



**PGI**® COMPILERS  
& TOOLS

# PGI UPDATES FOR 2019

Michael Wolfe  
GTC, March 2019



# AGENDA

## PGI Introduction

Platforms, Programming models, Performance, Interoperability

## What's new in PGI 2019

Features, Targets, CUDA Fortran, OpenACC, Containers and Clouds

## PGI Directions

OpenACC, OpenMP, Languages

# PGI — THE NVIDIA HPC SDK

Fortran/C/C++ Compilers and Tools

Multi-Platform Solution

X86-64 and OpenPOWER Multicore CPUs

NVIDIA Tesla GPUs

Supported on Linux, macOS, Windows

High Performance Parallelism

Optimization, SIMD Vectorization

OpenMP and OpenACC on CPU

OpenACC and CUDA Fortran on GPU

Interoperability

GNU compatible C++

Rogue Wave Totalview, Arm Allinea DDT, nvprof

OpenACC, CUDA Fortran, CUDA C++

# PGI®

## The Compilers & Tools for Supercomputing



19.4 Coming Soon!

# PGI COMPILERS FOR EVERYONE

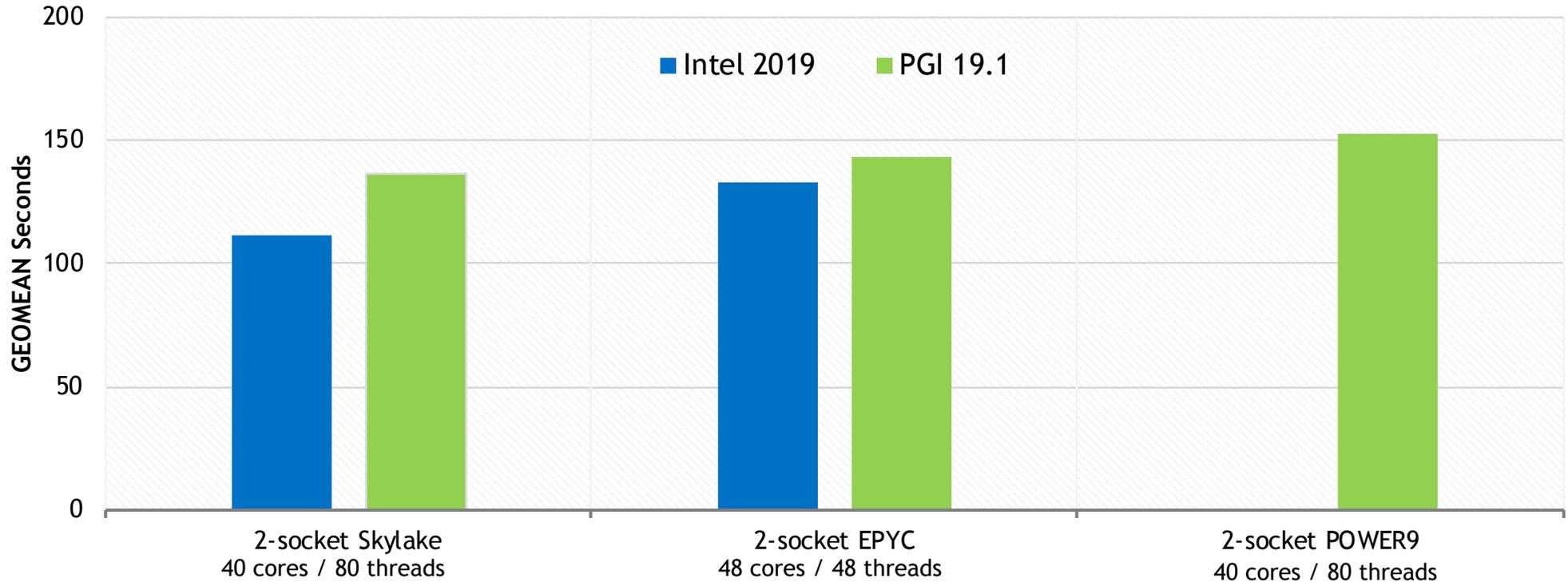
## The PGI 18.10 Community Edition

**FREE**

	<b>PGI</b> Community EDITION	<b>PGI</b> Professional EDITION	<b>PGI</b> Enterprise EDITION
<b>PROGRAMMING MODELS</b> OpenACC, CUDA Fortran, OpenMP, C/C++/Fortran Compilers and Tools	✓	✓	✓
<b>PLATFORMS</b> X86, OpenPOWER, NVIDIA GPU	✓	✓	✓
<b>UPDATES</b>	1-2 times a year	6-9 times a year	6-9 times a year
<b>SUPPORT</b>	User Forums	PGI Support	PGI Premier Services
<b>LICENSE</b>	Annual	Perpetual	Volume/Site

# SPEC CPU 2017 FP SPEED BENCHMARKS

OpenMP 3.1 Performance on Multicore CPUs - smaller is better



Performance measured January, 2019. Skylake: Two 20 core Intel Xeon Gold 6148 CPUs @ 2.4GHz w/ 376GB memory, hyperthreading enabled. EPYC: Two 24 core AMD EPYC 7451 CPUs @ 2.2GHz w/ 256GB memory. POWER 9: Two 20 core IBM POWER9 D2.2 @ 3.4GHz w 128GB memory.

Compilers: PGI 19.1, Intel 2019.

