

Chimera on Titan:

Exploiting new parallelism for supernova simulations

Presented by:

Michael Lentz (U. Tennessee/ORNL)

Jason Messer (OLCF)

Adam Lewkow (UConn)

Anne Parete-Koon (UT -> NCCS)

the Chimera Team:

Guenn (FAU), J. Blondin (NCSU), A. Chertkow (UT), E. Endeve (ORNL), W. R. H. (ORNL/UT), E. Lingerfelt (ORNL), P. Marronetti (FAU), K. Yakunin (FAU) and A.



Diagram of a stellar core showing concentric layers of chemical elements. From the center outwards, the layers are labeled: Fe (Iron), Si (Silicon), O (Oxygen), and He (Helium). The layers are represented by concentric circles in shades of purple, blue, green, and yellow.

Stellar collapse to explosion

Massive stars (~8x solar mass) build inert Fe-core that collapses until center reaches nuclear density (2×10^{14} g/cc)

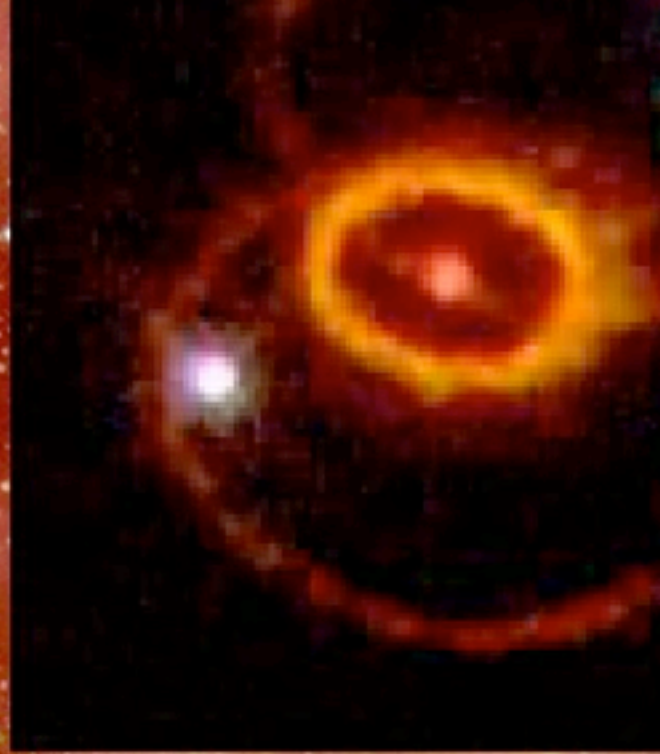
Compressed proto-Neutron Star (pNS) launches shock that “stalls” expanding falling layers, leaving a Standing Accretion Shock (SAS) that must be revived - explode star

Over ~ 1 second, the “mantle” between the pNS and shock is heated by neutrino emission and neutrino-driven convection and builds up enough energy to drive the shock through the rest of the star. Bang!

revived shock drives out through rest of star. Shock expells stellar envelope and triggers nuclear burning, reaching the surface with a flash after a few hours.

exploding star appears in sky as a supernova

Scientific Goals:



Understanding the core-collapse supernova mechanism

Determining the fates of massive stars (variation in explosions; neutron star/pulsar or black hole)

