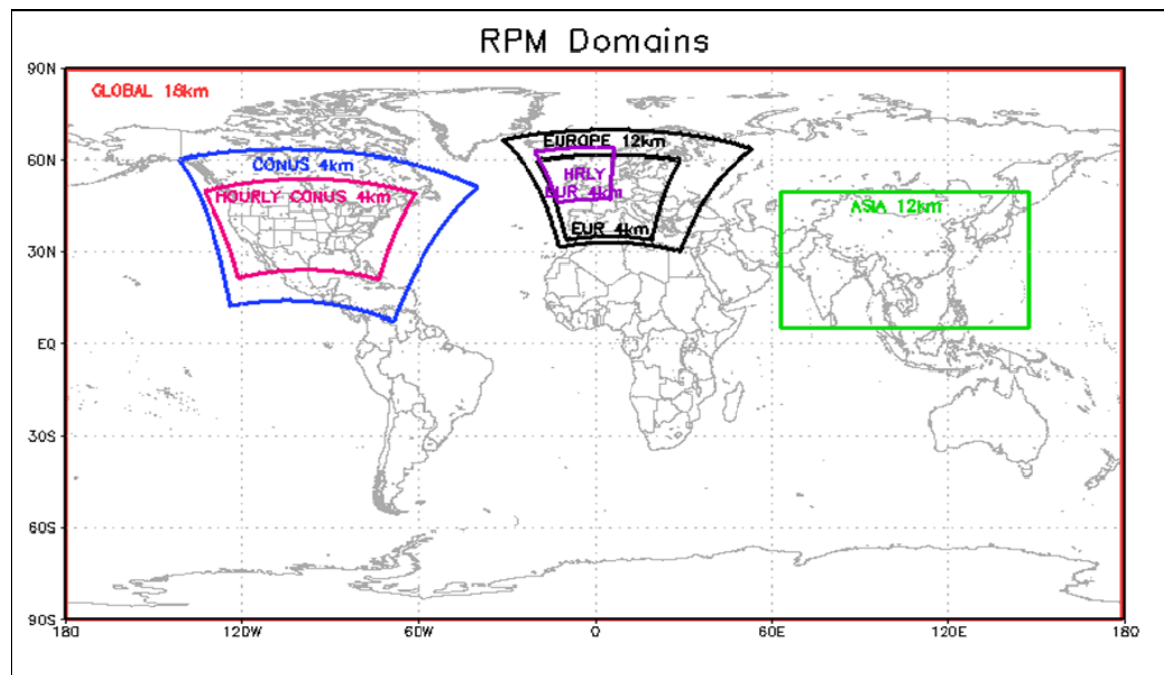# Performance analysis of an operational implementation of WRF

**Todd Hutchinson**

**WSI Corporation**

- **Analyze operational configuration of WRF to find potential for performance improvement**



**WSI°**

# Hardware and Software Used
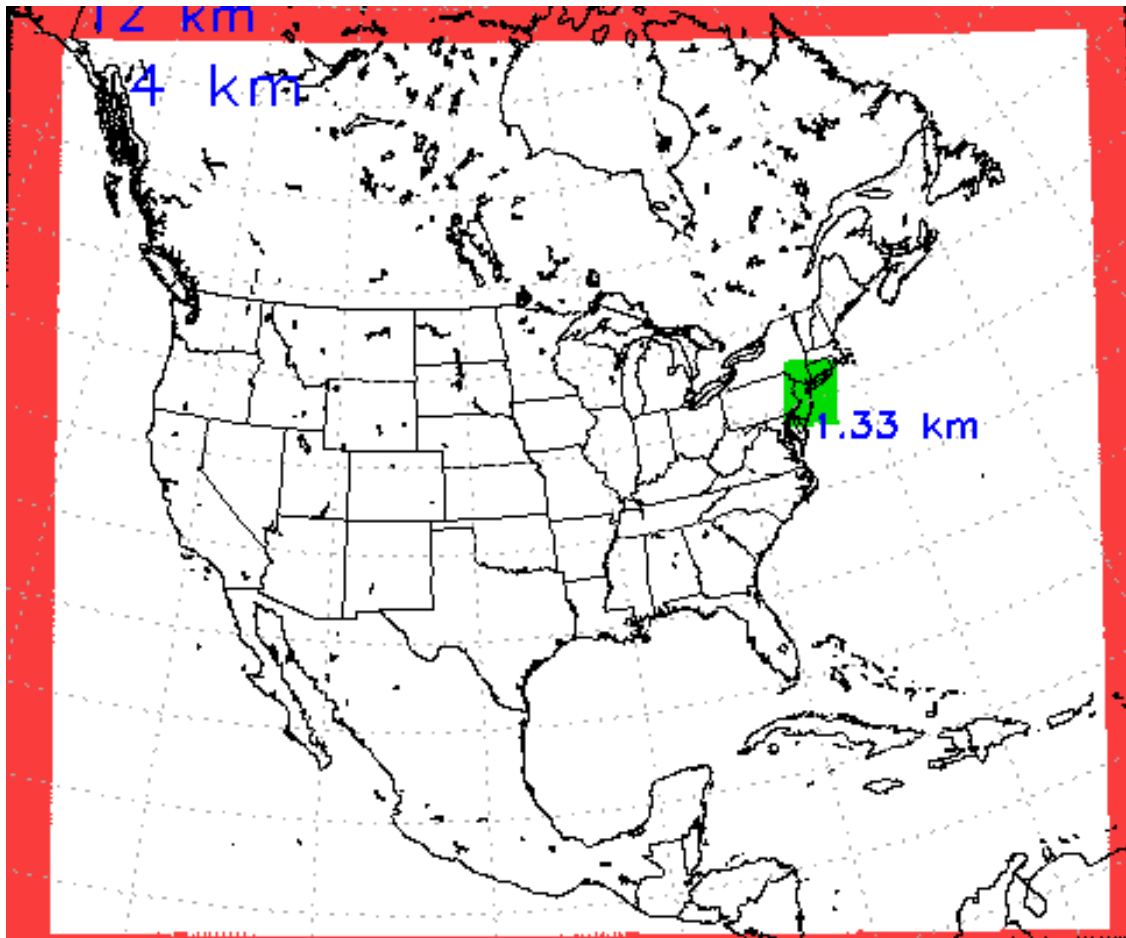
- **Hardware:**
  - Infiniband FDR (56 Gb/s, 0.7 µs)
  - 18 Compute Nodes, each with:
    - 2 Intel E5-2697 Processors = 24 cores/node
    - 64 GB memory
  - Total of 432 cores

