran, we should select the “CPU Details” tab



PROFILING SEQUENTIAL CODE

CPU Details

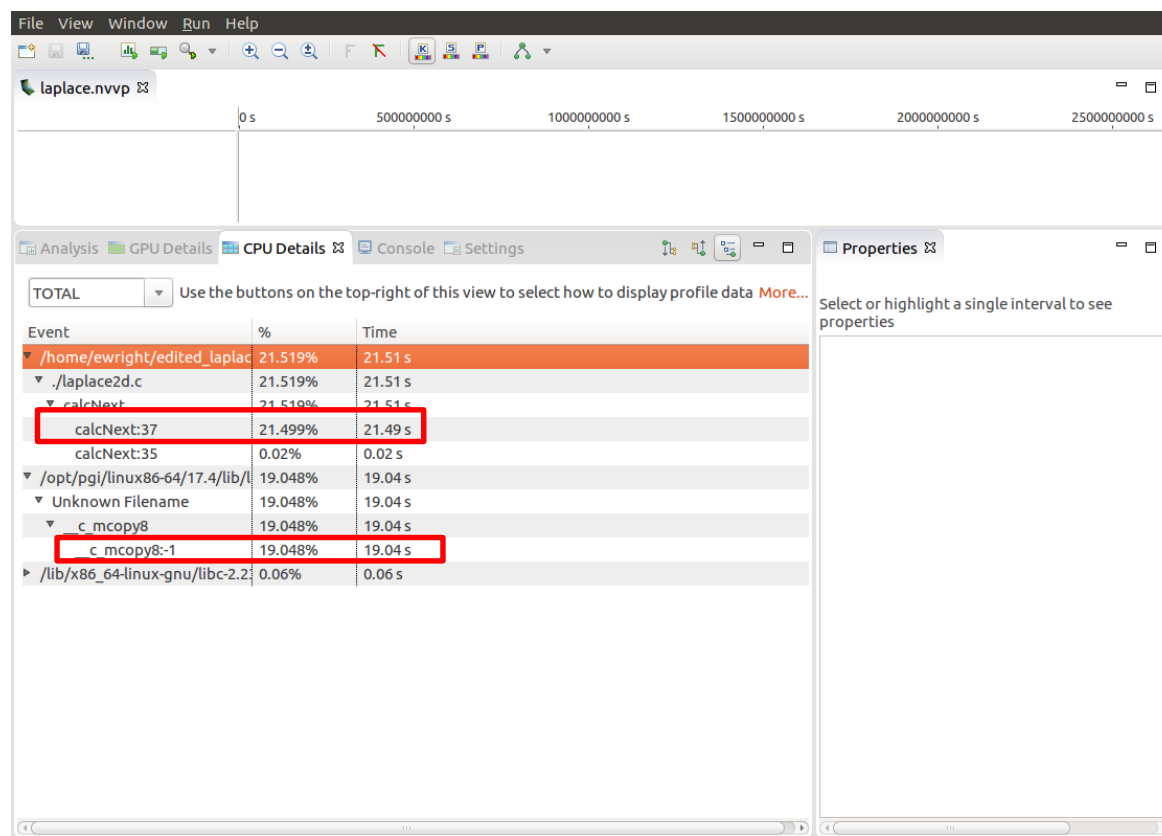
- Within the “CPU Details” tab, we can see the various parts of our code, and how long they took to run
- We can reorganize this info using the three options in the top-right portion of the tab
- We will expand this information, and see more details about our code



PROFILING SEQUENTIAL CODE

CPU Details

- We can see that there are two places that our code is spending most of its time
- 21.49 seconds in the “calcNext” function
- 19.04 seconds in a memcpy function
- The c_memcpy8 that we see is actually a compiler optimization that is being applied to our “swap” function



PROFILING SEQUENTIAL CODE

PGPROF

- We are also able to select the different elements in the CPU Details by double-clicking to open the associated source code
- Here we have selected the “calcNext:37” element, which opened up our code to show the exact line (line 37) that is associated with that element

The screenshot shows the PGPROF application interface. The top pane displays the source code for `laplace2d.c`, with line 37 highlighted. The bottom pane shows the CPU Details profile, with the `calcNext:37` event selected and highlighted in orange.

