

## Challenges at Petascale for Pseudo-spectral Methods on Spheres (A Last Hurrah?)

Tom Clune  
Software Integration and Visualization Office  
NASA GSFC



## Nice Mathematical Properties

- “Exponential” accuracy
  - Double the effective resolution compared to FD
  - Commonly used as baseline for comparison
- Fast inversion of elliptic operators
  - Diagonal or nearly diagonal matrices
  - Enables efficient implicit time stepping
- Natural boundary conditions

9/26/11

Pseudospectral - Clune



## Gone The Way of the Dinosaur?



9/26/11

Pseudospectral - Clune



## PS Weaknesses

- Poor quality near discontinuities – e.g. terrain
- Numerically expensive at high resolution –  $O(n^4)$
- Heavy data movement – limited by bandwidth of the interconnect
- Must respect symmetry. E.g. implicit coriolis

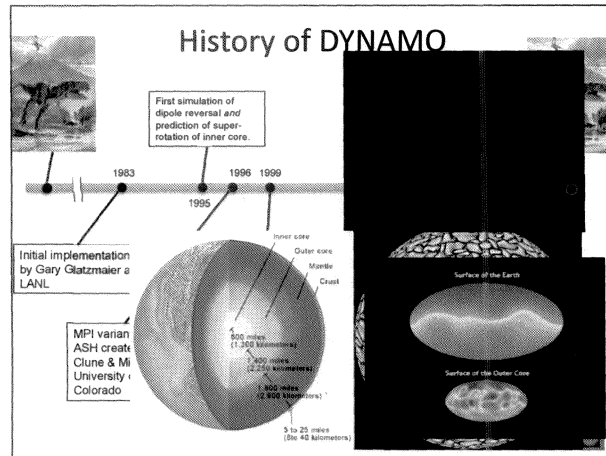


Not immediately  
obvious where the  
tradeoffs fall.

9/26/11

Pseudospectral - Clune





## Pseudospectral Refresher



## DYNAMO Goals

- Existing runs are for  $n \sim 500-1000$ 
  - Run on variety of clusters
  - Performance constrained by 1D decomposition
    - $O(500 \text{ nodes})$
    - Wasting lots of cores to get necessary memory
    - Limited to  $\sim 1 \text{ TF}$
- Would like to achieve  $\sim 10x$  resolution
  - Needs **at least** petascale platform
  - Need consistent 2D (or even 3D) domain decomp
- Stretch – implicit treatment of coriolis

9/26/11

Pseudospectral - CUB

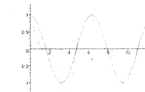


## Pseudospectral Methods

Complementary mathematical representations are used in different phases of calculations.

Spatial derivatives computed in frequency (spectral) domain:

$$\nabla^2 A(k_{ijk}) = |k_{ijk}|^2 A(k_{ijk})$$



Nonlinear products evaluated in spatial domain ("configuration space"):

$$f_{ijk} = b_{ijk} * u_{ijk}$$



Each type of calculation is extremely **simple** *and* **accurate** in the appropriate domain.

9/26/11

Pseudospectral - CUB



Spatial and spectral domains are connected via **transforms**:



Simple and efficient transforms are only possible in relatively small number of simple geometries.

- Periodic box: FFT  $O(n^3 \log n)$
- Separable coordinates:  $O(n^4)$
- Spherical shells are intermediate
- General case is  $O(n^6)$  - uncompetitive

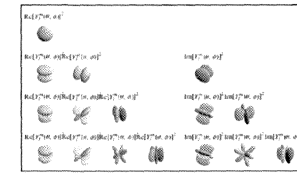
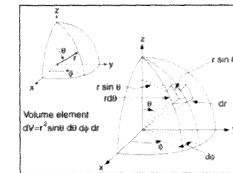
9/16/11

Penciltopical - C. D. Lee



## PS Methods on Spherical Shells

- Spectral expansion based upon spherical harmonics
  - $Y_{lm}(\theta, \phi)$  ( $m = 0, \dots, m_{\max}$ ,  $l = 0, \dots, l_{\max}(m)$ )
- Radial expansion based upon Chebyshev polynomials
  - $T_n(r)$  ( $n = 0, \dots, n_{\max}$ )
- Poisson operators block **diagonal** (nontrivial coupling in radial)