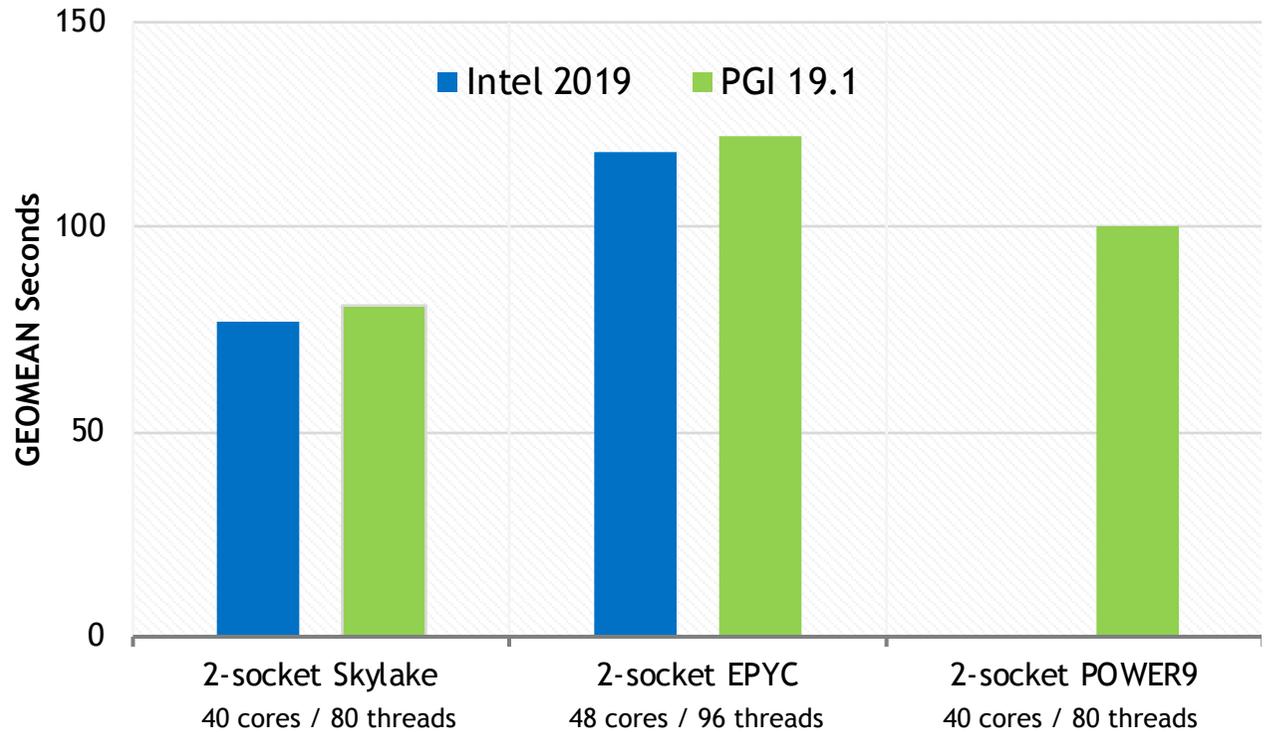
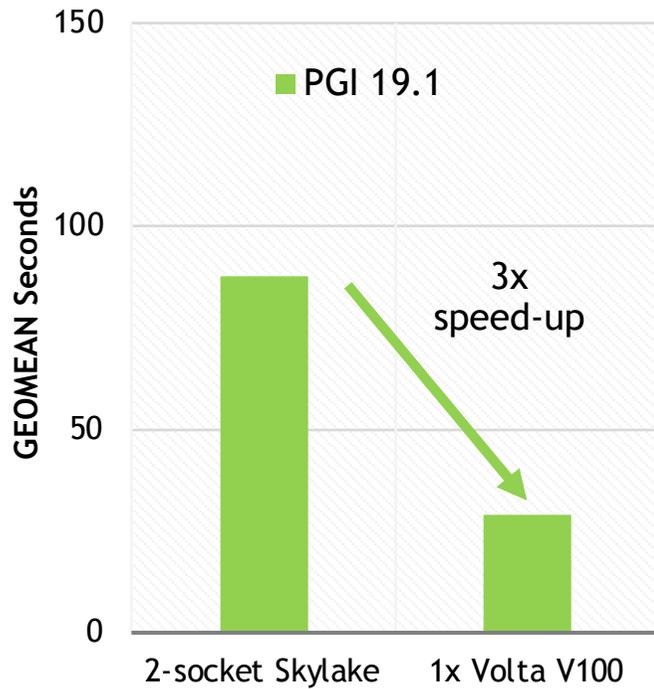
SPEC® is a registered trademark of the Standard Performance Evaluation Corporation ([www.spec.org](http://www.spec.org)).

# SPEC ACCEL 1.2 BENCHMARKS

## OpenACC vs OpenMP Performance on CPUs and GPUs

OpenACC

OpenMP 4.5



Performance measured January, 2019. Skylake: Two 20 core Intel Xeon Gold 6148 CPUs @ 2.4GHz w/ 376GB memory, hyperthreading enabled with two NVIDIA Tesla V100-PCIe-16GB GPU @ 1.53GHz. EPYC: Two 24 core AMD EPYC 7451 CPUs @ 2.2GHz w/ 256GB memory. POWER9: Two 20 core IBM POWER9 DD2.2 @ 3.4GHz w 128GB memory.

SPEC® is a registered trademark of the Standard Performance Evaluation Corporation ([www.spec.org](http://www.spec.org)).

# PGI | 19.1

## Fortran, C and C++ for the Tesla Platform

Full C++17 language

OpenMP 4.5 for CPUs

CUDA 10.x support

CUDA Fortran Tensor Core Support

OpenACC 2.6

OpenACC + printf, deep copy, C++ lambda

PCAST

PGI in the Cloud



**CUDA FORTRAN**

# CUDA FORTRAN

```
real, device, allocatable, dimension(:) :: Xd, Yd, Rd
real, pinned, allocatable, dimension(:) :: X, Y
real :: R

. . .

allocate (Xd(N), Yd(N), Rd((N+127)/128))
Xd = X(:)
Yd = Y(:)

call dsaxpy1<<<(N+127)/128, 128>>>(Yd,A,Xd,N)

call dsum1<<<(N+127)/128, 128>>>(Yd,Rd,N)

call dsum1<<<1, 128>>>(Rd,Rd,(N+127)/128)

Y(:) = Yd
R = Rd(1)
deallocate (Xd, Yd, Rd)

. . .
```

CPU Code

```
attributes(global) subroutine daxpy1(Y,A,X,N)
real :: Y(*), X(*)
integer, value :: N
real, value :: A
i = (blockidx%x-1) * blockdim%x + threadidx%x
if (i <= N) Y(i) = A * X(i) + Y(i)
end subroutine daxpy1
```

```
attributes(global) subroutine dsum1(Y,R,N)
real :: Y(*), R(*)
integer, value :: N
real, shared :: part(128)
s = 0
j = (blockidx%x-1) * blockdim%x + threadidx%x
do i = j, N, blockdim%x * griddim%x
    s = s + Y(i)
enddo
j = threadidx%x ; k = blockdim%x
part(j) = s
do while (k > 1)
    k = k / 2
    if (j <= k) part(j) = part(j) + part(j+k)
    call syncthreads()
enddo
R(blockidx%x) = part(1)
end subroutine dsum1
```

Tesla Code

# CUDA FORTRAN DIRECTIVES

```
real, device, allocatable, dimension(:) :: Xd, Yd, Rd  
real, pinned, allocatable, dimension(:) :: X, Y  
real :: R
```

...

```
allocate (Xd(N), Yd(N), Rd((N+127)/128))  
Xd = X(:)  
Yd = Y(:)  
R = 0.0  
!$cuf kernel do(1) <<<*,*>>>  
do i = 1, n  
  Yd(i) = A * Xd(i) + Yd(i)  
  R = R + Yd(i)  
enddo
```

```
Y(:) = Yd
```

```
deallocate (Xd, Yd)
```

...

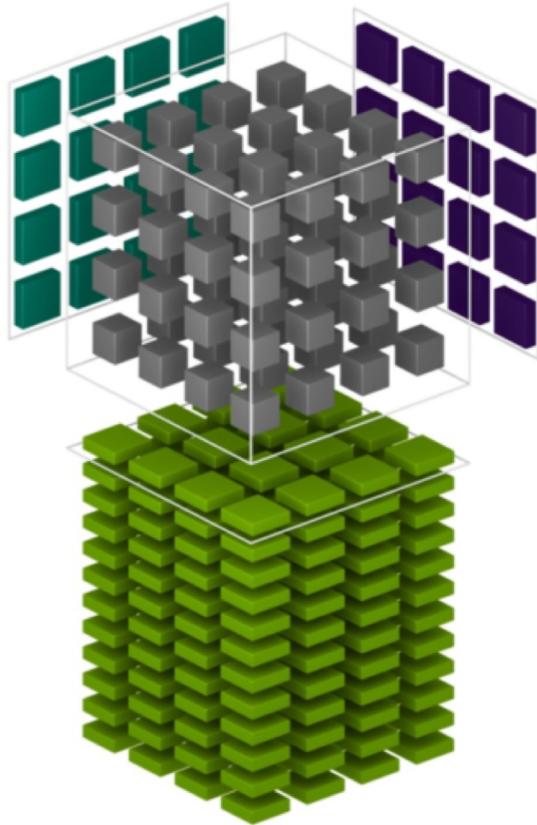
CPU Code

Compiler  
Generated!

