

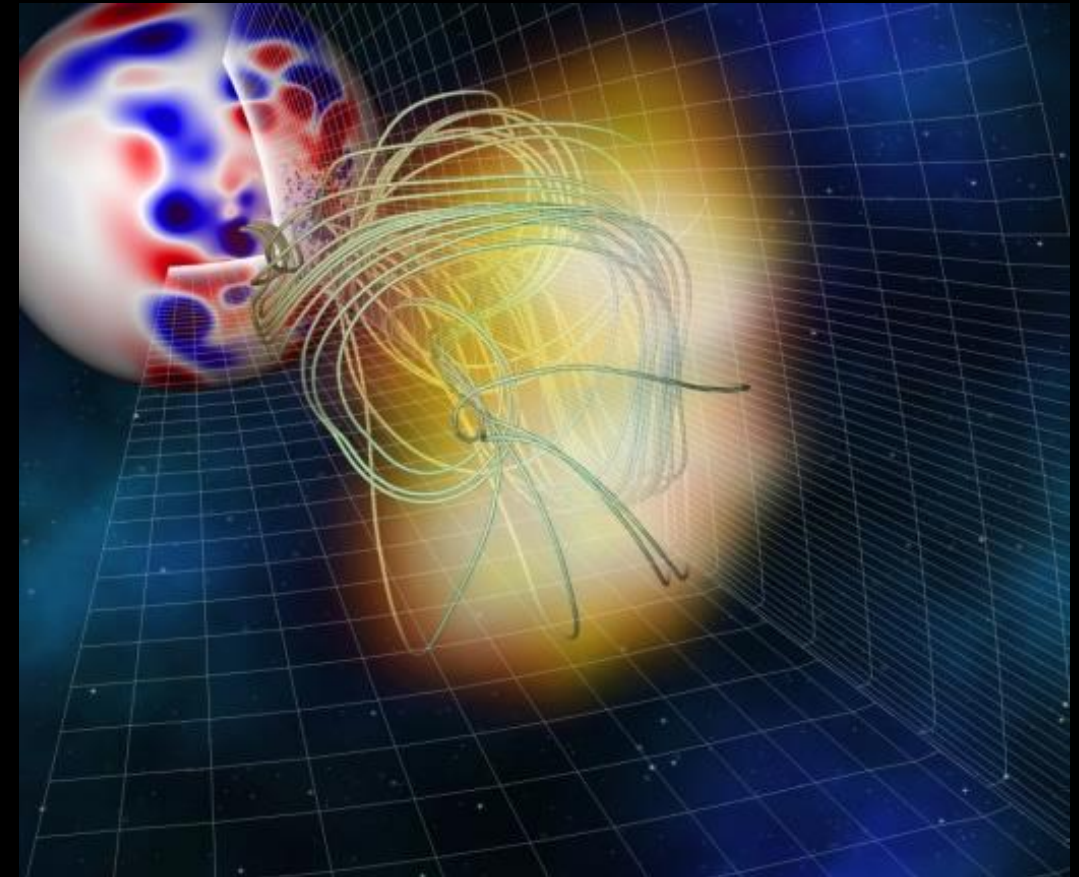
# Potential Field Solutions of the Solar Corona: Converting a PCG Solver from MPI to MPI+OpenACC

Ronald M. Caplan,  
Jon A. Linker, and Zoran Mikic  
Predictive Science Inc.  
[www.preds-ci.com](http://www.preds-ci.com)



**GPU** TECHNOLOGY  
CONFERENCE

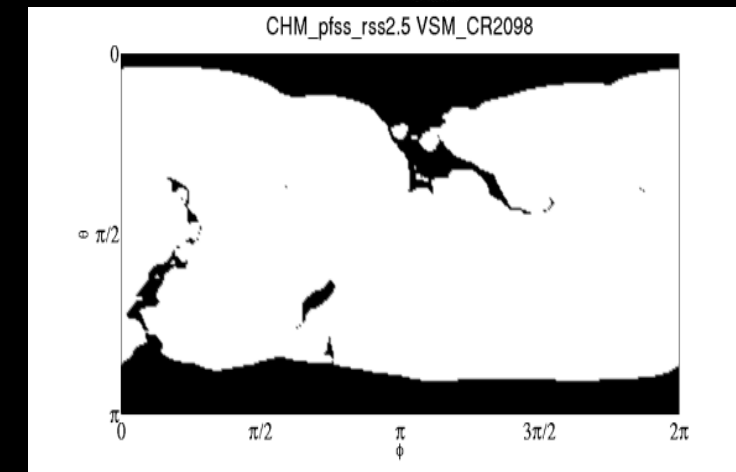
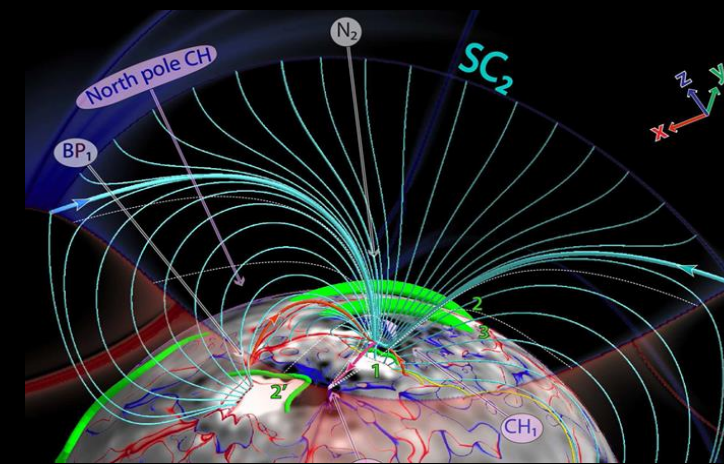
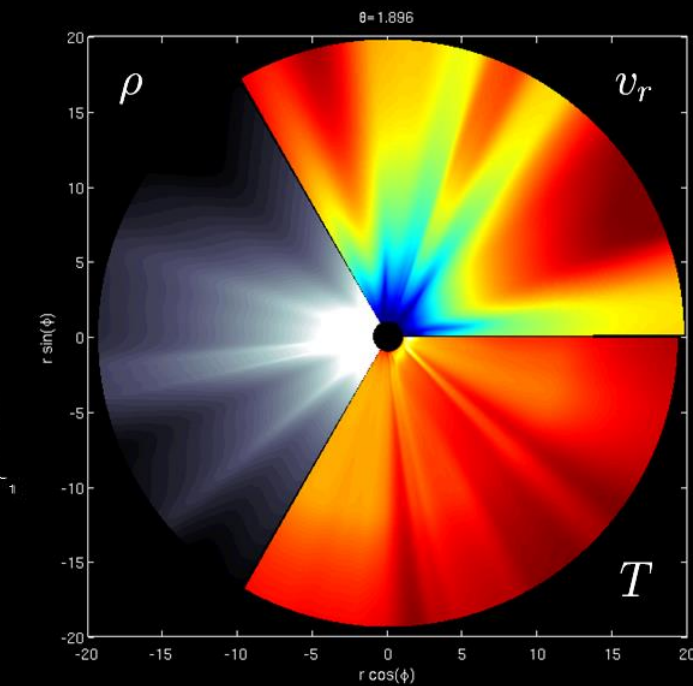
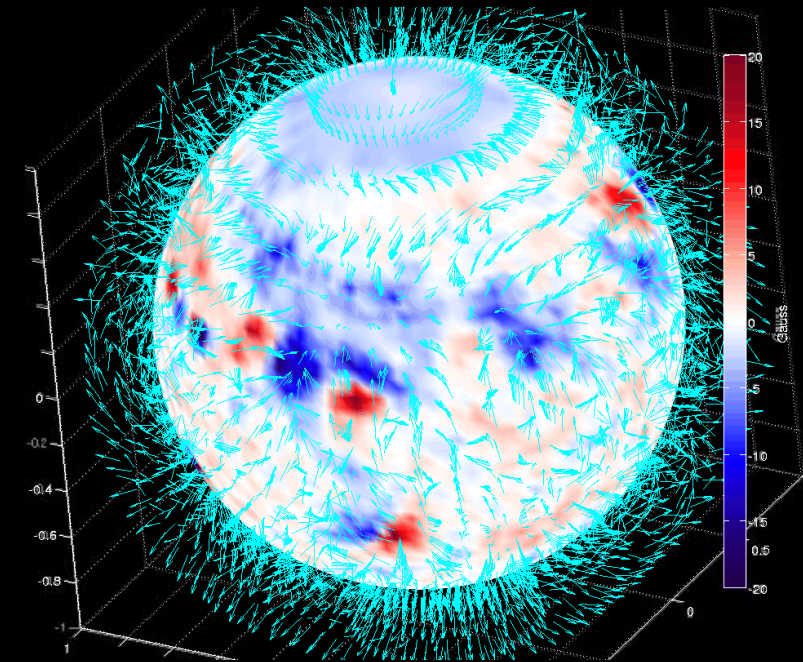
- ① Potential Field Solutions of the Solar Corona
- ② Preconditioned Conjugate Gradient
- ③ POT3D
- ④ MPI → MPI+OpenACC
- ⑤ Test Cases
- ⑥ Performance Results
- ⑦ Performance Portability





# Potential Field Solutions of the Solar Corona

- ☉ Studying the Sun's activity important for space weather forecasting
- ☉ Knowing the magnetic field is critical
- ☉ Available observations: Earth-facing 2D field at solar surface
- ☉ Need 3D global field
  - ☉ Magnetohydrodynamic simulations  
[SLOW! ~ days]
  - ☉ Potential field (PF) solution  
[FASTER! ~ minutes]
- ☉ PF solutions useful for:
  - ☉ Estimating magnetic topology and IMF levels
  - ☉ Long-time statistical analysis
  - ☉ Providing boundary/initial conditions for modeling solar storms with MHD simulations





# Potential Field Solutions of the Solar Corona

Maxwell equation:  $\nabla \times \vec{B} = \mu_0 \vec{J}$

Assume no current:  $\cancel{\vec{J}} \rightarrow \nabla \times \vec{B} = 0$

Solution form

$$\vec{B} = \nabla \Phi$$

Divergence-free condition

$$\nabla \cdot \vec{B} = 0$$

$$\nabla^2 \Phi = 0$$

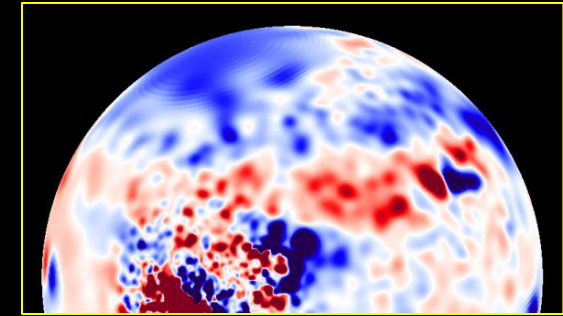
Laplace equation

## Boundary Conditions

$$r = R_{\odot}$$

Observed Surface

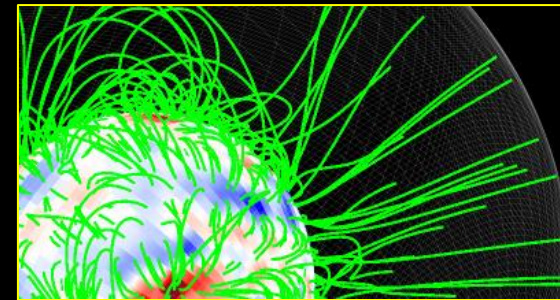
$$\left. \frac{\partial \Phi}{\partial r} \right|_{R_{\odot}} = \left. B_r \right|_{R_{\odot}}$$



$$r = R_1$$

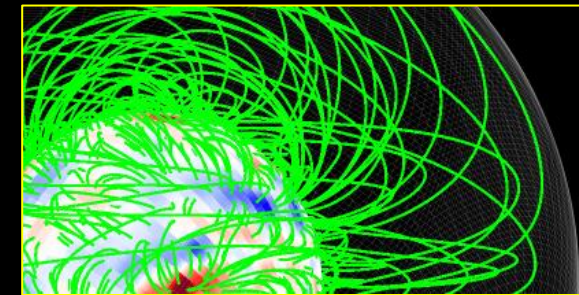
(a) Source Surface

$$\Phi|_{R_1} = 0$$



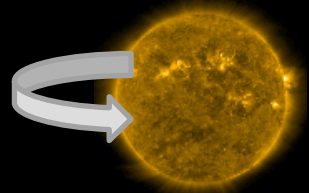
(b) Closed Field

$$\left. \frac{\partial \Phi}{\partial r} \right|_{R_1} = \left. B_r \right|_{R_1} = 0$$



$$\phi = 0, 2\pi$$

Periodic  $\Phi|_{\phi=0} = \Phi|_{\phi=2\pi}$



$$\theta = 0, \pi$$

Polar average

$$\Phi|_{\theta_0} = \frac{1}{2\pi} \int_{\phi=0}^{\phi=2\pi} \Phi|_{\theta_0 \pm \epsilon} d\phi$$



# Preconditioned Conjugate Gradient

- ⌘ PCG: common linear solver for large symmetric sparse matrices
- ⌘ A 'preconditioner' is applied each iteration to improve/ensure convergence

$$\begin{aligned} \mathbf{A} \vec{x} &= \vec{b} & (\mathbf{A} \Phi = 0) \\ \downarrow \\ \mathbf{P}^{-1} \mathbf{A} \vec{x} &= \mathbf{P}^{-1} \vec{b} & (\mathbf{P} \approx \mathbf{A}) \end{aligned}$$

- ⌘ We use two *communication free* preconditioning options:

## PC1

Point-Jacobi/Diagonal scaling  
(**DIAG**)

Cheap, not very effective

## PC2

Non-overlapping domain decomposition zero-fill  
incomplete LU factorization (**ILU0**)

Expensive, much more effective!