Source Code (laplace2d.c):

```
30 #define OFFSET(x, y, m) ((x)*(m)) + (y)
31
32 double calcNext(double *restrict A, double *restrict Anew, int m, int n)
33 {
34     double error = 0.0;
35     for( int j = 1; j < n-1; j++)
36     {
37         for( int i = 1; i < m-1; i++)
38         {
39             Anew[OFFSET(j, i, m)] = 0.25 * ( A[OFFSET(j, i+1, m)] + A[OFFSET(j, i-1, m)]
40             + A[OFFSET(j-1, i, m)] + A[OFFSET(j+1, i, m)]);
41             error = fmax( error, fabs(Anew[OFFSET(j, i, m)] - A[OFFSET(j, i, m)]));
42         }
43     }
44     return error;
45 }
46
47 void swap(double *restrict A, double *restrict Anew, int m, int n)
48 {
49     for( int j = 1; j < n-1; j++)
50     {
51         for( int i = 1; i < m-1; i++)
52         {
53             ...
54         }
55     }
56 }
```

CPU Details Profile:

Event	%	Time
TOTAL		
▼ /home/ewright/edited_laplace2d.c	21.519%	21.51 s
▼ ./laplace2d.c	21.519%	21.51 s
▼ calcNext	21.519%	21.51 s
calcNext:37	21.499%	21.49 s
calcNext:35	0.02%	0.02 s
▶ /opt/pgi/linux86-64/17.4/lib/lin	19.048%	19.04 s
▶ /lib/x86_64-linux-gnu/libc-2.2	0.06%	0.06 s

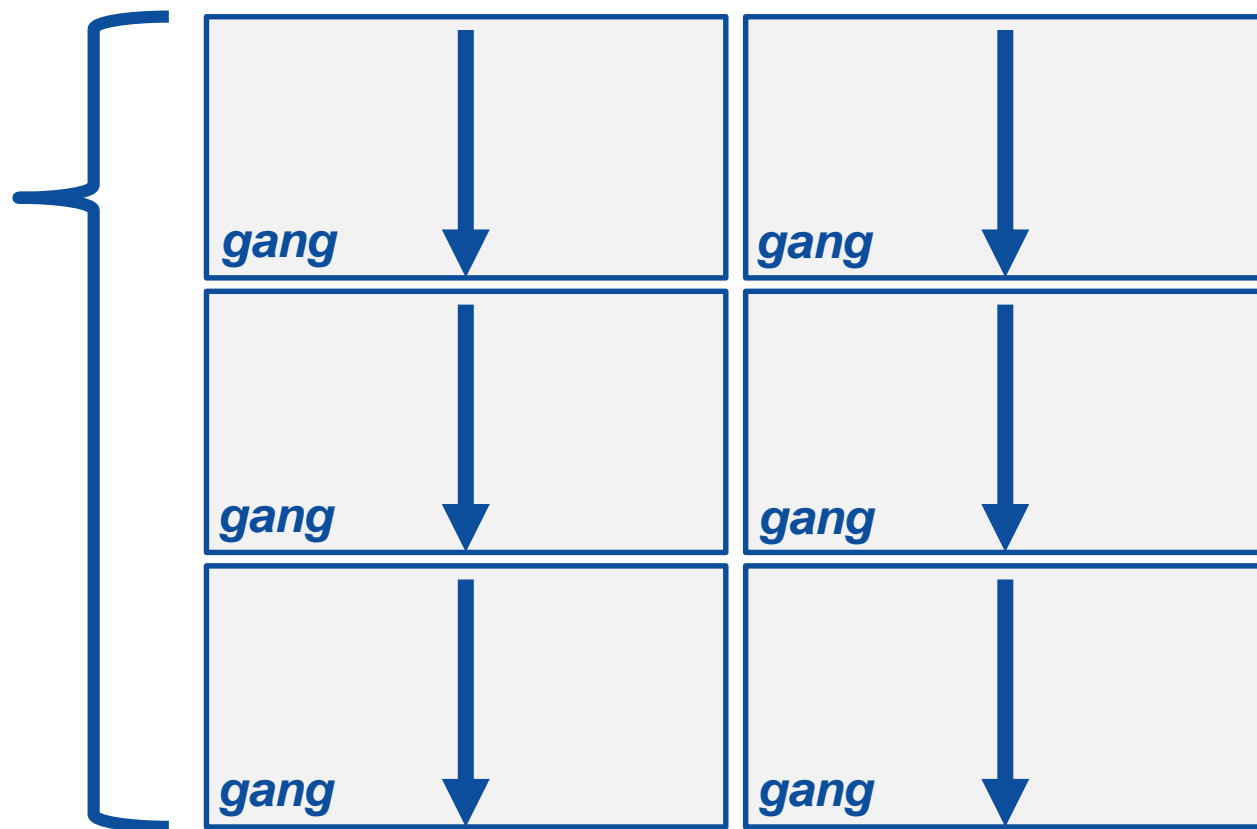
OPENACC PARALLEL LOOP DIRECTIVE

OPENACC PARALLEL DIRECTIVE

Expressing parallelism

```
#pragma acc parallel  
{
```

When encountering the ***parallel*** directive, the compiler will generate *1 or more parallel gangs*, which execute redundantly.



```
}
```