Tesla Code

# V100 TENSOR CORES

Fast Matrix Multiply Using CUDA Fortran



Fast FP16 matrix multiply and accumulation into FP16 or FP32

8x to 16x faster than pure FP32 or FP64 in the same power envelope

Up to 125 TFlops FP16 performance

# CUDA FORTRAN TENSOR CORES

```
attributes(global) subroutine ww(a, b, c, m, k, n)
  use wmma
  integer, value :: m, k, n
  CUFReal12, device :: a(m, k)
  CUFReal12, device :: b(k, n)
  CUFReal12, device :: c(m, n)
  WMMASubMatrix(WMMAMatrixA, 32, 8, 16, Real, WMMAColMajor):: sa
  WMMASubMatrix(WMMAMatrixB, 32, 8, 16, Real, WMMAColMajor):: sb
  WMMASubMatrix(WMMAMatrixC, 32, 8, 16, Real, WMMAColMajor):: sc

  call wmmaLoadMatrix(sa, a(1,1), m)
  call wmmaLoadMatrix(sb, b(1,1), k)
  call wmmaMatmul(sc, sa, sb, sc)
  call wmmaStoreMatrix(c(1,1), sc, m)
end subroutine
```

**OPENACC**

# OPENACC EXAMPLE

```
for (iter = 0; iter < 100; ++iter) {  
    for (i = 1; i < n-1; ++i) {  
        for (j = 1; j < n-1; ++j) {  
            NH[j+i*n] = 0.25 * (OH[(j-1)+i*n] + OH[(j+1)+i*n] +  
                                OH[j+(i-1)*n] + OH[j+(i+1)*n]);  
        }  
    }  
    tmp = OH; OH = NH; NH = tmp;  
}
```

# SUPPORT FOR PRINTF() IN OPENACC

Debugging / tracing GPU-accelerated code regions

```
#pragma acc kernels copy(x[0:256])
{
  for (int i = 0; i < 256; i++) {
    if (x[i] < small) {
      printf("x at loc %03d is small: %d\n",i,x[i]);
      x[i] = small;
    } else if (x[i] > large) {
      printf("x at loc %03d is large: %d\n",i,x[i]);
      x[i] = large;
    }
  }
}
```

Support for formatted output using printf() statements in OpenACC compute regions

Debug and trace OpenACC programs during development and tuning

Works with both CPU and GPU

# OPENACC TRUE DEEP COPY DIRECTIVES

```
typedef struct points {
    float* x; float* y; float* z;
    int n;
    float coef, direction;
    #pragma acc shape(x[0:n],y[0:n])
} points;

void sub ( int n, float* y ) {
    points p;

    p.n = n;
    p.x = ( float*) malloc ( sizeof ( float )*n );
    p.y = ( float*) malloc ( sizeof ( float )*n );
    p.z = ( float*) malloc ( sizeof ( float )*n );

    #pragma acc data copy(p)
    {
        #pragma acc parallel loop
        for ( i =0; i<p.n; ++i ) p.x[i] += p.y[i];
        . . .
    }
}
```

Still in definition by the  
OpenACC Committee  
Beta Implementation  
since PGI 18.7

# OPENACC TRUE DEEP COPY DIRECTIVES

```
typedef struct points {
    float* x; float* y; float* z;
    int n;
    float coef, direction;
    #pragma acc policy<dpmove> copy(x[0:n]) copyin(y[0:n])
} points;

void sub ( int n, float* y ) {
    points p;

    p.n = n;
    p.x = ( float*) malloc ( sizeof ( float )*n );
    p.y = ( float*) malloc ( sizeof ( float )*n );
    p.z = ( float*) malloc ( sizeof ( float )*n );

    #pragma acc data copy(p<dpmove>)
    {
        #pragma acc parallel loop
        for ( i =0; i<p.n; ++i ) p.x[i] += p.y[i];
        . . .
    }
}
```

Still in definition by the  
OpenACC Committee  
Beta Implementation  
since PGI 18.7

# C++14 LAMBIDAS AND OPENACC

## Lambdas with Capture by Value in OpenACC Regions

```
template <typename Execution_Policy, typename BODY>
double bench_forall ( int s, int e, BODY body ) {
    StartTimer ();
    if ( is_same<Execution_Policy, Serial> :: value ) {
        for ( int i = s; i < e; ++i )
            body ( i );
    } elseif ( is_same<Execution_Policy, OpenACC> :: value ) {
        #pragma acc parallel loop
        for ( int i = s; i < e; ++i )
            body ( i );
    } return EndTimer ( );
}
```

```
using T = double;
void do_bench_saxpy ( int N, T*a, T*b, Tx) {
    auto saxpy = [=]( int i ) /* Capture-by-Value */
    { b[i] += a[i] * x; };
}
```

```
double stime = bench_forall<Serial>(0, N, saxpy);
double time = bench_forall<OpenACC>(0, N, saxpy);
printf ( "OpenACC Speedup %f \n", stime / time );
}
```

# DIRECTIVE-BASED HPC PROGRAMMING

## Who's Using OpenACC?

3 OF TOP 5 HPC APPS



5 OF 13 CAAR CODES



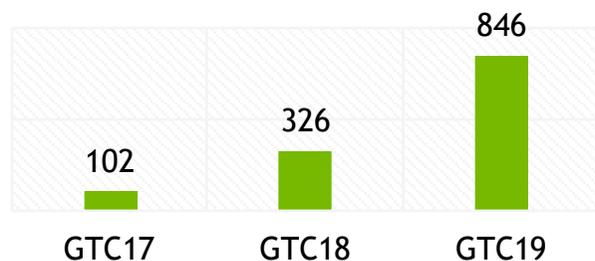
ACCELERATED APPS



725 TRAINED EXPERTS



SLACK MEMBERS



160,000+ DOWNLOADS

PGI®

Community  
EDITION

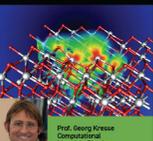
### GAUSSIAN



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-enhanced work. In the end, we could use the same code base for SMP cluster networks and GPU parallelism. PCs, compilers were essential to the success of our efforts. ☐

### VASP



Prof. Georg Kresse  
Lead Software Developer  
Materials Physics  
Vienna

For VASP, OpenACC is the easy forward for GPU acceleration. Performance is similar to CUDA. But OpenACC is generally easier for GPU developers and maintenance efforts. ☐