Why use/test **PC1** at all?

- ⌘ For "easy" solves, **PC1** can be faster than **PC2**

⌘  $N_{\text{processors}} \rightarrow N_{\text{grid}} \rightarrow \text{PC2} \rightarrow \text{PC1}$

# Preconditioned Conjugate Gradient

## Key components of the PCG algorithm

$$\mathbf{P} = \text{pc}()$$

### Initialization of PC

- **Expensive**, but only done once

$$\vec{y} = \mathbf{A} \vec{x}$$

### Matrix-Vector Product

- **Vectorization-friendly**
- Point-to-point MPI communication (**good scaling**)

$$\vec{z} = a \vec{x} + b \vec{y}$$

### Vector-Vector Operations

- **Vectorization-friendly**

$$r = \vec{x} \cdot \vec{y}$$

### Dot-Products

- **Vectorization friendly**
- Global MPI communication & synchronization (**reduced scaling**)

$$\vec{y} = \mathbf{P}^{-1} \vec{x}$$

### Application of PC

- **PC1**: Vector-Operation (**vectorization-friendly**)
- **PC2**: Triangular Solves (**vectorization-unfriendly**)

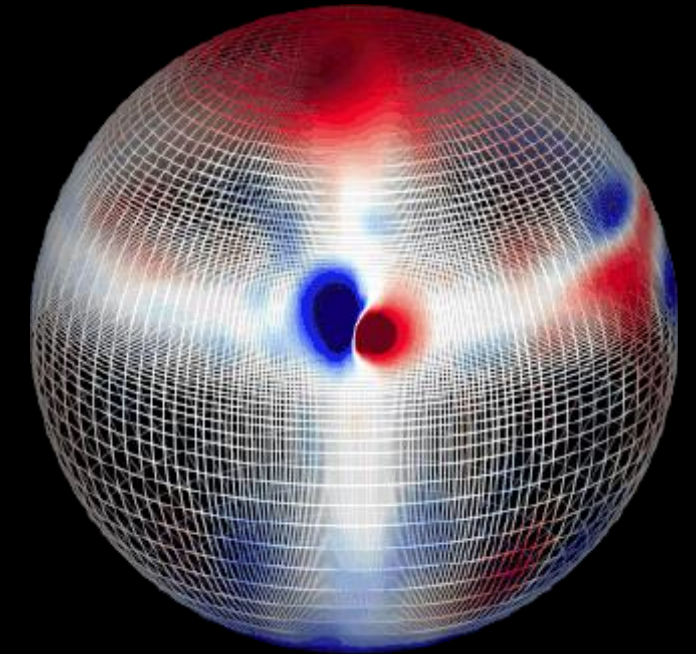


# POT3D Code

- Written in **FORTRAN90**, parallelized with **MPI**
- Finite difference on logically rectangular non-uniform spherical grid

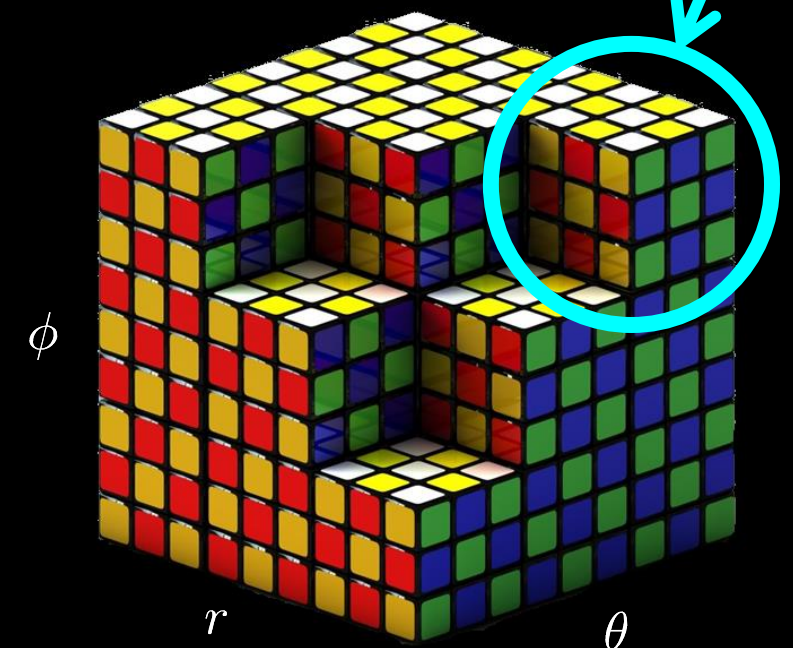
$$\begin{aligned}\nabla^2 \Phi_{i,j,k} \approx & \frac{1}{\Delta r_i} \left[ \frac{\Phi_{i+1,j,k} - \Phi_{i,j,k}}{\Delta r_{i+\frac{1}{2}}} - \frac{\Phi_{i,j,k} - \Phi_{i-1,j,k}}{\Delta r_{i-\frac{1}{2}}} \right] \\ & + \frac{1}{\sin \theta_j \Delta \theta_j} \left[ \sin \theta_{i,j+\frac{1}{2}} \frac{\Phi_{i,j+1,k} - \Phi_{i,j,k}}{\Delta \theta_{j+\frac{1}{2}}} - \sin \theta_{i,j-\frac{1}{2}} \frac{\Phi_{i,j,k} - \Phi_{i,j-1,k}}{\Delta \theta_{j-\frac{1}{2}}} \right] \\ & + \frac{1}{\sin^2 \theta_j \Delta \phi_k} \left[ \frac{\Phi_{i,j,k+1} - \Phi_{i,j,k}}{\Delta \phi_{k+\frac{1}{2}}} - \frac{\Phi_{i,j,k} - \Phi_{i,j,k-1}}{\Delta \phi_{k-\frac{1}{2}}} \right]\end{aligned}$$

- 3D Logical Domain Decomposition  
[MPI\_Cart\_create()]
- Operator matrix stored in **DIA** format  
(boundary conditions matrix-free)
- PC2** preconditioner stored in **CSR** format



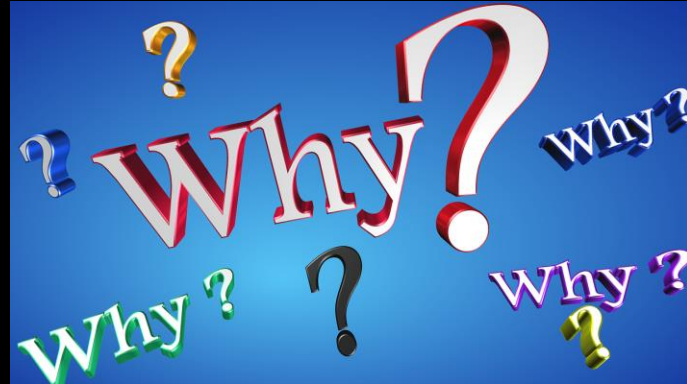
$(\Delta r, \Delta \theta, \Delta \phi)$

Proc #x





# MPI to MPI+OpenACC – Why?



- 1) Want to run “in-house”
  - a) Off-the-shelf desktops: 2-4 GPUs
  - b) HPC single-node rack: 8-16 GPUs



- 2) Use less HPC allocation \$\$\$
- 3) Performance Boost:  
Given allocation limits, can GPU code be faster?
- 4) Stepping stone:  
Our thermodynamic 3D MHD code  
has up to 80% run-time in PCG solvers

**Want to do all this in a *portable, single-source way!***

**-OK!**



As we go along, watch out for:



Difficulties encountered



Things to be careful of!





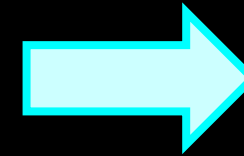
# MPI to MPI+OpenACC

Unstructured data regions

Essential\* and easy!

```
module fields  
  real(r_typ), dimension(:, :, :), allocatable :: phi  
end module
```

```
allocate ( phi (nr,nt,np) )  
phi = 0.  
[CPU sets initial values of phi...]  
!$acc enter data copyin (phi)
```



Compute with phi in  
OpenACC regions using  
“**present**” clause.....

```
!$acc update self (phi)  
[Areas where CPU needs phi (e.g. I/O)]
```

```
!$acc exit data delete(phi)  
deallocate (phi)
```

# MPI to MPI+OpenACC

## FORTRAN array operations

Need to use “kernels” or expand into loops

```
!$acc kernels present(x,p,ap,r)
  x=x+alpha*i*p
  r=r-alpha*i*ap
  ap=r
!$acc end kernels
```

```
4440, Generating present(ap(:),p(:),r(:),x(:))
4441, Loop is parallelizable
      Accelerator kernel generated
      Generating Tesla code
      4441, !$acc loop gang, vector(128) !
      blockidx%x threadidx%x
```



In current OpenACC 2.5, no control on things like vector length....

```
!$acc kernels present(x,br0)
  x( 1,2:ntm1,2:npm1)= x(2,2:ntm1,2:npm1)
  &      -vmask*br0(2:ntm1,2:npm1)*dr1
!$acc end kernels
```

```
6096, Generating present(br0(:,,:),x(:,,:))
6097, Loop is parallelizable
      Accelerator kernel generated
      Generating Tesla code
      6097, !$acc loop gang, vector(32) !
      blockidx%x threadidx%x
              !$acc loop gang, vector(4) !
      blockidx%y threadidx%y
```

OpenACC 2.6 draft:

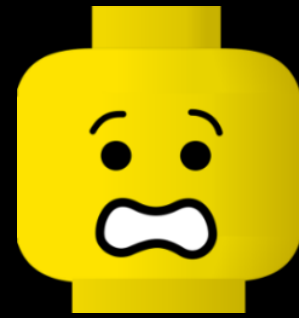
Add "num\_gangs", "num\_workers", and "vector\_length" to “kernels”!



# MPI to MPI+OpenACC

## Multiple devices

The simplest method:  
1 GPU per MPI rank  
How easy?



OpenACC with multiple devices is  
*NOT* that scary ...  
... even on a large cluster using MPI!

## How about in three lines!

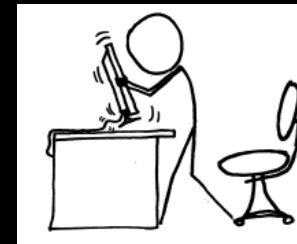
```
ngpus_per_node = 4  
igpu = MODULO(iprocw, ngpus_per_node)  
!$acc set device_num(igpu)
```



Make sure to submit job with  
-ntasks-per-node=ngpus\_per\_node  
(or equivalent)



Assumes linear affinity of MPI ranks!  
To be safe, should set custom host\_file  
or mpiexec CPU bindings manually....



PGI 16.10: ~~!\$acc set device\_num~~  
Had to use non-portable API calls

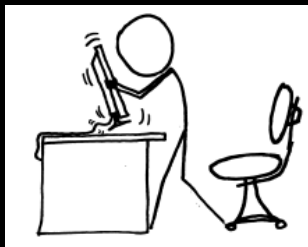
```
use openacc  
typ = acc_get_device_type  
call acc_init(typ)  
igpu = MODULO(iprocw, ngpus_per_node)  
call acc_set_device_num(igpu, typ)
```

PGI >17.1: !\$acc set device\_num supported!

# MPI to MPI+OpenACC

## CUDA-aware MPI

- ❖ MPI finite difference communication involves loading buffers to send data to neighboring processors
- ❖ Do **NOT** want buffers filled on CPU!
- ❖ **CUDA-aware MPI**: Send device pointers directly to MPI API and “it just works”! 😊
- ❖ This allows use of **GPUDirect/NVlink**!