- **Software:**
  - Intel compilers version 15.0.0.90
  - Intel mpi 5.0 Update 1
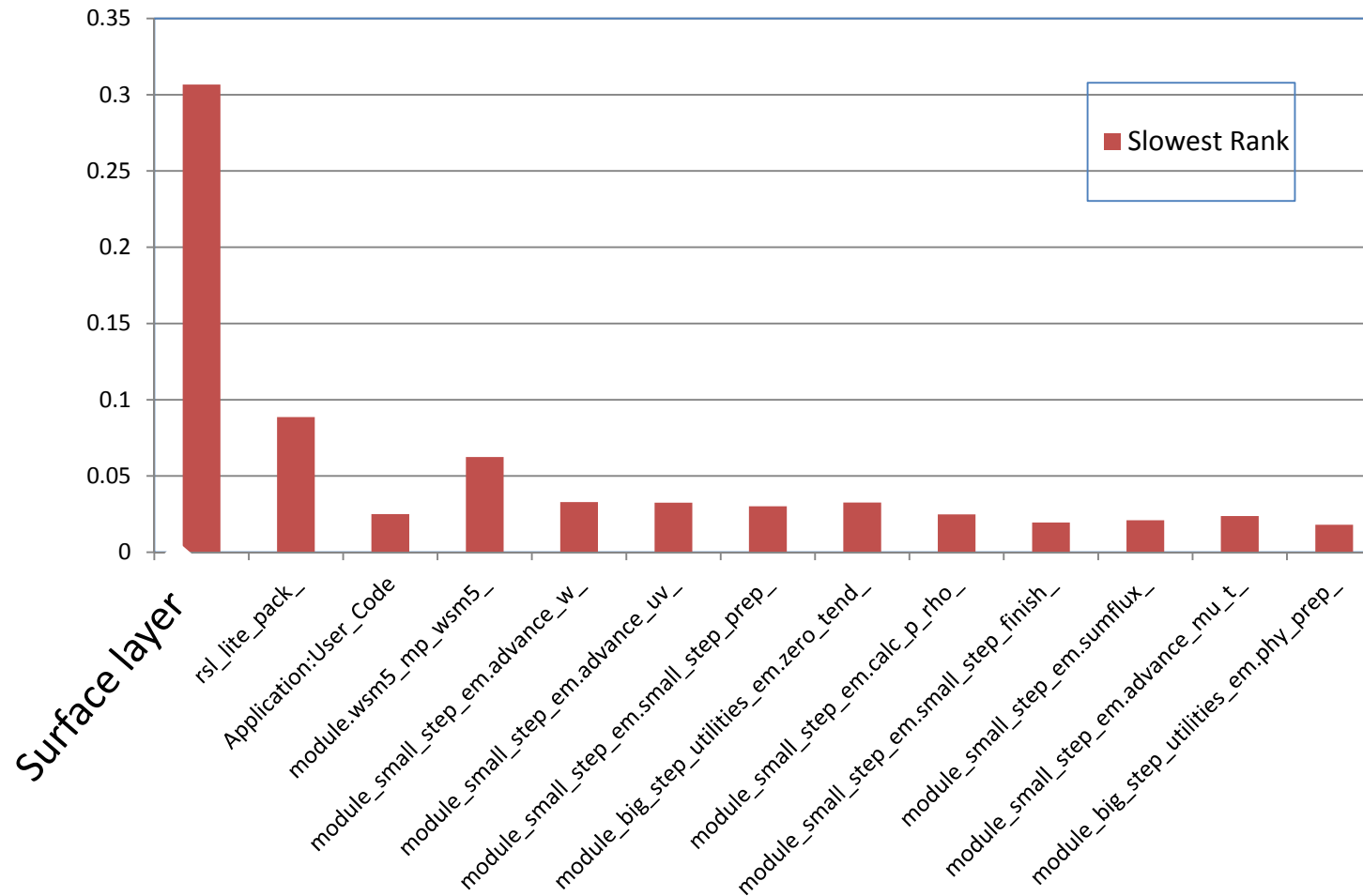  - Intel Trace Analyzer and Collector 9.0 Update 1

**WSI**°

# WRF Configuration



**Model Setup:**
- **WRF ARW v3.6.1**
- **12/4 km, 38 levels**
- **4k feeds back to 12k**
- **Physics:**
  - CU: KF on 12k
  - PBL: YSU
  - MP: WSM5
  - RA: RRTM/Goddard
- **MPI only (no OpenMP)**
- **51/72 hr model runs:**
  - We'll look at first 12 hrs here

**WSI**°

# A quick profile of one rank, one time step
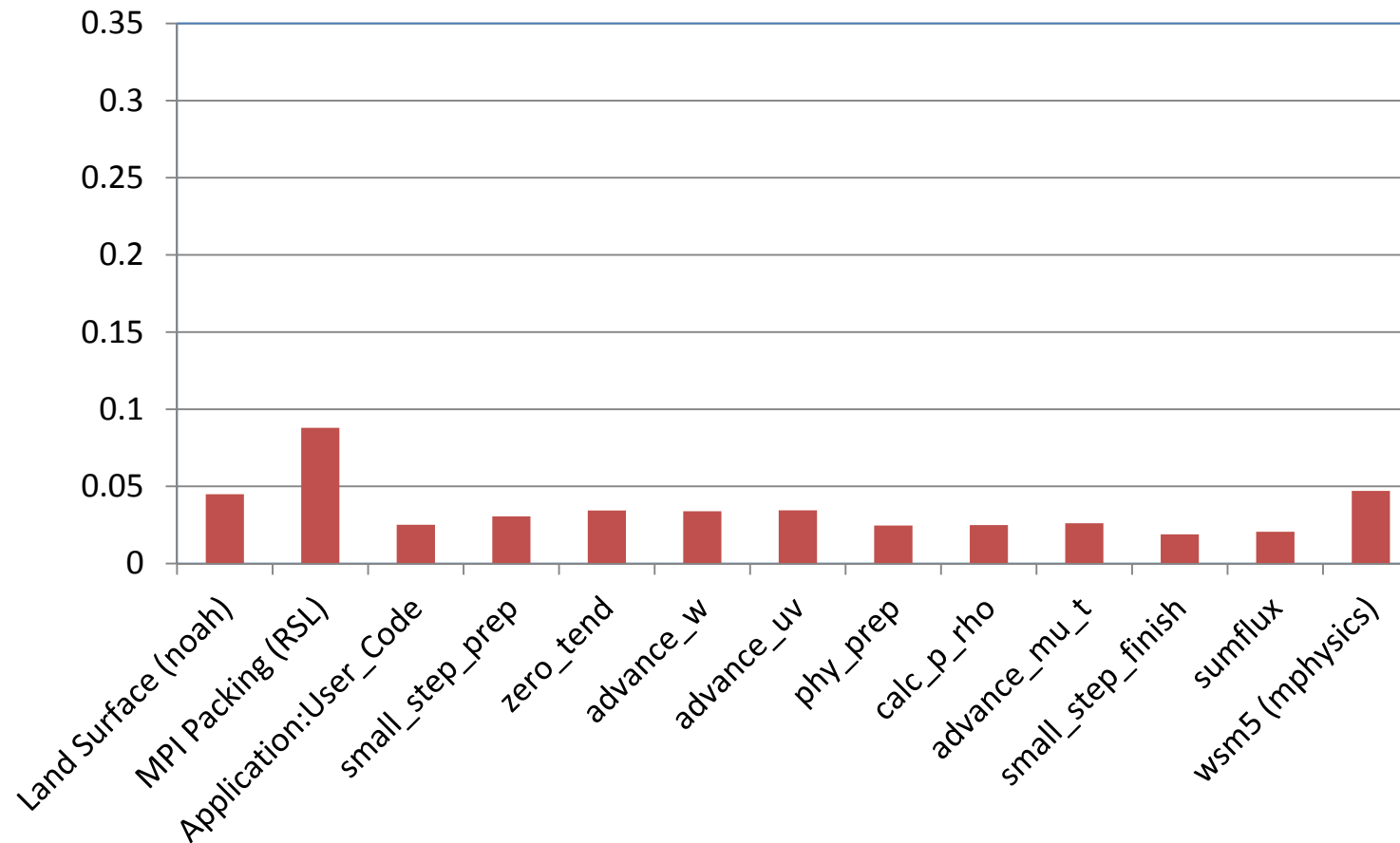


Time Spent (s) in Each Function

# The profile reveals:

- **The surface layer calculation (sf_sfclay=1) is taking up 40% of the run time!**