OPENACC PARALLEL DIRECTIVE

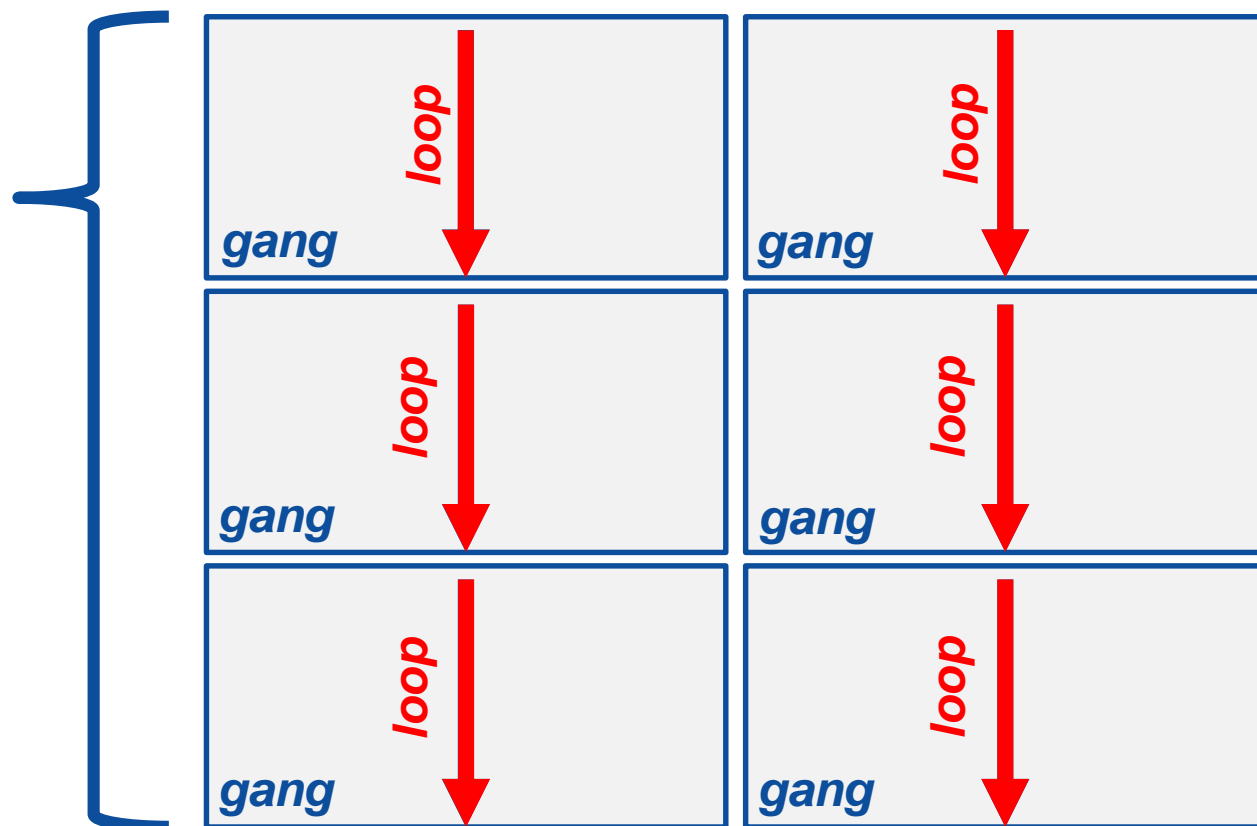
Expressing parallelism

```
#pragma acc parallel  
{
```

```
    loop  
    for(int i = 0; i < N; i++)  
    {  
        // Do Something  
    }
```

```
}
```