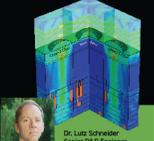
### ANSYS Fluent



Surel Sathe  
Lead Software Developer  
ANSYS Fluent  
Vienna

We recently started evaluating OpenACC for parallelizations on multicore CPUs. Based on some early results, one of our OpenACC based solvers is as fast as the OpenMP version on multicore CPUs and delivers speed-ups of 4x on a Tesla P100 compared to all the cores of a 3x8 server. ☐

### SYNOPTICS



Dr. Lutz Schneider  
Senior R&D Engineer  
SYNOPTICS  
Vienna

Using OpenACC, we have accelerated the Synoptics TAD Sensarius Direct CFD solver to support up to 100 simulations per image sensor by a factor of 10 on a single V100 GPU compared to a desktop of Broadwell server. GPUs are key to improving simulation throughput in the design of advanced image sensors. ☐

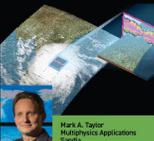
### COSMO



Dr. Oliver Fuhrer  
Senior Scientist  
Mentimeter  
Vienna

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 CMAR 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 invasive programming model approaches. ☐

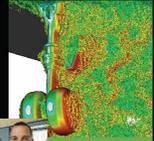
### VMD



John Stone  
Senior Research Programmer  
Bechman Institute

Due to AMD's ban, we need to port more parts of our code to the CPU and we're going to sweat it out. 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 self-starting because with the current algorithms, CPU performance is bandwidth-bound. ☐

### LAVA



Michael Brandt  
Research Aerospace Engineer  
NASA Ames Research Center

We used OpenACC to port our LAVA Laminar-Bottom-mix app to GPUs. Running a single block grid at 250^3, one V100 GPU with 50 Tensor Flow was 5x-6x faster than a 32-core Xeon E5-2600 v3. ☐

### SOMA



Ludwig Schneider  
PhD Student  
Georgia Institute of Technology

OpenACC enables us to compile a single code base for multiple architectures. This keeps the code maintainable and facilitates use for future accelerators. For our OpenACC accelerated "CMF" algorithm, a single V100 outperforms 26 CPU cores by roughly a factor of 10. ☐

### MAS



Reynald M. Caplan  
Computational Scientist  
Predictive Science Inc.

Adding OpenACC into MAS has given us the ability to integrate medium-sized simulations from a multi-node GPU 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 storm modeling. ☐

### SANJEEVINI



Ashish Jain  
Project Scientist  
Indian Institute of Technology  
New Delhi

In an academic environment, maintenance and debugging of existing codes is a tedious task. OpenACC provides a great platform for computational scientists to accomplish both tasks without involving a lot of effort or manpower in spending on the extra computational cost. ☐

### MPAS-A



Rafael Lott  
Director, Technology  
Development, NCAR

Our team has been evaluating OpenACC as a pathway to performance portability for the Model for Prediction Across Scales atmospheric model. Using this approach on the MPAS dynamical core, we have achieved performance on a single P100 GPU equivalent to 2.7 bank-scaled Intel Xeon nodes on our new Cheyenne supercomputer. ☐

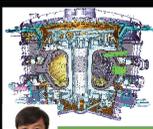
### GAMERA FOR GPU



Takuma Yamamoto, Koki Fujita,  
Tatsuya Ichimura, Masaru Mori, Lath Wilber  
The University of Tokyo

With OpenACC and compute nodes based on NVIDIA's Tesla P100 GPUs, we achieved more than a 1x speed-up over a K-Corner node running our earthquake disaster simulation code. ☐

### GTC



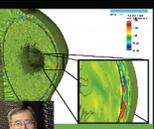
Zihong Lin  
Professor and Principal Investigator  
UC Irvine

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

### XGC



C.S. Chang  
Principal Investigator  
Plasma Physics Laboratory,  
Princeton University

Using a combination of CUDA and OpenACC for our recent compute intensive benchmarks, the GPU accelerated version of XGC delivers over 1x speed-up compared to CPU-only calculation when running at scale on 2048 nodes of ORNL's new Summit supercomputer. ☐

### HIPSTAR



Richard Sandberg  
Investigator  
University of Melbourne

For a relatively small production case of about 80 million grid points, our OpenACC-enabled version running on two GPU nodes with 4 P100s each was approximately 26.2x faster than two nodes with 2x Haswell cores each. We're looking forward to running a larger simulation in Summit early in 2019. ☐

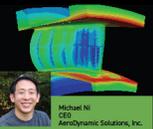
### CASTRO/MAESTRO



Adam Jacobs  
Ph.D. Candidate  
Stony Brook University

For scientific applications that run on several different supercomputing architectures and need to be usable for many generations of architecture, the cost of something like CUDA outweighs the profit. That's why we prefer OpenACC. ☐

### ADS CFD



Michael Ni  
CEO  
Aerodynamic Solutions, Inc.

Using OpenACC on a single Tesla V100 in the Amazon Cloud on GPU accelerated Code Law Flow solver runs 4 times faster at 30% less cost compared to runs using all 2x GPU instances on a Xeon Platinum C124 large instance. We feel this will revolutionize the aerospace design cycle, enabling the delivery of more affordable, more reliable and higher performing designs at reduced development cost. ☐

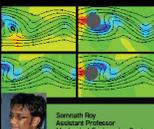
### FLASH



Brian Messer  
Senior Scientist  
Oak Ridge National Laboratory

We use OpenACC on Summit to accelerate our most compute-intensive kernels. We love OpenACC interoperability and how this allows us to use multiple methods to perform memory placement and movement. CPU/GPU performance is more than reasonable compared to something impossible to do on Titan, is 2x faster than CPU only. ☐

### IBM-CFD



Suresh Rav  
Associate Professor  
Mechanical Engineering Dept.  
Indian Institute of Technology  
Kharagpur

Using OpenACC to accelerate our in-house boundary integral velocity-CFD code, we were seeing an order of magnitude reduction in computing time and memory placement and movement. CPU/GPU performance is more than reasonable compared to something impossible to do on Titan, is 2x faster than CPU only. ☐

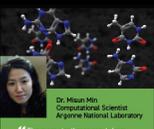
### HIFUN



Munirahna Nigam  
Chief Technology Officer  
SAI Engineering Solutions  
Pvt Ltd.