- **Further analysis showed that this was an issue introduced in WRF3.6**

- **WRF Developers at NCAR are investigating**

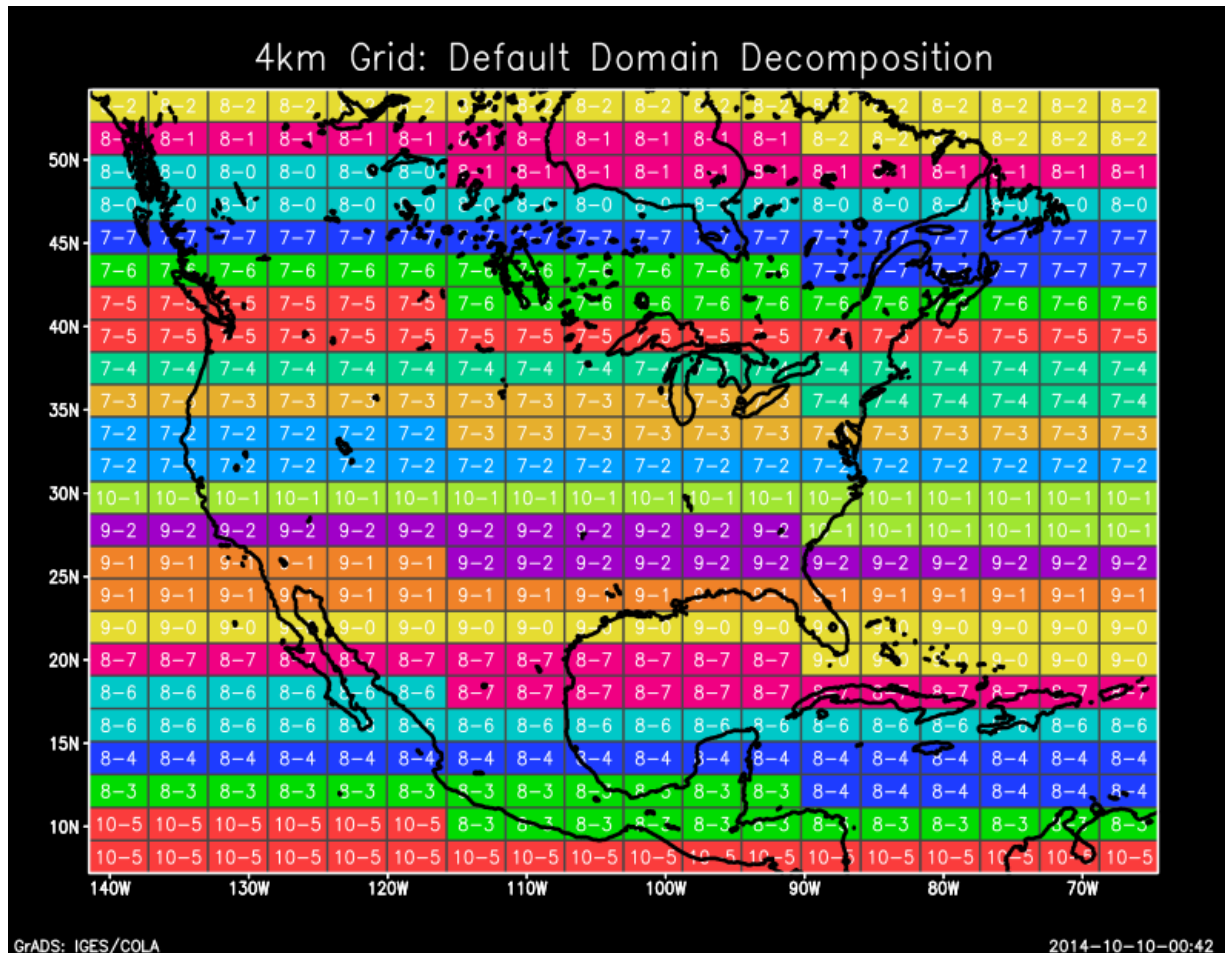- **The quick fix is to set sf_sfclay=91 (the older version of this surface layer scheme)**

**WSI**°

Time Spent (s) in Each Function

**Now, the surface layer doesn't make the top 13!**

# Default Decomposition of 4km Domain
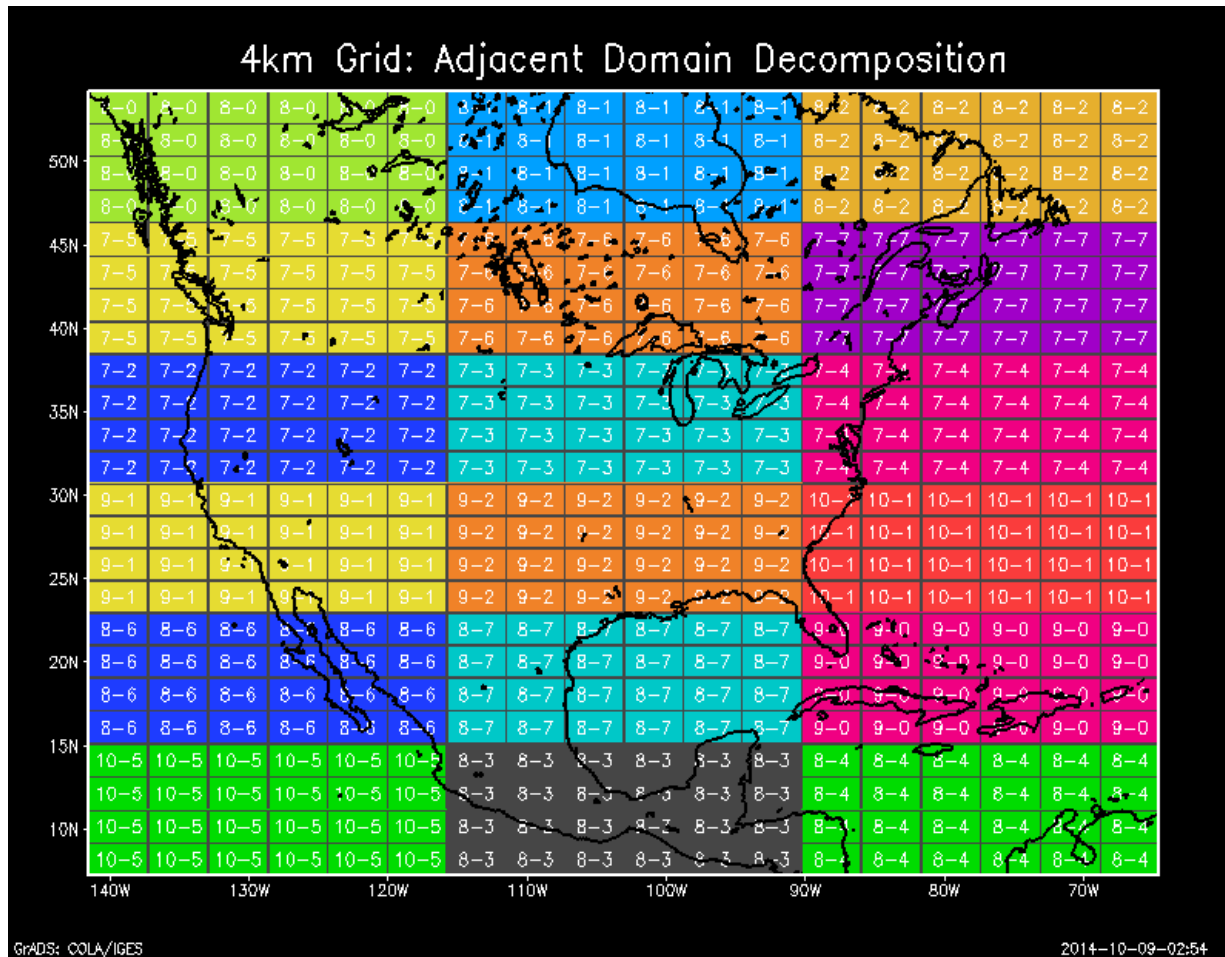


4km Grid: Default Domain Decomposition

- **Each Patch is a different cpu core**
- **Each color a different compute node**

