

IFS migrates from IBM to Cray 'CPU, Comms and I/O'

Deborah Salmond & Peter Towers
Research Department Computing Department

Thanks to Sylvie Malardel, Philippe Marguinaud, Alan Geer
& John Hague and many others ...

CCA and CCB - Cray XC30

2*19 Cabinets
with 2*3465
Compute Nodes

Nov 2013 -



C2A and C2B - IBM Power7
2*11 Frames with 2*768 Compute Nodes
June 2011 - Sept 2014



Comparison of IBM P7 and Cray XC30

	IBM	Cray
Processor	IBM Power7	Intel IvyBridge
Clock Speed	3.8 GHz	2.7 GHz
Switch	IBM HFI	Cray Aries
Nodes	768 *2	3465 *2
Cores per Node	32	24
Cores	24576 *2	83160 *2
Peak (Tflops)	754 *2	1796 *2
Memory per Node (GB)	64	64
Compiler	IBM XLF	Cray CCE
OS	AIX	CLE
Parallel File system	gpfs	lustre

Statistics for IFS 10-day Forecast

- Spectral truncation TL1279
- Horizontal Grid-Points = 2×10^6
- 16km grid-spacing
- Vertical levels = 137
- Timestep = 600 Seconds

- Floating-point ops = 12×10^{15}
- MPI Communications = 150 TB
- SL Halo-width=18

- 1 Million lines of Fortran + C
- Shared with Météo-France
- Bit-reproducible
- Vector length = NPROMA
- Hybrid MPI and OpenMP

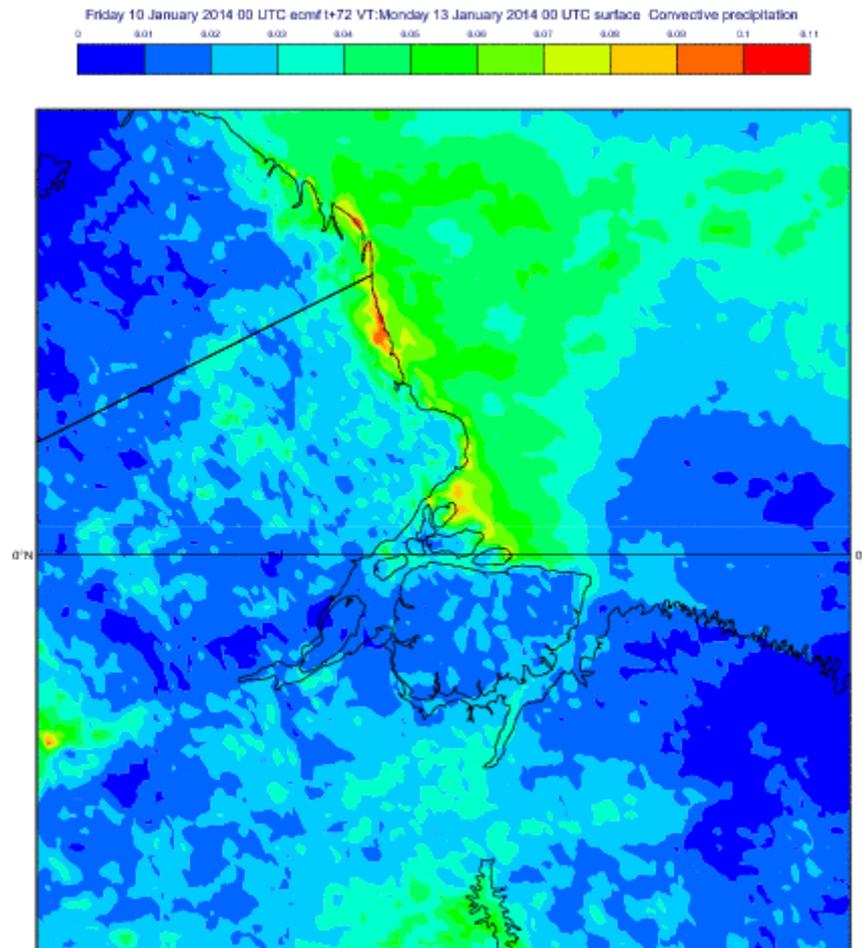
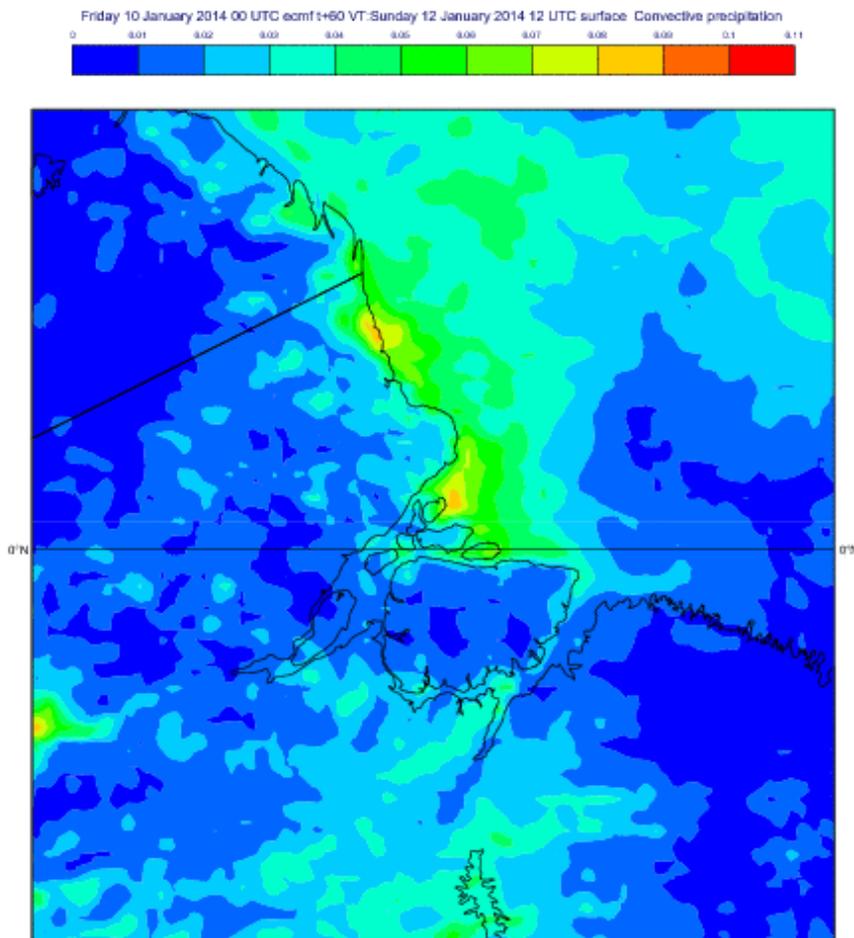
- Elapsed time ~3000 secs
- 4 Tflops

- IBM
60 Nodes = 1920 Cores (+SMT)
480 MPI x 8 OMP
6.8% Peak

- Cray
100 Nodes = 2400 Cores (+HT)
400 MPI x 12 OMP
7.7% Peak

* For RD config without full operational I/O

TL1279 (16km) and TC1279 (8km) Convective precipitation (accumulated over first 3 days of FC from 10 Jan 2014) in the Amazon delta