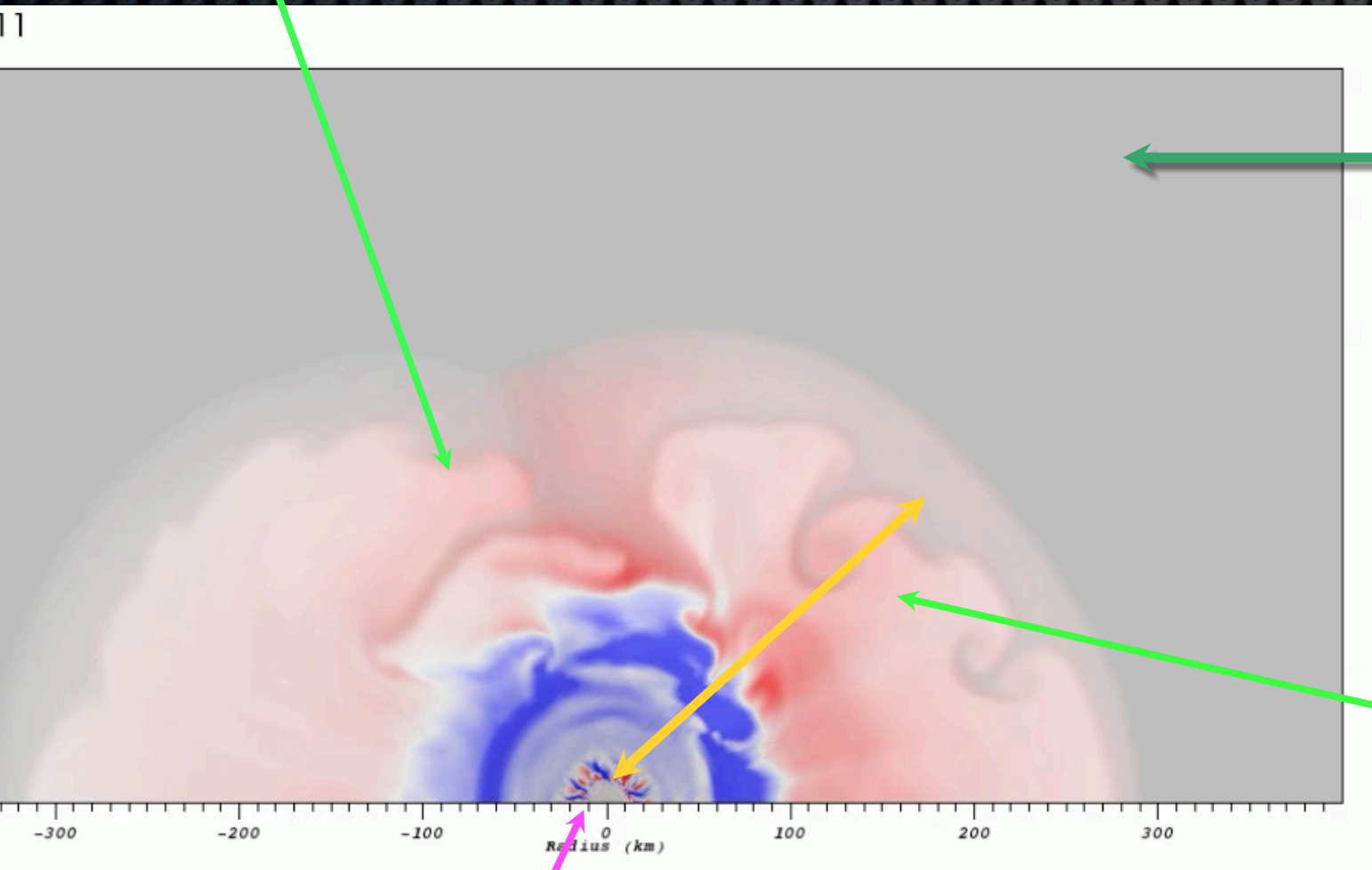
Optical properties of the ejecta (observed Supernovae)

Chemical element production in supernovae

Contribution to the formation of the interstellar medium

Nuclear Challenges

mic “Hot mantle”: neutrino
ed material inside shock in
nuclear equilibrium



$\approx O(10^9)$ sequential
• a few seconds p
 $\approx \sim 2000$ hours ru
• 2-/3-D Hydrody
• Radial neutrino t
• Nuclear Eq. of S
nuclear netw

Rxn. Network: I
material is not in
equilibrium. Use
network to comput
and nucleosynth
ejecta

Neutrino transport
energy from cooling
(blue) to heated
(Red) to energize

Small steps: Dense nuclear equation of state in

Chimera domain decomposition

3D dimensional splitting: (x,y,z) or: (r,θ,ϕ)

