#### Accelerating the Community Atmosphere Model – Spectral Element Method: The Challenges, the Science, and the Future Direction



#### **Rick Archibald**

Predictive Methods Group Computer Science and Mathematics Division Climate Change Science Institute Oak Ridge National Laboratory

Titan Users and Developers Workshop (East Coast) and Users Meeting Knoxville, TN February 19<sup>th</sup>, 2013



### **Teams and Projects**

Valentine Anantharaj<sup>1,3</sup> Rick Archibald<sup>1,2,3</sup> Chris Baker,<sup>2</sup> Ilene Carpenter<sup>1</sup> NREL Katherine Evans\*23 Jeffrey Larkin<sup>1</sup> Aaron Lott<sup>2</sup> LLNL Matthew Norman<sup>\*1</sup> Paulius Micikevicius<sup>1</sup> Mark Taylor<sup>23</sup> Carrol Woodward<sup>2</sup> LLNL

**ORNL: OLCF** ORNL: PMG/CCSI **ORNL: CEES/CASL** ORNL: CESG/CCSI  $Cray \longrightarrow NVIDIA$ **ORNL: OLCF NVIDIA SANDIA** 

\*Team/Project Leader

- 1. Center for Accelerated Application Readiness CAM-SE
- 2. Multiscale SCIDAC Process Integration
- 3. Ultra High Resolution Global Climate Simulation



### Motivation

OLCF 20

3



National Laboratory

TOP 10 - 06/2012

#### **Motivation: Scientific**

4



Variability of T85 Ensemble over 30 years of Simulation Global Mean Variation 3 Degrees



#### **Motivation: Scientific**

5



Variability of T341 Ensemble over 30 years of Simulation Global Mean Variation 3.5 Degrees



### Comparison





#### **Earth System Model Description**

#### A Climate Model closes the radiative and hydrologic cycles



An Earth System Model closes additional cycles as well



## State of the CESM

8





## What is CAM-SE?

9 OLCF

20

- Climate-scale atmospheric simulation for capability computing
- Comprised of (1) a dynamical core and (2) physics packages



## What is CAM-SE?

2.

- Climate-scale atmospheric simulation for capability computing
- Comprised of (1) a dynamical core and (2) physics packages



dcmip/jablonowski\_cubed\_sphere\_vorticity.png

10 OLCF 20

## **Dynamical Core**

- 1. "Dynamics": wind, energy, & mass
  - "Tracer" Transport: (H<sub>2</sub>O, CO<sub>2</sub>, O<sub>3</sub>, ...) Transport quantities not advanced by the dynamics



# What is CAM-SE?

- Climate-scale atmospheric simulation for capability computing
- Comprised of (1) a dynamical core and (2) physics packages



**Dynamical Core** 

- 1. "Dynamics": wind, energy, & mass
  - "Tracer" Transport: (H<sub>2</sub>O, CO<sub>2</sub>, O<sub>3</sub>, …) Transport quantities not advanced by the dynamics

http://esse.engin.umich.edu/groups/admg/ dcmip/jablonowski\_cubed\_sphere\_vorticity.png

## **Physics Packages**

2.

Resolve anything interesting not included in dynamical core (moist convection, radiation, chemistry, etc)



## **Gridding, Numerics, & Target Run**



- Cubed-Sphere + Spectral Element
- Each cube panel divided into elements





## **Gridding, Numerics, & Target Run**





## **Gridding, Numerics, & Target Run**

 $\square L C F$ 





### Scalability: CAM4 1/4 degree (3 tracers)





- Faster performance at high resolution on parallel computers
- Not cheaper cost in terms of core-hours is the same.

15 OLCF 20

• CAM-Eul the most efficient (cheapest), but with lowest peak SYPD.

'New capabilities and new dynamical cores in CAM'. Mark Taylor, 16<sup>th</sup> Annual CESM Workshop, 2011.