Not all HPC system's MPI libraries up-to-date and/or fully tested with CUDA-aware MPI

```
!$acc kernels present(a,sbuf11,sbuf12)
    sbuf11(:, :)=a(n1-1, :, :)
    sbuf12(:, :)=a( 2, :, :)
!$acc end kernels
```

← LOAD BUFFERS  
ON GPU

```
!$acc host_data use_device(sbuf11,sbuf12,rbuf11,rbuf12)
    call MPI_Isend (sbuf11,lbuf,ntype_real,iproc_rp,tag,
&                  comm_all,reqs(1),ierr)
    call MPI_Isend (sbuf12,lbuf,ntype_real,iproc_rm,tag,
&                  comm_all,reqs(2),ierr)
    call MPI_Irecv (rbuf11,lbuf,ntype_real,iproc_rm,tag,
&                  comm_all,reqs(3),ierr)
    call MPI_Irecv (rbuf12,lbuf,ntype_real,iproc_rp,tag,
&                  comm_all,reqs(4),ierr)
    call MPI_Waitall (4,reqs,MPI_STATUSES_IGNORE,ierr)
!$acc end host_data
```

```
    if (iproc_rm.ne.MPI_PROC_NULL) then
!$acc kernels present(a,rbuf11)
    a( 1, :, :)=rbuf11(:, :)
!$acc end kernels
    end if
    if (iproc_rp.ne.MPI_PROC_NULL) then
!$acc kernels present(a,rbuf12)
    a(n1, :, :)=rbuf12(:, :)
!$acc end kernels
    end if
```

← UNLOAD BUFFERS  
ON GPU



# MPI to MPI+OpenACC

## Calling cuSparse from FORTRAN (PC2 only!)

- ❖ PC2: Requires triangular solves (vector-unfriendly)
- ❖ Advanced algorithms difficult to implement... Use freely-available libraries! **cuSparse!**
- ❖ Portability: Call C code from FORTRAN (instead of FORTRAN bindings)
- ❖ For compiler-independent, “correct” way, see: *“Parallel Programming in OpenACC”* pg. 283

```
void lusol_cusparse(double* x, double* CSR_LU,
                  int* CSR_I, int* CSR_J, int N, int M)
{ // Forward solve (Lx=y)
  cusparseDcsrsv_solve(cusparseHandle,
    CUSPARSE_OPERATION_NON_TRANSPOSE, N, &one, L_described,
    CSR_LU, CSR_I, CSR_J, L_analyzed, x, x);
  // Backward solve (Uy=x)
  cusparseDcsrsv_solve(cusparseHandle,
    CUSPARSE_OPERATION_NON_TRANSPOSE, N, &one, U_described,
    CSR_LU, CSR_I, CSR_J, U_analyzed, x, x);
  cudaDeviceSynchronize();
}
```

**lusol.c**

```
module cusparse_interface
  interface
    subroutine lusol_cusparse(
      &          x, CSR_LU, CSR_AI, CSR_AJ, N, M)
      &          BIND(C, name="lusol_cusparse")
    use, intrinsic :: iso_c_binding
    integer(C_INT), value :: N
    type(C_PTR), value :: x, CSR_LU, CSR_AI, CSR_AJ
  end subroutine lusol_cusparse
end interface
end module
```

```
use, intrinsic :: iso_c_binding
use cusparse_interface
integer(c_int) :: cN, cM
```

```
!$acc host_data use_device(x, a_csr, a_csr_ja, ...)
  call lusol_cusparse(C_LOC(x(1)),
    &          C_LOC(a_csr(1)),
    &          C_LOC(a_csr_ia(1)),
    &          C_LOC(a_csr_ja(1)), cN, cM)
!$acc end host_data
```

**Compile!**

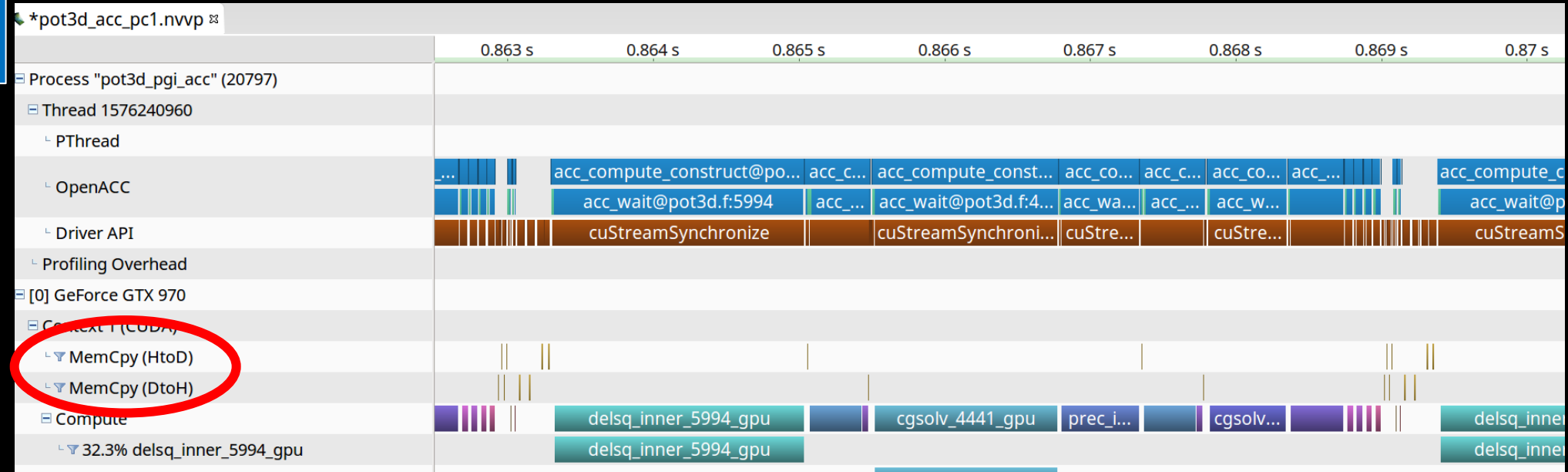
```
nvcc -c [FLAGS] lusol.c
pgf90 [FLAGS] lusol.o [LIBS] pot3d.f
```

# MPI to MPI+OpenACC - Profiling

PROFILE!

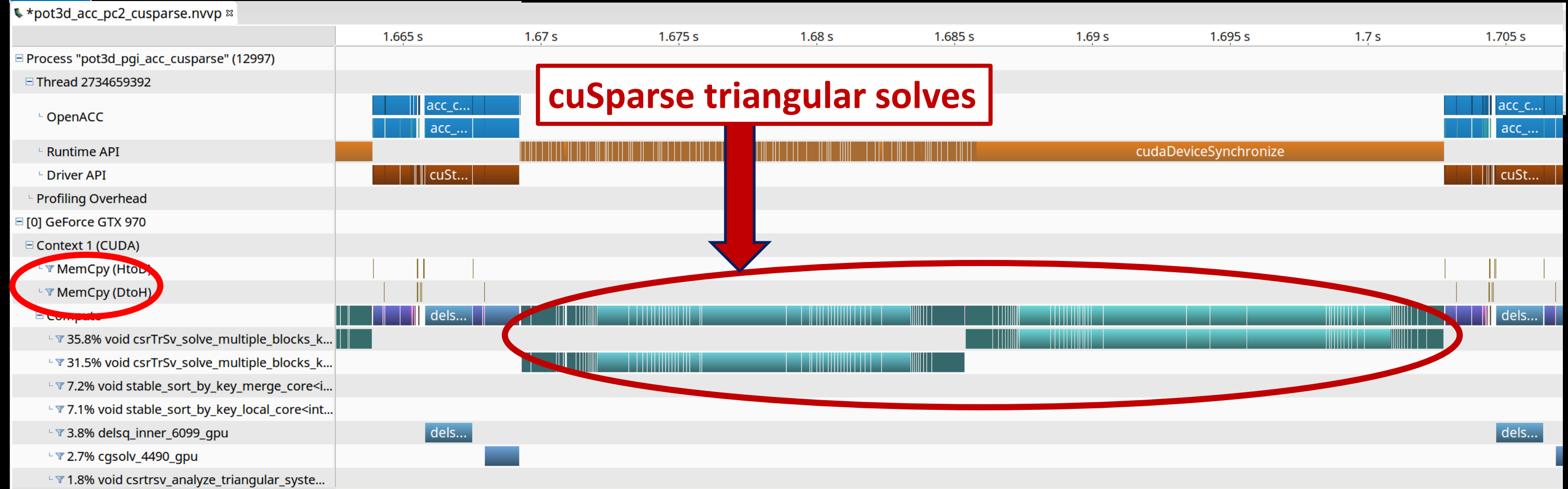
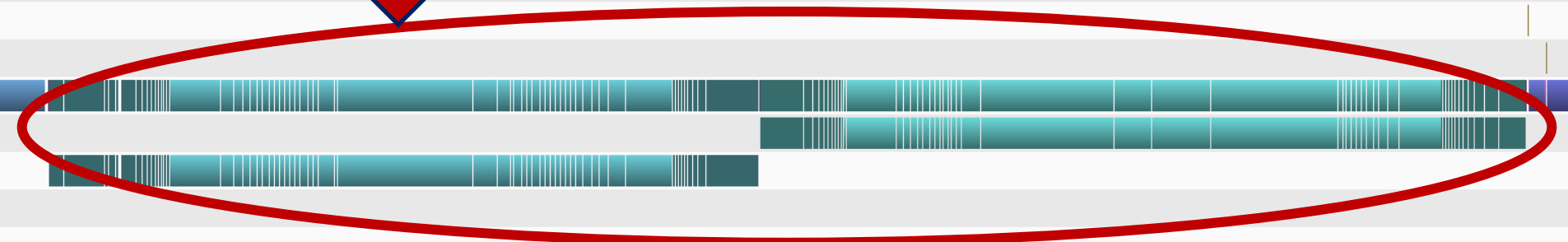
PGPROF

PC1



PC2

cuSparse triangular solves





# Summary of MPI to MPI+OpenACC

## Modification

## Portability

PC1

~1%

~90 lines of comments

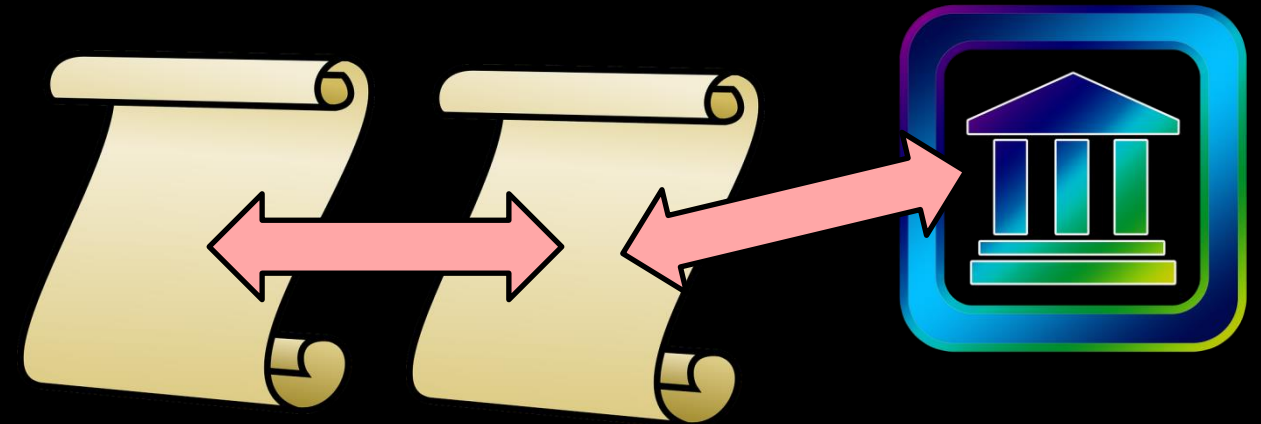


Single source for GPU and CPU!

PC2

~2%

~150 lines of comments  
and new code



New source code and library dependency  
(Could make 'portable' with 'ifdefs'...)

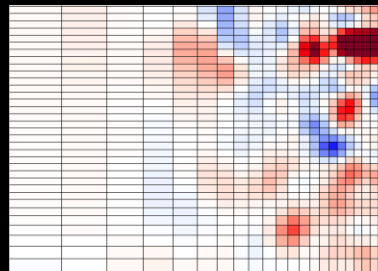
# Test Problem

- Real-world test problem
- Boundary magnetic map from HMI[CR2099+Daily(8/1/2010)]
- Used to study sympathetic solar storms [Titov, et al, Ap.J. 759:70 (2012)]

## Three problem sizes

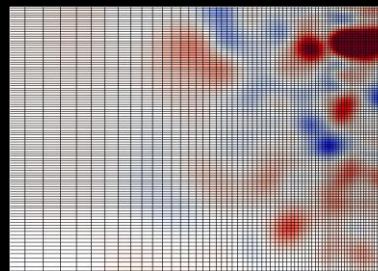
**Small**

**120x120x240**  
~ 3 million points



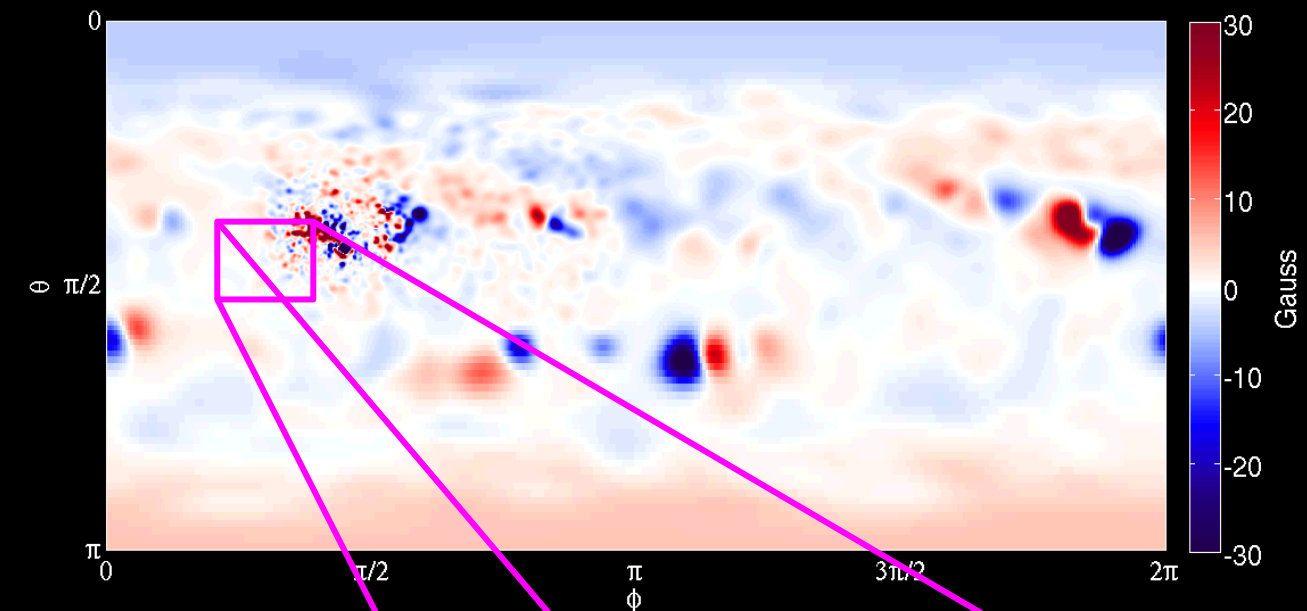
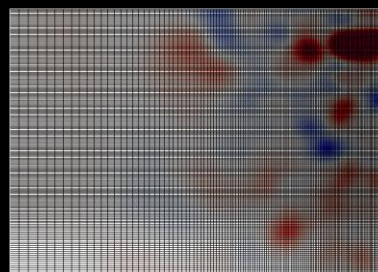
**Medium**