alternating sweeps: XYZ-ZYX, XZY-YZX

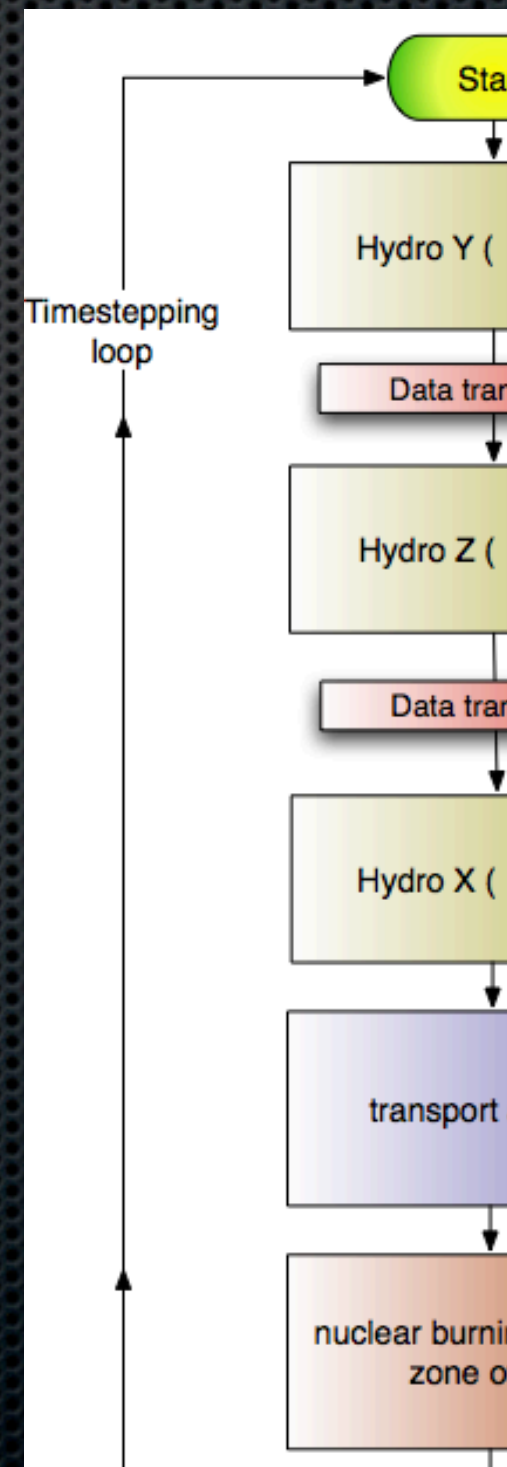
yz-rays: $1 \times 1 \times N$

bundles of several per MPI task

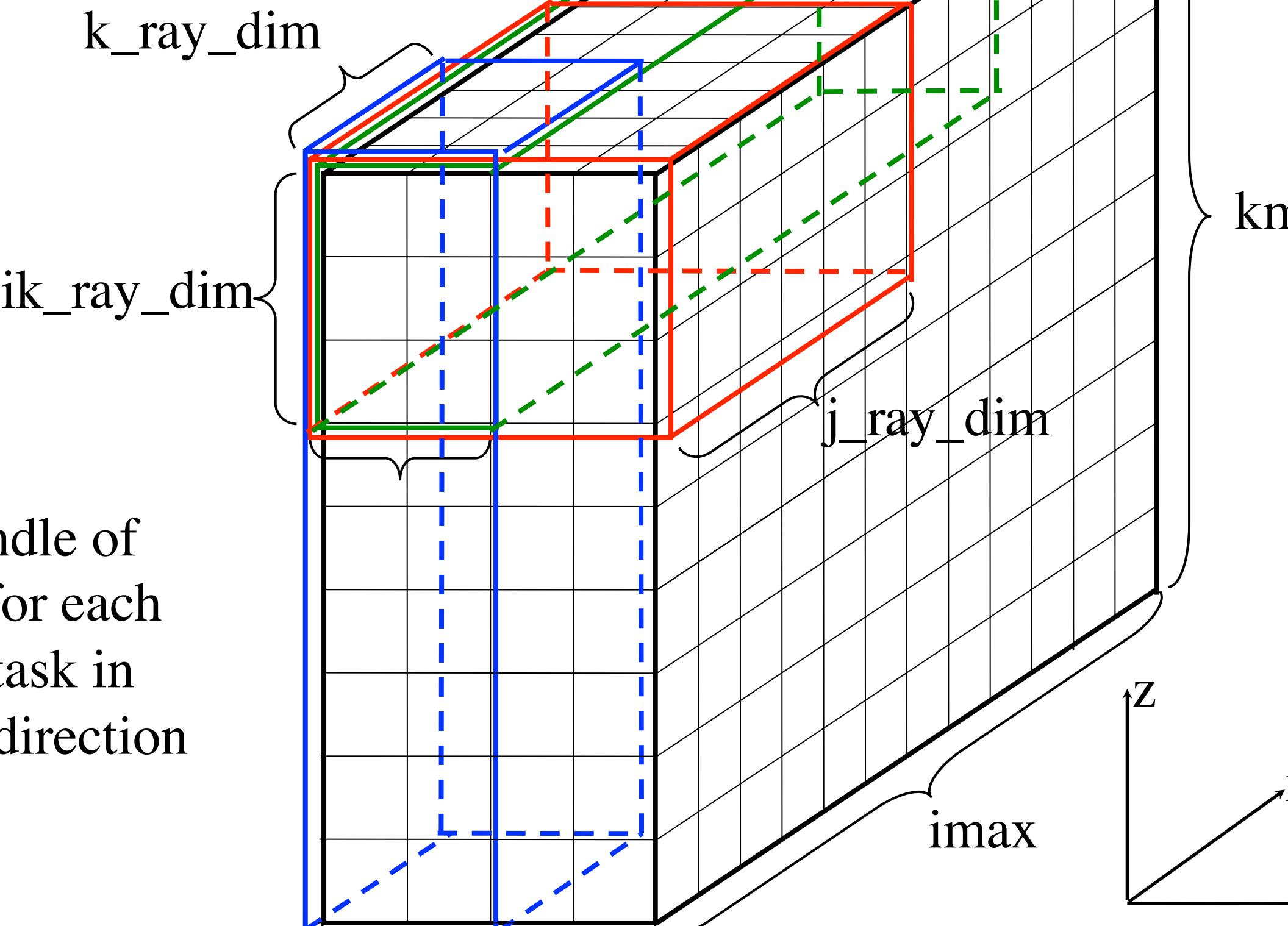
transpose data between sweeps (XY and XZ slabs") using MPI_alltoall on sub-communicators