• 16 billion degrees of freedom



- 16 billion degrees of freedom
  - 6 cube panels





- 16 billion degrees of freedom
  - 6 cube panels

18 OLCF 20

- 240 x 240 columns of elements per panel





- 16 billion degrees of freedom
  - 6 cube panels

- 240 x 240 columns of elements per panel
- 4 x 4 basis functions per element





- 16 billion degrees of freedom
  - 6 cube panels
  - 240 x 240 columns of elements per panel
  - 4 x 4 basis functions per element
  - 26 vertical levels





- 16 billion degrees of freedom
  - 6 cube panels
  - 240 x 240 columns of elements per panel
  - 4 x 4 basis functions per element
  - 26 vertical levels

21 OLCF 20

- 110 prognostic variables

 $\rho, \rho u, \rho v, p$ 

```
H_2O , CO_2 , O_3 , CH_4 , \ldots
```



- 16 billion degrees of freedom
  - 6 cube panels
  - 240 x 240 columns of elements per panel
  - 4 x 4 basis functions per element
  - 26 vertical levels
  - 110 prognostic variables
- Scaled to 14,400 XT5 nodes with 60% parallel efficiency





- 16 billion degrees of freedom
  - 6 cube panels
  - 240 x 240 columns of elements per panel
  - 4 x 4 basis functions per element
  - 26 vertical levels
  - 110 prognostic variables
- Scaled to 14,400 XT5 nodes with 60% parallel efficiency
- Must simulate 1-2 thousand times faster than real time
- With 10 second CAM-SE time step, need  $\leq$  10 ms per time step
  - 32-64 columns of elements per node, 5-10 thousand nodes





# CAM-SE Profile (Cray XT5, 14K Nodes)

- Original CAM-SE used 3 tracers (20% difficult to port)
- Mozart chemistry provides 106 tracers (7% difficult to port)
  - Centralizes port to tracers with mostly data-parallel routines





## **The Routines of Tracer Transport**

#### **Three Main Routines**

• Euler\_step (3x): Performs the actual transport

25 OLCF ZO

- Advance\_hypervis: Stabilizes the solution with  $\nabla^4$  diffusion
- Vertical\_Remap: Re-grids the data vertically to a reference grid

#### **Three Helper Routines**

- Edge\_pack: Place all element edges into process-wide buffer
   Edge\_unpack: Sum data at element edges
  - <u>Once</u> per euler\_step
     <u>Twice</u> per advance\_hypervis
- Limiter2d\_zero: Remove negative coefficients & conserve mass
  - <u>Once</u> per euler\_step
     <u>Once</u> per advance\_hypervis



## **Call Tree Structure**

prim\_main

26 OLCF 20

Main driver routines

prim_run_subcycle		Time stepping routines
Г	prim_advec_tracers	Tracer advection routines
	<pre>euler_step edgevpack bndry_exchangev }=pack-unpack edgevunpack laplace_sphere_wk pack-unpack remap_velocity</pre>	This routine called multiple times depending upon the type of timestep (ie, leapfrog, RK, etc).



## **Profile of a Tracer Time Step**

- We wrote a new remapping algorithm to reduce runtime
- Pack/exchange/unpack is ≈ half of euler\_step & advance\_hypervis





### **GPU CAM-SE Goal**









http://regmedia.co.uk/2011/05/22/cray-xk6\_super-blade.jpg



29 **DLCF ZD** 



#### **GPU Kernels**



http://regmedia.co.uk/2011/05/22/cray-xk6\_super-blade.jpg

30 **DLCF** 

20





http://regmedia.co.uk/2011/05/22/cray-xk6\_super-blade.jpg

31  $\square \square \square \square \square \square \square$ 

20

#### GPU Kernels PCI-e D2H





GPU Kernels PCI-e D2H PCI-e H2D



http://regmedia.co.uk/2011/05/22/cray-xk6\_super-blade.jpg

 $32 \square \square \square \square \square \square \square$ 

20





GPU Kernels PCI-e D2H PCI-e H2D MPI



http://regmedia.co.uk/2011/05/22/cray-xk6\_super-blade.jpg

33

OLCF

20





GPU Kernels PCI-e D2H PCI-e H2D MPI



The Free Lunch Is Over: A Fundamental Turn Toward Concurrency in Software









If you were plowing a field, which would you rather use: 2 strong oxen or 1024 chickens?

- Seymour Cray





*If you were plowing a field, which would you rather use: 12,160 bipolar oxen*<sup>1</sup> *or 786,432 chickens*<sup>2</sup>*?* 

- Not Seymour Cray

<sup>1</sup> A bipolar oxen is one AMD Bulldozer and one NVIDIA Fermi.
<sup>2</sup> A chicken is a Blue Gene/Q core.