Thanks to Sylvie Malardel

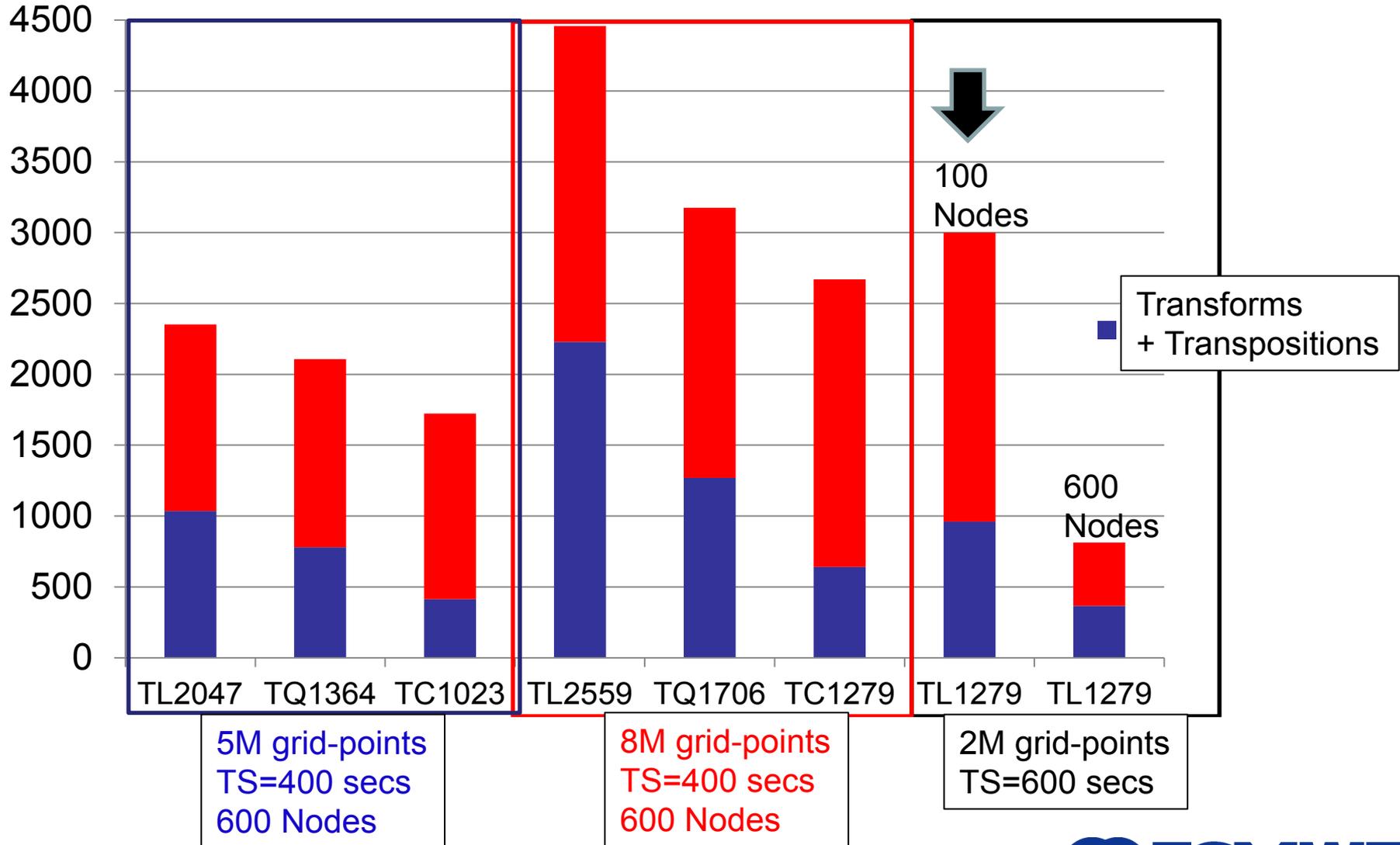
Choices for Higher resolution upgrade in 2015

Different Wavenumbers - Grid-point matches

Horizontal grid-points	Linear	Quadratic	Cubic
2140702 (16km)	TL1279		
5447118 (10km)	TL2047	TQ1364	TC1023
8505906 (8km)	TL2559	TQ1706	TC1279