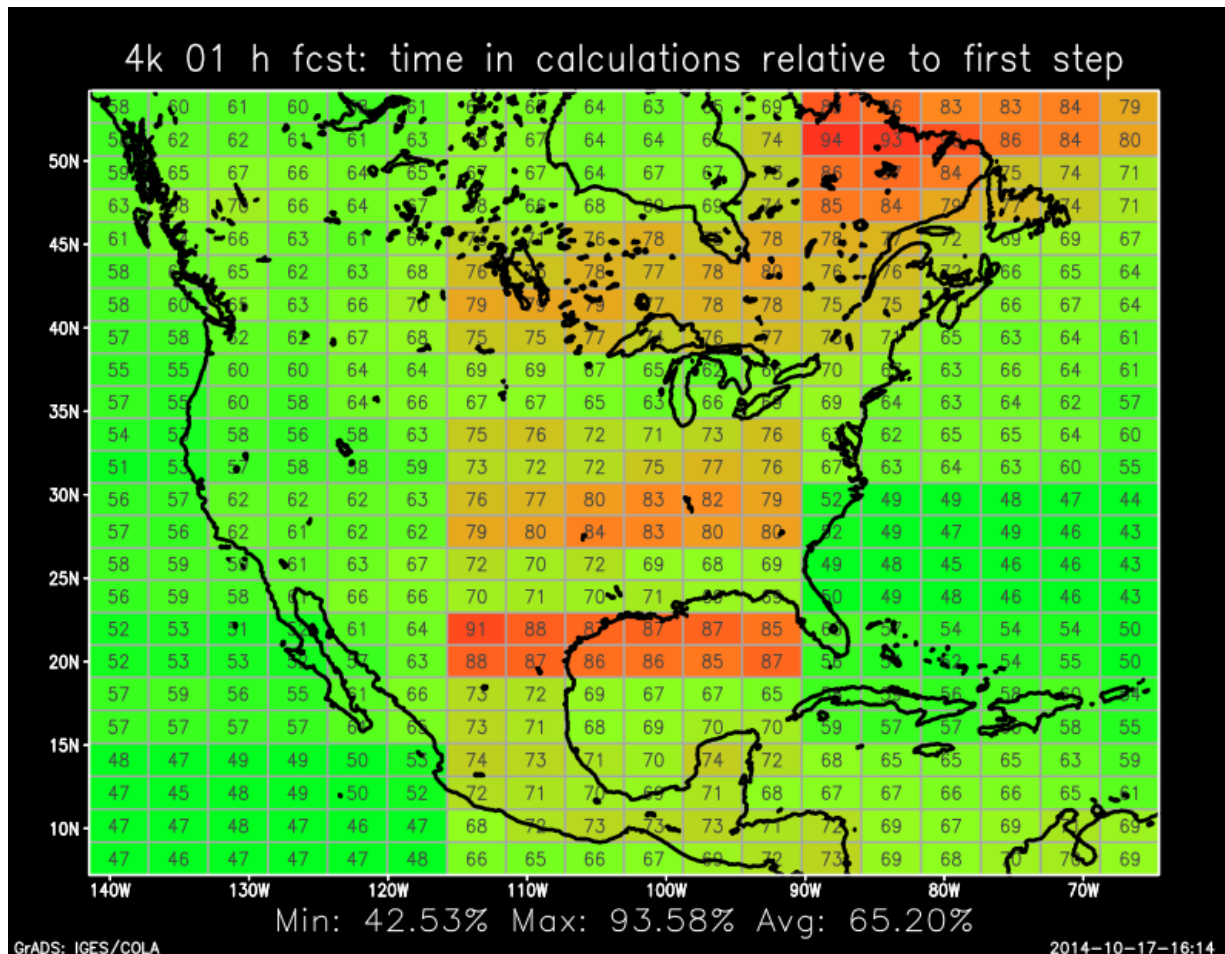# Default Decomposition: Relative Run Time for patches



- **Relative (%) time that is spent computing**
- **Green: Low compute time**
- **Red: High compute time**