OpenACC allowed us to port an in-house CFD solvers to hybrid CPU-GPU platforms in a form that is readable and maintainable, and enabled us to extend GPUs in our HIFUN MP-CFD solver in very little time. ☐

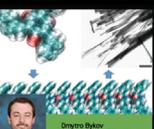
### NEKCEM



Dr. Mian Min  
Computational Scientist  
Argonne National Laboratory

The most significant result from our performance studies is that our computation with less energy consumption compared with our CPU-only runs. The GPU required only 20 percent of the energy needed for 16 GPUs to do the same computation. That OpenACC is an open standard was an important factor in our decision to use it for our research. ☐

### LSDALTON



Srinivas Balakrishna  
Computational Scientist  
Oak Ridge National Laboratory

Using OpenACC, we see large performance gains with very little when GPU acceleration runs over the course of our simulations due to differing program sizes, but is typically 3x-5x. On Summit we can now do simulations of several thousand atoms, compared to maybe 800 on Titan. ☐

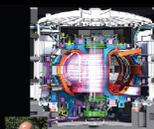
### FINE/Open



David Gutwiler  
Lead Software Developer  
NVIDIA

Porting our unstructured OpenACC 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. ☐

### CGYRO



Igor Okrop  
HPC Software Developer  
General Atomics

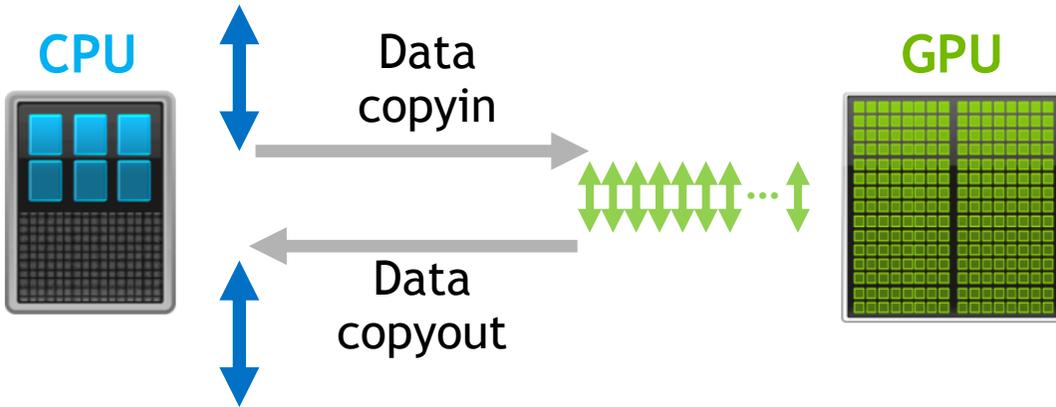
We planned to spend a month porting our OpenACC code from X86 CPUs to POWER9 CPUs. Using the PGI compilers and our standard build environment, we were working in an afternoon. It just worked. ☐

**PCAST**

**PGI COMPILER-ASSISTED SOFTWARE TESTING**

# OPENACC AUTO-COMPARE

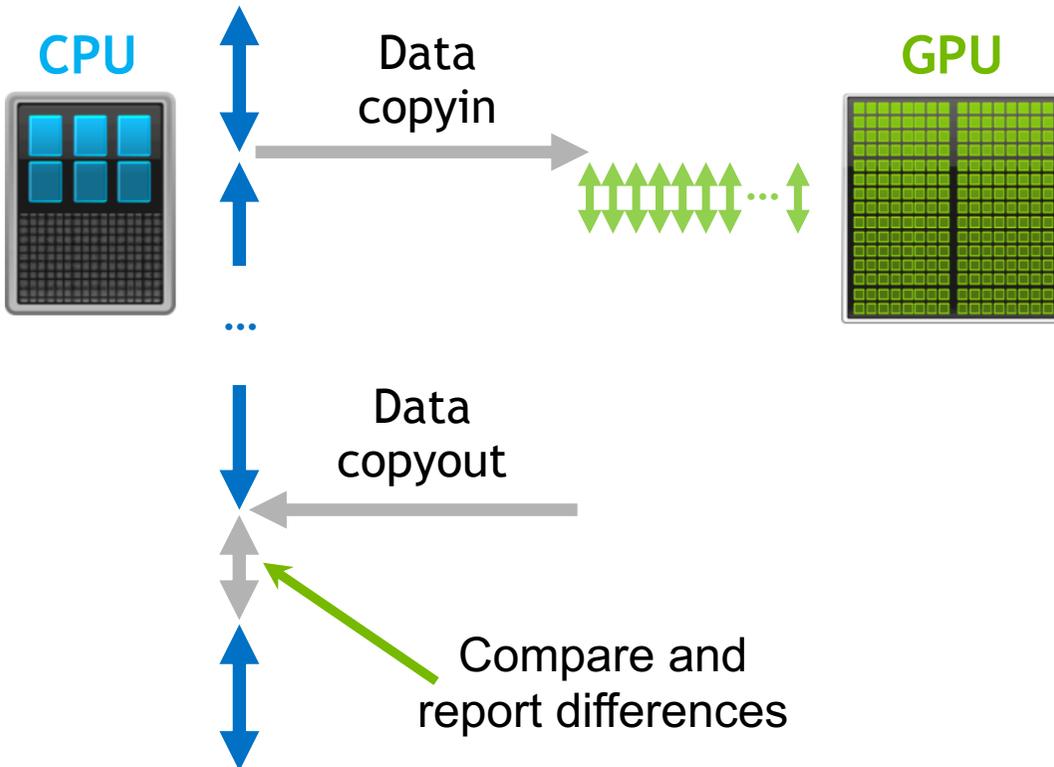
Find where CPU and GPU numerical results diverge



Normal OpenACC execution

# OPENACC AUTO-COMPARE

Find where CPU and GPU numerical results diverge



-ta=tesla:autocompare

Compute regions run redundantly on CPU and GPU

Results compared when data copied from GPU to CPU

[pgicompilers.com/pcast](http://pgicompilers.com/pcast)

# AUTO-COMPARE OVERVIEW

## PGI Compiler-Assisted Software Testing

```
$ gcc -ta=tesla:autocompare -o a.out example.c
```

```
$ PGI_COMPARE=summary,compare,abs=1 ./a.out
PCAST a1 comparison-label:0 Float
  idx: 0 FAIL ABS act: 8.40187728e-01
exp: 1.00000000e+00 tol: 1.00000001e-01
  idx: 1 FAIL ABS act: 3.94382924e-01
exp: 1.00000000e+00 tol: 1.00000001e-01
  idx: 2 FAIL ABS act: 7.83099234e-01
exp: 1.00000000e+00 tol: 1.00000001e-01
  idx: 3 FAIL ABS act: 7.98440039e-01
exp: 1.00000000e+00 tol: 1.00000001e-01
```

Test for program correctness, and determine points of divergence

Detect diverging results between CPU & GPU or multiple CPU code versions

Flexible calling and instrumentation options

[pgicompilers.com/pcast](http://pgicompilers.com/pcast)

# PCAST SUMMARY

OpenACC, GPU to CPU comparison

`-ta=tesla,autocompare`, compare at copyout or update host

`-ta=tesla,redundant`, redundant execution

`acc_compare()` to dynamically compare device-to-host data

CPU to CPU comparison

`pgi_compare()` to save data to file, then to compare to saved data



# PGI in the Cloud

Develop | Test | Benchmark | Deploy

# PGI ON AMAZON WEB SERVICES

**PGI Community Edition**  
By: Nvidia Latest Version: 18.10.1  
PGI Community Edition Fortran, C and C++ compilers with OpenACC and CUDA Fortran for heterogeneous CPU+GPU high-performance computing.  
Linux/Unix (0) **Free Tier**

**Continue to Subscribe**  
**Save to List**  
Typical Total Price  
**\$3.060/hr**  
Total pricing per instance for services hosted on p3.2xlarge in US East (N. Virginia). [View Details](#)