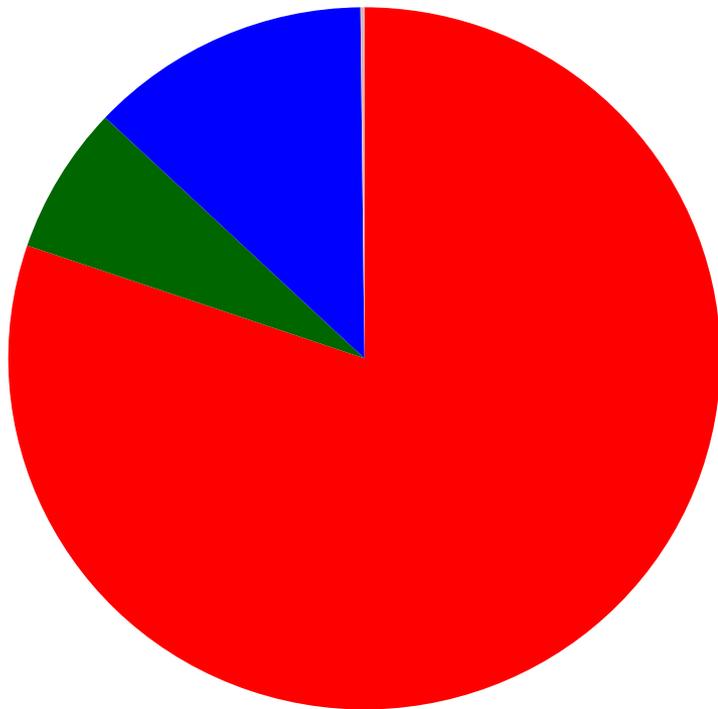
Costs of Different Resolutions

Seconds

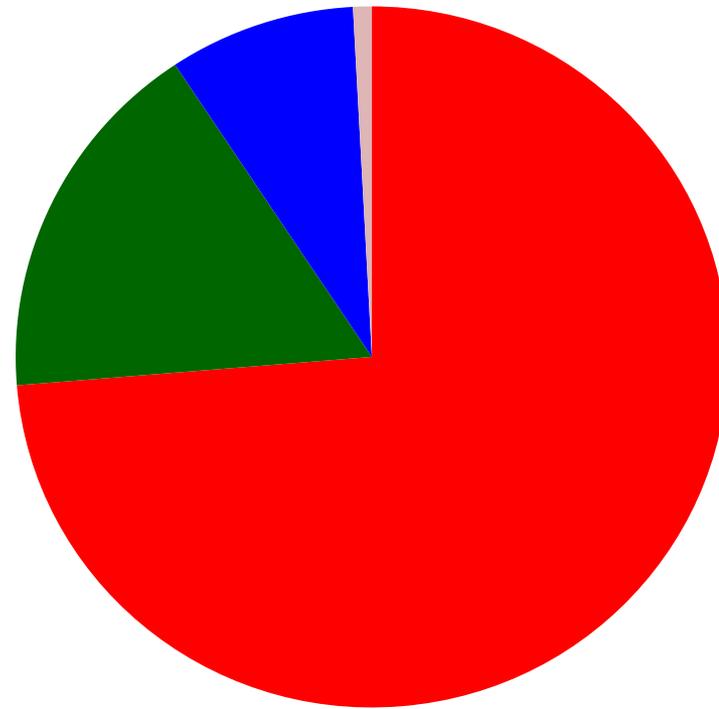


Comparison between IBM and Cray for IFS T1279 10-day Forecast

IBM - 60 Nodes



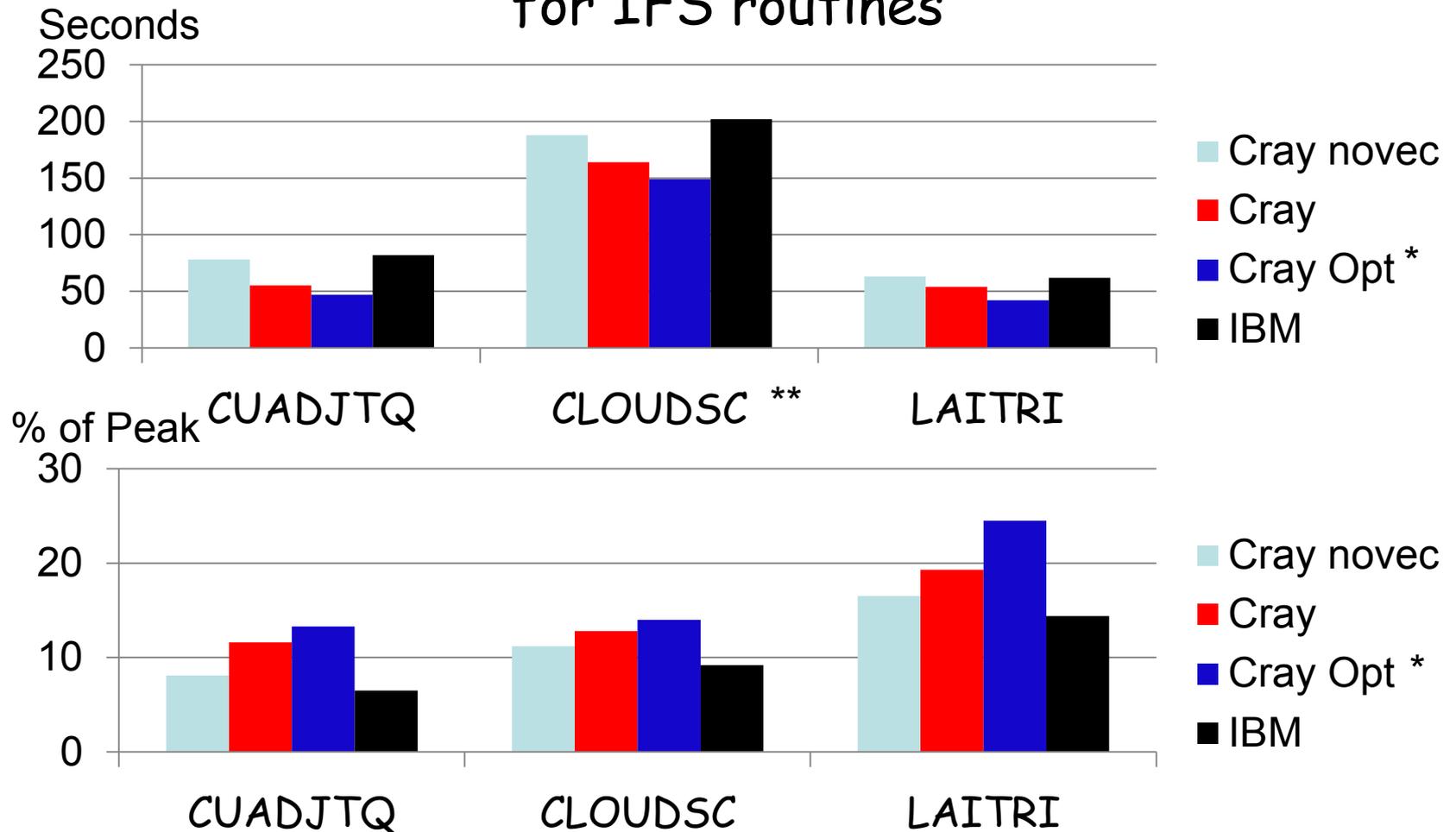
CRAY - 100 Nodes



- CPU
- Comms
- Barrier
- Serial

CCE compiler able to vectorise IFS source → Compute relatively faster
Aries has fewer hub chips per node than HFI → Comms relatively slower
Cray has light-weight kernel on application nodes → Less jitter

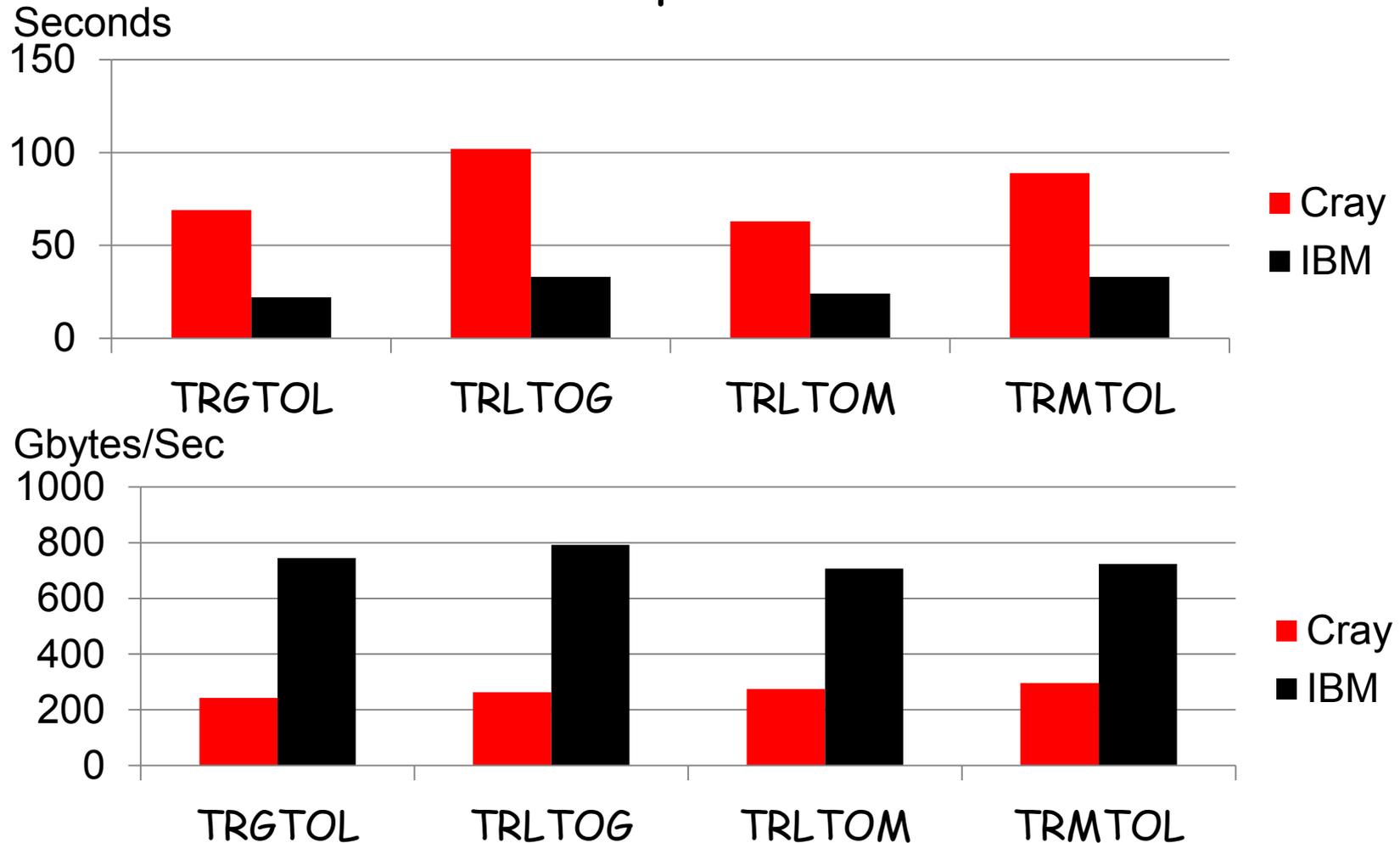
CPU Comparison: Cray (100 Nodes) and IBM (60 Nodes) for IFS routines



* Opt by improving vectorisation see John Hague's talk

** For more on CLOUDSC see Sami Saarinen's talk

Comms Comparison: Cray (100 Nodes) and IBM (60 Nodes) for IFS transposition routines