Herding millions of chickens: Programming models for Blue Gene/Q and beyond

Jeff Hammond

Leadership Computing Facility Argonne National Laboratory




#### **Communication Between Elements**







#### **Communication Between Elements**



Physically occupy the same location, Spectral Element requires them to be equal

Edges are averaged, and the average replaces both edges





#### **Communication Between Elements**



Physically occupy the same location, Spectral Element requires them to be equal

Edges are averaged, and the average replaces both edges

#### Implementation

Edge\_pack: pack <u>all</u> element edges into process-wide buffer. Data sent over MPI are contiguous in buffer.

Bndry\_exchange: Send & receive data at domain decomposition boundaries

Edge\_unpack: Perform a weighted sum for data at <u>all</u> element edges.



- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"





- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI Cycle 1



- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI
- Original pack/exchange/unpack

20

42





- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI
- Original pack/exchange/unpack
  - Pack all edges in a GPU Kernel





- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI
- Original pack/exchange/unpack
  - Pack all edges in a GPU Kernel
  - For each "send cycle"

- Send cycle over PCI-e (D2H)
- MPI\_Isend the cycle





- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI
- Original pack/exchange/unpack
  - Pack all edges in a GPU Kernel
  - For each "send cycle"

- Send cycle over PCI-e (D2H)
- MPI\_Isend the cycle





- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI
- Original pack/exchange/unpack
  - Pack all edges in a GPU Kernel
  - For each "send cycle"

- Send cycle over PCI-e (D2H)
- MPI\_Isend the cycle





- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI
- Original pack/exchange/unpack
  - Pack all edges in a GPU Kernel
  - For each "send cycle"

- Send cycle over PCI-e (D2H)
- MPI\_Isend the cycle





- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI
- Original pack/exchange/unpack
  - Pack all edges in a GPU Kernel
  - For each "send cycle"
    - Send cycle over PCI-e (D2H)
    - MPI\_Isend the cycle
  - For each "receive cycle"

- MPI\_Wait for the data
- Send cycle over PCI-e (H2D)





- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI
- Original pack/exchange/unpack
  - Pack all edges in a GPU Kernel
  - For each "send cycle"
    - Send cycle over PCI-e (D2H)
    - MPI\_Isend the cycle
  - For each "receive cycle"

OLCF 20

49

- MPI\_Wait for the data
- Send cycle over PCI-e (H2D)





- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI
- Original pack/exchange/unpack
  - Pack all edges in a GPU Kernel
  - For each "send cycle"
    - Send cycle over PCI-e (D2H)
    - MPI\_Isend the cycle
  - For each "receive cycle"

- MPI\_Wait for the data
- Send cycle over PCI-e (H2D)





- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI
- Original pack/exchange/unpack
  - Pack all edges in a GPU Kernel
  - For each "send cycle"
    - Send cycle over PCI-e (D2H)
    - MPI\_Isend the cycle
  - For each "receive cycle"

- MPI\_Wait for the data
- Send cycle over PCI-e (H2D)





- Edge\_pack ensures data for MPI is contiguous in buffer
- MPI communication occurs in "cycles"
- A cycle contains a contiguous data region for MPI
- Original pack/exchange/unpack
  - Pack all edges in a GPU Kernel
  - For each "send cycle"
    - Send cycle over PCI-e (D2H)
    - MPI\_Isend the cycle
  - For each "receive cycle"

52 OLCF ZO

- MPI\_Wait for the data
- Send cycle over PCI-e (H2D)
- Unpack all edges in a GPU Kernel





# **Optimizing Pack/Exchange/Unpack**

- For a cycle, PCI-e D2H depends only on packing that cycle
  - <u>Divide</u> edge\_pack into equal-sized cycles
    - 1. Find only the elements directly involved in each separate cycle
    - 2. Evenly divide remaining elements among the cycles
  - Associate each cycle with a unique CUDA stream
  - Launch each pack in its stream

53 **DLCF 2D** 

- After a cycle is packed, call async. PCI-e D2H in its Stream
- Edge\_unpack at MPI boundaries requires all MPI to be finished
- However, internal unpacks can be done directly after packing



- For each cycle
  - Launch edge\_pack kernel for the cycle in a unique stream
  - Call a cudaEventRecord for the stream's packing event