**Overview** Pricing Usage Support Reviews

### Product Overview

PGI Community Edition Fortran, C and C++ compilers and tools for multicore CPUs and NVIDIA GPUs, including all OpenACC, OpenMP and CUDA Fortran features. Enables development of performance-portable HPC applications with uniform source code across the most widely used parallel processors and systems. For more information see [pgicompilers.com/community](https://www.pgroup.com/community) or <https://www.pgroup.com/community>

This PGI Community Edition AMI version 18.10.1 also includes:

- Ubuntu: 16.04 LTS
- NVIDIA Driver: 410.79
- Docker CE: 18.09.0
- NVIDIA Container Runtime for Docker: (nvidia-docker2) v2.0.3

Version	18.10.1
By	Nvidia
Categories	Application Development High Performance Computing Education & Research
Operating System	Linux/Unix, Ubuntu 16.04
Delivery Methods	Amazon Machine Image

### Highlights

- Parallelizing Fortran, C and C++ compilers for x86-64 CPU + GPU high-performance computing.
- PGI compilers deliver the performance you need on CPUs, with OpenACC and CUDA Fortran for HPC applications development on GPU-accelerated systems.
- Includes support for NVIDIA Tesla V100 GPUs. OpenACC and CUDA programs can run several times faster on a single Tesla V100 GPU compared to all the cores of a dual-socket server, and interoperate with MPI and OpenMP to deliver the full power of today's multi-GPU servers.

Easy access to the latest NVIDIA GPU hardware and PGI Community Edition Fortran, C and C++ compilers

[pgicompilers.com/cloud](https://pgicompilers.com/cloud)

# PGI DOCKER IMAGE ON NVIDIA GPU CLOUD

Run PGI compilers on any Docker-enabled system

Use Containers to reliably run software when moving from one computing environment to another.



[ngc.nvidia.com](https://ngc.nvidia.com)

The screenshot shows the NVIDIA GPU Cloud interface for the 'PGI Compilers' container image. The page includes a header with 'NVIDIA GPU CLOUD', navigation links for 'Sign In', 'Create an Account', and 'Terms Of Use', and a 'CATALOG' sidebar. The main content area displays the 'PGI Compilers' image details, including a 'Pull ce' button, a table with columns for 'Publisher', 'Built By', 'Latest Tag', 'Modified', and 'Size', a description, labels for 'HPC', 'High Performance Computing', and 'Toolkit', and a 'Pull Command' field containing the command: `docker pull nvcr.io/hpc/pgi-compilers:ce`. Below this, there are tabs for 'Overview', 'Tags', and 'Layers'. The 'Overview' tab is active, showing a license agreement link, the version 'PGI Community Edition 18.10', and detailed text about the image's capabilities and supported hardware. A footer contains links for 'Documentation', 'User Forum', and 'Collapse'.

# Future Directions

```
#include <opencv.h>
// my managed vector datatype
template<typename element> class managed_vector
{
public:
    size_t size;
    bool iscopy;

    managed_vector( size_t size )
    {
        size = size;
        data = new element[size];
        iscopy = false;
    }
    managed_vector( const managed_vector &other )
    {
        size = other.size;
        data = copyof.data;
        iscopy = true;
    }
    void attach( void* data )
    {
        #pragma acc update dev( data )
    }
    void updatehost()
    {
        #pragma acc update dev( data )
    }
    void updatedev()
    {
        #pragma acc update dev( data )
    }
    ~managed_vector()
    {
        delete data;
    }
};
```



# NEW FEATURES IN OPENACC 2.7

Spec Finalized Nov 2018

Multicore CPU targets  
Partially-shared memory  
Array / struct reductions  
Multi-GPU programming



# UPCOMING FEATURES IN OPENACC 3.0

True Deep Copy directives, using feedback for PGI initial implementation

Support for C++ Lambdas, directives within/above Lambdas

Better support for multiple devices (device-to-device copy, ...)

Error handler

Better support for Fortran pointers, Fortran interfaces to all API routines

Update base languages, Fortran 2018, C++17, C18

Many smaller fixes

# PCAST IMPROVEMENTS

## From User Feedback

More control over compare frequency

Select function, file, variable names

Options to test single file or function

```
#pragma acc compare(a[0:n])
```

```
#pragma pgi compare(b[0:n])
```

# REAL(2) IN CUDA FORTRAN

```
real(2), allocatable, device :: x(:)
real(2), allocatable :: hx(:)

allocate(x(n), hx(n))

call fill<<<n/128,128>>>(x)

hx(:) = x(:)
print *, hx(1)
```

# SCALABLE PARALLELISM WITH OPENACC...

```
% pgfortran -fast -ta=tesla -Minfo -c PdV_kernel.f90
pdv_kernel:
...
77, Loop is parallelizable
79, Loop is parallelizable
Accelerator kernel generated
Generating Tesla code
77, !$acc loop gang, vector(4) ! blockidx%y
! threadidx%y
79, !$acc loop gang, vector(32)! blockidx%x
! threadidx%x
...
```

```
75 !$ACC KERNELS
76 !$ACC LOOP INDEPENDENT
77   DO k=y_min,y_max
78 !$ACC LOOP INDEPENDENT PRIVATE(right_flux,left_flux,top_flux,bottom_flux,total_flux,
min_cell_volume,energy_change,recip_volume)
79   DO j=x_min,x_max
80
81     left_flux= (xarea(j ,k )*(xvel0(j ,k )+xvel0(j ,k+1)           &
82                +xvel0(j ,k )+xvel0(j ,k+1)))*0.25_8*dt*0.5
83     right_flux= (xarea(j+1,k )*(xvel0(j+1,k )+xvel0(j+1,k+1)       &
84                 +xvel0(j+1,k )+xvel0(j+1,k+1)))*0.25_8*dt*0.5
85     bottom_flux=(yarea(j ,k )*(yvel0(j ,k )+yvel0(j+1,k )       &
86                 +yvel0(j ,k )+yvel0(j+1,k )))*0.25_8*dt*0.5
87     top_flux=   (yarea(j ,k+1)*(yvel0(j ,k+1)+yvel0(j+1,k+1)     &
88                 +yvel0(j ,k+1)+yvel0(j+1,k+1)))*0.25_8*dt*0.5
89     total_flux=right_flux-left_flux+top_flux-bottom_flux
90
91     volume_change(j,k)=volume(j,k)/(volume(j,k)+total_flux)
92
93     min_cell_volume=MIN(volume(j,k)+right_flux-left_flux+top_flux &
94                        ,volume(j,k)+right_flux-left_flux         &
95                        ,volume(j,k)+top_flux-bottom_flux)
96
97     recip_volume=1.0/volume(j,k)
98
99     energy_change=(pressure(j,k)/density0(j,k)+viscosity(j,k)/density0(j,k))*...
101    energy1(j,k)=energy0(j,k)-energy_change
103    density1(j,k)=density0(j,k)*volume_change(j,k)
105    ENDDO
106  ENDDO
107 !$ACC END KERNELS
```