**100x360x720**  
~25 million points



**Large**

**200x720x1440**  
~207 million points



PCG Iterations

PC1	3883
PC2	0697 $\rightarrow$ 0874

PCG Iterations

PC1	8719
PC2	1273 $\rightarrow$ 1871

PCG Iterations

PC1	19995
PC2	2563 $\rightarrow$ 3221

	N	min $\Delta$	max $\Delta$	max $\Delta\Delta$
$\theta$	480	0.1°	2.6°	$\approx 15\%$
$\phi$	1151	0.1°	0.9°	$\approx 4\%$

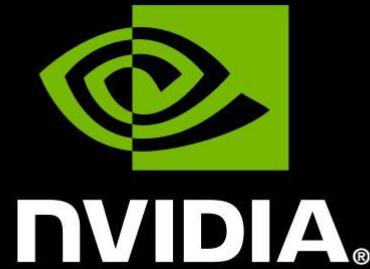


# HPC Environment

## SYSTEMS



Allocation provided by:



NVIDIA PSG Cluster

Allocation provided by:  
NVIDIA Cooperation



Local Desktop

## SOFTWARE



2015.2.164



**MVAPICH**

for infiniband v2.1



16.10 + CUDA 8.0



**OpenMPI**

1.10.2

## HARDWARE SPECS

	COMET NODE	DESKTOP
CPU Model	(2x) Intel Xeon E5-2680v3	Intel Xeon E5-1650v4
Base Clock	2500 MHz	3600 MHz
# Cores	24	6
DP FLOPS	$\approx 0.96$ TFlop/s	$\approx 0.35$ TFlop/s
RAM	128 GB	32 GB
Mem Band	68 GB/sec	76.8 GB/s

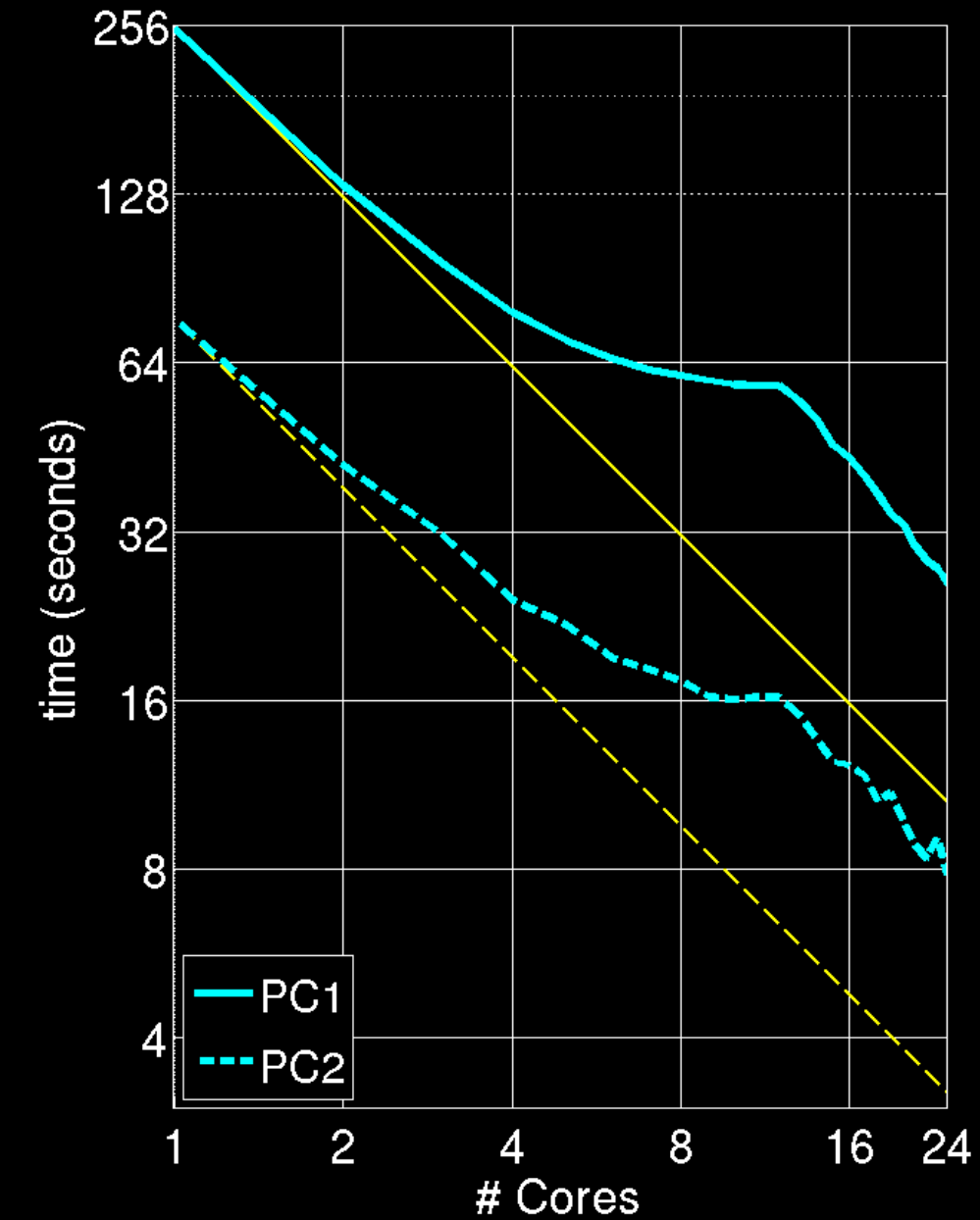
	COMET NODE	NVPSG NODE
GPU Model	NVIDIA Tesla K80 (x1)	NVIDIA Tesla P100 (PCIe)
Base Clock	560 MHz	1328 MHz
# DP Cores	832	1792
DP FLOPS	$\approx 0.93$ TFlops	$\approx 4.7$ TFlops
RAM	12 GB	16 GB
Mem Band	240 GB/sec	549 GB/sec



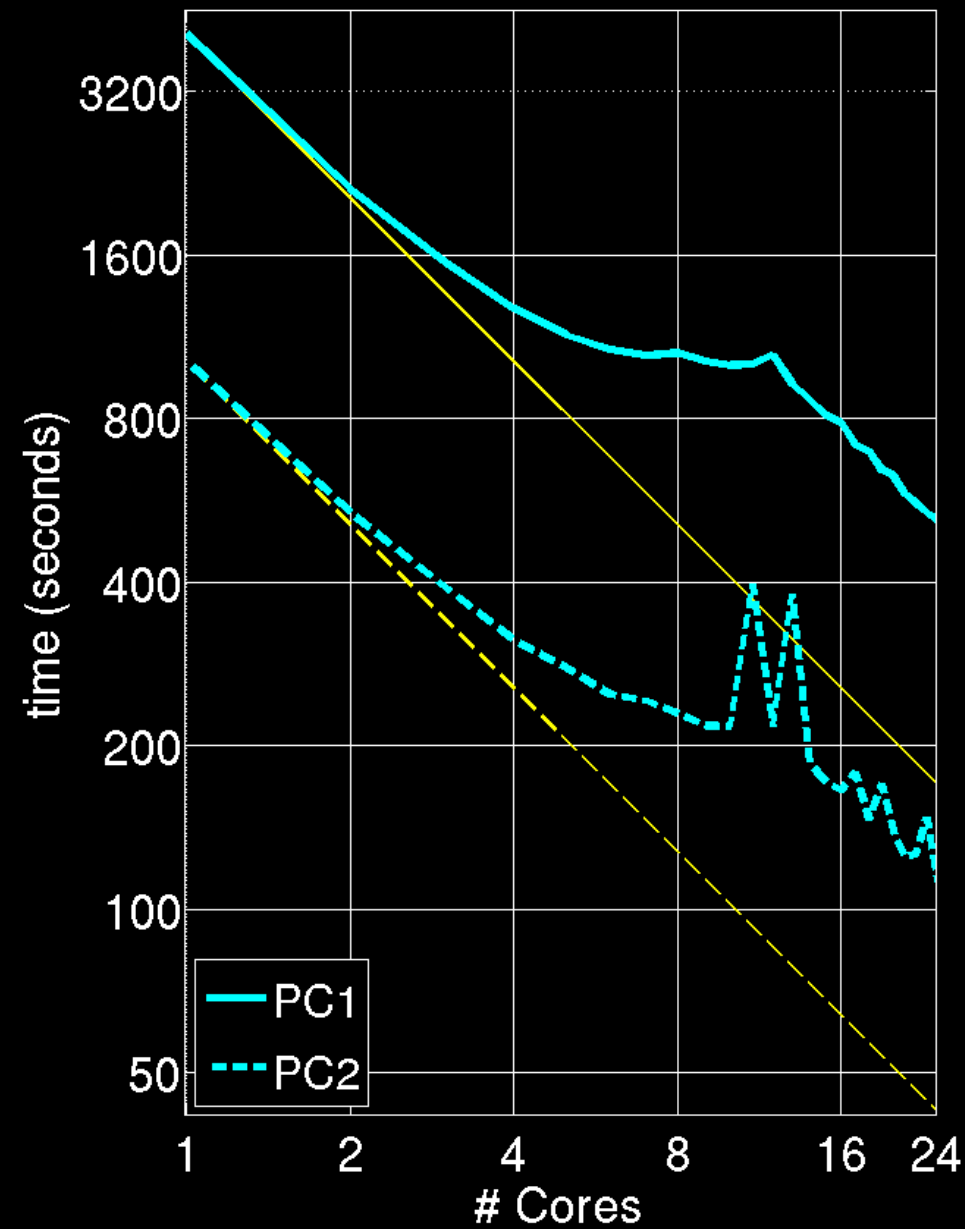
All timing results made with  
MPI timers for solve and setup  
(i.e. excludes I/O)

# Timings - Original CPU Scaling

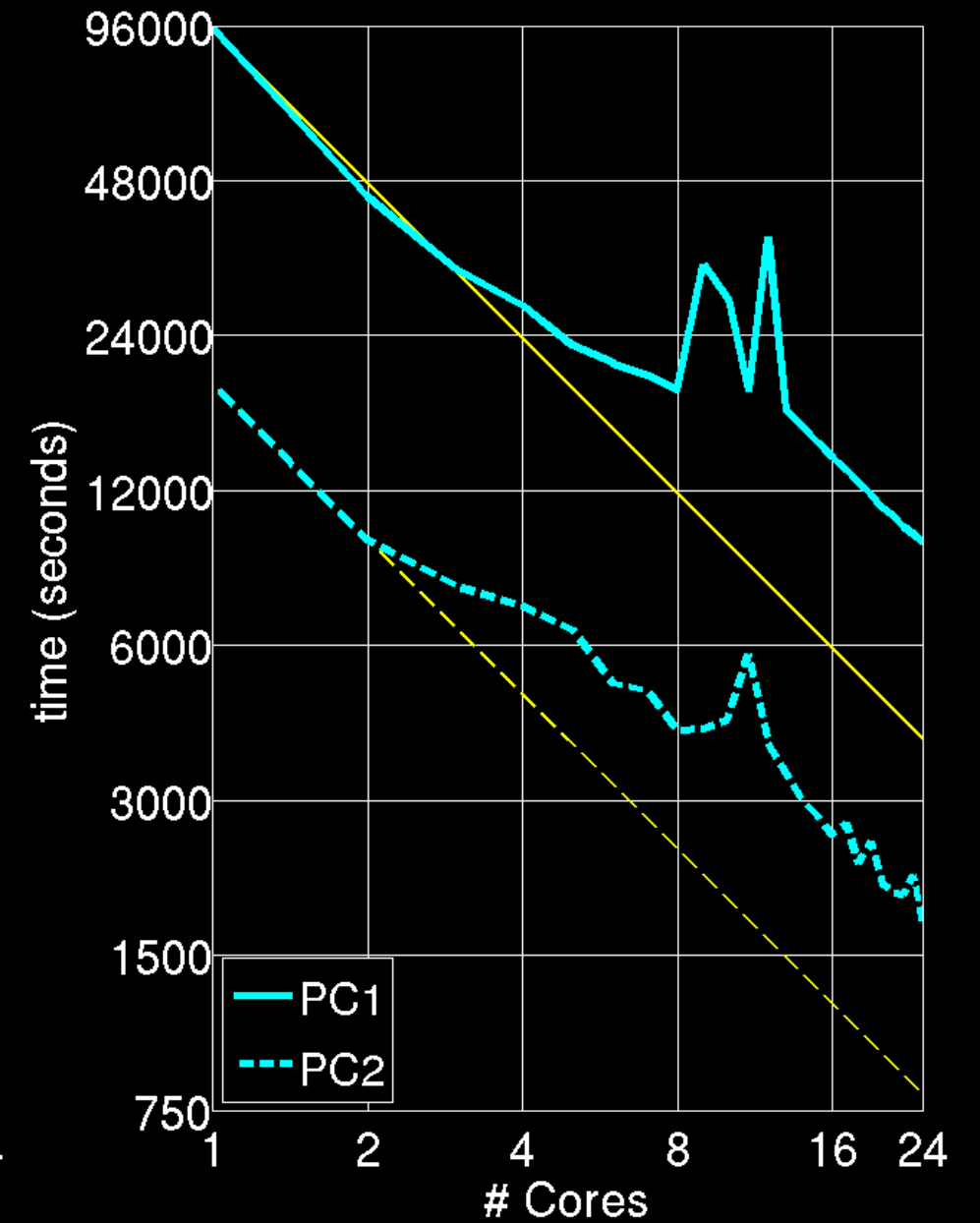
SMALL (SINGLE-NODE)