9/16/11

Penciltopical - C. D. Lee



## Interesting Software Aspects

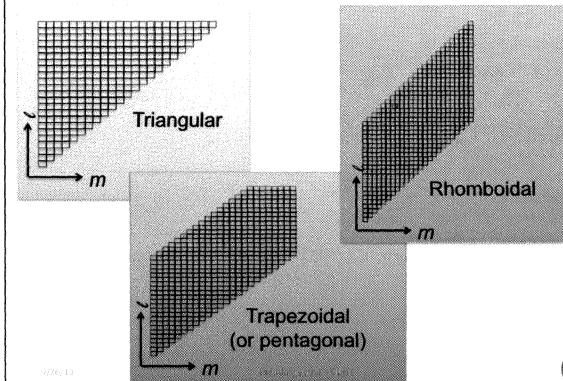
- Elegant program structure
  - Sequence of transforms coupled by memory transposes
  - Software *infrastructure* plays major role
    - Non-rectangular domains
    - Non-trivial domain decomposition
    - Non-obvious data layout
- Unique performance aspects
  - Different scaling properties: transpose vs. halo fill
  - Nothing to optimize!
    - FP workload largely in optimized libraries (FFT, DGEMM, ...)
    - All-to-all is part of HPC benchmarks

9/16/11

Penciltopical - C. D. Lee



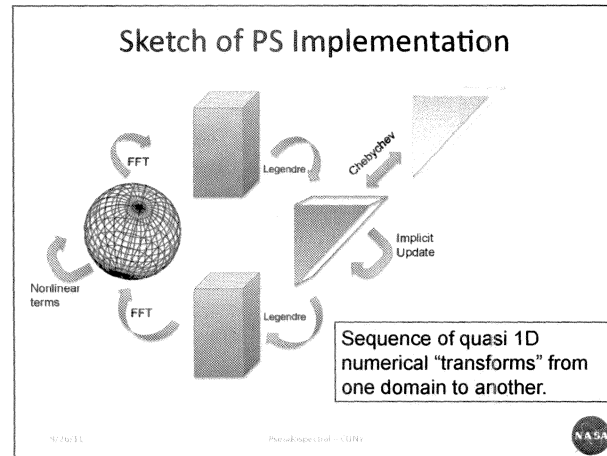
## Spherical Harmonic Truncations



9/16/11

Penciltopical - C. D. Lee





### Domain Decomposition Strategy For a Given Transform

- **Co-locate transform axis**
  - Perform operation "in processor"
  - Leverage serial implementation
  - Often available in optimized library (FFT, DGEMM)
- **Distribute remaining coordinates**
  - No computational dependencies
  - Effective 2D decomposition
- **Balance/Optimize**
  - Memory/performance optimization by grouping along "independent" axis.

10/4/11 Pseudospectral - CURU NASA

### Constraints on Decomposition

- **General anatomy of a transform**
  - Acts on 1 axis of domain  $Q(i,j,k) \rightarrow Q'(m,j,k)$
  - **Independent** of at least one axis
  - Possibly **parameterized** by another axis
- **Example: Legendre Transform**
  - Acts on meridional coordinate/degree
  - Independent of radial coordinate
  - Parameterized by wavenumber coordinate

10/4/11 Pseudospectral - CURU NASA

### Unfortunately ...

- **No** decomposition scheme works for all cases
  - Usually need separate layout for *each* transform!
  - Sometimes can do 2 transforms on one layout if using a 1D decomposition
- Load balance is nontrivial for non-rectangular domains.

10/4/11 Pseudospectral - CURU NASA

### Example: Legendre Transform

- Transforms degree  $\ell$  to/from  $\theta$  (DGEMV)
  - Independent of radial coordinate  $r$
  - Parameterized by azimuthal wavenumber  $m$
- Decomposition constraints:
  - Keep  $\ell$  and  $\theta$  in processor
  - Distribute wavenumbers across processors.
  - Group  $r$  on processor to improve cache reuse (DGEMV  $\rightarrow$  DGEMM)
    - Split blocks of  $r$  across processors to balance scalability against serial performance.

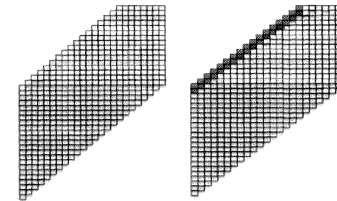
9/26/11

ParasolSpectral - CUIR

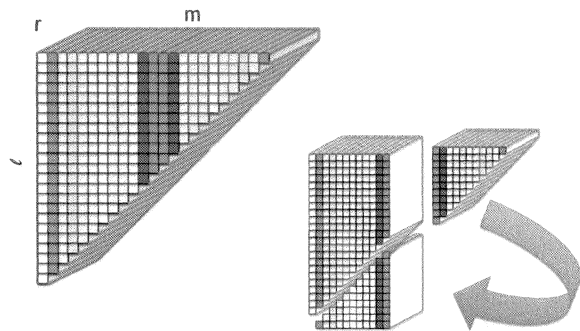


### Implicit Update

- Couples radial coordinate (DGETRS)
  - Mostly independent of azimuthal wavenumber
  - Dependent on degree  $\lambda$ .
- Strategy
  - Group **wavenumber** for BLAS2  $\rightarrow$  BLAS3



### Load Balance and Triangular Domains



9/26/11

ParasolSpectral - CUIR



### Spectral Framework

SPF



## Abstractions

**Axis:** Label(s) and coordinate indices

R
0
1
2
...