- For each cycle
  - Launch edge\_pack kernel for the cycle in a unique stream
  - Call a cudaEventRecord for the stream's packing event







- For each cycle
  - Launch edge\_pack kernel for the cycle in a unique stream
  - Call a cudaEventRecord for the stream's packing event







- For each cycle
  - Launch edge\_pack kernel for the cycle in a unique stream
  - Call a cudaEventRecord for the stream's packing event







- For each cycle
  - Launch edge\_pack kernel for the cycle in a unique stream
  - Call a cudaEventRecord for the stream's packing event







• Memory coalescing in kernels

59 OLCF 20

Know how threads are accessing GPU DRAM



#### **CPU Code**

```
do ie=1,nelemd
  do q=1,qsize
    do k=1,nlev
    do j=1,np
    do i=1,np
    coefs(1,i,j,k,q,ie) = ...
    coefs(2,i,j,k,q,ie) = ...
    coefs(3,i,j,k,q,ie) = ...
```

#### **GPU Code**

- ie = blockidx%y
- q = blockidx%x
- k = threadidx%z
- j = threadidx%y
- i = threadidx%x
- coefs(1,i,j,k,q,ie) = ... coefs(2,i,j,k,q,ie) = ...
- coefs(3,i,j,k,q,ie) = ...





```
RIDGE
National Laboratory
```



#### **CPU Code**

```
do ie=1,nelemd
  do q=1,qsize
    do k=1,nlev
    do j=1,np
    do i=1,np
    coefs(1,i,j,k,q,ie) = ...
    coefs(2,i,j,k,q,ie) = ...
    coefs(3,i,j,k,q,ie) = ...
```



- Memory coalescing in kernels
  - Know how threads are accessing GPU DRAM
- Use of shared memory

- New vertical remap does nearly all computations on shared memory
- Only accesses to DRAM are at start and end of a very large kernel
- Hence, a 5.5x speed-up over CPU code for vertical remap
- Watch out for banking conflicts



- Memory coalescing in kernels
  - Know how threads are accessing GPU DRAM
- Use of shared memory
  - New vertical remap does nearly all computations on shared memory
  - Only accesses to DRAM are at start and end of a very large kernel
  - Hence, a 5.5x speed-up over CPU code for vertical remap
  - Watch out for banking conflicts
- Overlapping kernels with CPU code execution
  - Any non-dependent sections of code overlapped with kernels
  - Computed & PCI-e copied data for future kernels during prior kernels



- Memory coalescing in kernels
  - Know how threads are accessing GPU DRAM
- Use of shared memory

- New vertical remap does nearly all computations on shared memory
- Only accesses to DRAM are at start and end of a very large kernel
- Hence, a 5.5x speed-up over CPU code for vertical remap
- Watch out for banking conflicts
- Overlapping kernels with CPU code execution
  - Any non-dependent sections of code overlapped with kernels
  - Computed & PCI-e copied data for future kernels during prior kernels
- PCI-e copies: consolidate if small, break up & pipeline if large



#### Speed-Up: Fermi GPU vs 1 Interlagos / Node

- Benchmarks performed on XK6 using end-to-end wall timers
- All PCI-e and MPI communication included





#### **Computational Time-step barrier**

- Without it, climate simulation will hit the time-step barrier
  - Little benefit from higher scales of computing
  - Must choose between higher resolution or tolerable throughput
- Dramatically improve accuracy of highly-coupled models
- Dramatically accelerate spin-up of new simulations
- Potentially revolutionize climate simulation and other application areas with long time integration
- Crucial for climate simulation, promising for other fields





- *Explicit* time integration
  - Directly compute future state using derivatives at current or past times
  - Time step limited by numerical stability
- *Implicit* time integration

- Define future state using derivates at that time, resulting in a system of (nonlinear) equations
- Use (nonlinear) solver to determine future state
- Numerically stable for large time steps
- Time step limited by accuracy
- Performance limited by solver



# **Accelerating Implicit solutions**



Merging these separate efforts into the trunk of CAM-SE will allow accelerated time-stepping methods on hybrid architecture that keeps pass with community algorithmic development



#### **Why Third Party Software?**





#### **Modules available to CAM-SE**

#### derived\_type\_mod

- type(derived\_type) :: object (HOMME variables)
- subroutine initialize(object, HOMME vars)