transport, nuclear network, I/O performed in radial direction. This provides a natural load balancing

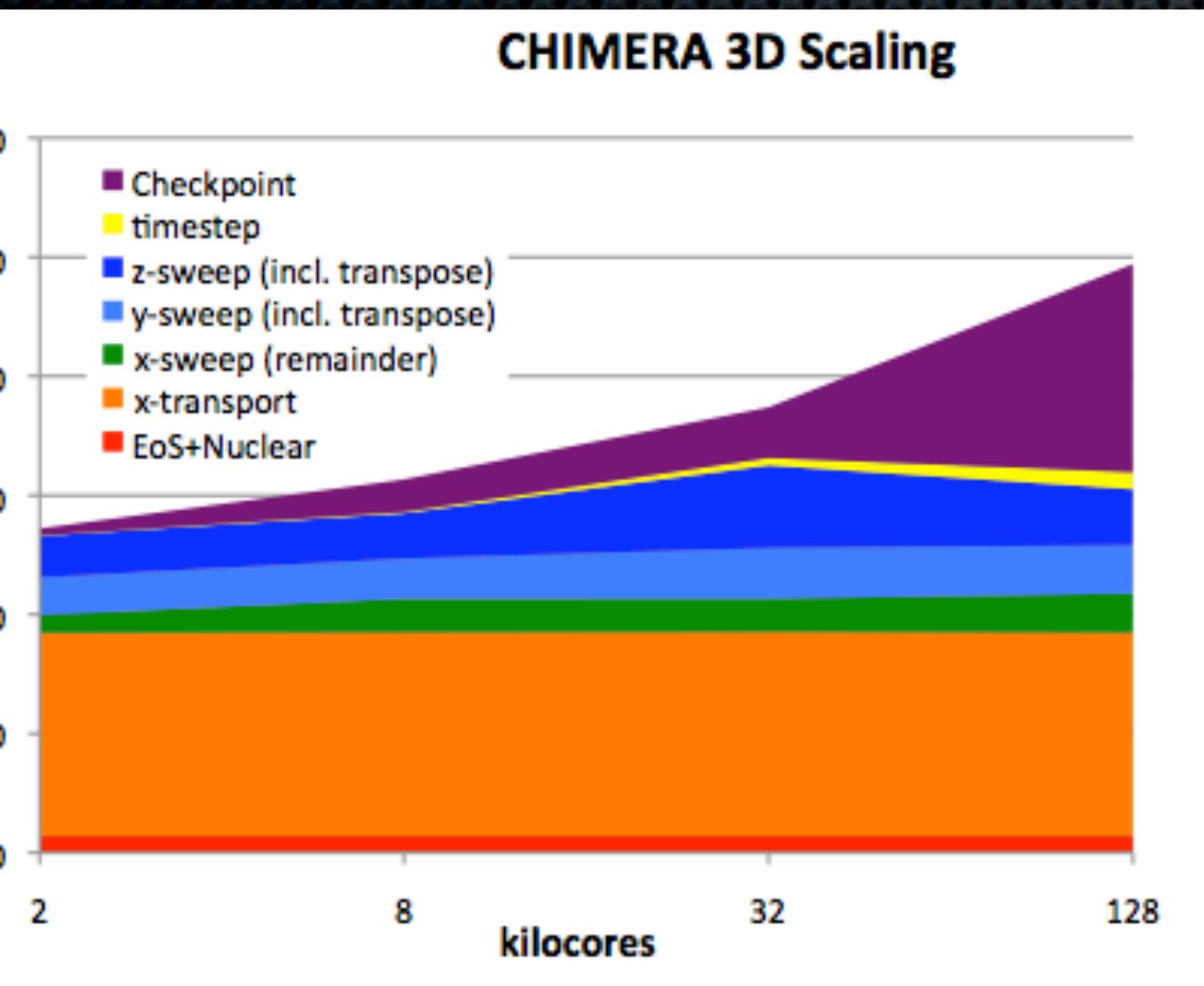
all communication is collective, all-to-all, each MPI



