

Code optimization and the accumulated impact on scientific throughput of an HPC center

John Dennis

Rory Kelly, Jim Edwards, Brian Dobbins, Chris Kerr,
Youngsung Kim, Raghu Kumar, Sheri Mickelson,
Rich Loft, Mariana Vertenstein, Sean Santos

May 16, 2018

Two related code optimization investments

- Application Scalability And Performance (ASAP) led effort
 - An effort to explore the use of accelerator and other future technology (KNL) on existing weather and climate model codes
- Strategic Parallel Optimization and Optimization Computing (SPOC)
 - An NCAR-wide effort to increase the performance and efficiency of NCAR community does on CESM, WRF, and MPAS

Approach: Incrementally improve existing codebase

Related/Collaborative Activities

- Funding from Intel Parallel Computing Center (IPCC-WACS)
- NESAP (NERSC Exascale Science Application Program)
 - Bi-weekly: NERSC-Cray-NCAR telecon on CESM & HOMME performance (Feb 2015)
- Weekly Intel-TACC-NREL-NERSC-NCAR telecon
 - Concall focused on CESM/HOMME KNC performance

Optimized the following pieces of CESM (>1% reduction in cost)

- Aerosol wet deposition
- Morrison Gettelman micro-physics
- Rapid radiative transport model
- Planetary boundary layer
- Heterogenous freezing in the Modal Aerosol Model
- Random number generator
- Implicit chemical solver
- Spectral element dynamical core (SE-dycore)
- CSLAM advection algorithm in the SE-dycore
- CICE/POP boundary exchange
- Better load balance of CESM

Optimized the following pieces of CESM =>1% reduction in cost

- Aerosol wet deposition
- Morrison Gettelman micro-physics
- **Rapid radiative transport model**
- Planetary boundary layer
- **Heterogenous freezing in the Modal Aerosol Model**
- Random number generator
- Implicit chemical solver
- **Spectral element dynamical core (SE-dycore)**
- CSLAM advection algorithm in the SE-dycore
- **CICE/POP boundary exchange**
- Better load balance of CESM

Four sub-projects

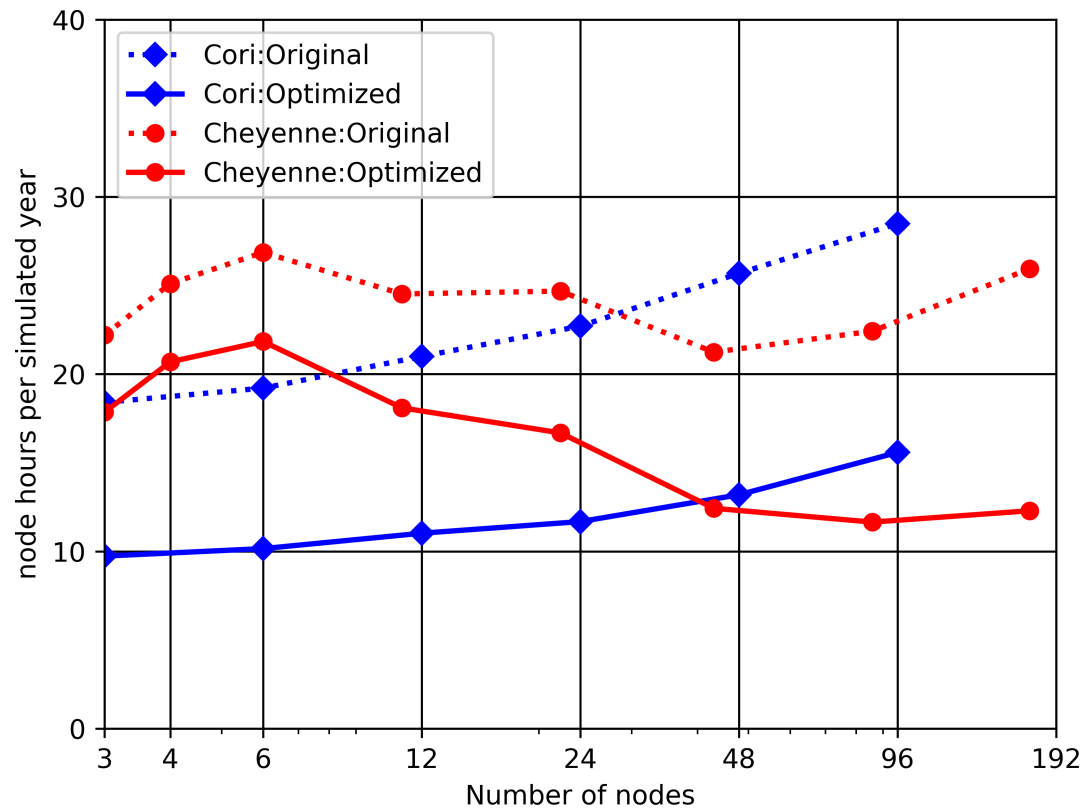
- Rapid radiative transport model
 - 6 FTE-month
 - 1% overall impact
 - ROI: 0.5 %/FTE-year
- Spectral element dynamical core (SE-dycore):
 - 18 FTE-month
 - 10% overall impact (high-resolution)
 - ROI: 7%/FTE-year
- Heterogenous freezing in the Modal Aerosol Model
 - 4 FTE-hour
 - 1% overall impact
 - ROI: 50%/FTE-year
- CICE/POP boundary exchange:
 - 1 FTE-month
 - 20% overall impact (ultra-high-resolution)
 - ROI: 120%/FTE-year