Comms developments: Overlap MPI with CPU

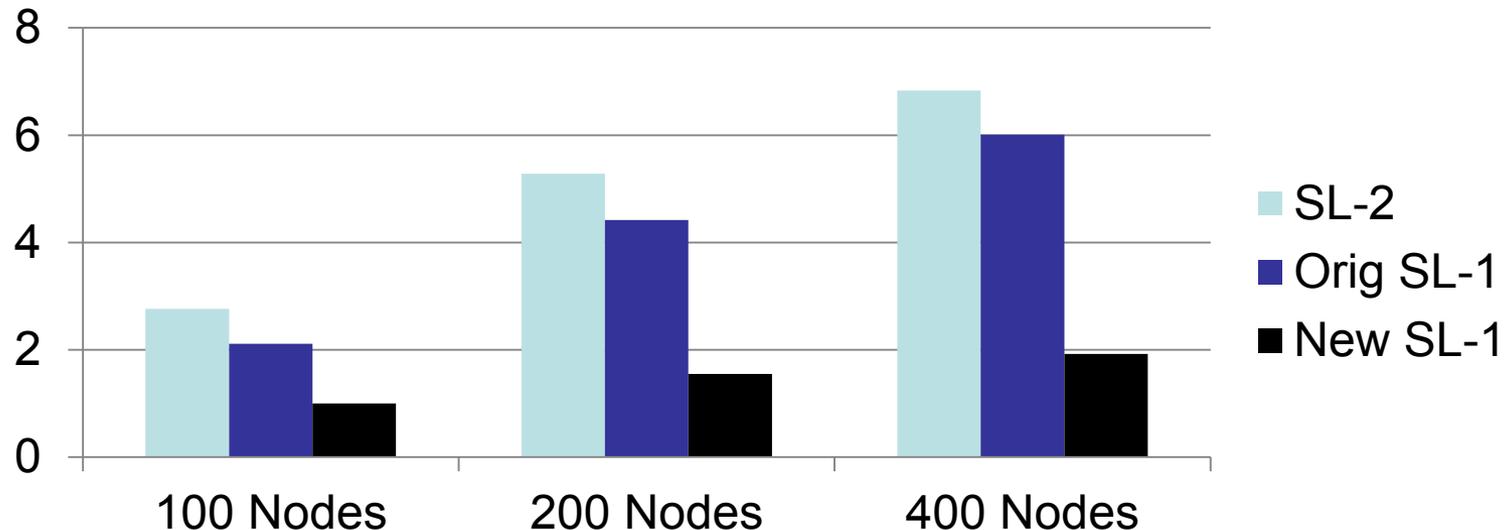
- Investigation of overlapping communications in direct and inverse Legendre transform by Philippe Marguinaud (Météo-France)
- MPI_alltoallv → MPI_Ialltoallv in transpositions
- Gives improvements for current and future resolutions

	T2047 (1024 nodes)	T1198 (64 nodes)
LT-INV: Orig	180 secs	121 secs
LT-INV: New	153 secs	113 secs

Comms developments: More on-demand SL-Comms using MPI

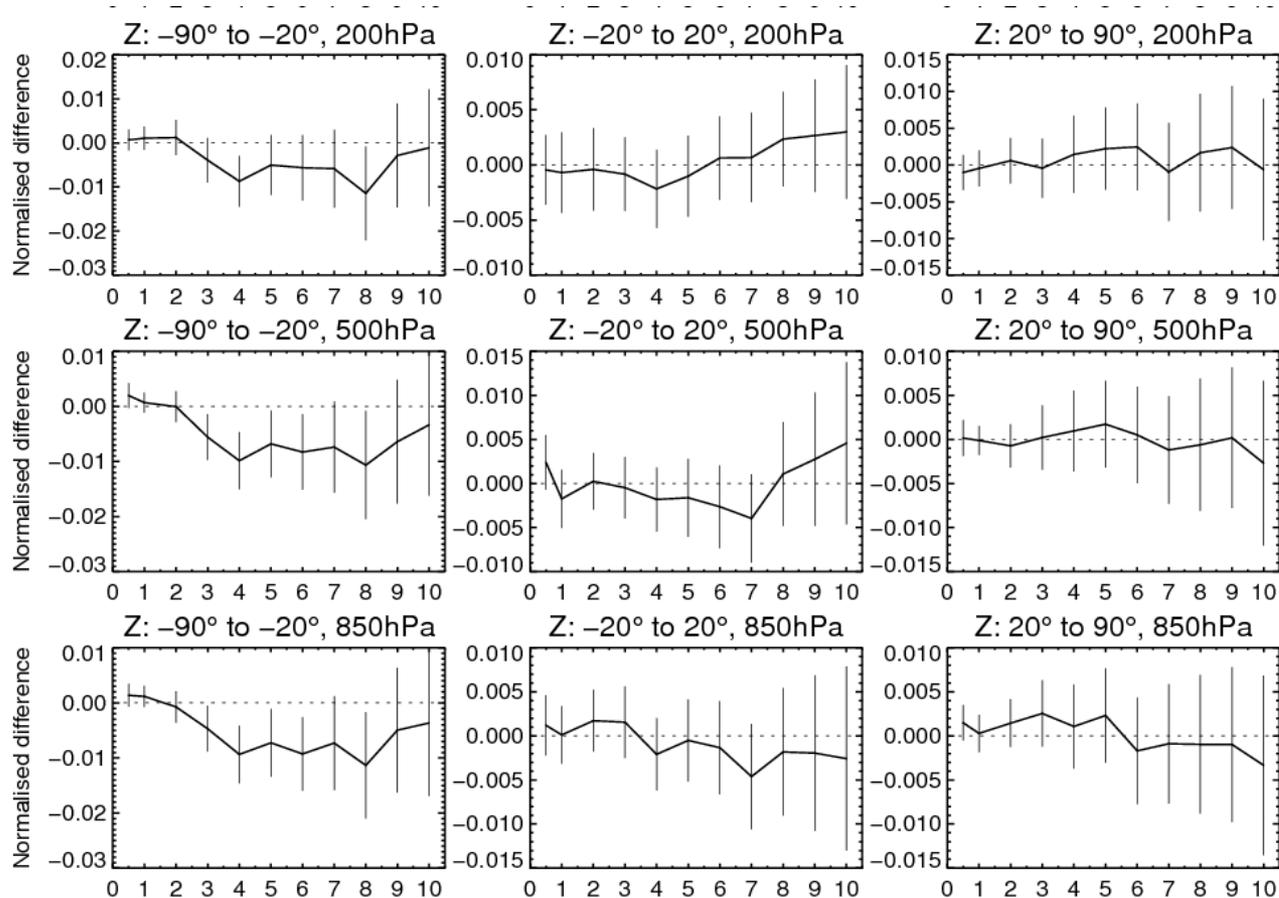
- 'SL-Comms part 2' already on-demand
- 'SL-Comms part 1' made on-demand to reduce volume of communications

% of total time



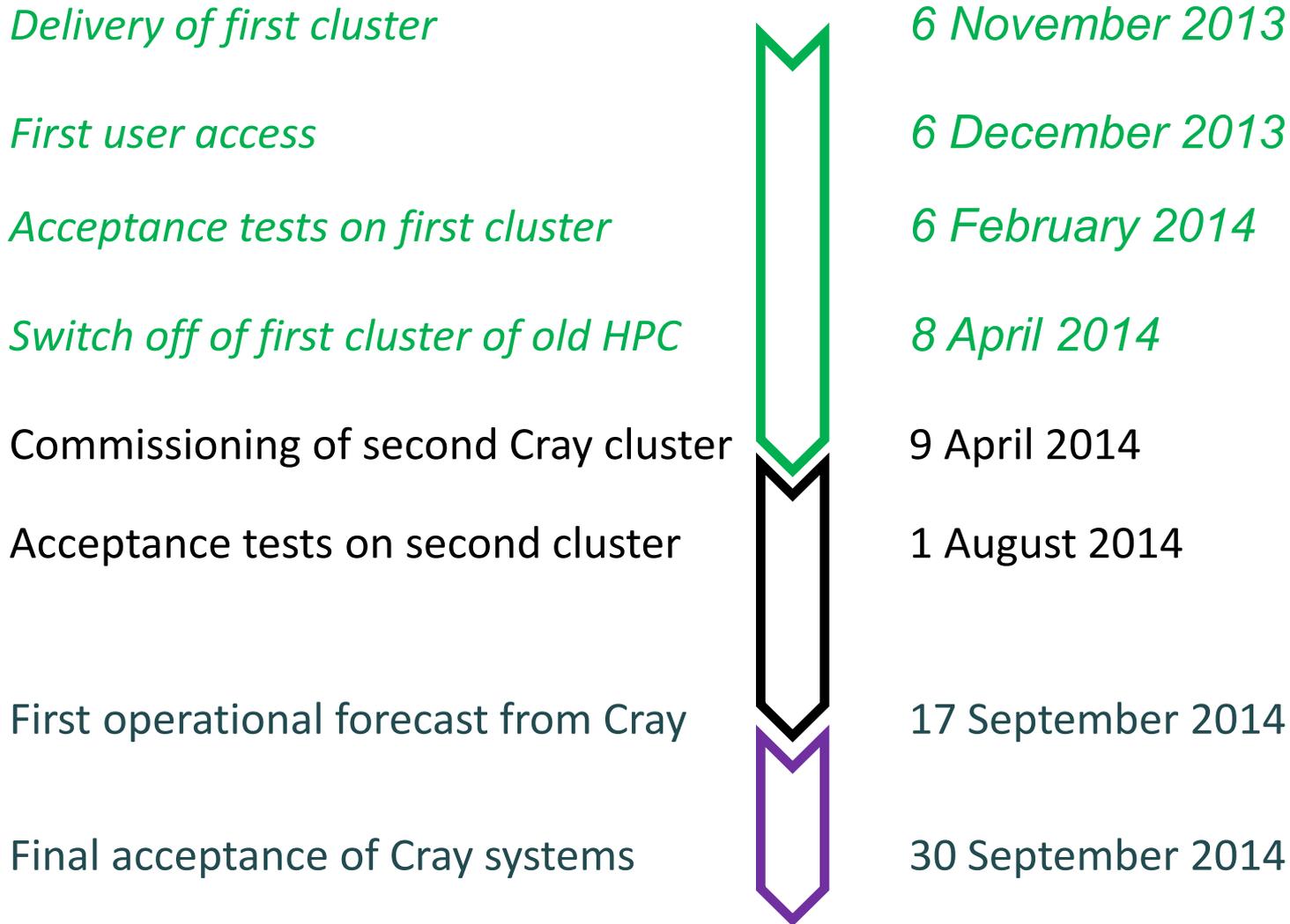
Improvement to scores with fixes to cce compiler difference of RMS error between runs with 8.2.7 & 8.2.0