# Adjacent Decomposition

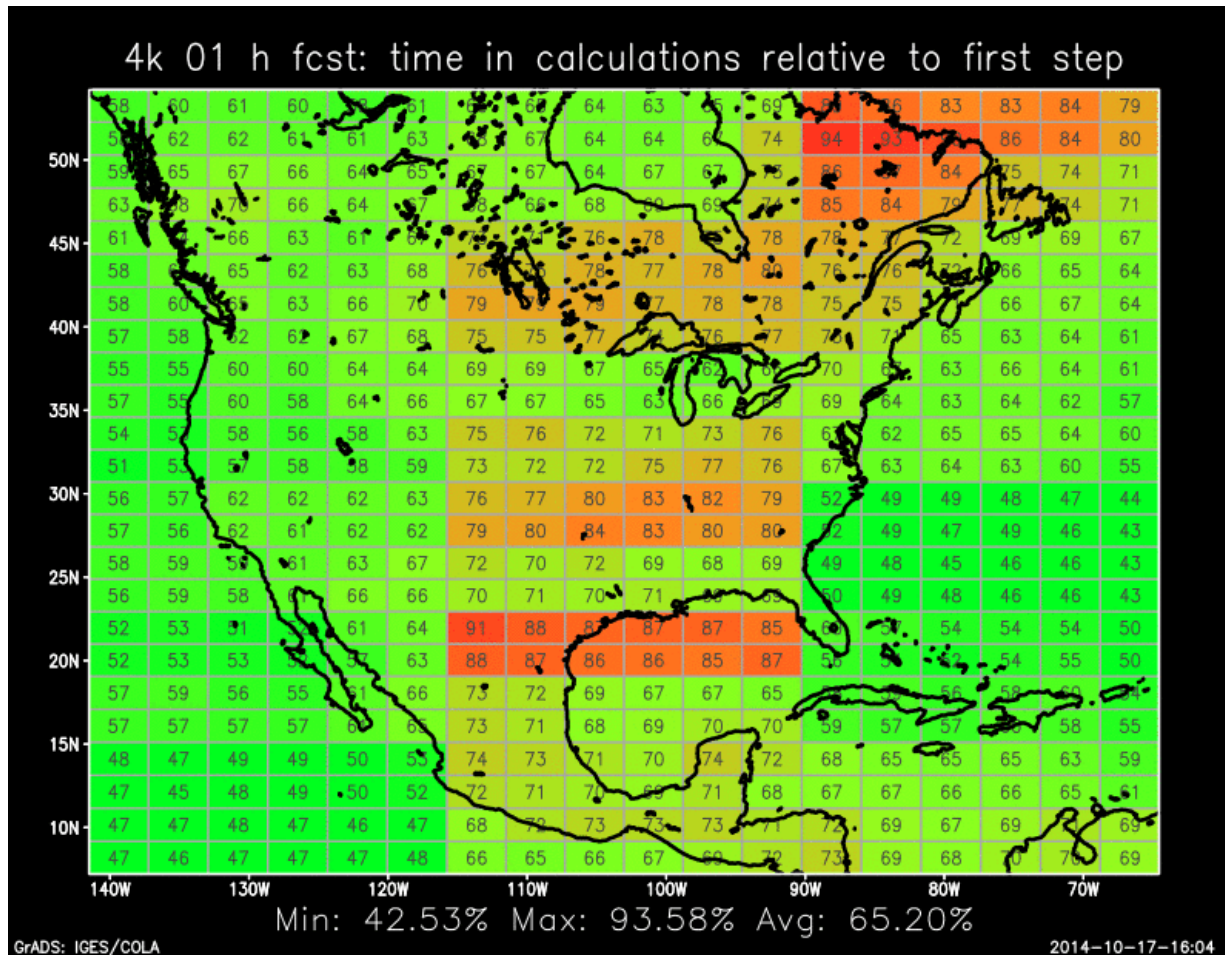

4km Grid: Adjacent Domain Decomposition

- **That's easier to see!**

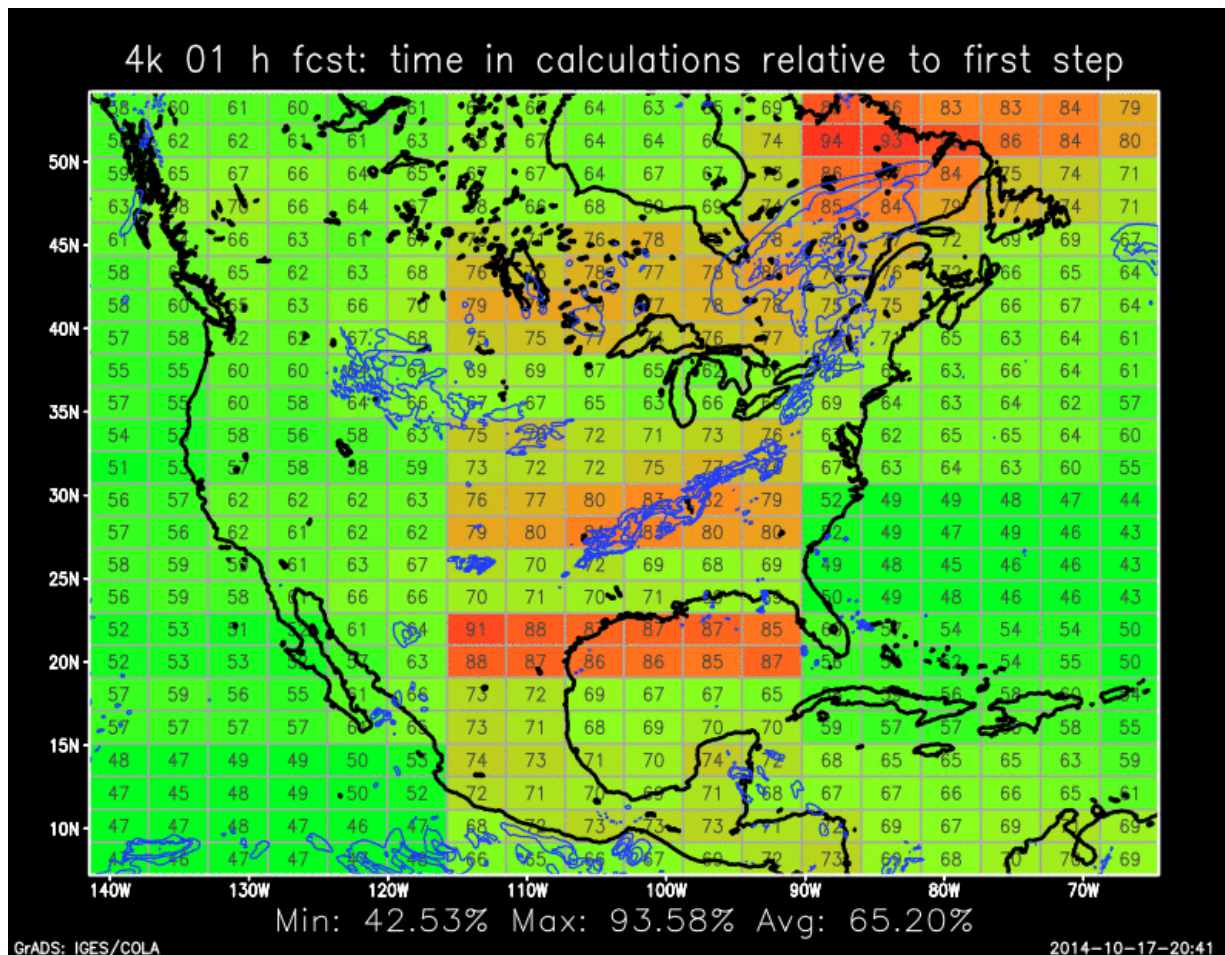# Adjacent Decomposition Timing



- **And now we see some dependence on processor**
- **Total run-time same as before, average decreased by 2%**