Spectral element dynamical core

Optimization phases of the SE-dycore

1. Threading memory copy in boundary exchange [Jamroz]
2. Restructure data-structures for vectorization [Vadlamani & Dennis]
3. Rewrite message passing library/ specialized comm ops [Dennis]
4. Rearrange calculations in euler_step for cache reuse [Dennis]
5. Reduced # of divides [Dennis]
6. Restructured/alignment for better vectorization [Kerr]
7. Rewrote and optimized limiter [Demeshko & Kerr]
8. Redesign of OpenMP threading [Kerr & Dennis]
9. Flexible MPI message passing back-ends [Dennis]
 1. MPI_Put/Get (MPI3)
 2. MPI neighborhood collectives (MPI3)
10. Replaced all functions with subroutines [Kerr & Dennis]
11. Custom OpenMP barrier [Dobbins]

Simulation cost for HOMME on Xeon and Xeon Phi @ 100 km



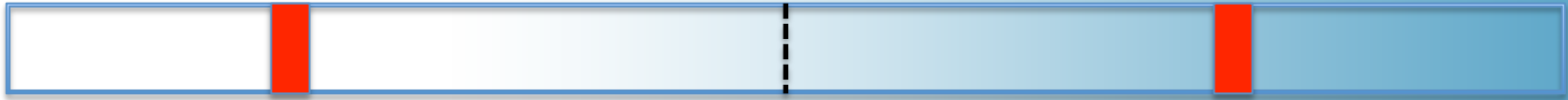
Sea-ice (CICE) and Ocean (POP) model boundary exchange optimizations

CICE/POP boundary exchange optimizations



Tripole message buffer (0.1° grid)

3600 words



Reduce ~20-30x



Specific impact of CICE/POP boundary exchange

- CICE @ 0.1 degree on ~20K cores:
 - 56% reduction in CICE cost
 - 10% reduction in ultra-high-resolution CESM
- POP @ 0.1 degree on ~7K cores
 - ~30% reduction POP cost
 - 10% reduction in ultra-high-resolution CESM

Estimate Overall impact

- What impact did this investment have on scientific throughput?
- Challenging because CESM code base has changed from both a scientific and code optimization perspective
- Approach
 - Detailed measurements of execution time of CESM2 on Cheyenne
 - Adjust execution time of segments of the code based historical timing information
 - I.e. reduced execution time of short-wave length radiation by 33% on Yellowstone....