This loop will be
executed redundantly
on each gang



OPENACC PARALLEL DIRECTIVE

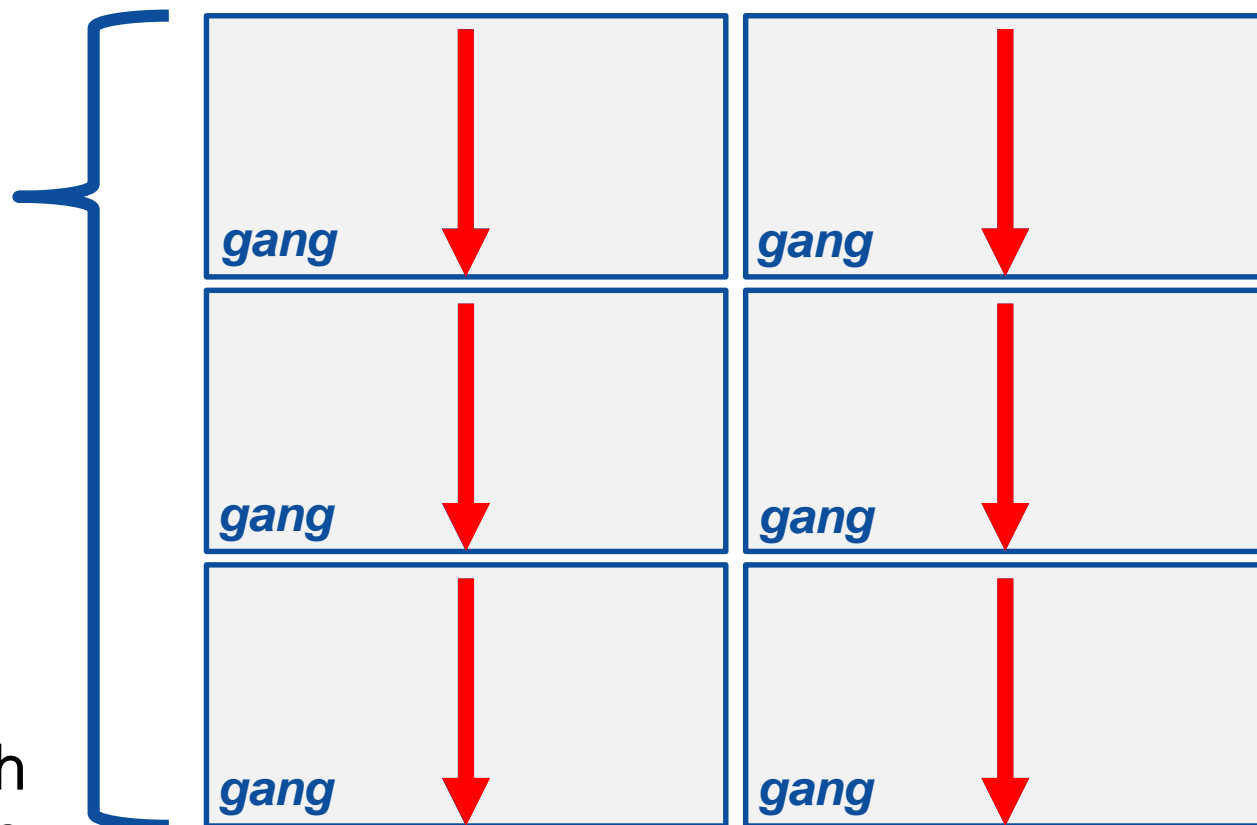
Expressing parallelism

```
#pragma acc parallel  
{
```

```
  for(int i = 0; i < N; i++)  
  {  
    // Do Something  
  }
```

```
}
```

This means that each **gang** will execute the entire loop



OPENACC PARALLEL DIRECTIVE

Parallelizing a single loop

C/C++

```
#pragma acc parallel
{
    #pragma acc loop
    for(int i = 0; i < N; i++)
        a[i] = 0;
}
```

Fortran

```
!$acc parallel
!$acc loop
do i = 1, N
    a(i) = 0
end do
!$acc end parallel
```

- Use a **parallel** directive to mark a region of code where you want parallel execution to occur
- This parallel region is marked by curly braces in C/C++ or a start and end directive in Fortran
- The **loop** directive is used to instruct the compiler to parallelize the iterations of the next loop to run across the parallel gangs

OPENACC PARALLEL DIRECTIVE

Parallelizing a single loop

C/C++

```
#pragma acc parallel loop
for(int i = 0; i < N; i++)
    a[i] = 0;
```

Fortran

```
!$acc parallel loop
do i = 1, N
    a(i) = 0
end do
```

- This pattern is so common that you can do all of this in a single line of code
- In this example, the parallel loop directive applies to the next loop
- This directive both marks the region for parallel execution and distributes the iterations of the loop.
- When applied to a loop with a data dependency, parallel loop may produce incorrect results