# SCALABLE PARALLELISM WITH OPENMP..

```
% pgfortran -fast -mp=tesla -Minfo -c PdV_kernel.f90
pdv_kernel:
...
77, Tesla device code generated
77, OMP teams, thread, simd parallelism
79, collapsed
...
```

```
75 !$OMP TARGET TEAMS DISTRIBUTE PARALLEL DO SIMD COLLAPSE(2)
76 !$OMP& PRIVATE(right_flux, left_flux, top_flux, bottom_flux, total_flux,
77 !$OMP& min_cell_volume, energy, change, recip_volume)
78 DO k=y_min,y_max
79 DO j=x_min,x_max
80
81 left_flux= (xarea(j ,k )*(xvel0(j ,k )+xvel0(j ,k+1) &
82 +xvel0(j ,k )+xvel0(j ,k+1)))*0.25_8*dt*0.5
83 right_flux= (xarea(j+1,k )*(xvel0(j+1,k )+xvel0(j+1,k+1) &
84 +xvel0(j+1,k )+xvel0(j+1,k+1)))*0.25_8*dt*0.5
85 bottom_flux=(yarea(j ,k )*(yvel0(j ,k )+yvel0(j+1,k ) &
86 +yvel0(j ,k )+yvel0(j+1,k )))*0.25_8*dt*0.5
87 top_flux= (yarea(j ,k+1)*(yvel0(j ,k+1)+yvel0(j+1,k+1) &
88 +yvel0(j ,k+1)+yvel0(j+1,k+1)))*0.25_8*dt*0.5
89 total_flux=right_flux-left_flux+top_flux-bottom_flux
90
91 volume_change(j,k)=volume(j,k)/(volume(j,k)+total_flux)
92
93 min_cell_volume=MIN(volume(j,k)+right_flux-left_flux+top_flux-bottom_flux &
94 ,volume(j,k)+right_flux-left_flux &
95 ,volume(j,k)+top_flux-bottom_flux)
96
97 recip_volume=1.0/volume(j,k)
98
99 energy_change=(pressure(j,k)/density0(j,k)+viscosity(j,k)/density0(j,k))*...
100
101 energy1(j,k)=energy0(j,k)-energy_change
102
103 density1(j,k)=density0(j,k)*volume_change(j,k)
104
105 ENDDO
106 ENDDO
107
```

# PARALLEL FEATURES IN FORTRAN AND C++

## Fortran 2018

Array syntax (F90)

FORALL (F95)

Co-arrays (Fo8, F18)

DO CONCURRENT (Fo8, F18)

## C++17

Threads (C++11)

pSTL Parallel Algorithms (C++17)

# CLOVERLEAF V1.3



AWE Hydrodynamics mini-app

6500+ lines, !\$acc kernels

OpenACC or OpenMP

Source on GitHub

*<http://uk-mac.github.io/CloverLeaf>*

# FORTRAN WITH OPENACC DIRECTIVES

```
% pgfortran -fast -ta=tesla:managed -Minfo -c
PdV_kernel.f90
pdv_kernel:
...
77, Loop is parallelizable
79, Loop is parallelizable
Accelerator kernel generated
Generating Tesla code
77, !$acc loop gang, vector(4) ! blockidx%y
! threadidx%y
79, !$acc loop gang, vector(32)! blockidx%x
! threadidx%x
...
```

```
75 !$ACC KERNELS
76 !$ACC LOOP INDEPENDENT
77   DO k=y_min,y_max
78 !$ACC LOOP INDEPENDENT PRIVATE(right_flux,left_flux,top_flux,bottom_flux,total_flux,
min_cell_volume,energy_change,recip_volume)
79   DO j=x_min,x_max
80
81     left_flux= (xarea(j ,k )*(xvel0(j ,k )+xvel0(j ,k+1)           &
82                +xvel0(j ,k )+xvel0(j ,k+1)))*0.25_8*dt*0.5
83     right_flux= (xarea(j+1,k )*(xvel0(j+1,k )+xvel0(j+1,k+1)       &
84                +xvel0(j+1,k )+xvel0(j+1,k+1)))*0.25_8*dt*0.5
85     bottom_flux=(yarea(j ,k )*(yvel0(j ,k )+yvel0(j+1,k )       &
86                +yvel0(j ,k )+yvel0(j+1,k )))*0.25_8*dt*0.5
87     top_flux=   (yarea(j ,k+1)*(yvel0(j ,k+1)+yvel0(j+1,k+1)     &
88                +yvel0(j ,k+1)+yvel0(j+1,k+1)))*0.25_8*dt*0.5
89     total_flux=right_flux-left_flux+top_flux-bottom_flux
90
91     volume_change(j,k)=volume(j,k)/(volume(j,k)+total_flux)
92
93     min_cell_volume=MIN(volume(j,k)+right_flux-left_flux+top_flux &
94                        ,volume(j,k)+right_flux-left_flux       &
95                        ,volume(j,k)+top_flux-bottom_flux)
96
97     recip_volume=1.0/volume(j,k)
98
99     energy_change=(pressure(j,k)/density0(j,k)+viscosity(j,k)/density0(j,k))*...
101    energy1(j,k)=energy0(j,k)-energy_change
103    density1(j,k)=density0(j,k)*volume_change(j,k)
105    ENDDO
106  ENDDO
107 !$ACC END KERNELS
```