2-Aug-2013 to 31-Mar-2014 from 322 to 360 samples. Confidence range 95%. Verified against own-analysis.

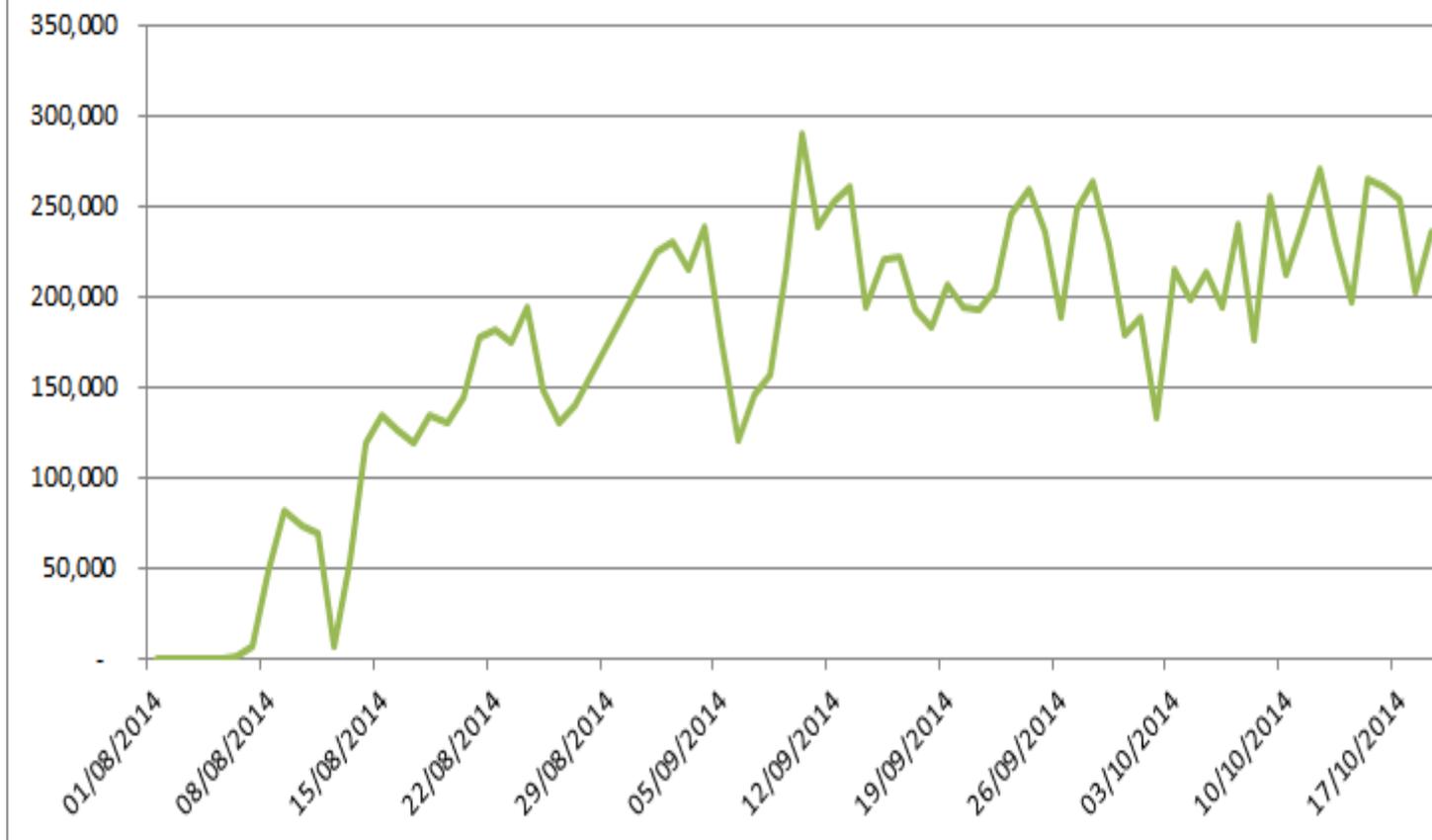


Thanks to Alan Geer

Migration Timeline



Jobs per day on one system



Operations on Cray

- Operational from mid September
- Many thanks to Cray staff
 - A huge effort by a large team of people
 - Both local and in the USA

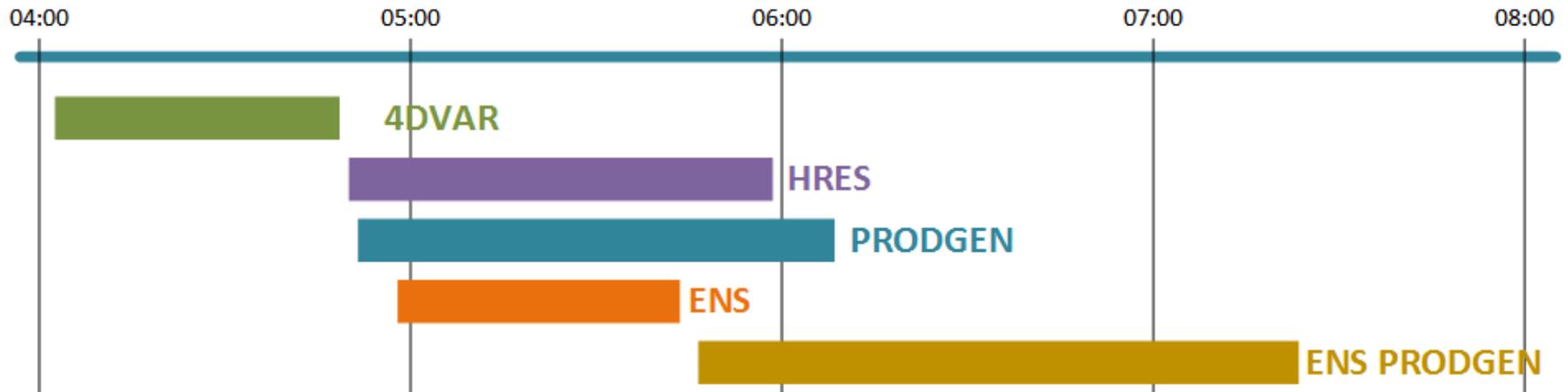
Operations on Cray

- Operational from mid September
- Many thanks to Cray staff
 - A huge effort by a large team of people
 - Both local and in the USA
- But.....