LM	L	M
0	0	0
1	0	1
2	1	1
...		

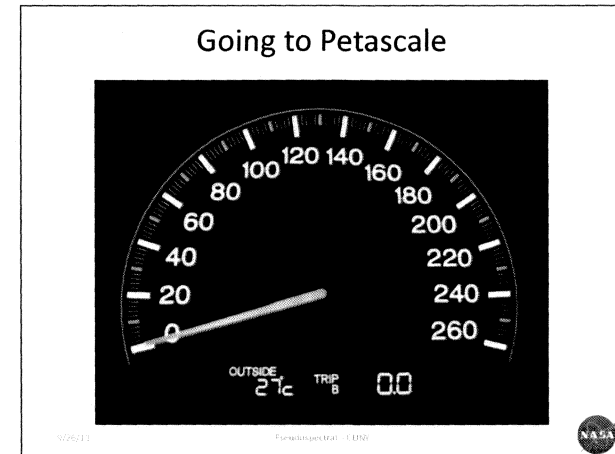
Triangular domain  
treated as 1D

**Domain:** Product of axes

- Note, collection of fields can be an "axis"
- DistributedDomain: subclass with "map"

**Field:** Domain reference with array of values

7/26/11
Pseudospectral - CUBS



## Abstractions (cont'd)

**Transformer:**

- Produces: Decomposition constraints, Cost function
- Field<sub>in</sub> → Field<sub>out</sub>

**Balancer:**

- Ingest: Transformer, Auxiliary Coords, Communicator
- Produces: In/Out Domains, Instantiated Transformer

**Transposer:**

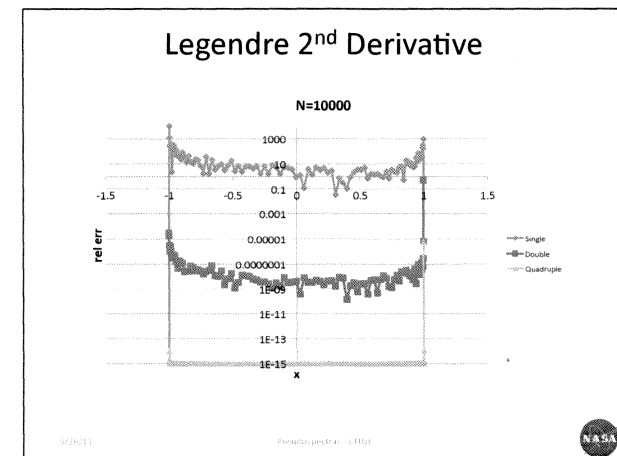
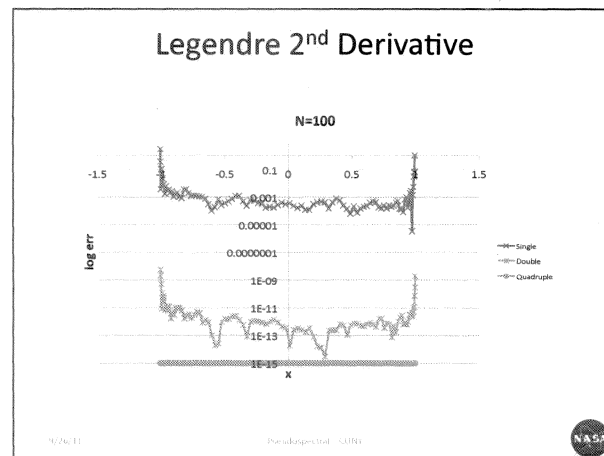
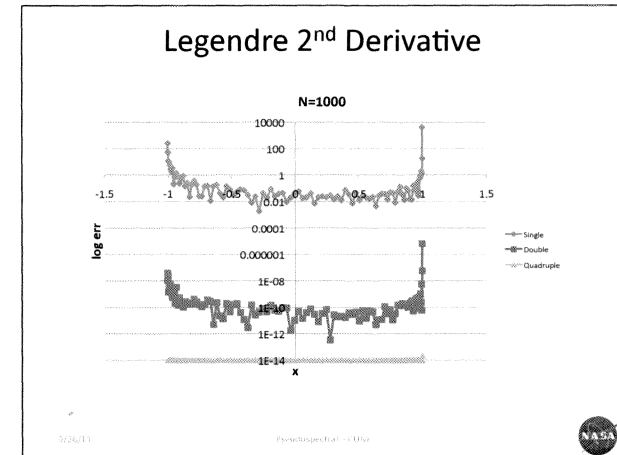
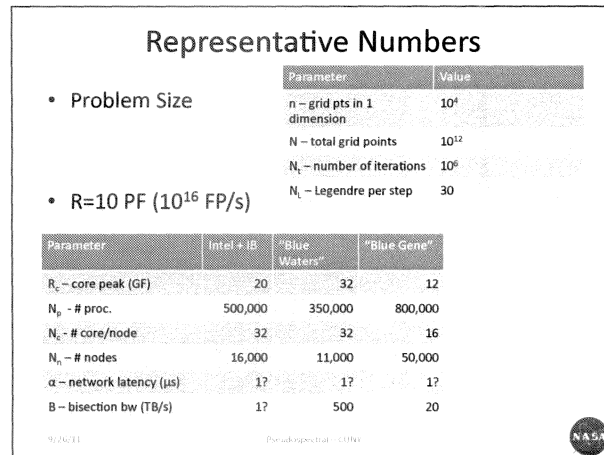
- Initialization: {Domain<sub>A</sub>, Domain<sub>B</sub>}
- Field<sub>in</sub> → Field<sub>out</sub>