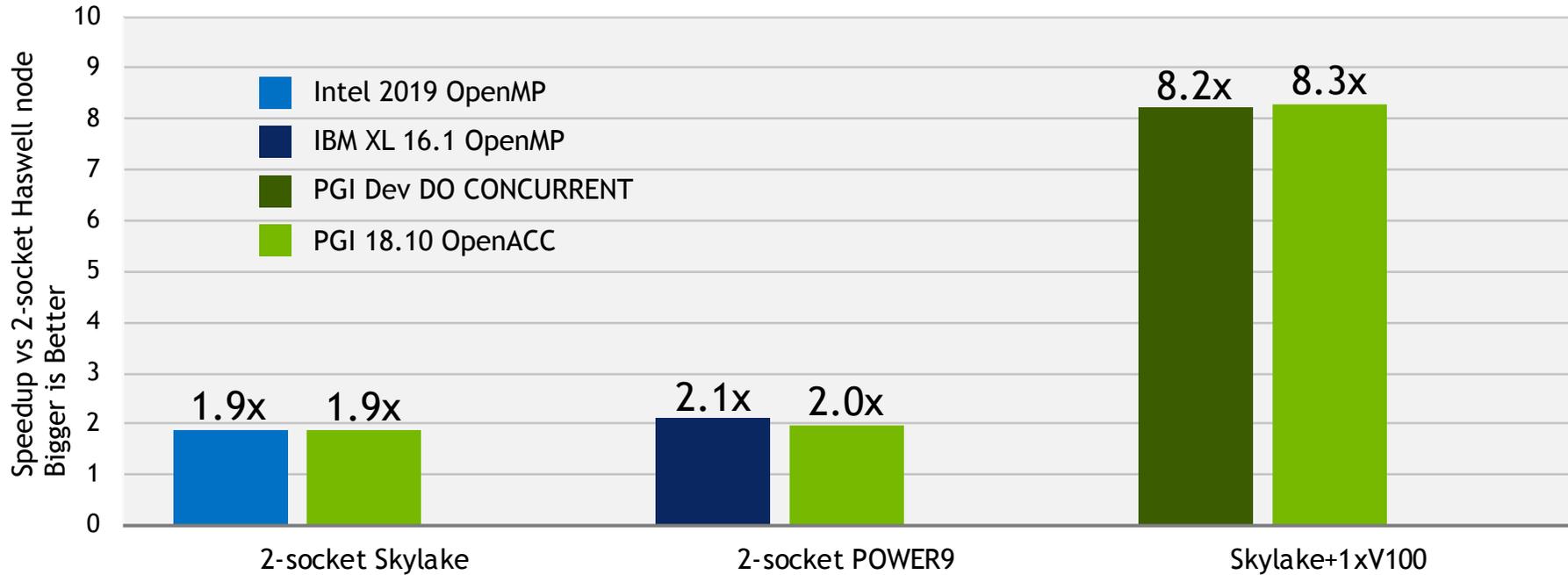
# FORTRAN 2018 DO CONCURRENT

```
% pgfortran -fast -ta=tesla:managed -Minfo -c
PdV_kernel.f90
pdv_kernel:
...
77, Do concurrent is parallelizable
Accelerator kernel generated
Generating Tesla code
77, !$acc loop gang, vector(4) ! blockidx%y
! threadidx%y
!$acc loop gang, vector(32)! blockidx%x
! threadidx%x
...
```

```
75
76
77 DO CONCURRENT (k=y_min:y_max, j=x_min:x_max) &
78 LOCAL (right_flux,left_flux,top_flux,bottom_flux,total_flux, &
min_cell_volume,energy_change,recip_volume)
79
80
81 left_flux= (xarea(j ,k )*(xvel0(j ,k )+xvel0(j ,k+1) &
82 +xvel0(j ,k )+xvel0(j ,k+1)))*0.25_8*dt*0.5
83 right_flux= (xarea(j+1,k )*(xvel0(j+1,k )+xvel0(j+1,k+1) &
84 +xvel0(j+1,k )+xvel0(j+1,k+1)))*0.25_8*dt*0.5
85 bottom_flux=(yarea(j ,k )*(yvel0(j ,k )+yvel0(j+1,k ) &
86 +yvel0(j ,k )+yvel0(j+1,k )))*0.25_8*dt*0.5
87 top_flux= (yarea(j ,k+1)*(yvel0(j ,k+1)+yvel0(j+1,k+1) &
88 +yvel0(j ,k+1)+yvel0(j+1,k+1)))*0.25_8*dt*0.5
89 total_flux=right_flux-left_flux+top_flux-bottom_flux
90
91 volume_change(j,k)=volume(j,k)/(volume(j,k)+total_flux)
92
93 min_cell_volume=MIN(volume(j,k)+right_flux-left_flux+top_flux-bottom_flux &
94 ,volume(j,k)+right_flux-left_flux &
95 ,volume(j,k)+top_flux-bottom_flux)
96
97 recip_volume=1.0/volume(j,k)
99 energy_change=(pressure(j,k)/density0(j,k)+viscosity(j,k)/density0(j,k))*...
101 energy1(j,k)=energy0(j,k)-energy_change
103 density1(j,k)=density0(j,k)*volume_change(j,k)
105
106 ENDDO
107
```

# FORTRAN 2018 DO CONCURRENT FOR V100

CloverLeaf AWE Hydrodynamics mini-App, bm32 data set



Systems: Skylake 2x20 core Xeon Gold server (sky-4), POWER9 2x20 core POWER9 DD2.1 server (wsn1)

Compilers: Intel 2019.0.045, PGI 18.10, IBM XL 16.1

Benchmark: CloverLeaf v1.3 OpenACC, OpenMP and DoConcurrent versions downloaded from <https://github.com/UoB-HPC> the week of October 29, 2018

[https://github.com/UoB-HPC/CloverLeaf\\_doconcurrent](https://github.com/UoB-HPC/CloverLeaf_doconcurrent)

# C++17 PARALLEL ALGORITHMS (PSTL)

## Parallel and Vector Concurrency for the Standard Algorithms

Execution Policies for most STL algorithms

`std::execution::seq, par, par_unseq`

A few new parallel algorithms

`std::transform_reduce`

`std::for_each_n`

...

*C++17 Parallel Algorithms and Beyond*, Bryce Lebach, CppCon 2016

# C++17 PARALLEL ALGORITHM: FOR\_EACH\_N()

```
% pgc++ -fast -ta=tesla:managed,pstl -Minfo -c lulesh.cc  
...
```

- No compiler feedback
- The parallel STL is a library, not a language
- Library calls are opaque to the compiler
- See S9770 David Olsen Talk C++17 Parallel Algorithms for NVIDIA GPUs with PGI C++

```
2100 static inline  
2101 void CalcVelocityForNodes(Domain &domain, const Real_t dt,  
2102                             const Real_t u_cut, Index_t numNode)  
2103 {  
2104  
2105  
2106     std::for_each_n (POLICY_PAR, idx(0), numNode, [=](Index_t i)  
2107     {  
2108         Real_t xdtmp, ydtmp, zdtmp ;  
2109  
2110         xdtmp = domain.xd(i) + domain.xdd(i) * dt ;  
2111         if( FABS(xdtmp) < u_cut ) xdtmp = Real_t(0.0);  
2112         domain.xd(i) = xdtmp ;  
2113  
2114         ydtmp = domain.yd(i) + domain.ydd(i) * dt ;  
2115         if( FABS(ydtmp) < u_cut ) ydtmp = Real_t(0.0);  
2116         domain.yd(i) = ydtmp ;  
2117  
2118         zdtmp = domain.zd(i) + domain.zdd(i) * dt ;  
2119         if( FABS(zdtmp) < u_cut ) zdtmp = Real_t(0.0);  
2120         domain.zd(i) = zdtmp ;  
2121     });  
2122 }
```

# THE FUTURE OF GPU PROGRAMMING

## Standard Languages | Directives | CUDA

```
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
```

GPU Accelerated  
C++17 and Fortran 2018

```
#pragma acc data copy(x,y) {  
  ...  
  std::transform(par, x, x+n, y, y,  
 [=] (float x, float y) {  
     return y + a*x;  
 });  
  ...  
}
```

Incremental Performance  
Optimization with OpenACC

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

Maximize GPU Performance  
with CUDA C++/Fortran

# PGI C++ AND FORTRAN COMPILERS & TOOLS

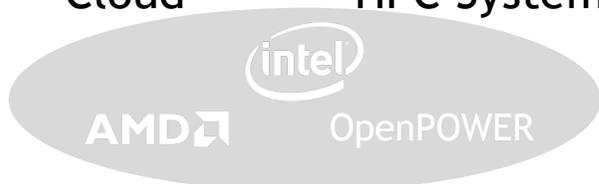
A cross-platform HPC SDK and Easy On-ramp to GPU Computing



Cloud



HPC Systems



Workstations



World-class HPC compilers for  
multicore CPUs

Enabling  
Parallel Computing



From Sequential, to Parallel,  
to Accelerated Computing



Cloud



HPC Systems



Workstations



GPU Acceleration with CUDA,  
Directives, native C++17/F18