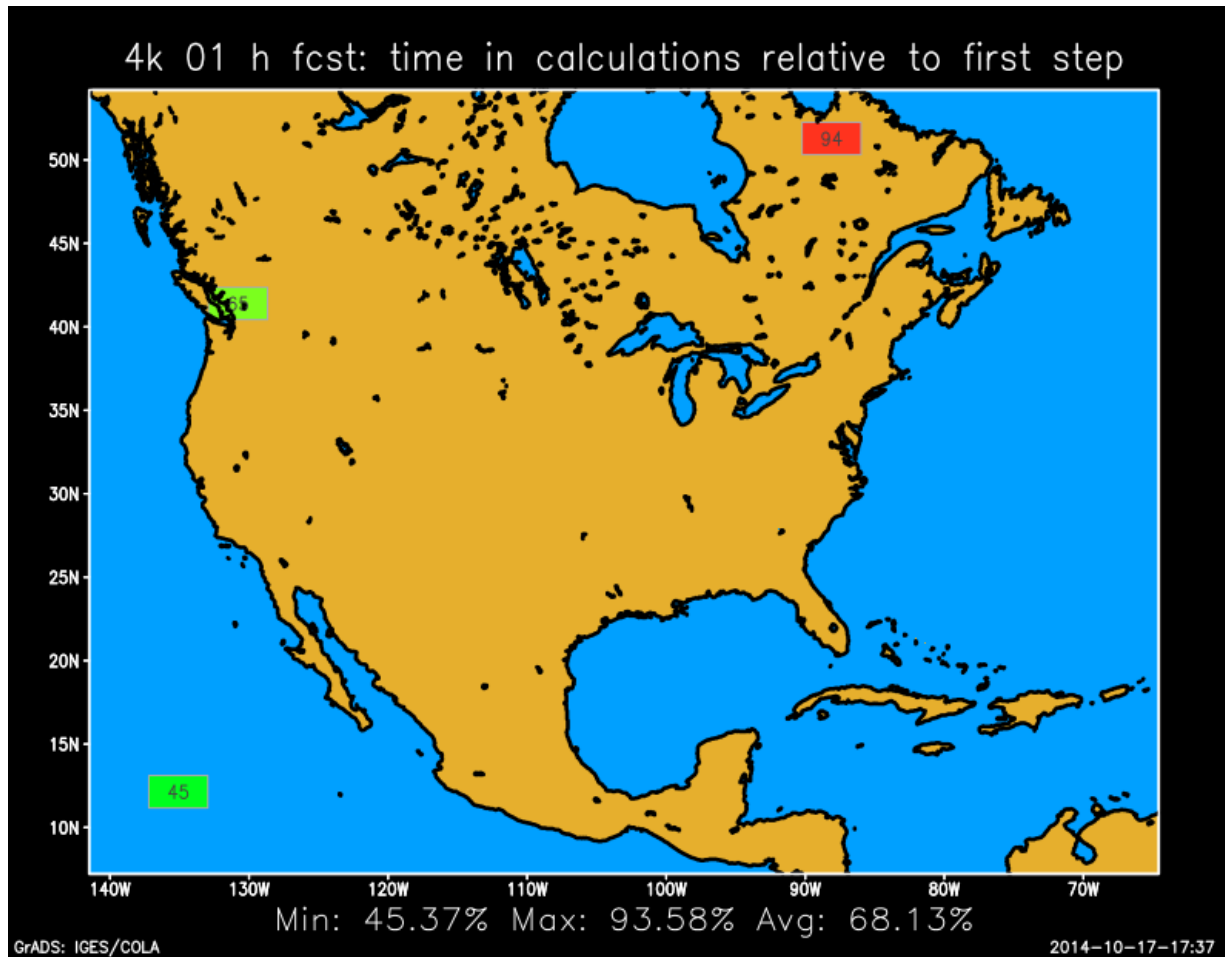# Patch timing every hour



4k 01 h fcst: time in calculations relative to first step

Min: 42.53%  Max: 93.58%  Avg: 65.20%

GrADS: IGES/COLA                                    2014-10-17-16:04

- **Patch timing relative to first 1 hour step time**

# Patch timing with precip



4k 01 h fcst: time in calculations relative to first step

Min: 42.53% Max: 93.58% Avg: 65.20%

GrADS: IGES/COLA                    2014-10-17-20:41

- **60-minute precip accumulation overlaid**

**WSI°**

# Slowest, Median, Fastest Patches

# Traces over 1 time step



Application Processing

MPI Wait

Fast

Median

Slow

Data Packing

**WSI**
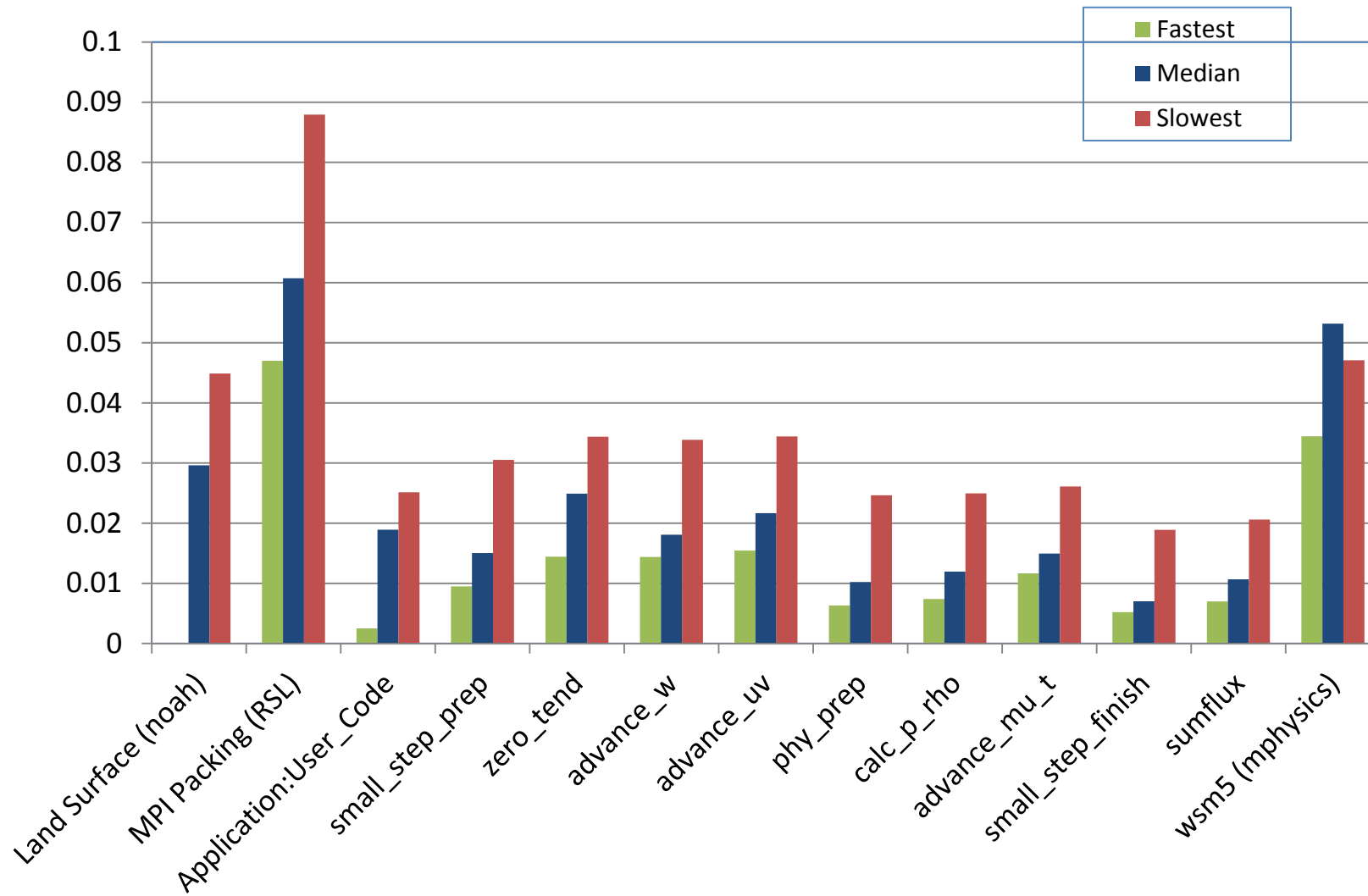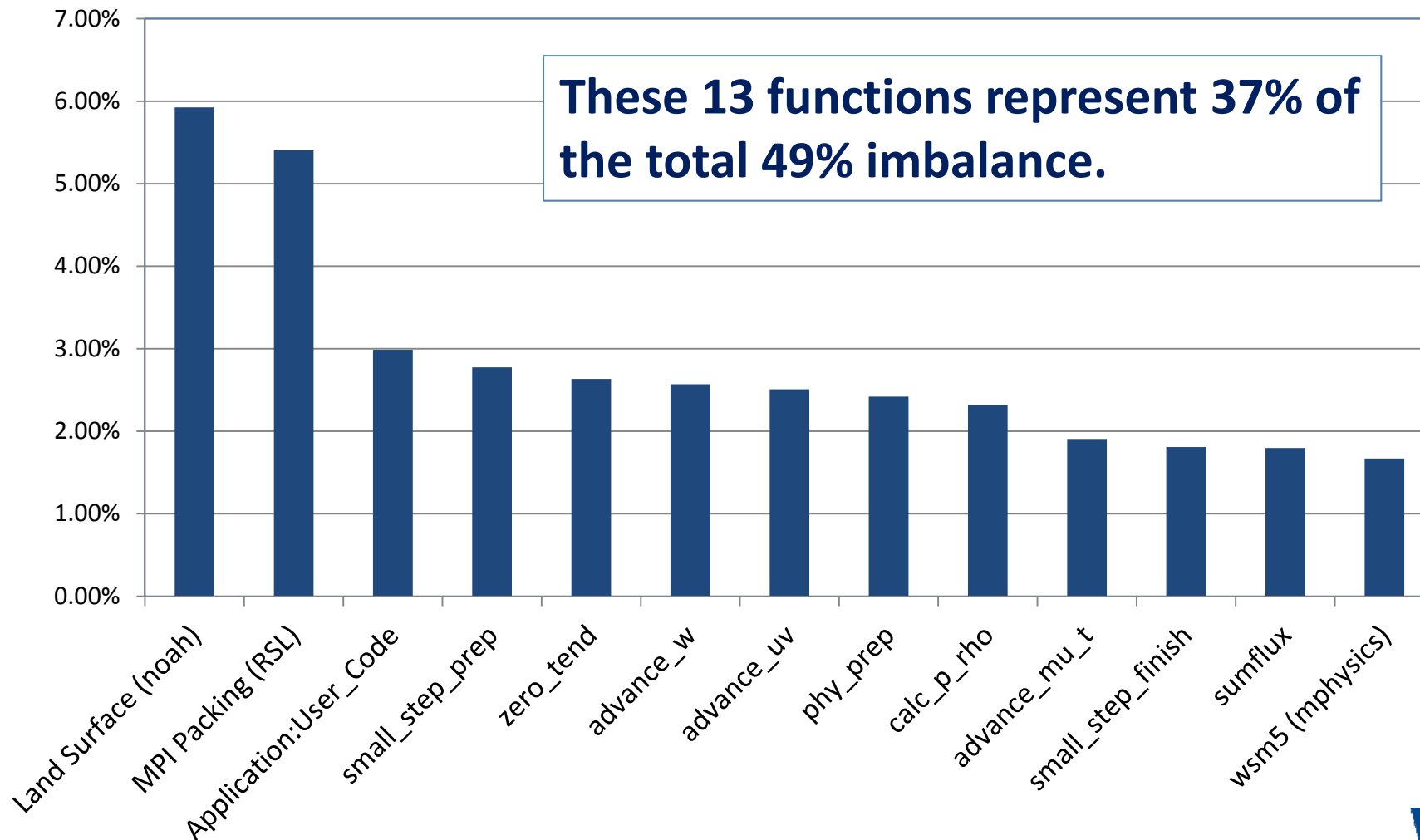
# Profile of the 3 traces



Time Spent (s) in Each Function