OPENACC PARALLEL DIRECTIVE

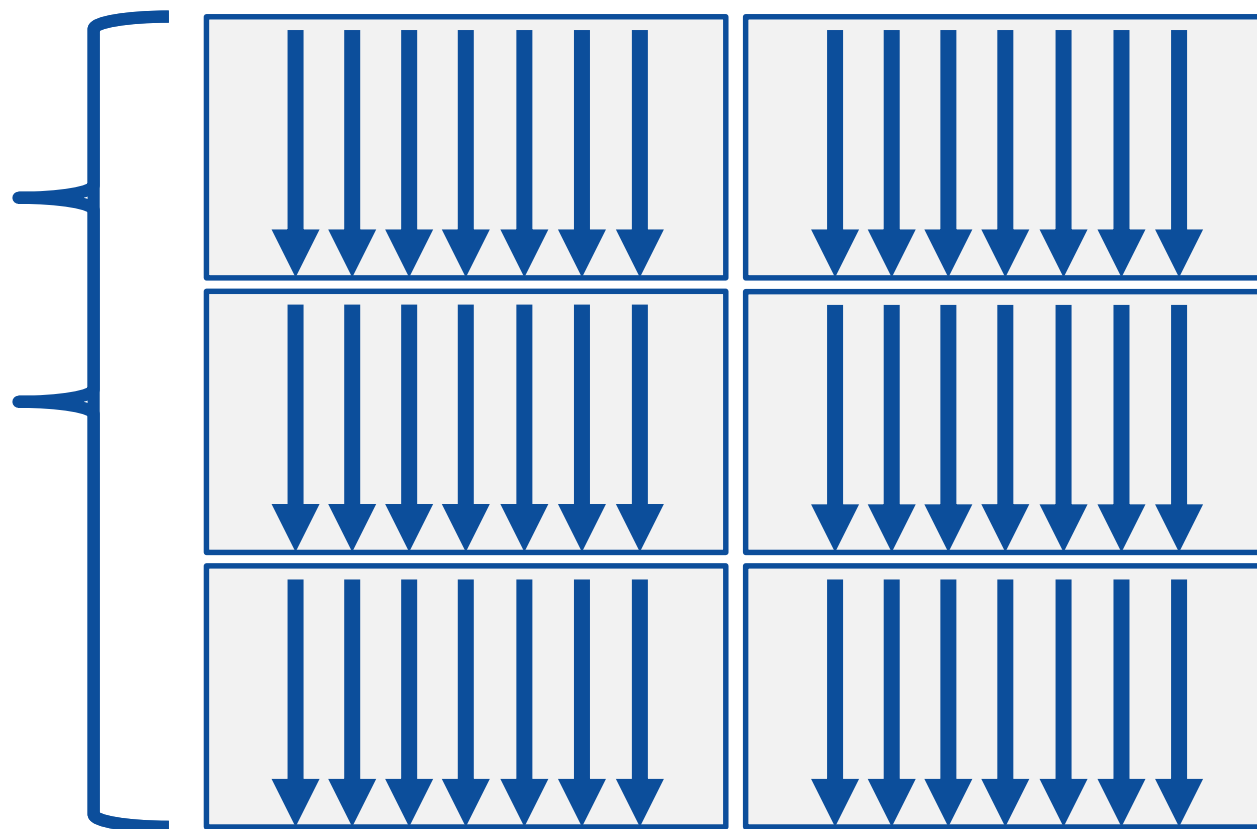
Expressing parallelism

```
#pragma acc parallel  
{
```

```
    #pragma acc loop  
    for(int i = 0; i < N; i++)  
    {  
        // Do Something  
    }
```

```
}
```

The *loop* directive informs the compiler which loops to parallelize.



OPENACC PARALLEL LOOP DIRECTIVE

Parallelizing many loops

```
#pragma acc parallel loop
for(int i = 0; i < N; i++)
    a[i] = 0;

#pragma acc parallel loop
for(int j = 0; j < M; j++)
    b[j] = 0;
```

- To parallelize multiple loops, each loop should be accompanied by a parallel directive
- Each parallel loop can have different loop boundaries and loop optimizations
- Each parallel loop can be parallelized in a different way
- This is the recommended way to parallelize multiple loops. Attempting to parallelize multiple loops within the same parallel region may give performance issues or unexpected results

PARALLELIZE WITH OPENACC PARALLEL LOOP

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

    #pragma acc parallel loop reduction(max:err)
    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 parallel loop
    for( int j = 1; j < n-1; j++) {
        for( int i = 1; i < m-1; i++ ) {
            A[j][i] = Anew[j][i];
        }
    }
    iter++;
}
```

Parallelize first loop nest,
max *reduction* required.

Parallelize second loop.

We didn't detail *how* to
parallelize the loops, just *which*
loops to parallelize.

REDUCTION CLAUSE

- The **reduction** clause takes many values and “reduces” them to a single value, such as in a sum or maximum
- Each thread calculates its part
- The compiler will perform a final reduction to produce a **single global result** using the specified operation