3-D Decomposition in 3-D



3D Car Scaling



with increased
resolution

- Jaguar-XT5 ea
access
- 4 checkpoint
outputs, up to
1800 sec
- rest scale nice

Improved collectives in MPI/O drop HDF5 time to ~100-200

Time step size goes down with each resolution increase

Strong scaling

loops over rays in each dimension

Current configuration

- 1 “x”-ray, or radial ray per MPI task

- typically 512 radial zones per task (Equivalent to 8

- 1 MPI task per core

Reached limit of Strong scaling for pure MPI code

Next: more cores than radial rays + GPUs

How to go faster...

Fewer time steps!

Improve single thread performance (striding)

Add OpenMP threading for multiple cores/task (x-ray)

Add GPU acceleration to computationally intense regions (transport, nuclear network)

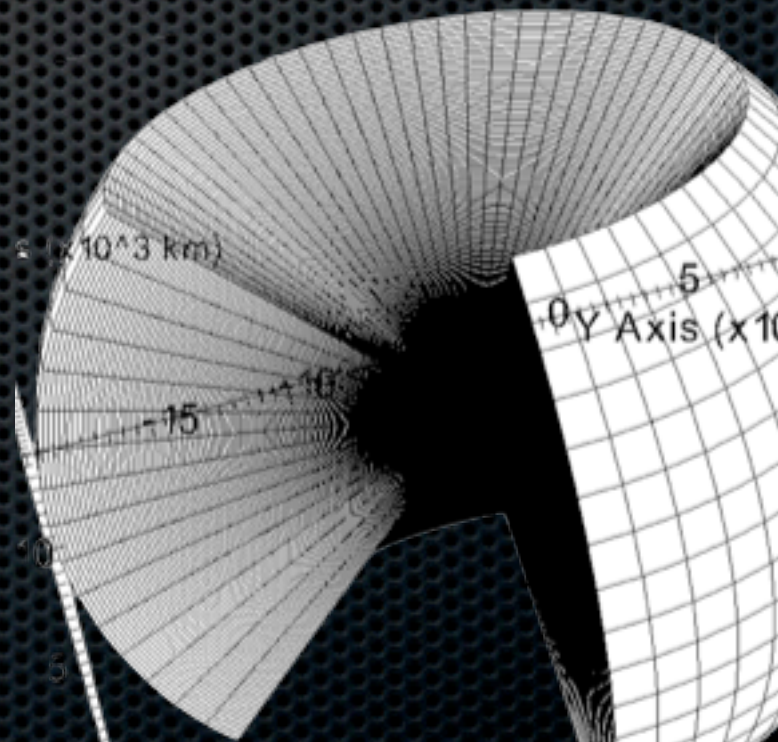
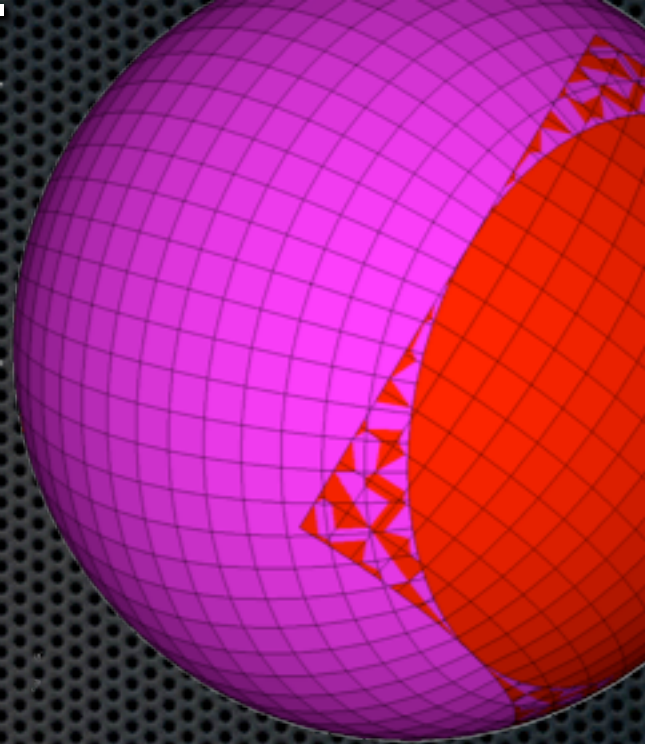
More efficient inter-process communication

One-step control = new grid

Lin/Yang or “Baseball” grid

Two spherical grids with poles
removed and rotated together.