# Imbalance across the three Patches

**Imbalance across Ranks as a Percentage of Time Step Time**

These 13 functions represent 37% of the total 49% imbalance.



WSI°

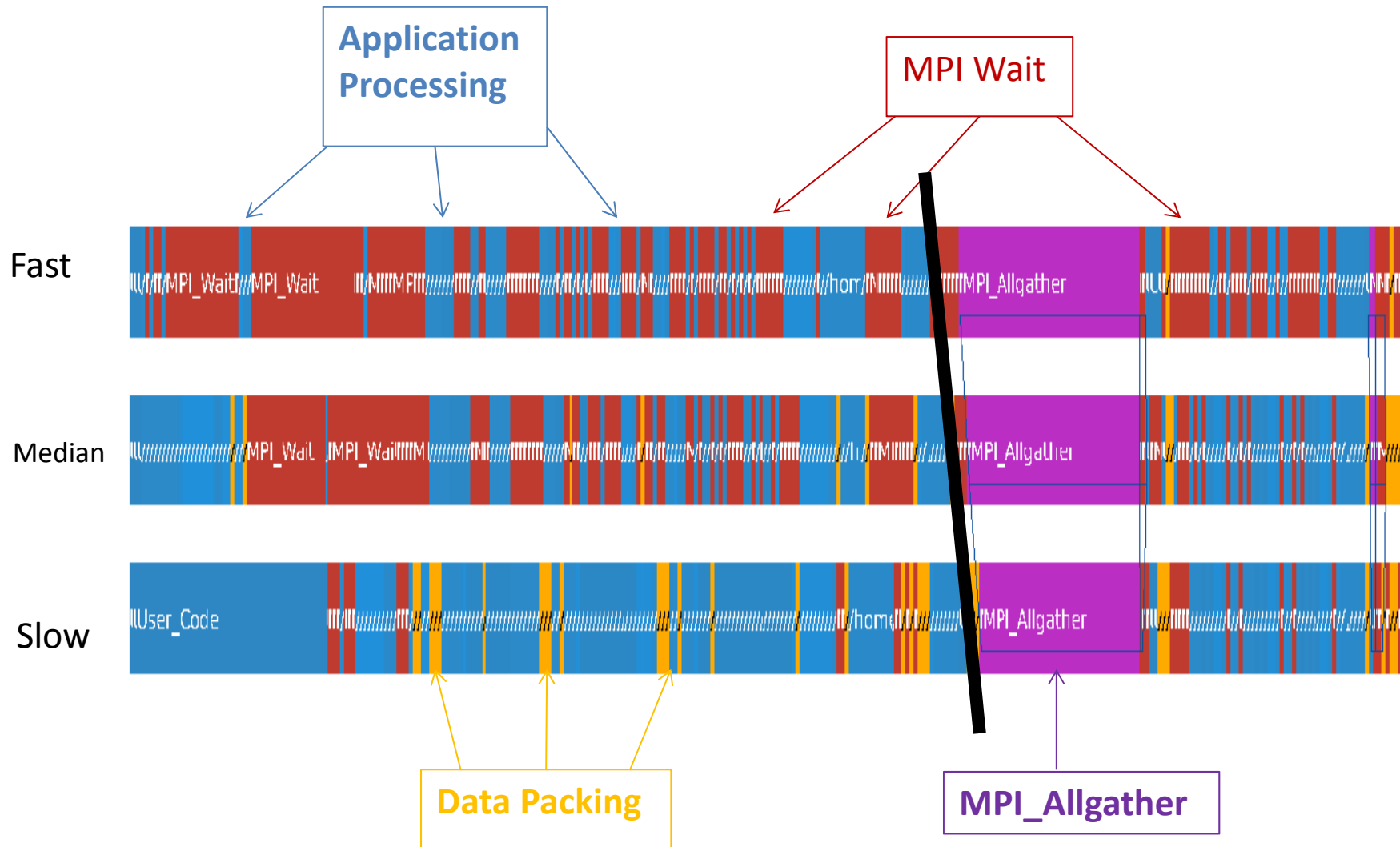# WRF Configuration



**Model Setup:**

- **WRF ARW v3.6.1**
- **12/4/1.33 km, 38 levels**
- **4k feeds back to 12k**
- **Physics:**
  - CU: KF on 12k
  - PBL: YSU
  - MP: WSM5
  - RA: RRTM/Goddard
  - 38 levels
- **MPI only (no OpenMP)**
- **51/72 hr model runs:**
  - We'll look at first 12 hrs

**WSI**°

# Time of a Single Parent Step

|  | Time (s) | Increase Size x # steps | Cumulative increase in grid points | Cumulative % increase in time |
|---|---|---|---|---|
| 12k | 0.131 | 555x473 x 1 | -- | -- |
| 12k + 4k | 2.941 | 1525x1309 x 3 | 22.8x | 22.5x |
| 12k + 4k + 1.3k | **5.123** | 229x280 x 9 | 1.092x | **1.74x** |