MEDIUM (SINGLE-NODE)



LARGE (SINGLE-NODE)



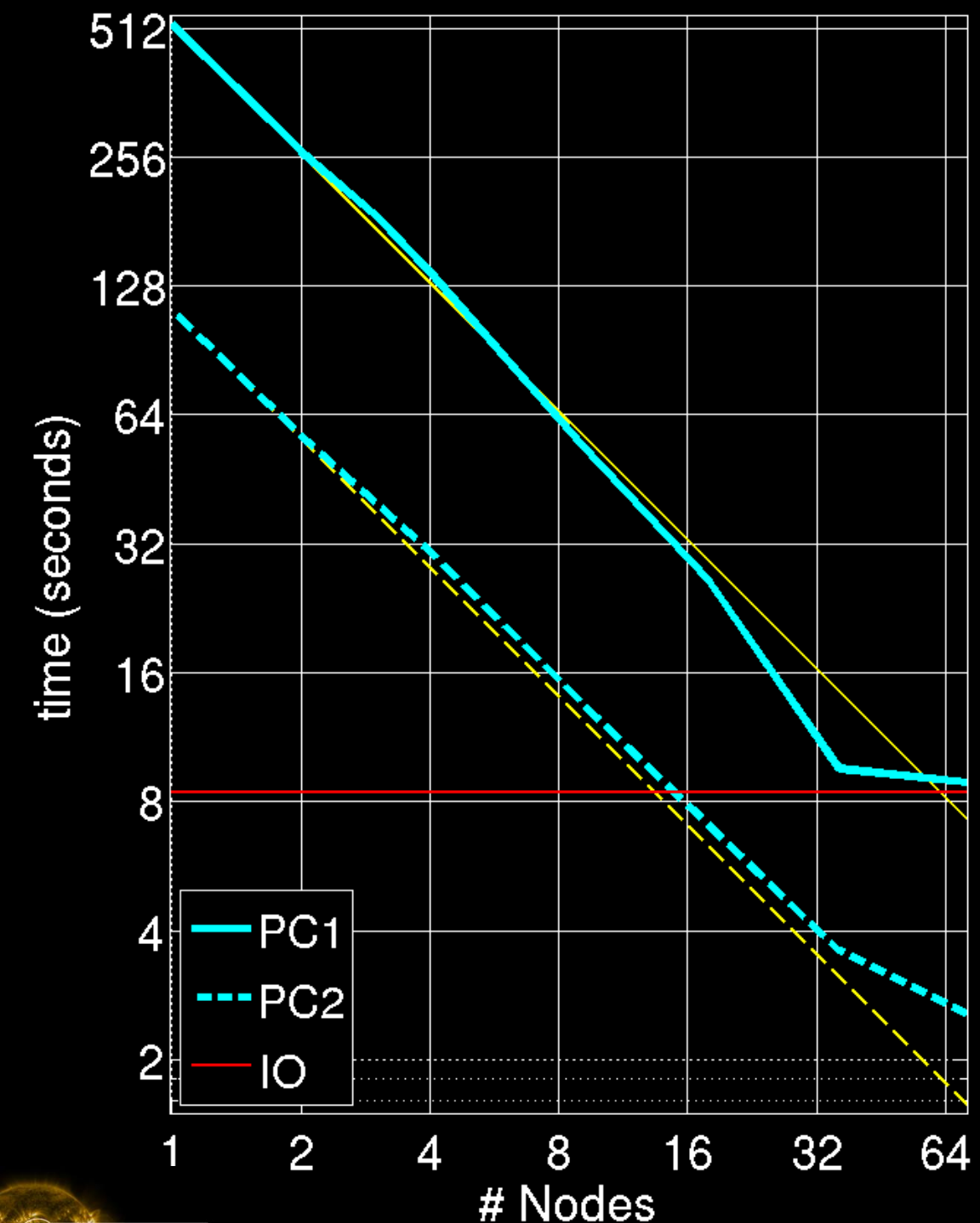
☹ PC2 more than twice as fast as PC1!

☹ Strange scaling behavior

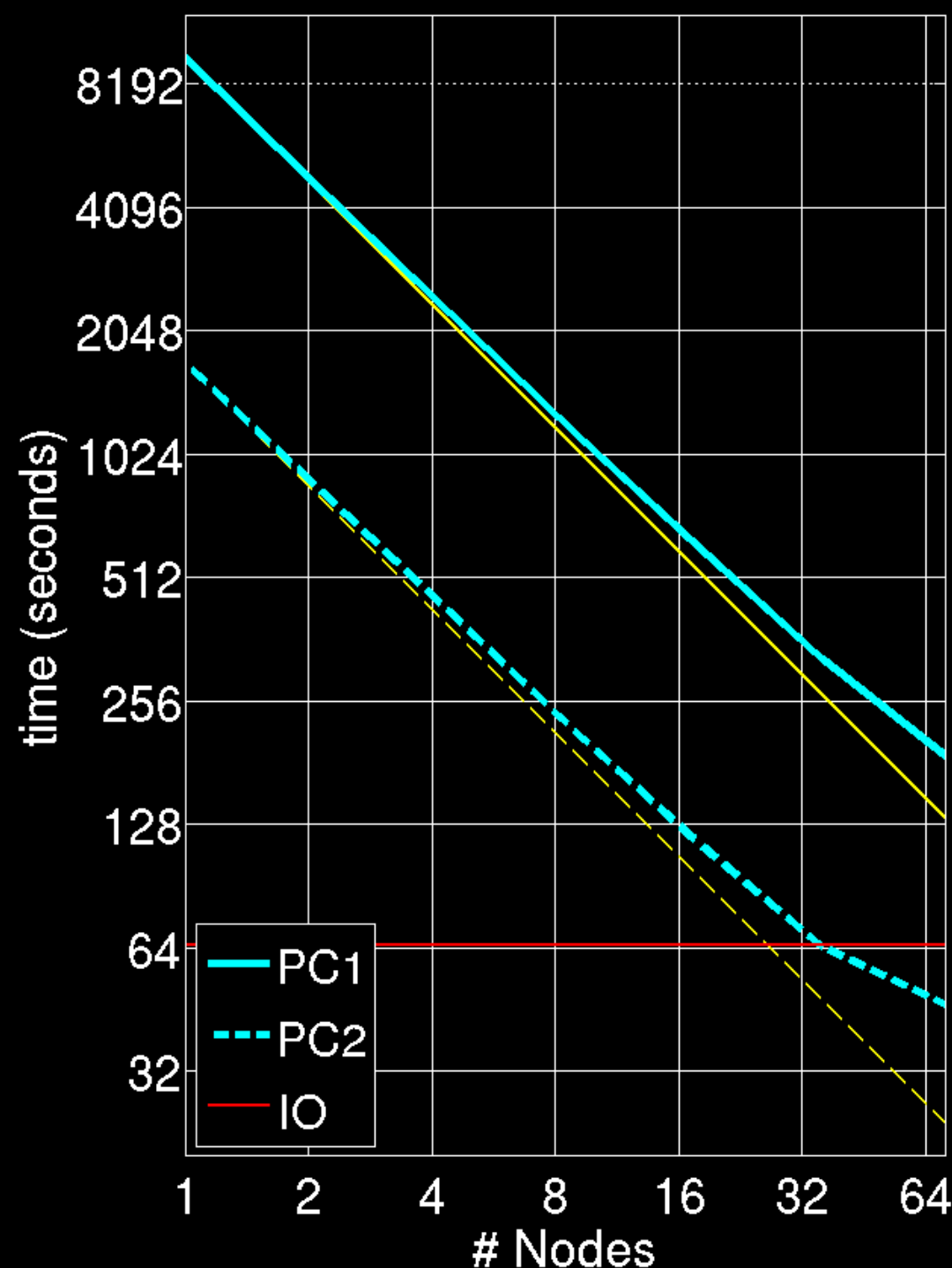


# Timings - Original CPU Scaling

## MEDIUM (MULTI-NODE)



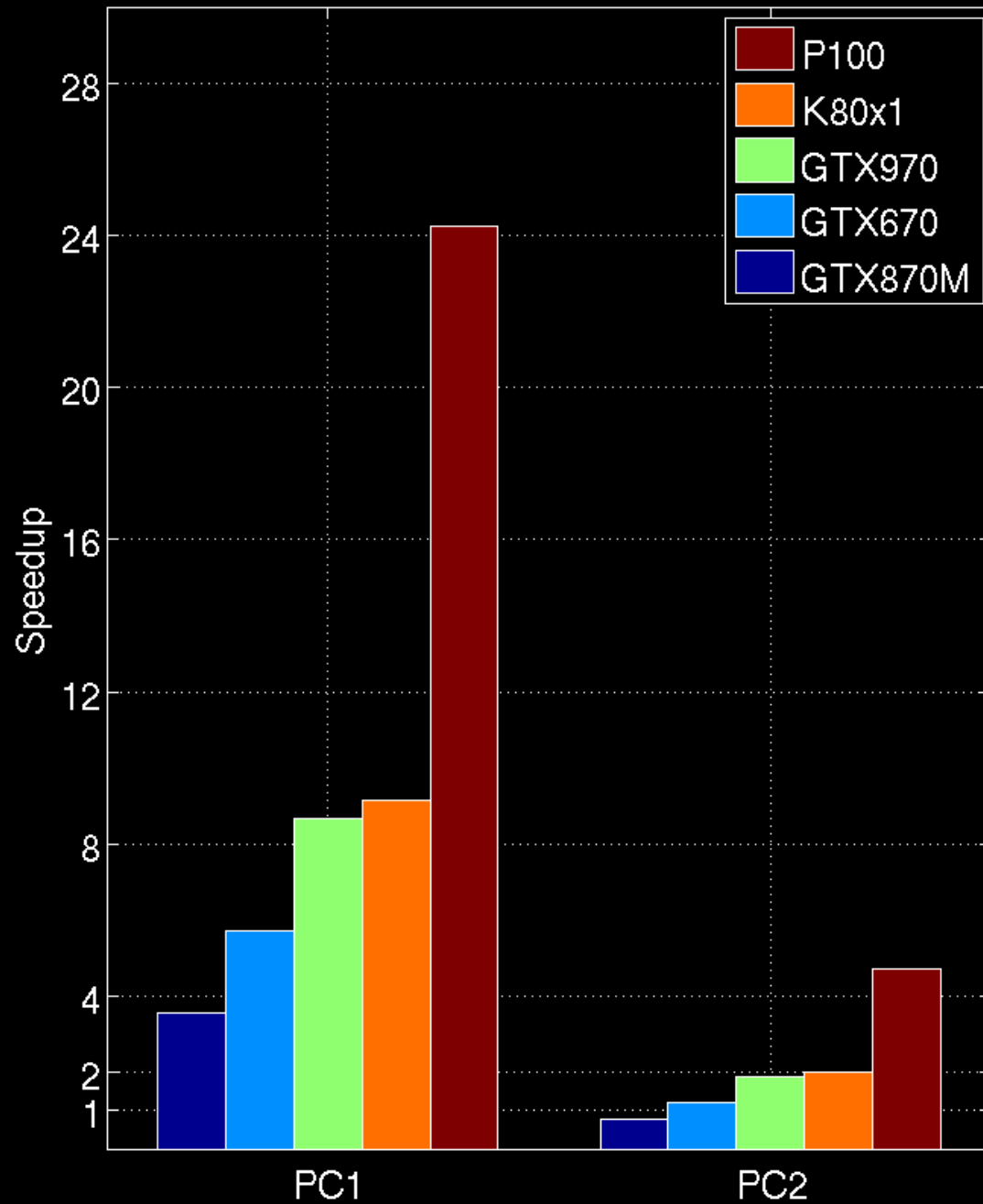
## LARGE (MULTI-NODE)



- ⌘ Multi-node scaling much better than single-node
- ⌘ Scaling starts to weaken near the 40-50 node level
- ⌘ Algorithm hits I/O time for largest runs