Eliminates short Courant times
at poles and numerical artifacts.



Improving scalar performance

Improving array storage order

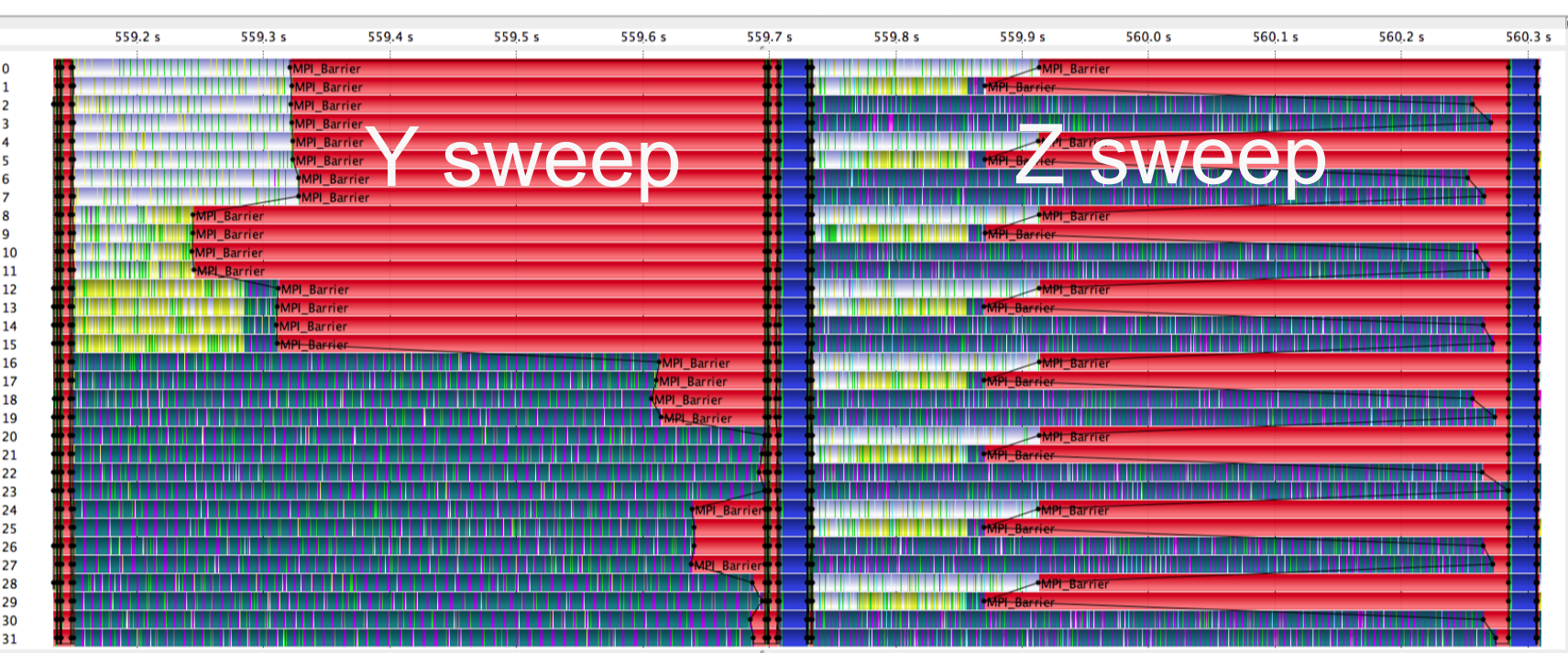
Removing duplicate memory structures

Pushing loops inside subroutines

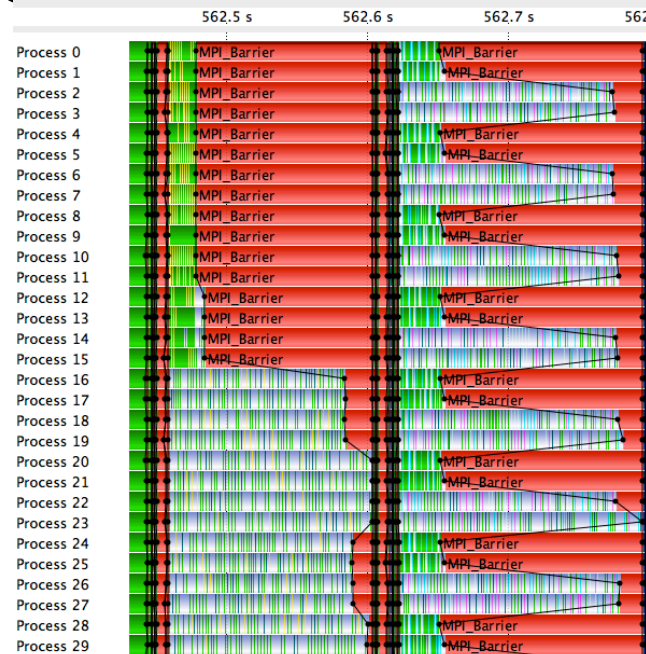
Example: Trimming cost of “lateral” computations