7/26/11
Pseudospectral - CUBS

## Specific Challenges

- Accuracy
- Transposes
- Initialization

7/26/11
Pseudospectral - CUBS

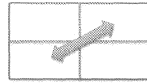


## Back of the Envelope

- Compute time dominated by Legendre DGEMM:

$$T_{comp} \approx \frac{N_t N_i n^4}{R_{eff} p}$$

- Communication dominated by global transposes
  - Each transpose moves 1/4 of data to other half of machine



- Each process sends packet to  $p^{1/D}$  other processes (D=1,2,3)

$$T_{comm} \approx N_t N_i \left( \alpha^{eff} p^{1/D} + \frac{4n^3}{B^{eff}} \right)$$

9/26/11

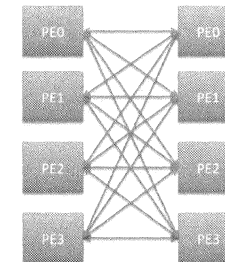
Pseudospectral - CUBS



## Generalized Redistribution

Distribution 1

Distribution 2



Naive Initialization:  
Each processor exchanges all metadata  
with all others and determine overlapping  
elements. Complexity is  $O(N^2)$

Consider a petascale problem:  
 $N=10^{12}$  ( $10^4$  pts each axis)  
10 peta-op platform

At 10 ops/comparison and 10%  
efficiency, initialization requires  $10^{10}$  sec.  
or 300 years!

9/26/11

Pseudospectral - CUBS



## Desirable Characteristics

- Domination by computations:  $n > \frac{4Rp}{B}$ 
  - 8 - 4000
- Bandwidth dominates latency:  $n < \left( \frac{\alpha B p^{1/D}}{4} \right)^{1/3}$ 
  - 250 - 125000
- Complete in  $T_{sim}$  seconds:  $n < \left( \frac{T_{sim} R_{eff} p}{N_t N_i} \right)^{1/4}$ 
  - 4000 for 1 week turnaround
- Efficient level 3 BLAS:  $n \gg p^{1/D}$ 
  - ~10000

9/26/11

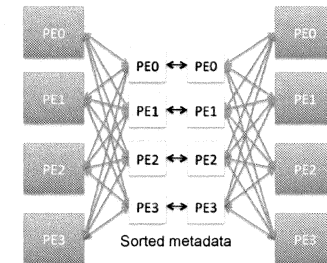
Pseudospectral - CUBS



## Fast Initialization

Distribution 1

Distribution 2



1. Global (sample) sort is used to order metadata.
  - Complexity  $\sim O((N/p) \log_2 p)$
  - Use same pivots for both distributions.
  - Data tagged with original PE
2. Each processor compares sorted metadata to determine overlap.
3. Overlap data is then "unsorted".

Petascale initialization  
should require ~10 minutes.

9/26/11

Pseudospectral - CUBS





## PGAS (CAF)

- Efficient implementations of realistic global transposes are intricate and tedious in MPI.
- PS at petascale requires exploration of a variety of strategies for spreading local and remote communications.
- PGAS allows far simpler implementation and thus rapid exploration of variants.

10/26/11

Pseudo-spectral - CSHR



## Conclusions

- Proper software abstractions should enable rapid-exploration of platform-specific optimizations/tradeoffs.
- Pseudo-spectral methods are marginally viable for at least some classes of petascale problems.
  - A GPU based machine with good bisection would be best.
- *Scalability* at exascale is possible, but the necessary resolution will make algorithm prohibitively expensive.

10/26/11

Pseudo-spectral - CSHR