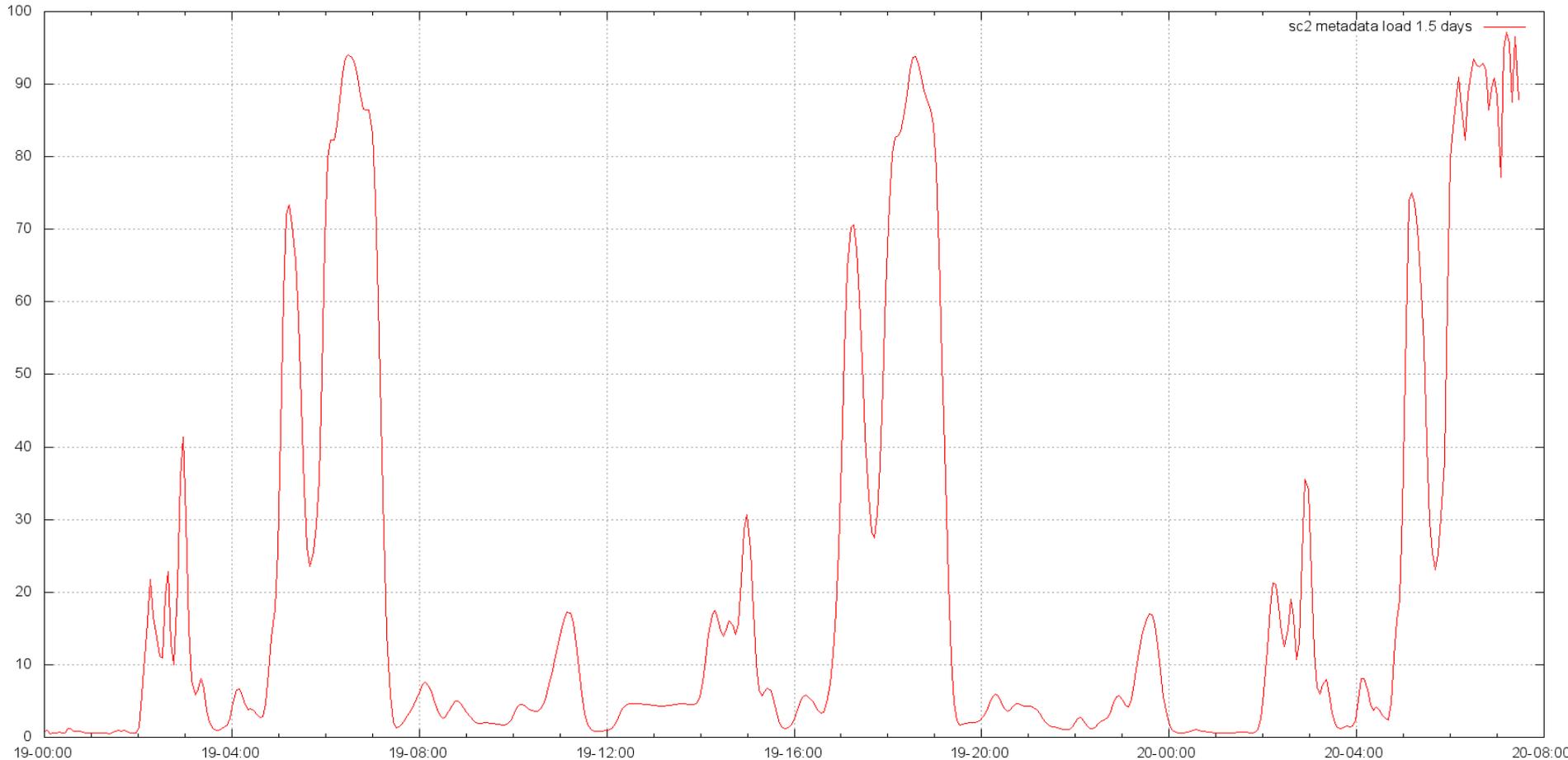
We have a Tuning Challenge

- Lustre is NOT GPFS
 - Different performance characteristics
 - Seeing delays due to IO jitter
- Need to streamline
 - Both workflow and IO load
- Tuning efforts started
- Cray to provide additional expertise