Typical problem size: $64(\theta) \times 128(\phi) \times 512(r)$ zones on
 $64 \times 128 = 8192$ MPI tasks (cores)

y-&z-sweeps each 5% of run time + 15% transpose
time. The “low cost” lateral components are 1/4 of
run seemingly out of proportion to expectations.



$$\Delta t = 0.36 \text{ sec}$$



Fixes:

Refactoring subroutines

NSE: replace direct EoS call with interpolation

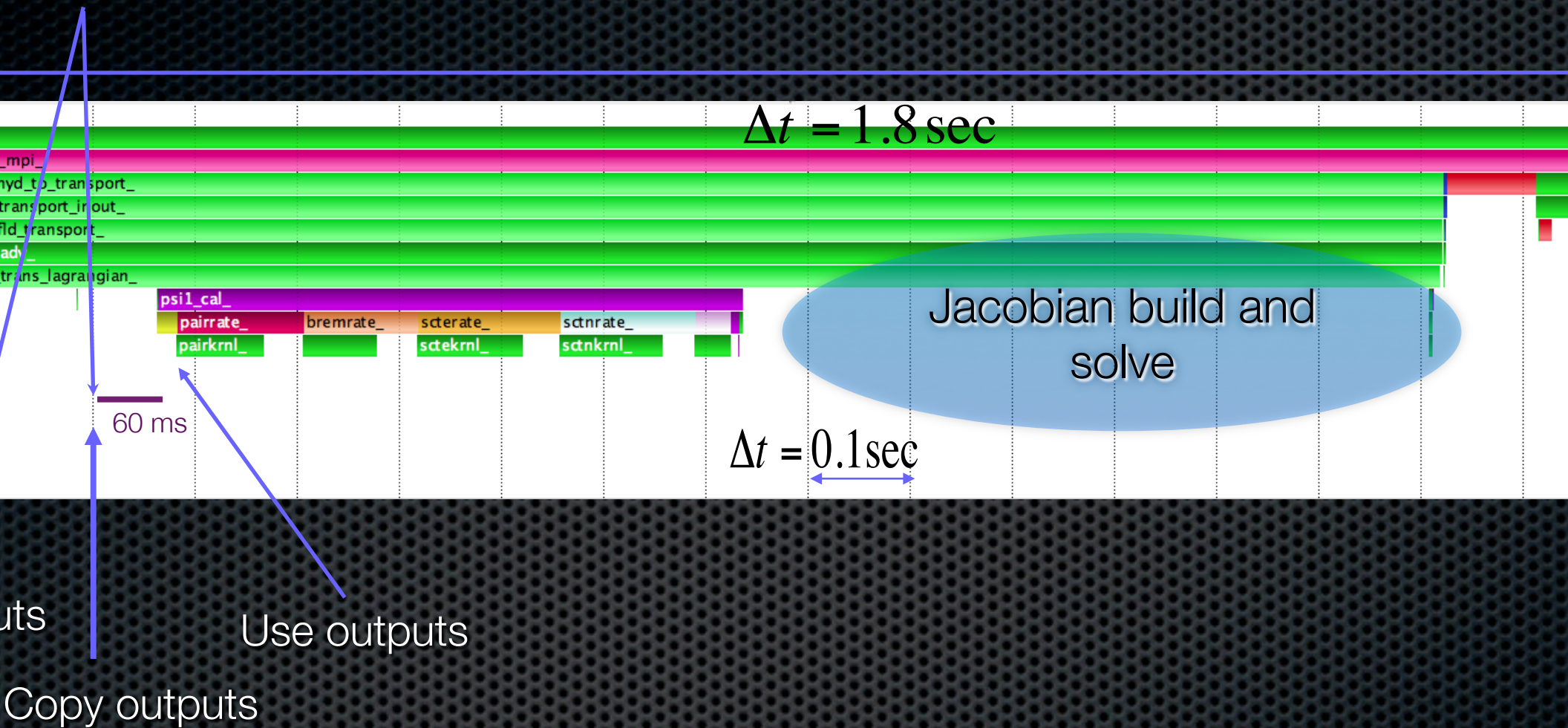
Network region: discovered broadcast to skip calls to unneeded part of