```
for( i = 0; i < size; i++ )  
    for( j = 0; j < size; j++ )  
        for( k = 0; k < size; k++ )  
            c[i][j] += a[i][k] * b[k][j];
```

```
for( i = 0; i < size; i++ )  
    for( j = 0; j < size; j++ )  
        double tmp = 0.0f;  
        #pragma parallel acc loop \  
            reduction(+:tmp)  
        for( k = 0; k < size; k++ )  
            tmp += a[i][k] * b[k][j];  
        c[i][j] = tmp;
```

REDUCTION CLAUSE OPERATORS

Operator	Description	Example
+	Addition/Summation	<code>reduction(+:sum)</code>
*	Multiplication/Product	<code>reduction(*:product)</code>
max	Maximum value	<code>reduction(max:maximum)</code>
min	Minimum value	<code>reduction(min:minimum)</code>
&	Bitwise and	<code>reduction(&:val)</code>
 	Bitwise or	<code>reduction(:val)</code>
&&	Logical and	<code>reduction(&&:val)</code>
 	Logical or	<code>reduction(:val)</code>

BUILD AND RUN THE CODE

PGI COMPILER BASICS

pgcc, pgc++ and pgfortran

- The command to compile C code is 'pgcc'
- The command to compile C++ code is 'pgc++'
- The command to compile Fortran code is 'pgfortran'
- The **-fast** flag instructs the compiler to optimize the code to the best of its abilities

```
$ pgcc -fast main.c  
$ pgc++ -fast main.cpp  
$ pgfortran -fast main.F90
```


PGI COMPILER BASICS

-Minfo flag

- The -Minfo flag will instruct the compiler to print feedback about the compiled code
- **-Minfo=accel** will give us information about what parts of the code were accelerated via OpenACC
- **-Minfo=opt** will give information about all code optimizations
- **-Minfo=all** will give all code feedback, whether positive or negative

```
$ pgcc -fast -Minfo=all main.c  
$ pgc++ -fast -Minfo=all main.cpp  
$ pgfortran -fast -Minfo=all main.f90
```

PGI COMPILER BASICS

-ta flag

- The -ta flag enables building OpenACC code for a “Target Accelerator” (TA)
- -ta=multicore – Build the code to run across threads on a multicore CPU
- -ta=tesla:managed – Build the code for an NVIDIA (Tesla) GPU and manage the data movement for me (more next week)

```
$ pgcc -fast -Minfo=accel -ta=tesla:managed main.c  
$ pgc++ -fast -Minfo=accel -ta=tesla:managed main.cpp  
$ pgfortran -fast -Minfo=accel -ta=tesla:managed main.f90
```

BUILDING THE CODE (MULTICORE)