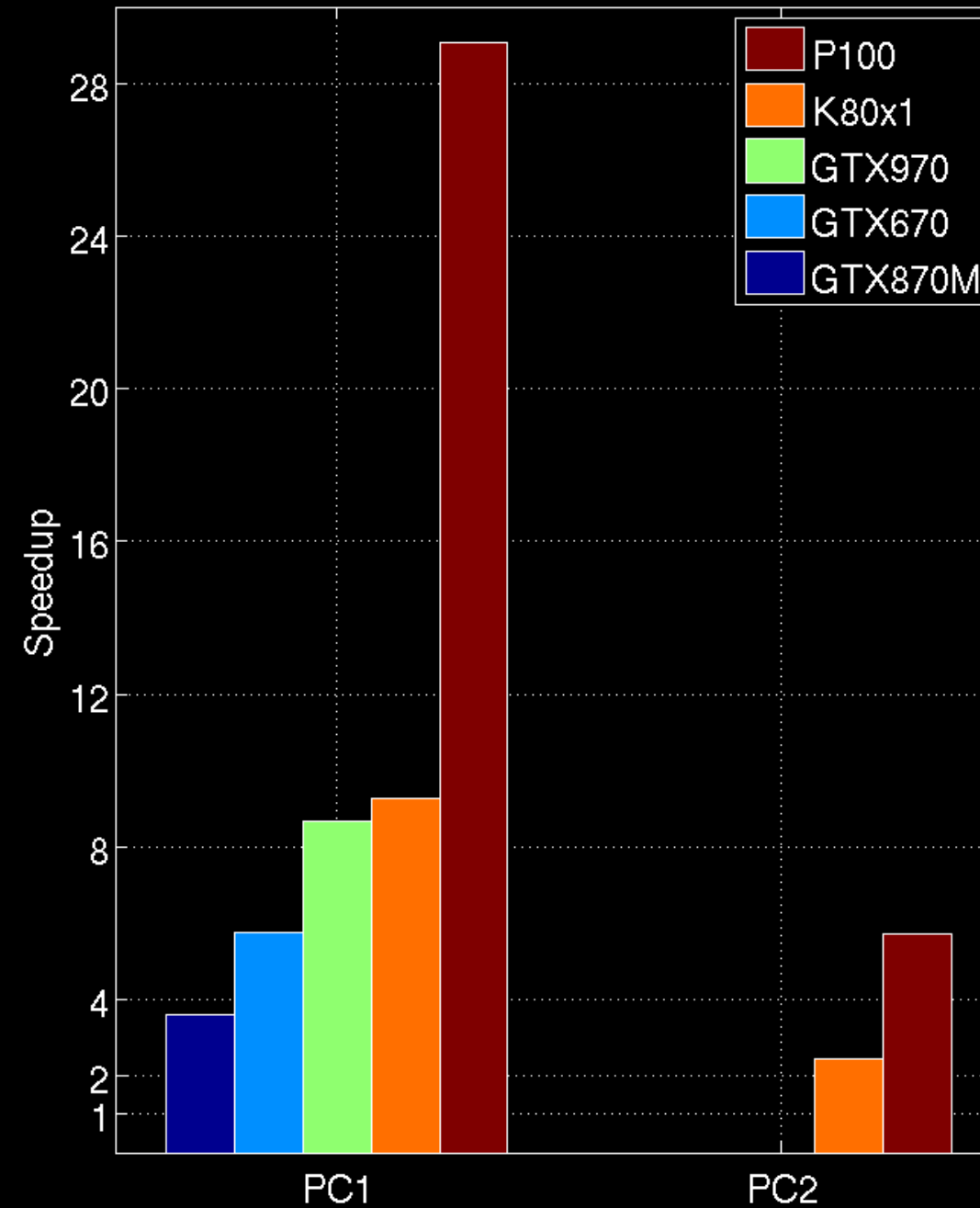
# Timings - Single GPU vs Serial CPU

Single GPU speedup over 1 core of Xeon E5-1650v4 3.6 GHz



**SMALL**

Single GPU speedup over 1 core of Xeon E5-1650v4 3.6 GHz



**MEDIUM**

⚡ **PC1** speedup  
much better  
than **PC2**

⚡ Pascal much  
faster than  
Kepler

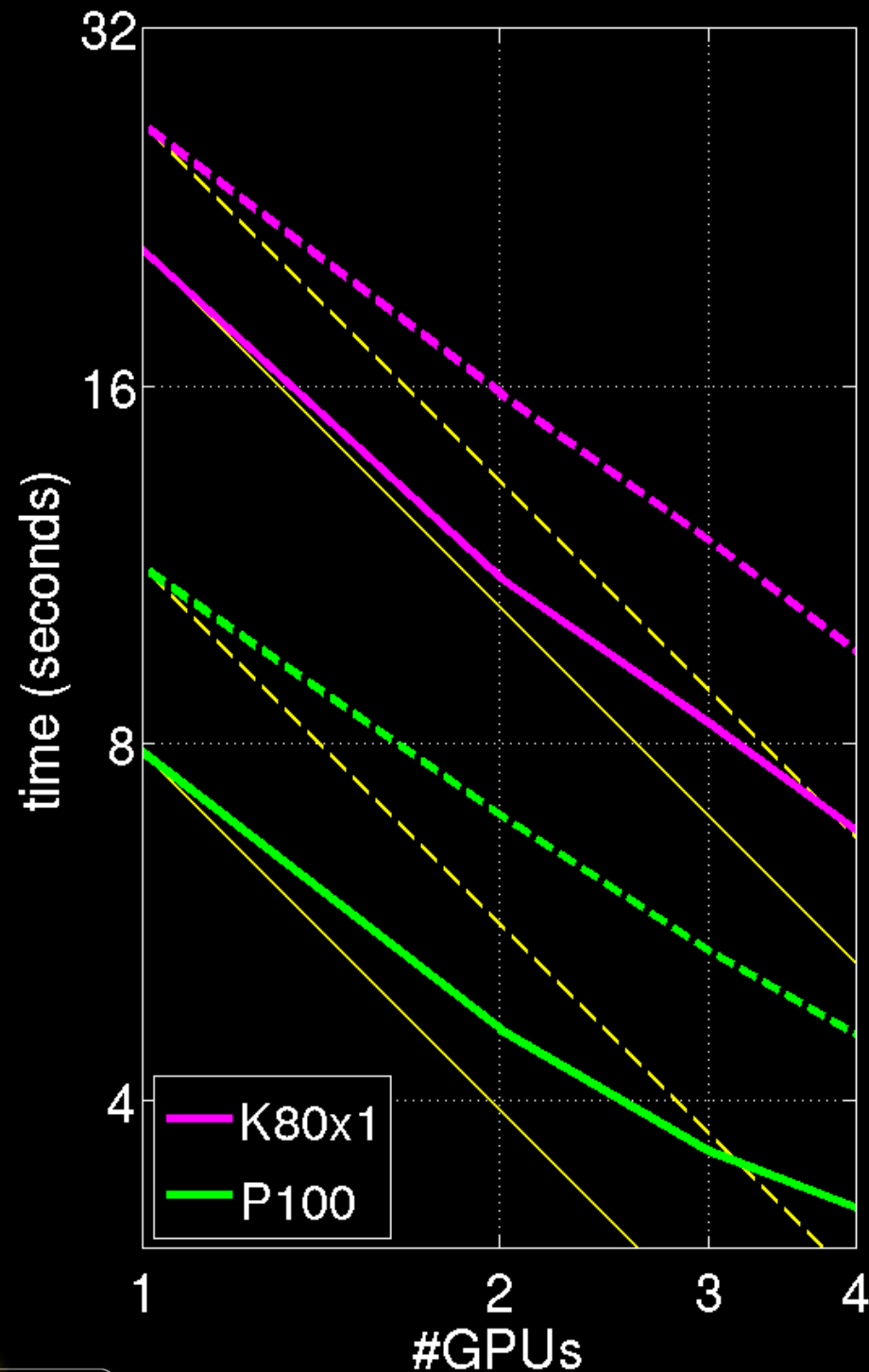


**Problem size  
important!  
(must fit on GPU!)**

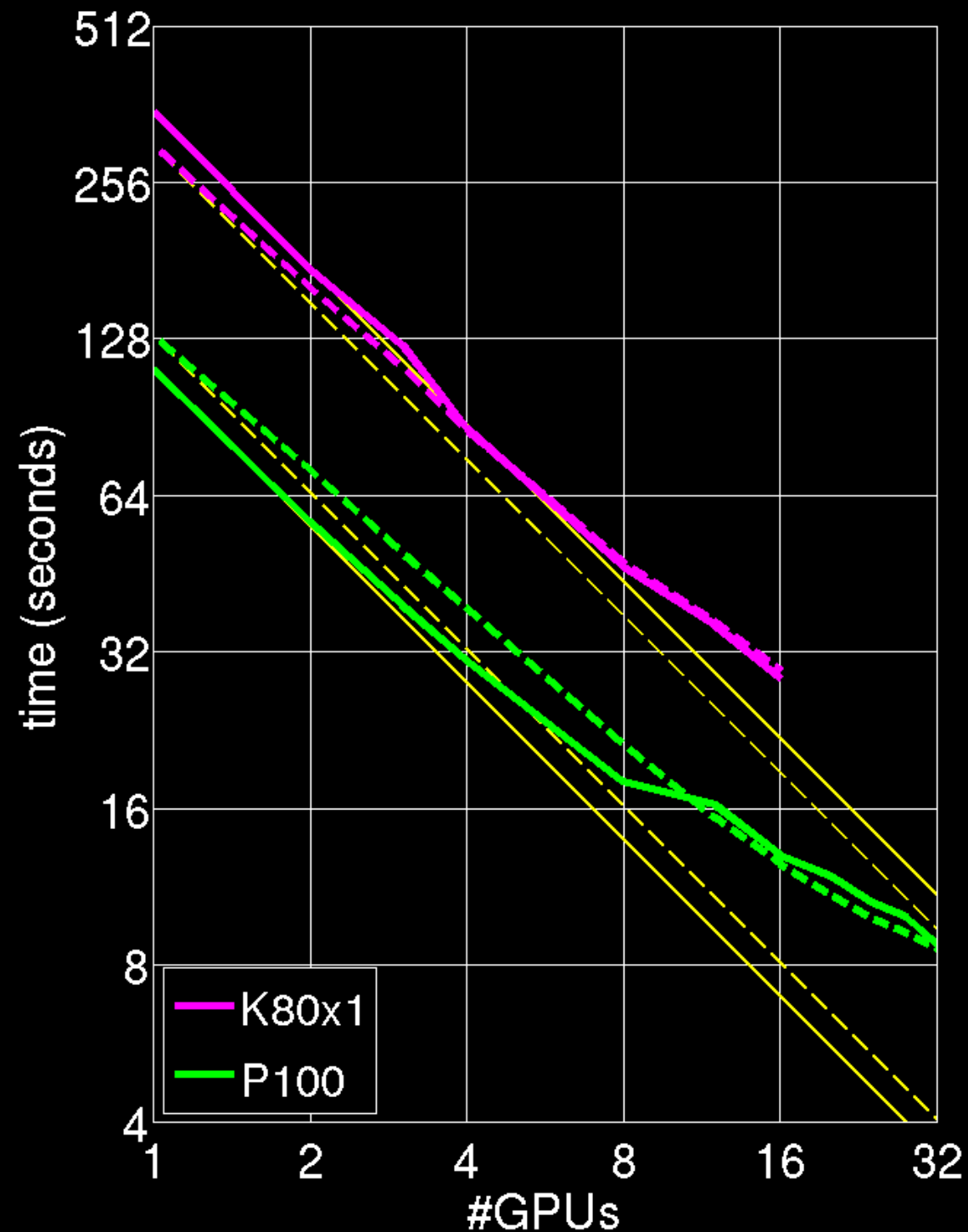


# Timings - GPU Scaling

## SMALL (MULTI-GPU)



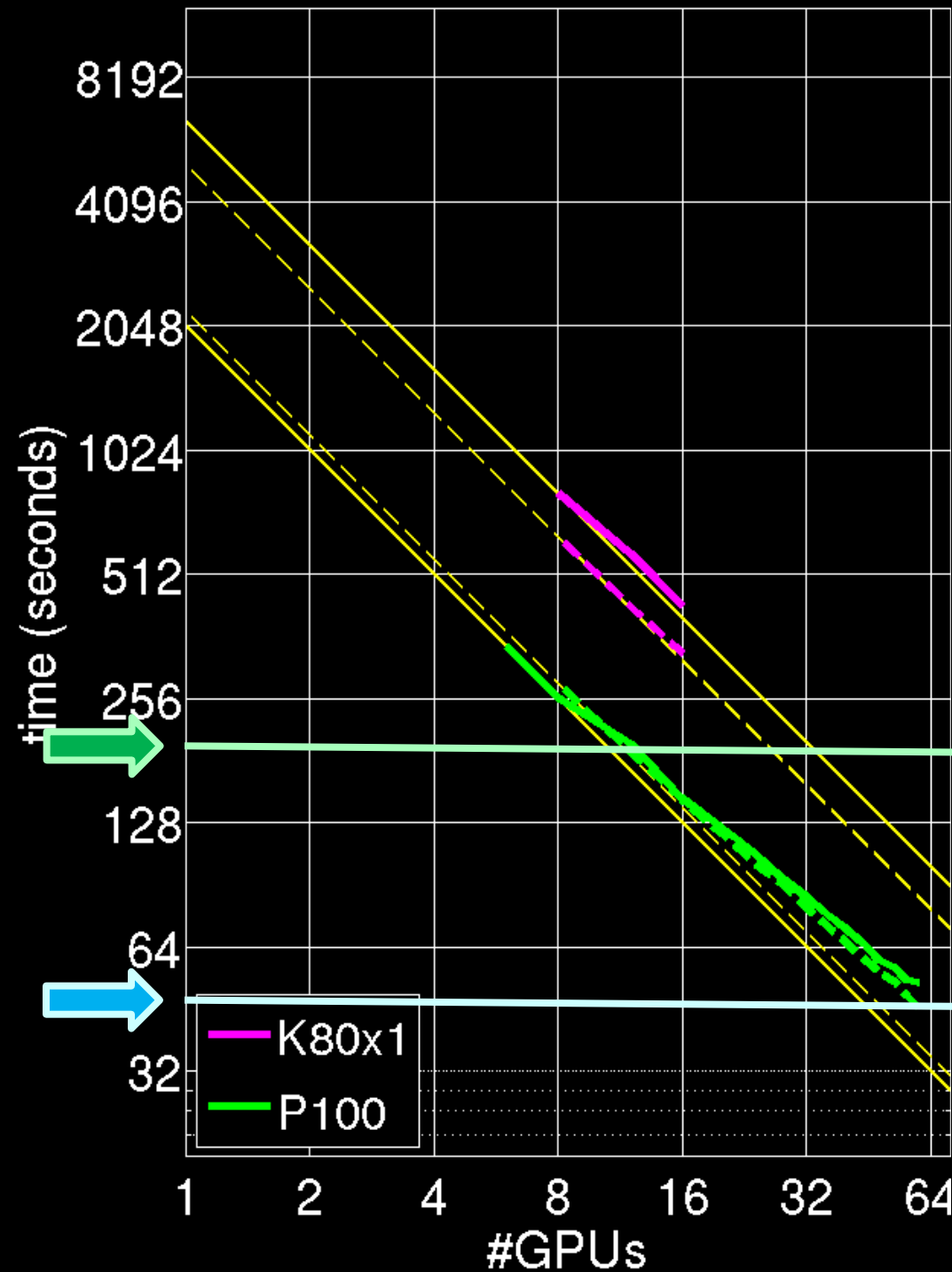
## MEDIUM (MULTI-GPU)