What was achieved ?

Higher efficiency = more science

CESM configuration	Atmos Resolution (km)	Ocean Resolution (km)	Speedup
Low-res IPCC	100	100	13%
Low-res WACCM chemistry	100	100	20%
High-res IPCC	25	100	25%
Ultra-high Ocean eddy permitting	25	10	45%

What was achieved ?

Higher efficiency = more science

CESM configuration	Atmos Resolution (km)	Ocean Resolution (km)	Speedup
Low-res IPCC	100	100	13%
Low-res WACCM chemistry	100	100	20%
High-res IPCC	25	100	25%
Ultra-high Ocean eddy permitting	25	10	45%

- CCSM/CESM consumes 57% of all Cheyenne
- The TCO to provision 1% more climate computing is \$285K over the 4-year life of Cheyenne
- Investment has enabled between **\$3.7M and \$12.6M of additional science throughput** on Cheyenne. Since CESM is a community model the valuation is larger.

What was achieved ?

Simulation rate & cost @ 100 km

NCAR System	Intel Xeon Processor	CESM Version	Capability (sim yr/day)	Cost (node-hrs per sim yr)
Cheyenne	18c Broadwell	CESM2	30	97
Yellowstone	8c Sandybridge	CESM2	19.6	323
Yellowstone	8c Sandybridge	CESM1	10.6	95

What was achieved ?

Simulation rate & cost @ 100 km

NCAR System	Intel Xeon Processor	CESM Version	Capability (sim yr/day)	Cost (node-hrs per sim yr)
Cheyenne	18c Broadwell	CESM2	30	97
Yellowstone	8c Sandybridge	CESM2	19.6	323
Yellowstone	8c Sandybridge	CESM1	10.6	95

- CESM2 on Cheyenne can deliver **2.8x the capability** compared to CESM1 on Yellowstone

Challenges

- New code being added **30x** quicker than it can be optimized
 - CESM1: 1.3M lines of code
 - CESM2: 1.6M lines of code
 - ~10K lines of code has been optimized
- Scientific evolution of codebase is unpredictable
 - SE dynamical core
 - Cloud Layer Unified by Bi-normals (CLUBB)
- 250x difference in ROI for optimization effort → choose wisely
- How to choose wisely?
 - What is expensive? [Trivial]
 - What is inefficient? [See POP CoE]
 - What will have longevity in code base? [Domain scientist input]
- Most optimization efforts performed twice



Conclusions/Future work

- Concerted/sustained effort reduced cost of CESM on Xeon and Xeon Phi
 - 1.8x speedup on current platform
- Investment in code optimization increased scientific throughput of Cheyenne by \$3.7M to \$12.6M
- Huge range of ROI
- Need to improve process efficiency
 - Optimizing code after insertion not a viable long term approach
 - Teaching code optimization: RRTMGP (Next generation radiation model) & Robert Pincus
- Incremental approach does not address transformative architecture changes

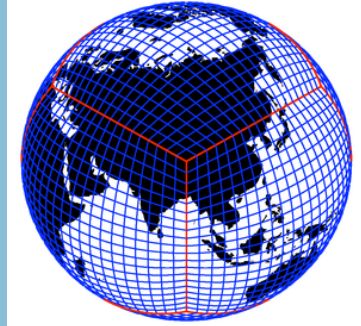
Questions?