```
$ pgcc -fast -ta=multicore -Minfo=accel laplace2d_uvm.c
```

```
main:
```

```
63, Generating Multicore code
```

```
64, #pragma acc loop gang
```

```
64, Accelerator restriction: size of the GPU copy of Anew,A is unknown  
Generating reduction(max:error)
```

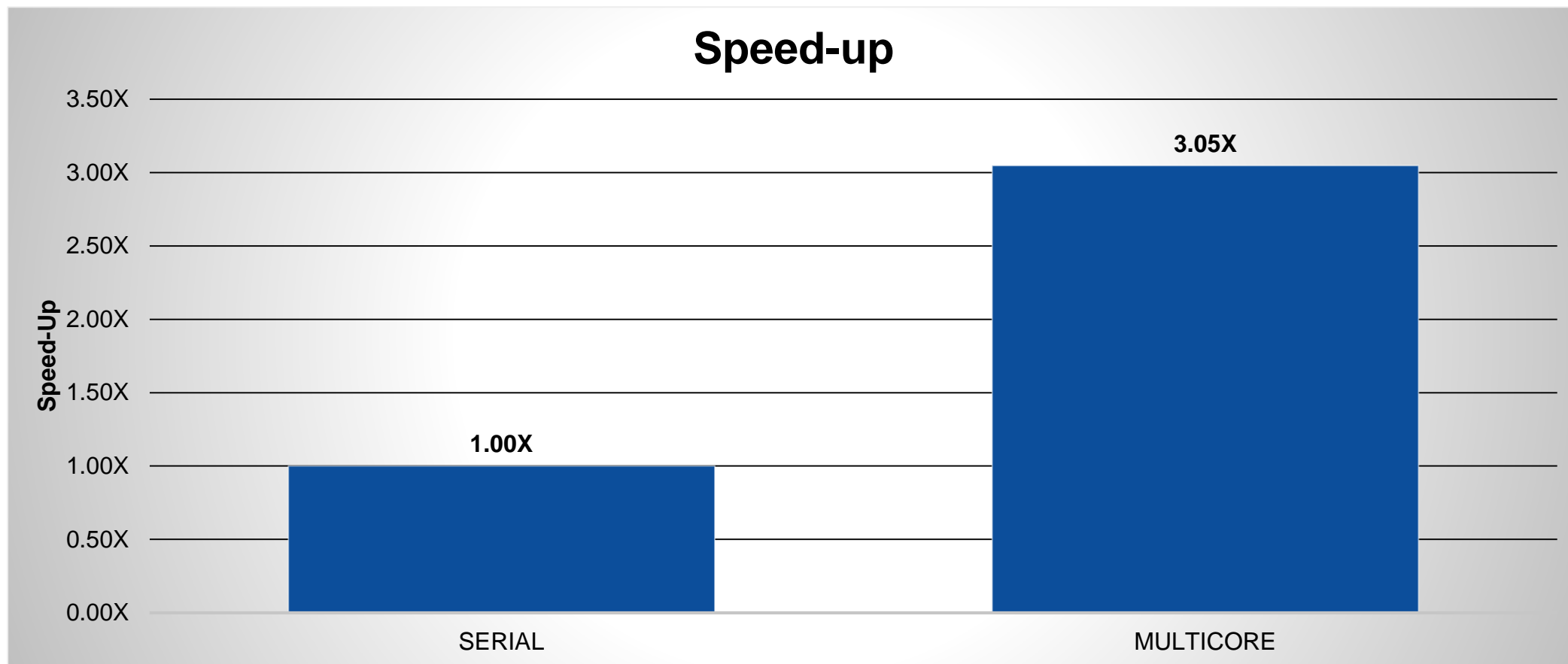
```
66, Loop is parallelizable
```

```
74, Generating Multicore code
```

```
75, #pragma acc loop gang
```

```
75, Accelerator restriction: size of the GPU copy of Anew,A is unknown  
77, Loop is parallelizable
```

OPENACC SPEED-UP

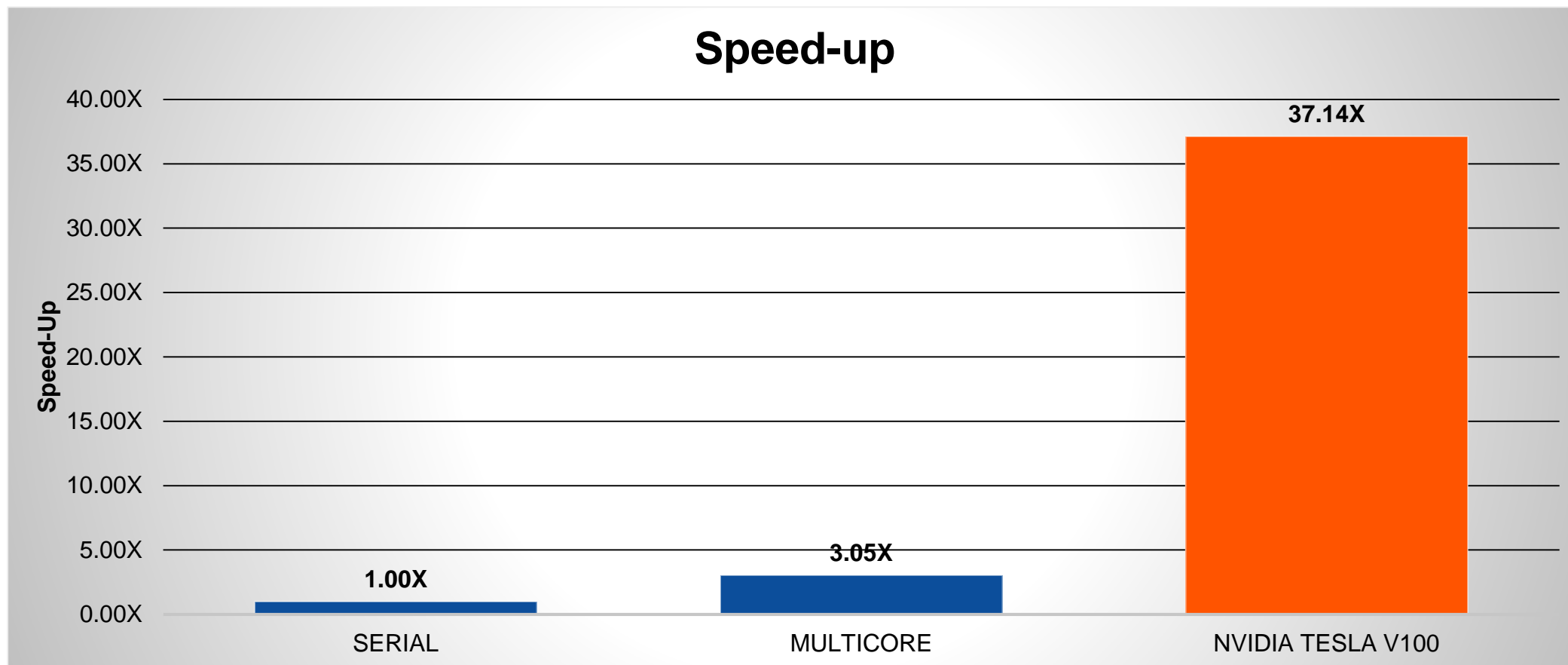


BUILDING THE CODE (GPU)