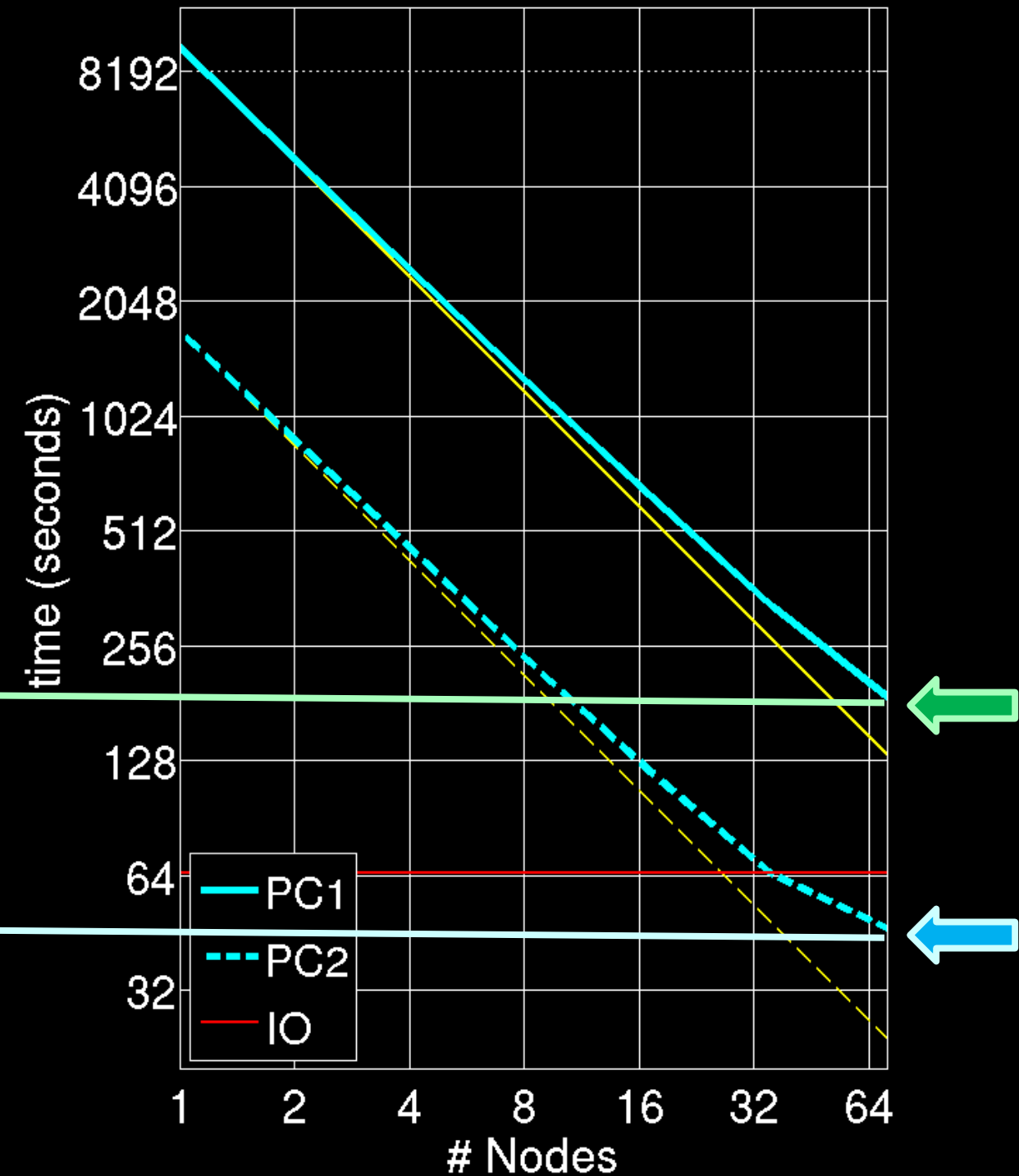
- Ⓢ Pascal more than twice as fast as Kepler
- Ⓢ PC2 slower than PC1 for small run!
- Ⓢ PC1~PC2 for medium run

# Timings - GPU Scaling

LARGE (MULTI-GPU)



LARGE (MULTI-NODE)



- Scaling has improved
- P100: PC1~PC2
- PC1: 14-GPU run faster than CPU run on max 72 nodes
- PC2: 60-GPU run faster than CPU run on max 72 nodes

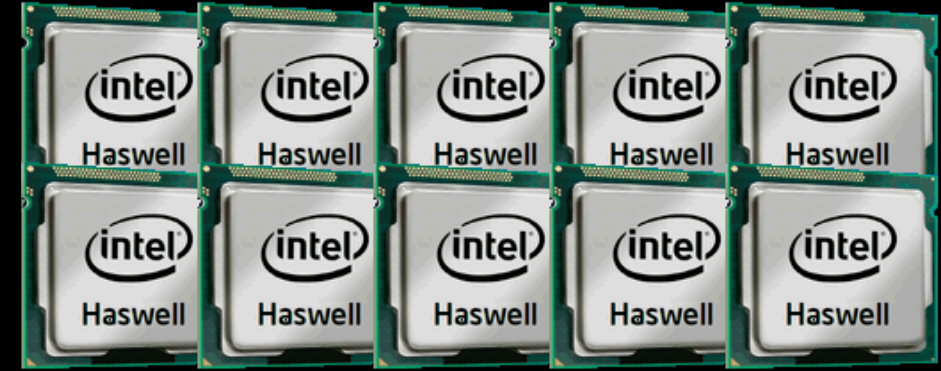


# Summary of GPU Performance (Large Run)

PC1



1x P100

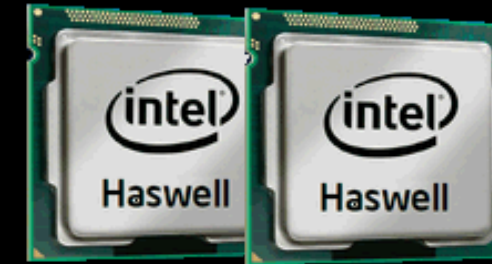


5x 2x12-core Haswell

PC2



1x P100



1x 2x12-core Haswell

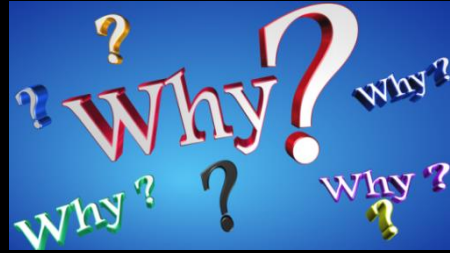
- ⌘ Triangular solve performance on GPUs (in this CSR case) needs improvement
- ⌘ Are there vectorizable preconditioners with comparable performance as ILU?

**Is a simple diagonal scaling the best preconditioner for conjugate gradients on supercomputers?**

Giorgio Pini and Giuseppe Gambolati

*[Adv. Water Resources, 1990, Vol. 13, No. 3]*

# Summary of MPI+OpenACC POT3D on NVIDIA GPUs



## *How did we do?*

PC1

PC2

### 1) Wanted to run "in-house"

Desktop with 4x P100s ~ 20x/04x HPC 24-core nodes for PC1/PC2  
HPC Node with 16x P100s ~ 80x/16x HPC 24-core nodes for PC1/PC2



### 2) Wanted to use less HPC allocation

With 4x P100s per node, allocation savings of  
~ 20x/4x for PC1/PC2



### 3) Wanted a performance boost

PC1: 60x P100s yield 4x speedup over max 72-nodes CPU  
PC2: 60x P100s yield only slight speedup over max 72-nodes CPU



**Wanted to do all this in a portable, single-source way...**





# Performance Portability

*One code to run on them all... efficiently...*

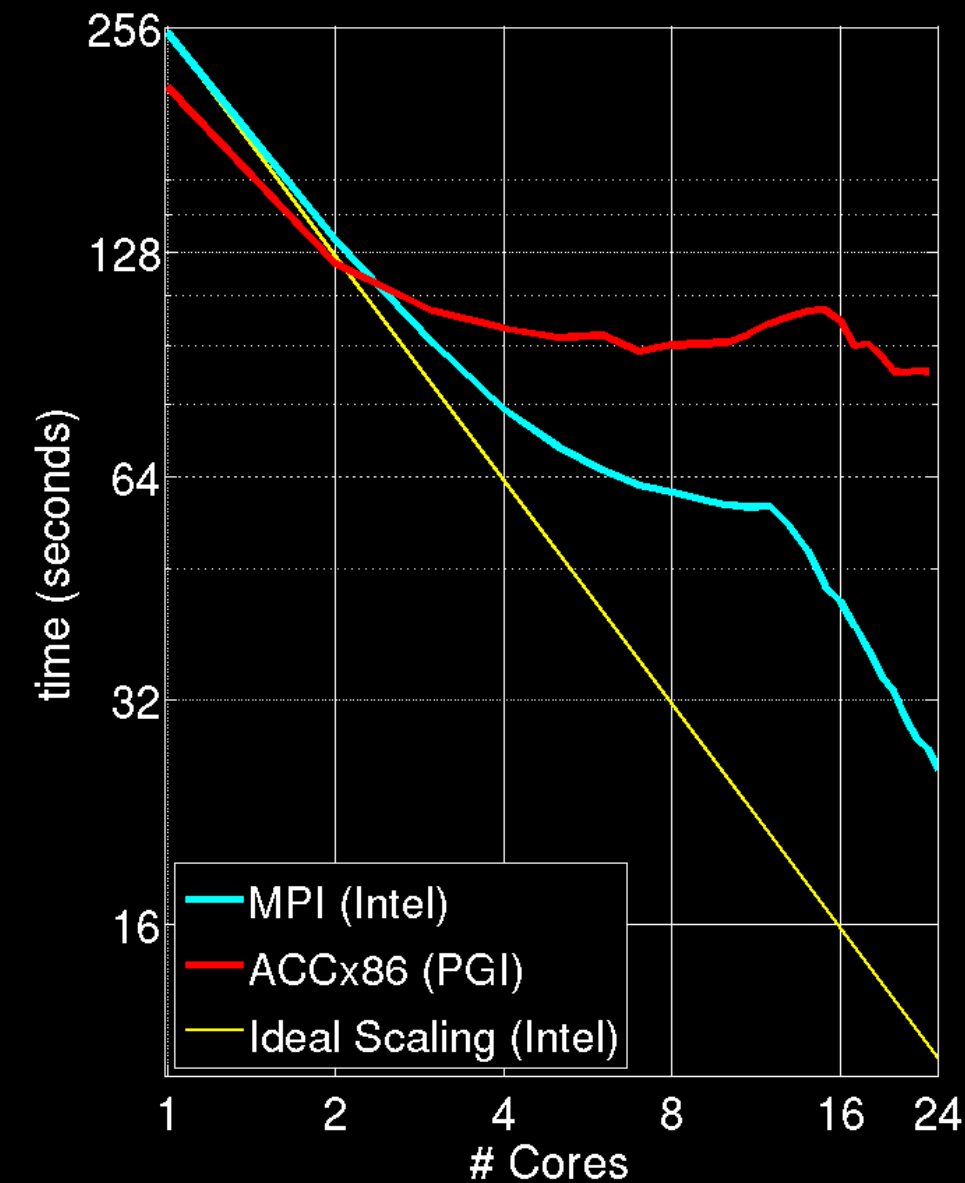


- ❧ OpenACC for x86 Multicore
- ❧ OpenACC for NVIDIA GPUs
- ❧ OpenACC for AMD GPUs
- ❧ OpenACC for OpenPower
- ❧ OpenACC for Intel KNL