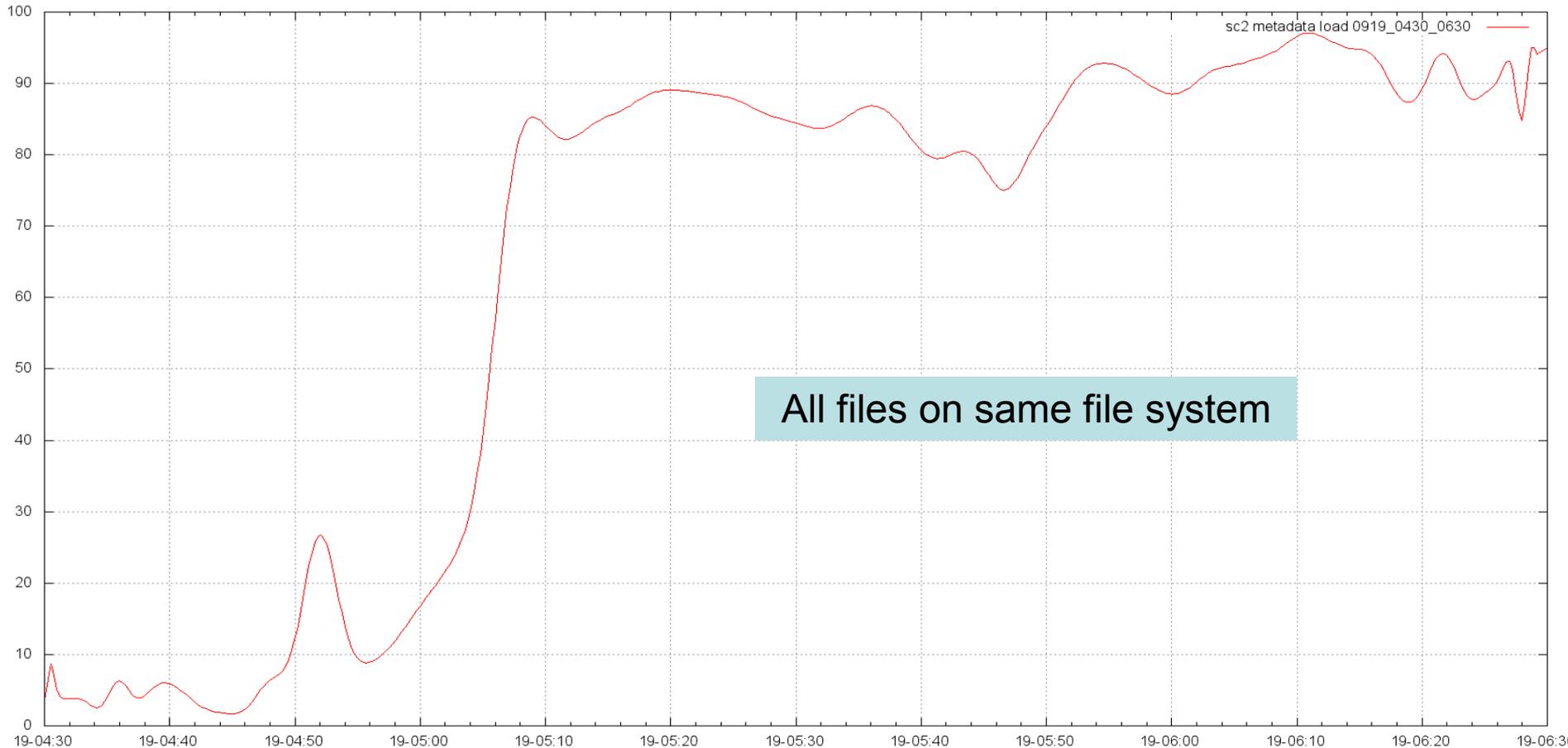
00Z Operational Schedule



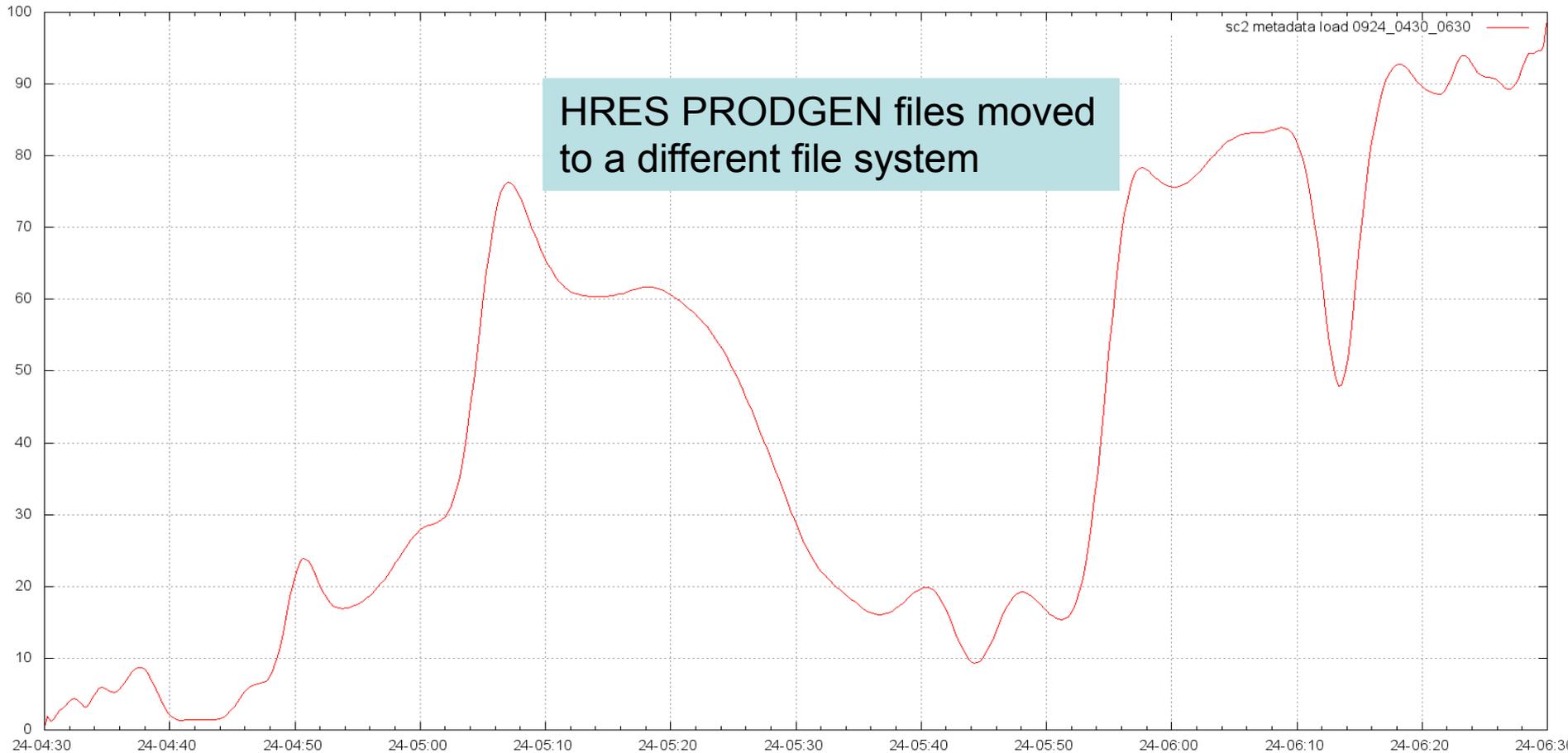
Operational File System on CCB: CPU Load on Meta Data Server on Oct 19



Operational File System on CCB: CPU Load on Meta Data Server on Sept 19



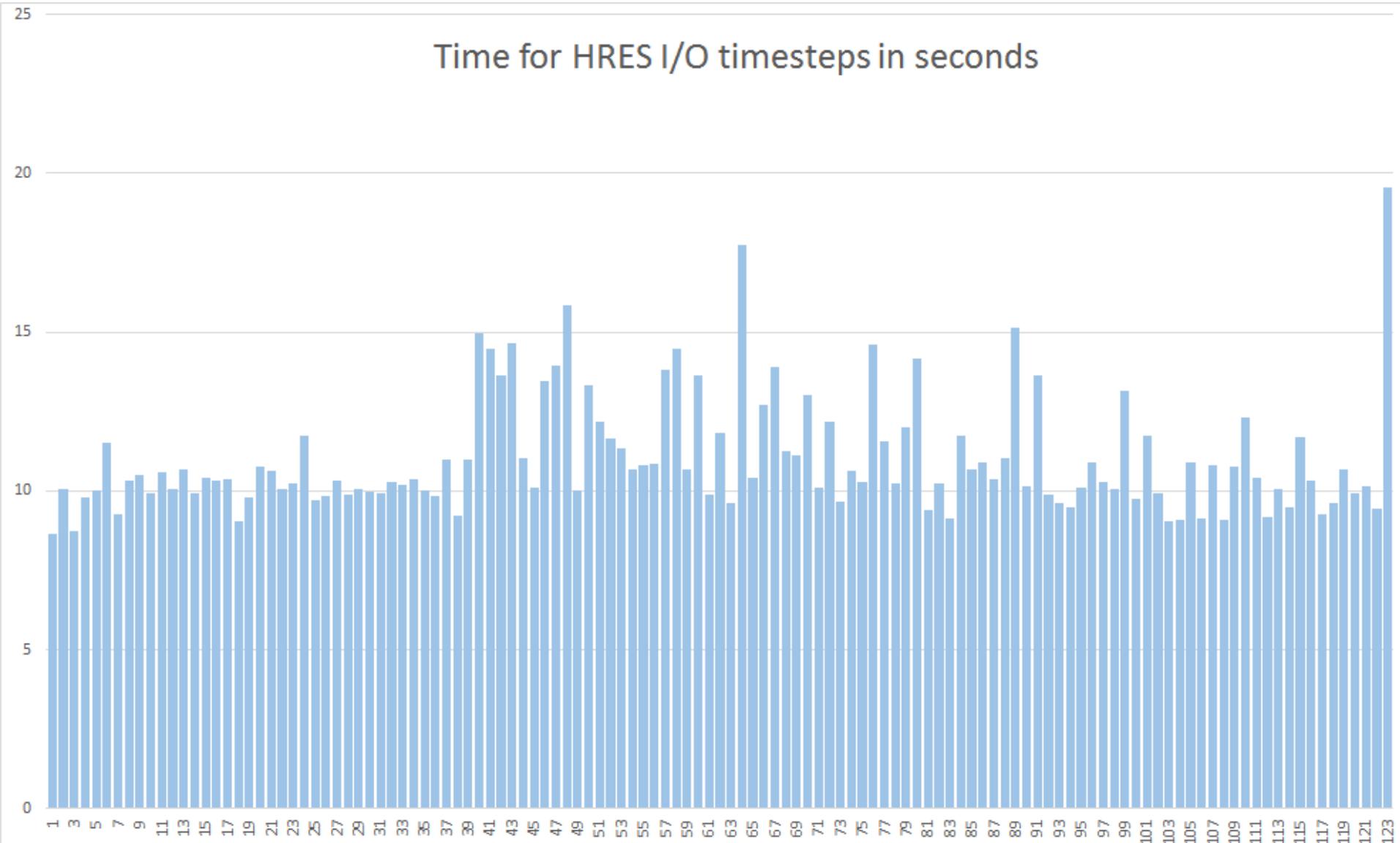
Operational File System on CCB: CPU Load on Meta Data Server on Sept 24