```
$ pgcc -fast -ta=tesla:managed -Minfo=accel laplace2d_uvm.c  
main:
```

```
63, Accelerator kernel generated  
Generating Tesla code  
64, #pragma acc loop gang /* blockIdx.x */  
Generating reduction(max:error)  
66, #pragma acc loop vector(128) /* threadIdx.x */  
63, Generating implicit copyin(A[:])  
Generating implicit copyout(Anew[:])  
Generating implicit copy(error)  
66, Loop is parallelizable  
74, Accelerator kernel generated  
Generating Tesla code  
75, #pragma acc loop gang /* blockIdx.x */  
77, #pragma acc loop vector(128) /* threadIdx.x */  
74, Generating implicit copyin(Anew[:])  
Generating implicit copyout(A[:])  
77, Loop is parallelizable
```

OPENACC SPEED-UP



CLOSING REMARKS

KEY CONCEPTS

This week we discussed...

- What is OpenACC
- How profile-driven programming helps you write better code
- How to parallelize loops using OpenACC's **parallel loop** directive to improve time to solution

Next Week:

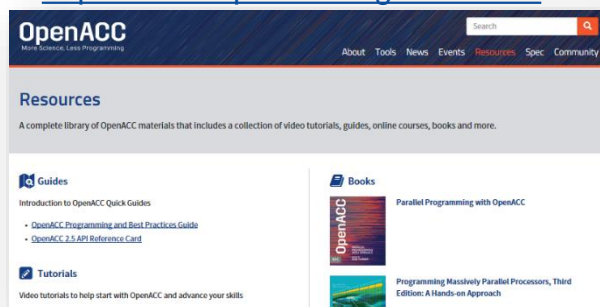
- Managing your data with OpenACC

OPENACC RESOURCES

Guides • Talks • Tutorials • Videos • Books • Spec • Code Samples • Teaching Materials • Events • Success Stories • Courses • Slack • Stack Overflow

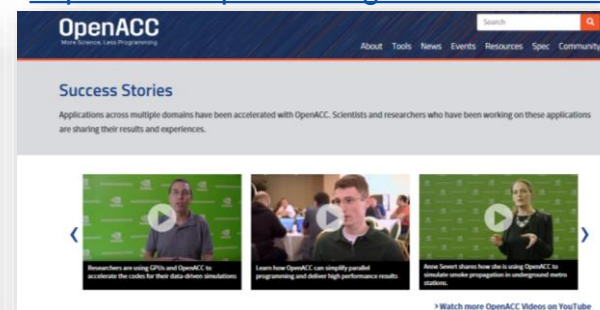
Resources

<https://www.openacc.org/resources>



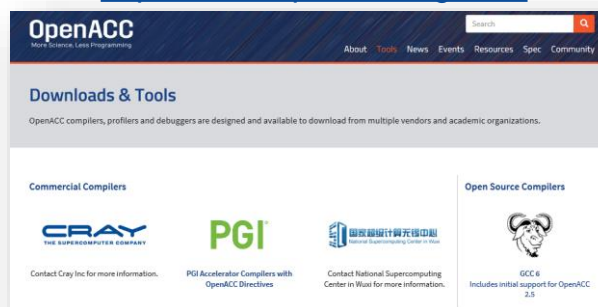
Success Stories

<https://www.openacc.org/success-stories>



Compilers and Tools

<https://www.openacc.org/tools>



Events

<https://www.openacc.org/events>