dennis@ucar.edu

Aquaplanet @ 100 km

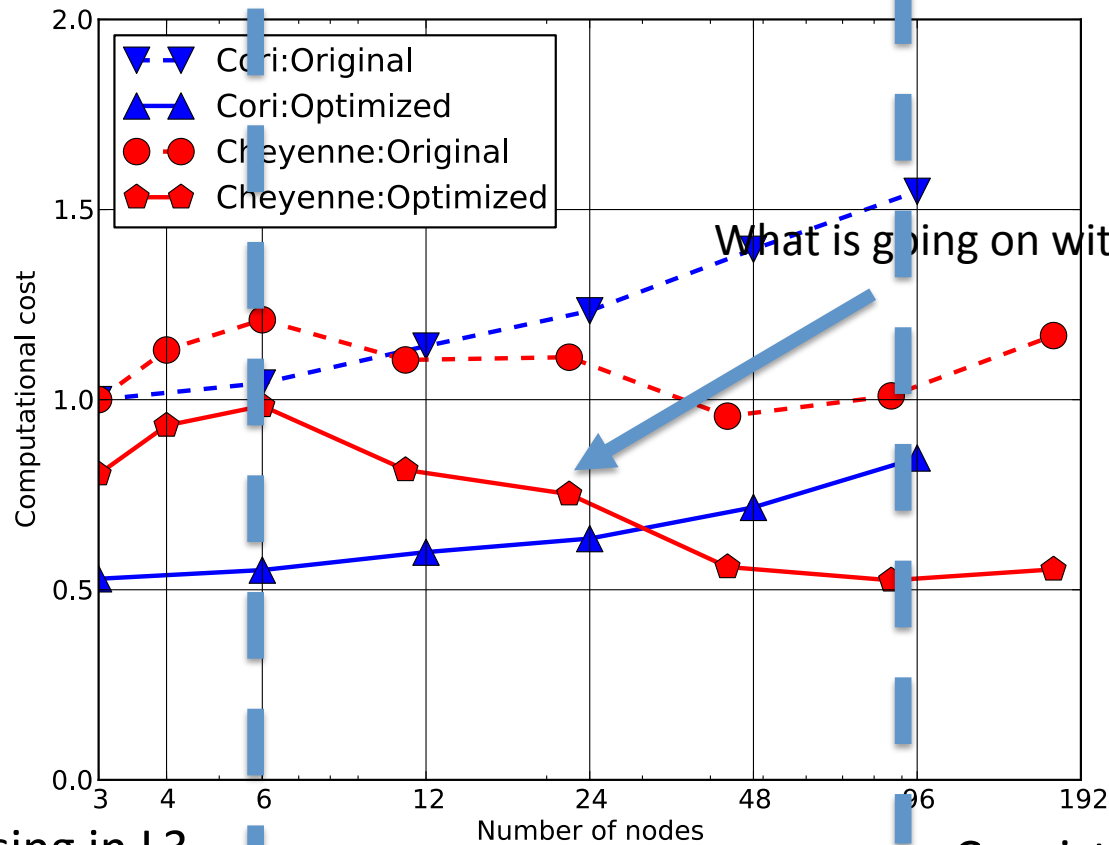
Dynamical core + advection algorithm	MPI rank x OpenMP threads	Capability (sim yr/day)	Cost (core-hrs per sim yr)	Increase/decrease relative to CAM-fv @ 1152x3
CAM-fv	1152x3	40.4	57	0.0%
CAM-SE/ eulerian	2700x1	33.6	54	-6.0%
	5400x1	58.8	61	7.3%
CAM-SE/ CSLAM	2700x1	29.8	60	5.8%

Motivation of HOMME optimization effort



- Atmosphere dynamical core (HOMME)
 - CAM: 35% of time (vert levels=32, # of tracers=25)
- Much easier to optimize than physics 😊
- Benchmark code
 - CORAL (CAM-SE)
 - NSF625
- Useful for evaluating full system performance

Simulation cost for HOMME on Xeon and Xeon Phi @ 100 km



What is going on with the cost curve?