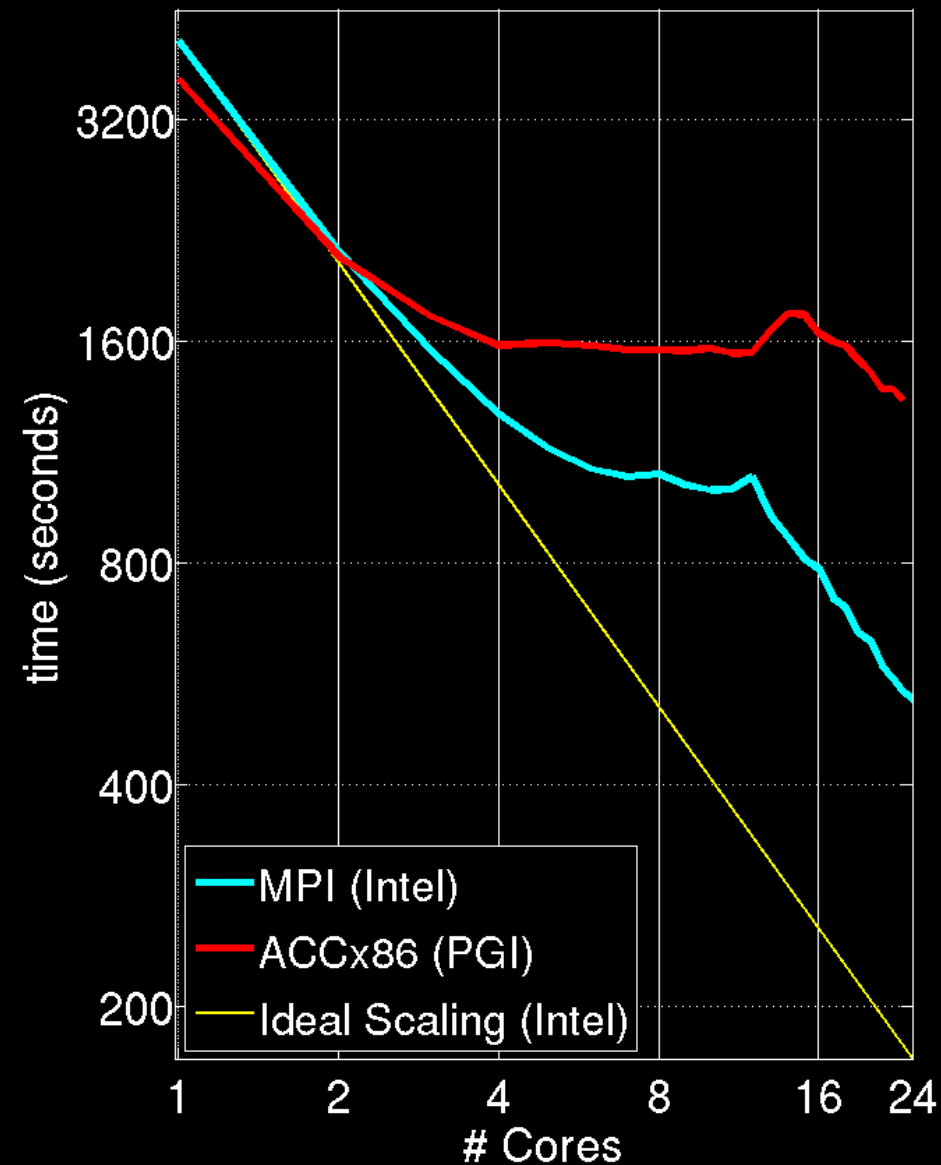
# Performance Portability (x86 multicore)

For **PC1**, with ***no code modifications***, can compile for multi-core x86 CPU (-ta=multicore)

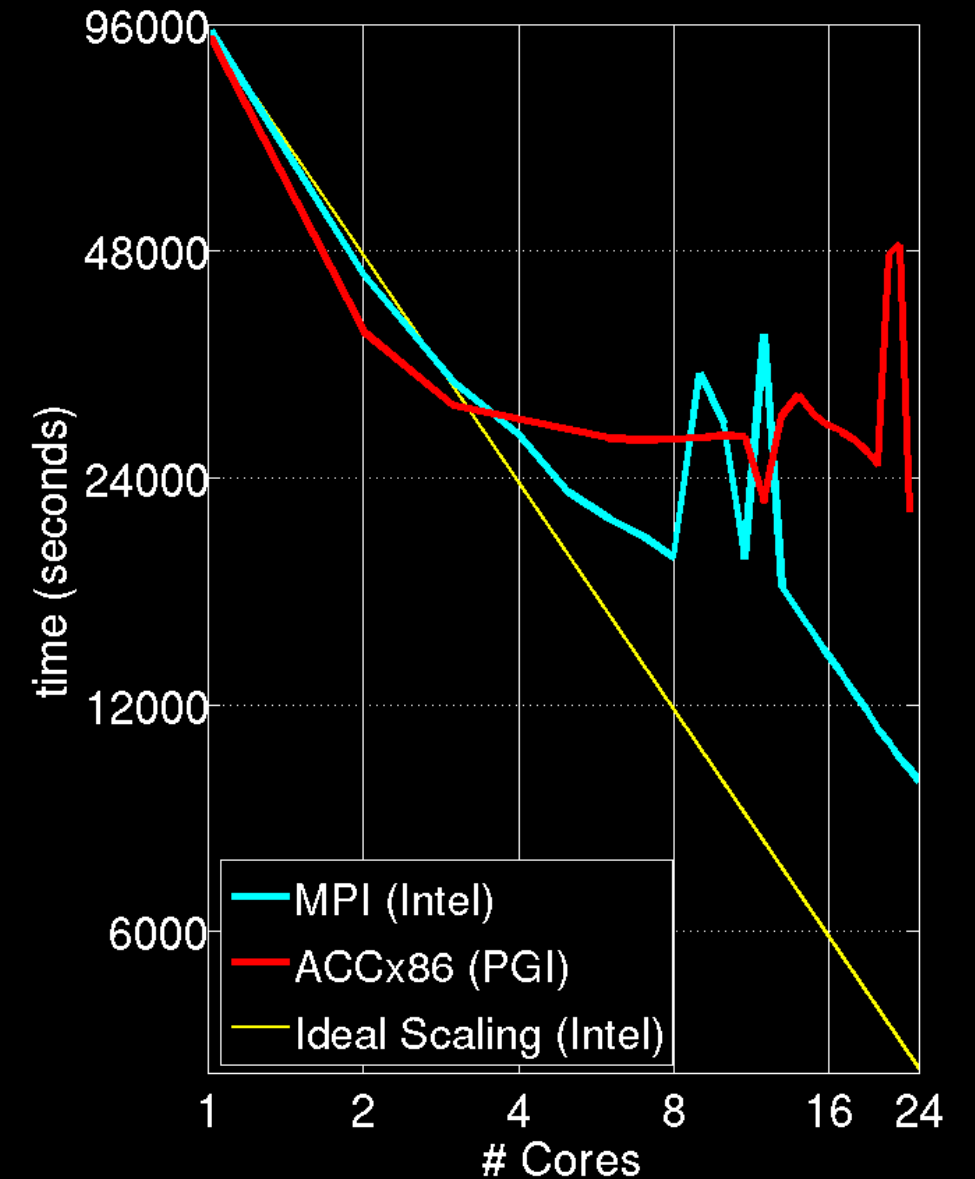
SMALL (SINGLE-NODE) PC1



MEDIUM (SINGLE-NODE) PC1



LARGE (SINGLE-NODE) PC1

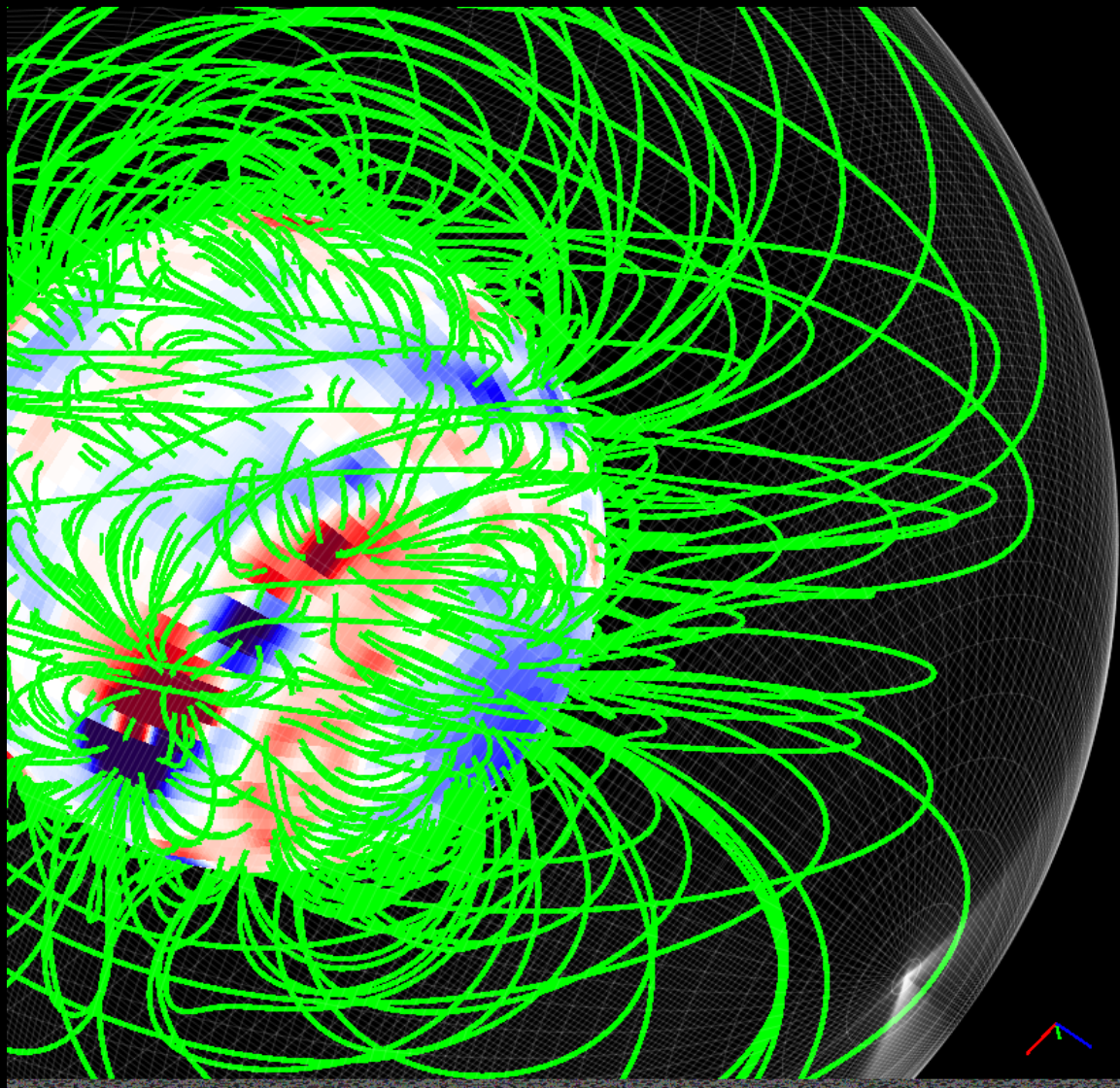


❶ OpenACC x86 scaling poorer than MPI-only (expected to match OpenMP)



# Future Work

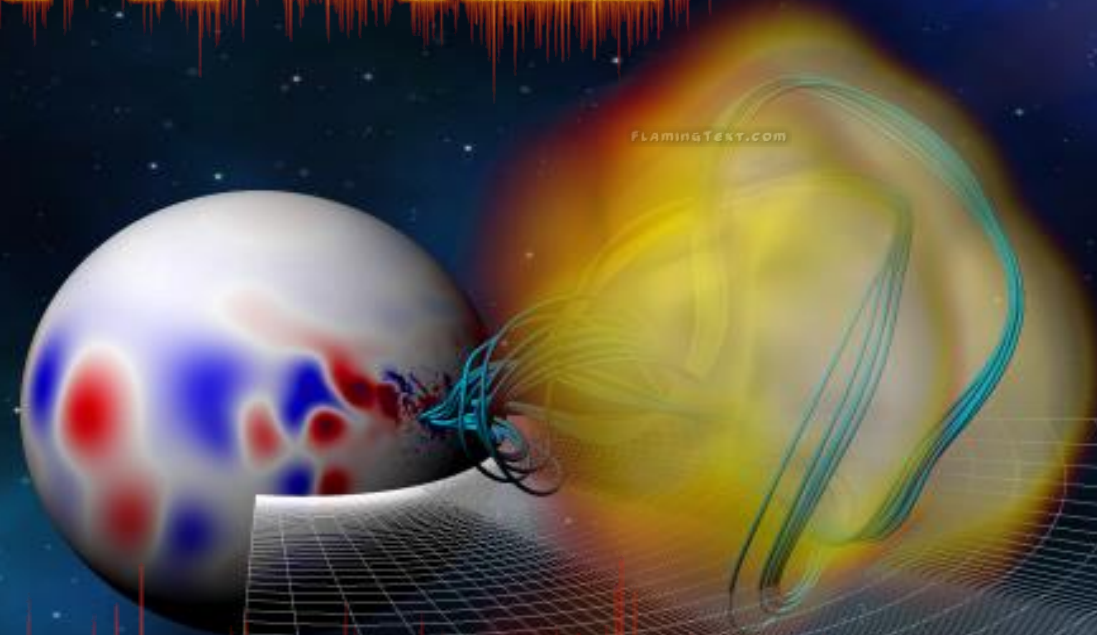
- ❖ OpenACC POT3D will be used in computing over 20,000 PF solutions covering over 4 years of solar data



- ❖ Our global MHD simulation code uses similar PCG solvers, so look for an initial OpenACC implementation at next GTC!

## MAS

MAGNETOHYDRODYNAMIC  
ALGORITHM  
OUTSIDE A SPHERE





# Questions?



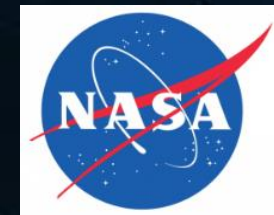
## OpenACC User Group

Directives for Accelerators —> More Science. Less programming

Twitter @OpenACCCorg  
Facebook @OpenACCCorg  
LinkedIn OpenACC Developers

*This work was supported by*

- NSF's Frontiers in Earth System Dynamics program
- NASA's Living with a Star program ([lws.gsfc.nasa.gov](http://lws.gsfc.nasa.gov))
- Air Force Office of Scientific Research



*We gratefully acknowledge NVIDIA Cooperation for donating allocation use of their PSG Cluster for the P100 timings.*



**GPU** TECHNOLOGY  
CONFERENCE



**Predictive Science Inc.**