# Having a small nest within a large parent appears to be expensive!

WSI°

# Trace of 4km and 1.3km steps



Application Processing

MPI Wait

Fast

Median

Slow

Data Packing

MPI_Allgather

# Analysis of 1.33km step



Overhead    1km steps    Overhead

1    2    3

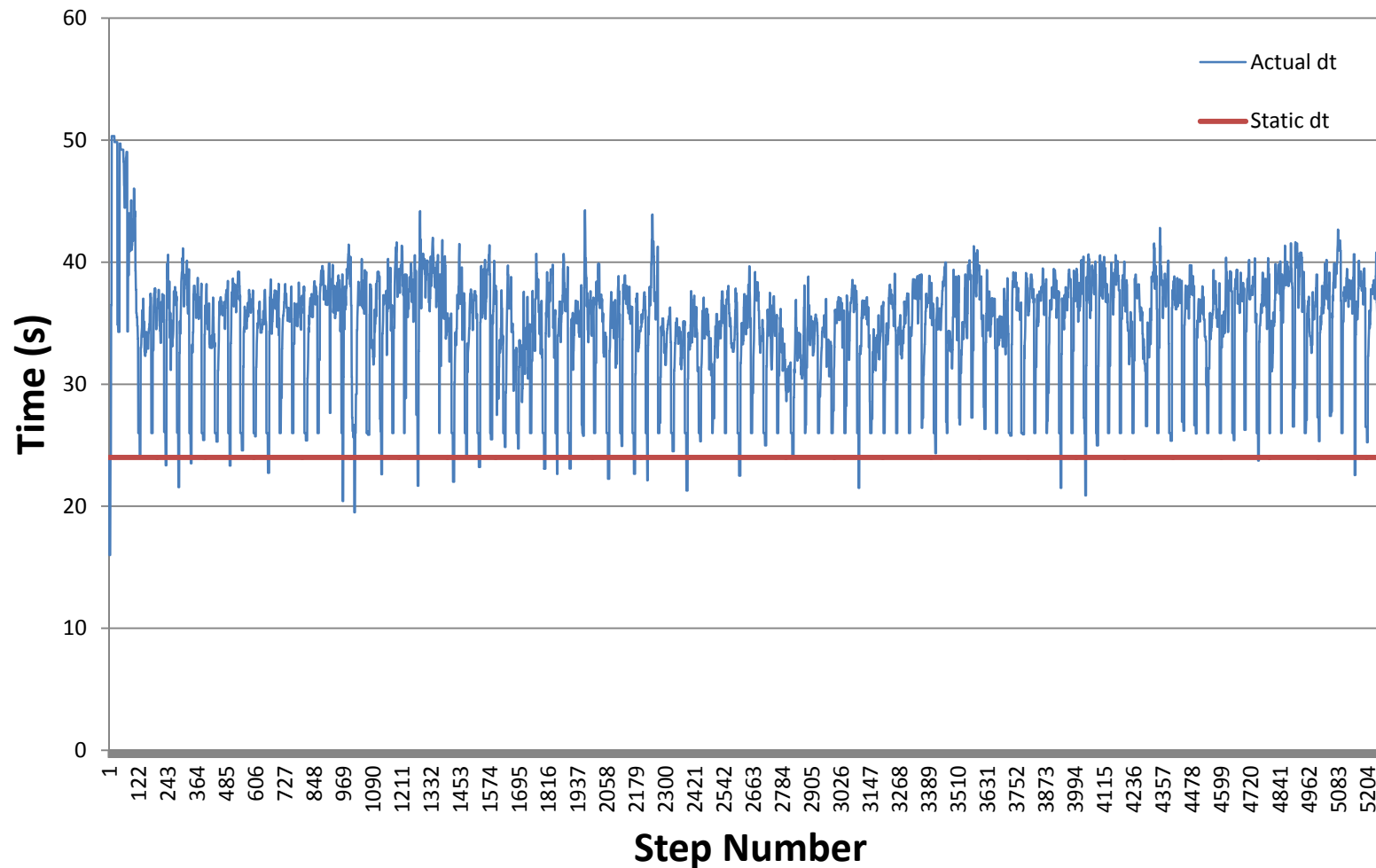MPI_Allgather

Memory Allocation

WSI

# Adaptive Time-Step

- **Automatically adapt the time-step to support maximum horizontal and vertical motions**

- **Adaptation assures stability**

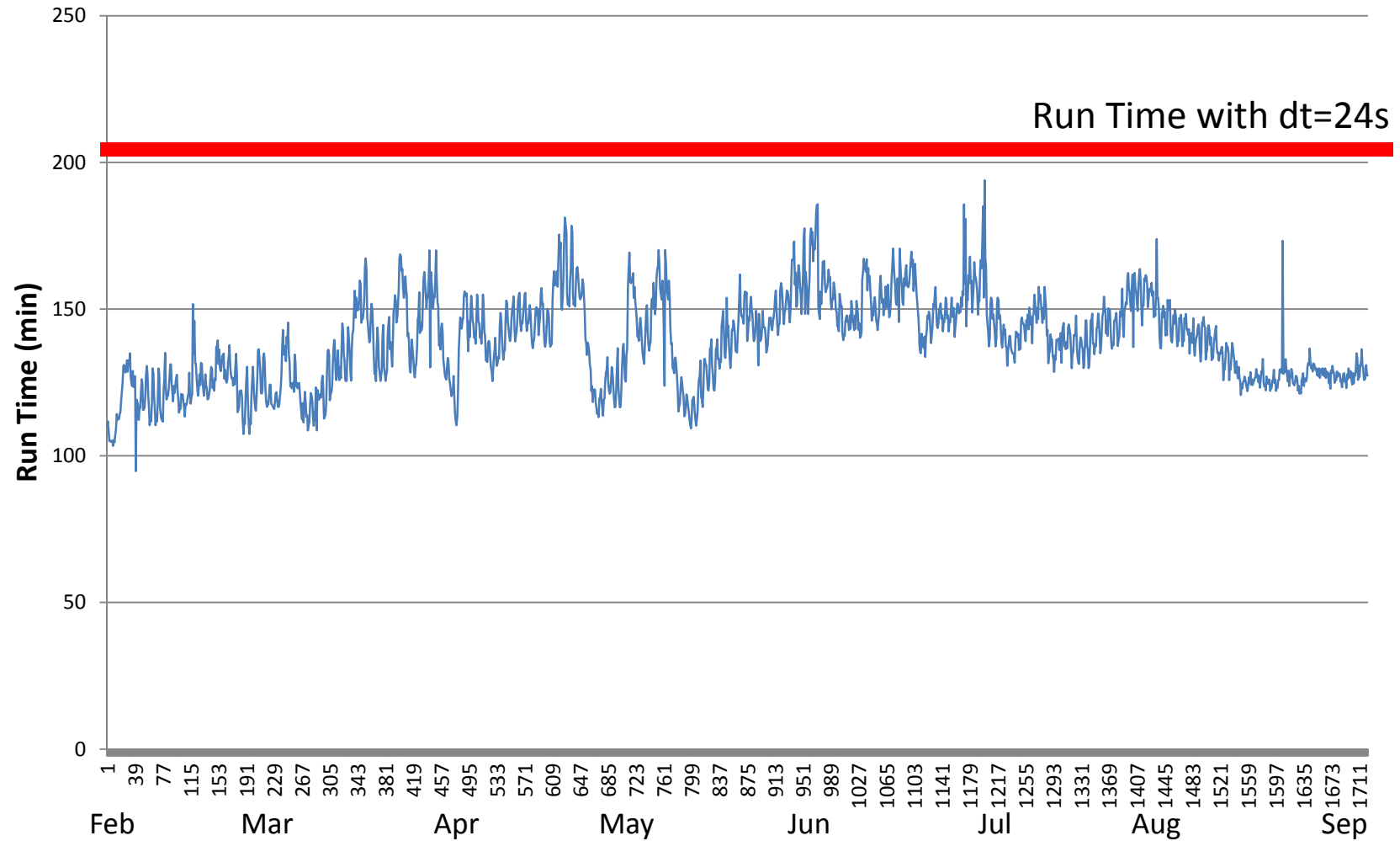- **Total Run-time reduced if average time-step exceeds static time-step**

WSI

# Time Step: Operational Run 26 Oct 2014



Time Step: 51 hr 4km North American Domain

# Run Time



Run-Time: 72 hour North Amer. 12/4/1.3k

# Implementation Details/Challenges

- **Look-ahead and adjust time-step to fall upon "history" and boundary input times**
  - Look ahead 2 steps to prevent very short steps
- **Nesting: Assure that nest and parent are sync'ed:**
  - Parent step adjusted to be a multiple of nest step
  - Let most expensive domain control the time-step
- **Very steep topography (Mt. St. Elias...):**
  - Have seen crash (on order of 1-2x per year out of 3000 model runs)
  - Recent adjustments to "look-ahead" with nesting may have resolved this

**WSI**°

# Summary/Findings

- **Performance:**
  - WRF runs faster over oceans than land
  - There's a lot of "noise" across the domain that needs to be better understood
  - Nesting Can be Expensive!
- **Adaptive time-step can be used to optimize performance**

**WSI**