A Look at HRES IO

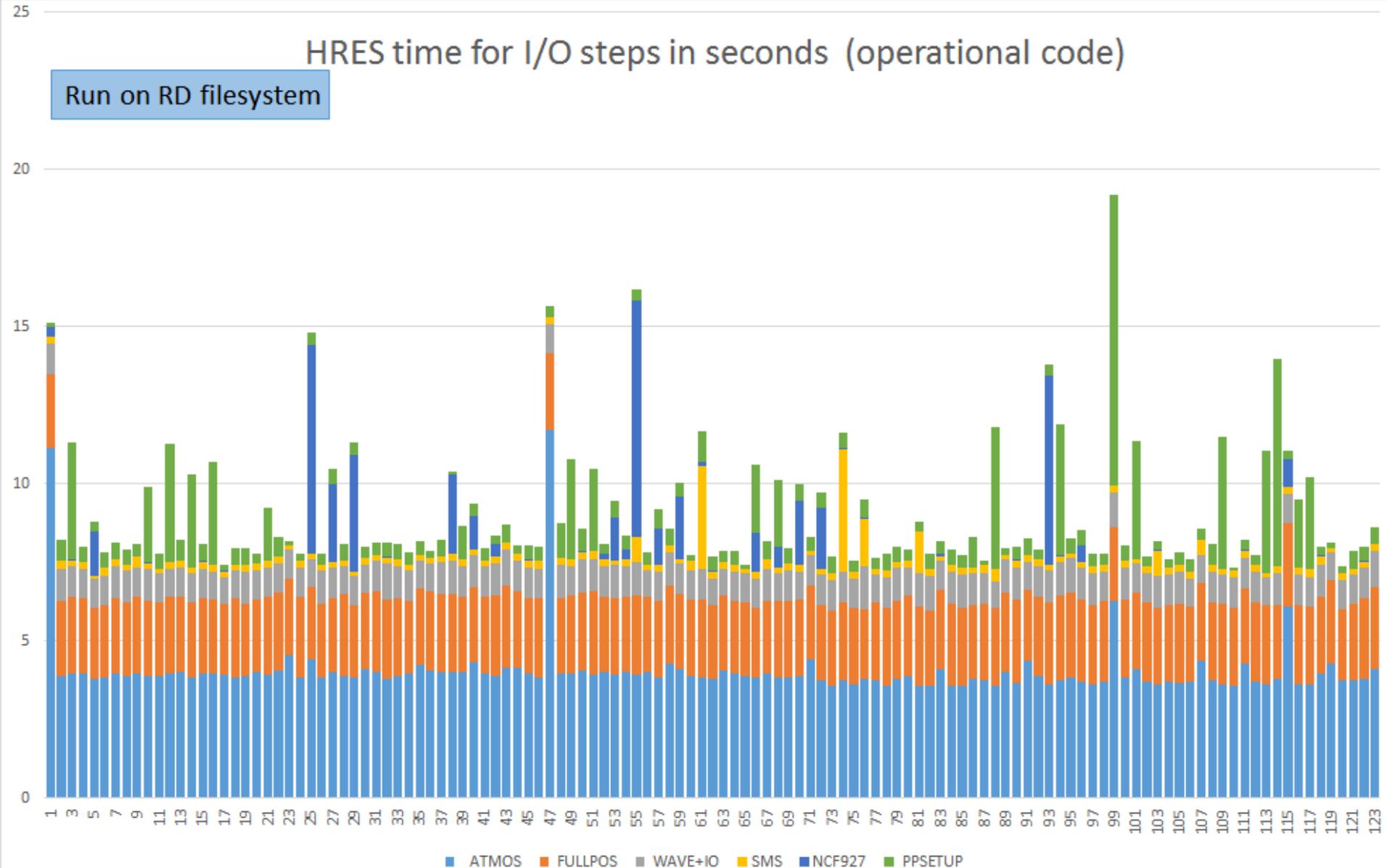
- 1.2TB total over 125 Post Processing IO steps
- Output to 120 files
- Bulk IO goes well (200+MB/s per file)
 - some delays due to jitter
- But we have IO jitter elsewhere
- Gets worse on a heavily loaded file system
- Tracked down to other IO operations every PP step
 - Open/read/close a file called dirlist by every task
 - Open/read/close 2 name list files by every task
 - Open/write/close a file called NCF927 by the master task

Time for HRES I/O timesteps in seconds



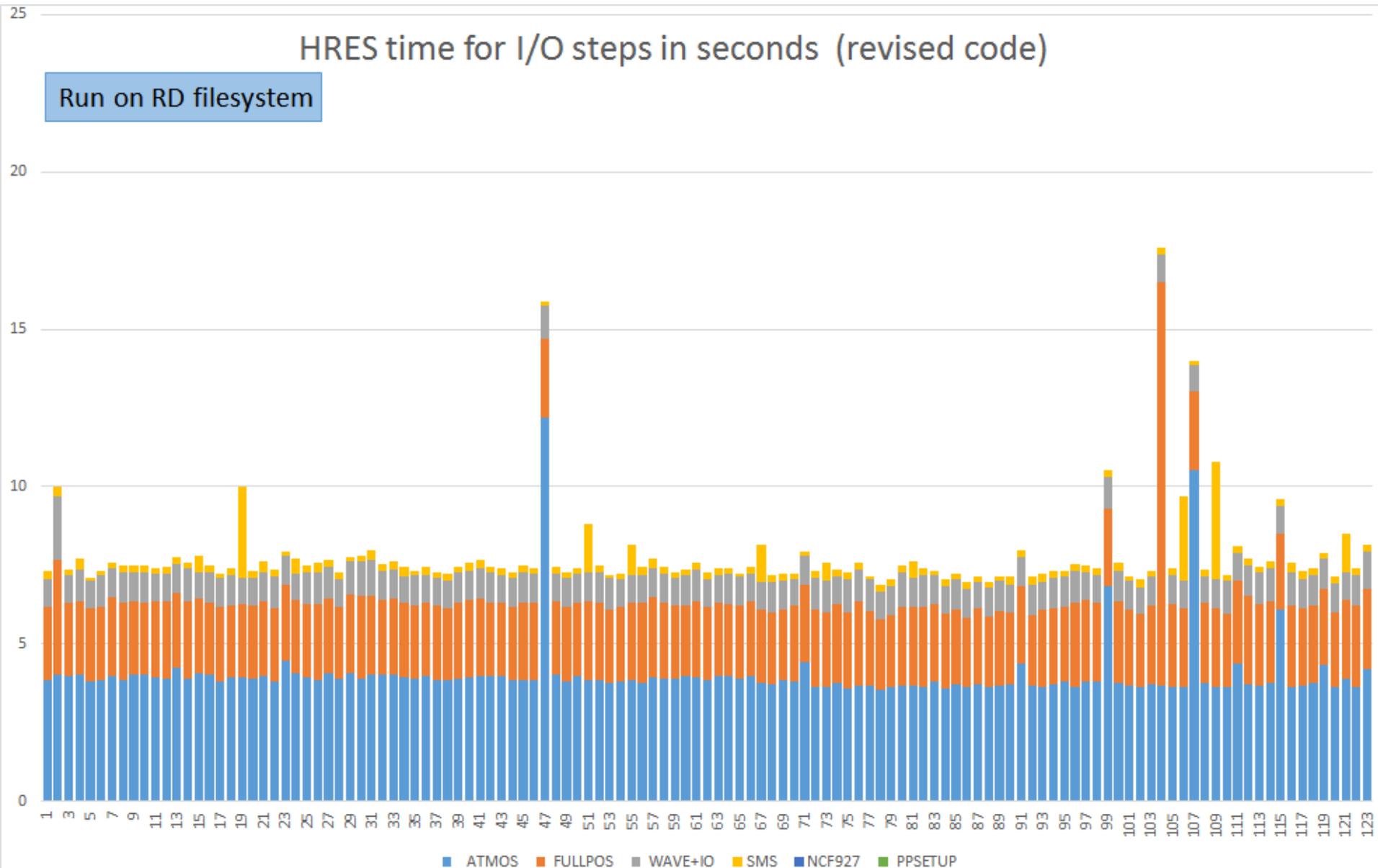
HRES time for I/O steps in seconds (operational code)

Run on RD filesystem



HRES time for I/O steps in seconds (revised code)

Run on RD filesystem



IO that needs to be removed

HRES - 180K sets of open/read or write/close

ENS - 1.2M sets of open/read or write/close

HRES PRODGEN - 930K temporary files

- 2.8 M sets of open/read or write/close

ENS PRODGEN - 880K temporary files

- 2.6M sets of open/read or write/close

Lessons Learned

- Ever increasing complexity
- More resources required
 - To meet applications and systems challenges
- Must allow more time for performance testing of the operational suite