1 XY
before
red: mp

1/3 time
mb. % up
mb. Time down
size problem,
sweeps are much
than transpose

Interpolation kernel

GPU Capacity Interpolation



out most intensive interpolation kernels for 1 or more opacities
inputs in asynchronously when updating table
polation can begin after temperature updated by nuclear burnin
cores do other preparation work while GPU interpolates
builds and solves Jacobian. GPU can start on next ray.

Net: Nuclear Network component

outer regions, nuclei not in equilibrium–Solve nuclear reaction network

Current usage: 14-species “alpha” network ~5% run

Newton-Raphson solution, Jacobian build and solve

scales N^{2-3} with network size

future: 150-species or more network

est. 200-300% time cost of current run w/o accel.
(numerical or physical)

OpenMP threaded, GPU work underway...

ANET. CULA Tools and the GPU

CULA is a GPU accelerated linear algebra library that works on CUDA enabled NVIDIA GPUs.

Host interface is very simple if you have the right setup.

A slow PCI bus offsets the speed of the GPU.

Table Solver Speed for CULA.

matrix order	LAPACK			CULA LAPACK		
	matrix time (s)	matrix time (%)	Step count	matrix time (s)	matrix time (%)	Step count
150	2.82	76	953	25.0	79	953
300	11.7	89	557	25.3	91	557
1072	724	87	664	224	67	668
2184	2837	92	547	500	66	547

OpenMP strategies

n- or low-level threads

each sweep there is a loop over rays, but...

computationally intensive radial orientation want $N_{\text{cores}} > N_{\text{rays}}$

radial rays have lots of work internally to thread by zone, or task

Code is not yet thread-safe

threading individual radial rays

network thread by zones

threading building of opacity table, interpolation

threading of Jacobian build/ solve

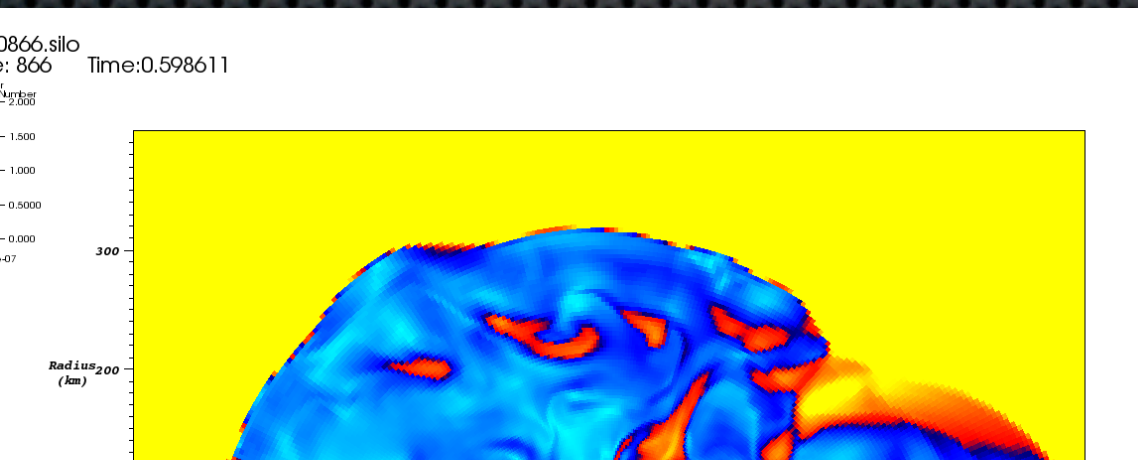
I/O - Analysis - Visualization

Checkpoint files written every 10-20 min, ~10-15 seconds each
for restart, analysis, and visualization

Automated plotting with Bellerophon automation tool

Checkpoint files are ~1% of memory footprint. More efficient to
move analysis offline.

Bellerophon also performs nightly automated code compilation
and testing. Future usage could include offline analysis, archiving
&V, and performance testing.



Chimera on Titan



Current work on new grid and scalar performance should benefit immediately with transition to Bulldozer Gemini configuration of 2012 transition system.

Adding OpenMP + GPU acceleration will allow Chimera to take full advantage of final XK6 system with bigger simulations that finish sooner.

Financial support: DOE, NASA, NSF

Computational support: OLC E NICS, LLN