#### implicit\_mod

- subroutine imp\_solver(HOMME variables)
  - Set c\_ptr -> f\_ptr -> object
- In c++, c\_funptr points to calc\_f in Fortran
- subroutine calc\_f(x(l), f(l), l, c\_ptr ) bind(C,name='user\_sub\_name')
- noxlocainterface.cpp calls calc\_f()
- noxlocainterface.hpp



#### 72 **DLCF 20**
# **Next Steps – Implicit Framework**

73 OLCF 20

```
do n=1, nvar
     do ie=nets,nete
     do k=1, nlev
       do j=1,np
          do i=1.np
          lx = lx+1
          if (n==1) xstate(lx) = elem(ie)%state%v(i,j,1,k,np1)
          if (n==2) xstate(lx) = elem(ie)%state%v(i,j,2,k,np1)
          if (n==3) xstate(lx) = elem(ie)%state%p(i,j,k,np1)
          end do !np
       end do !np
     end do !nlev
    end do !ie
     end do !nvar
call noxsolve(size(xstate), xstate, c_ptr_to_object, c_ptr_to_pre)
     lx = 0
    do n=1, nvar
     do ie=nets,nete
     do k=1, nlev
       do j=1,np
          do i=1, np
          lx = lx+1
          if (n==1) elem(ie)%state%v(i,j,1,k,np1)=xstate(lx)
          if (n==2) elem(ie)%state%v(i,j,2,k,np1)=xstate(lx)
           if (n==3) elem(ie)%state%p(i,j,k,np1)=xstate(lx)
          end do !np
        end do !np
     end do !nlev
    end do !ie
     end do !nvar
subroutine residual(xstate, fx, nelemd, c_ptr_to_object) bind(C,name='calc_f')
 use , intrinsic :: iso c binding
 use kinds, only : real_kind
 use dimensions_mod, only : np, nlev, nvar, nelem
```



# **Next Steps – GPU Framework**

subroutine remap\_velocityQ\_launcher(n0,np1,dt,elem,hvcoord,nets,nete,compute\_diagnostics,rkstage)
use hybrid\_mod, only: hybrid\_t

```
! Setup Thread Partitioning
blockdim = dim3(np,np, EUL BLK )
griddim = dim3(qsize_d, int(ceiling(dble(nete-nets+1)/dble(_EUL_BLK_))),1)
if(compute_mean_flux==1 .and. prescribed_wind==0) use_mean_flux=.true.
do ie=nets,nete
 do i=1, np
    do j=1,np
      do k=1, nlev
          dp_np1_h(i,j,k,ie) = (hvcoord%hyai(k+1) - hvcoord%hyai(k))*hvcoord%ps0 + &
                               (hvcoord%hybi(k+1) - hvcoord%hybi(k))*elem(ie)%state%ps v(i,j,np1)
        endif
      enddo
    enddo
  enddo
enddo
ierr = cudaMemcpy(dp_star_d,dp_star_h,size(dp_star_h),cudaMemcpyHostToDevice)
ierr = cudaMemcpy(dp_np1_d ,dp_np1_h ,size(dp_np1_h ),cudaMemcpyHostToDevice)
ierr = cudaDeviceSetCacheConfig(cudaFuncCachePreferL1)
call remap_velocityQ_kernel_new<<<griddim,blockdim>>>(n0,np1,dt,device_arrays%Qdp,nets,nete,comput
ierr = cudaThreadSynchronize()
call t_stopf('remap_velocityQ')
```

```
end subroutine remap_velocityQ_launcher
```

#### 74 OLCF 20



### **Next Steps – Linked Framework**



Merge framework so that array state is manipulated once and 'calc\_f' executes on the GPU

75 OLCF 20



### **Future Efforts**

# $\frac{\partial \psi}{\partial t} = D\left(\psi\right) + P\left(\psi\right),$

 $\psi$  - prognostic variable,

- $\boldsymbol{D}$  dynamical core component,
- ${\cal P}$  physical parameterization suite.

Change coupling between land and atmosphere from bulk-average of prognostic variables to distributions by adaptive sampling of ensembles.

 $\psi$ 



 $\{\psi_1,\ldots,\psi_n\}$ 