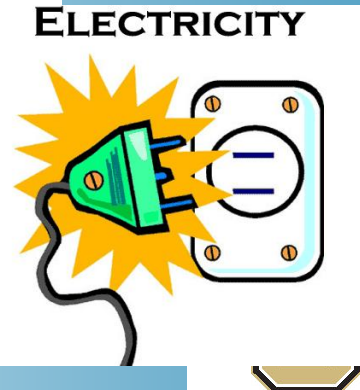
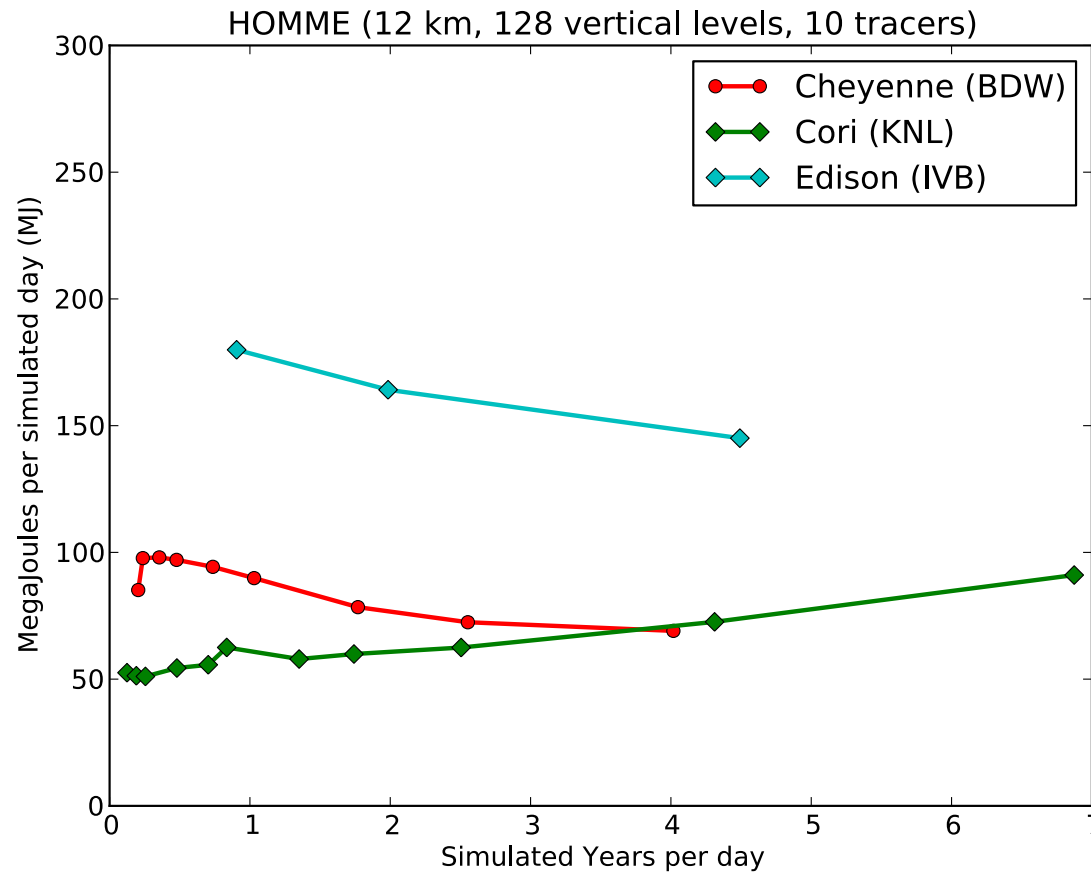
Missing in L3

Consistently hitting in L3

Group/Team

- Rich Loft, Division Director (NCAR)
- John Dennis, Scientist (NCAR)
- Chris Kerr, Software Engineer, contractor
- Youngsung Kim, Software Engineer (NCAR) /Graduate Student (CU)
- Brian Dobbins, Software Engineer (NCAR)
- Raghu Raj Prasanna Kumar, Associate Scientist (NCAR)
- Sheri Mickelson, Software Engineer (NCAR) / Graduate Student (CSU)
- Ravi Nanjundiah, Professor (IISc)

Energy usage for HOMME (NGGPS-like) on Xeon and Xeon Phi @ 12 km



Cost of Aquaplanet @ 100 km for several different dynamical cores

Simulation rate for HOMME on Xeon and KNL

