# I/O and post-processing for High-resolution climate data

Dr. John Dennis
dennis@ucar.edu

# Motivation

- High-resolution climate generates a large amount of data!

# Outline

- PIO update and Lustre optimizations
- How do we analyze high-resolution climate data faster?
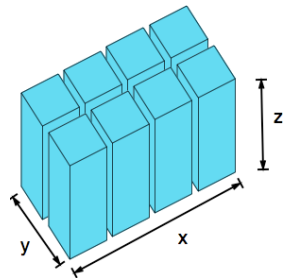- Wavelet compression for climate data

# PIO update and Lustre optimizations

# Parallel I/O library (PIO)

- Goals:
  - Reduce memory usage
  - Improve performance
- Writing a single file from parallel application
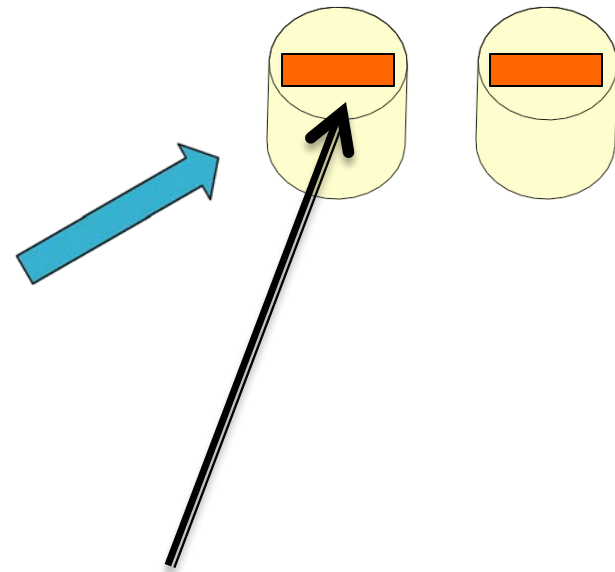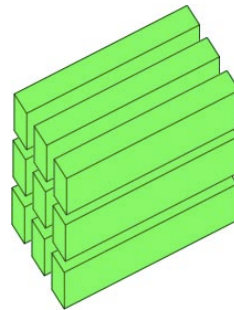  - Flexibility in  I/O libraries
  - MPI-IO,NetCDF3, NetCDF4, pNetCDF

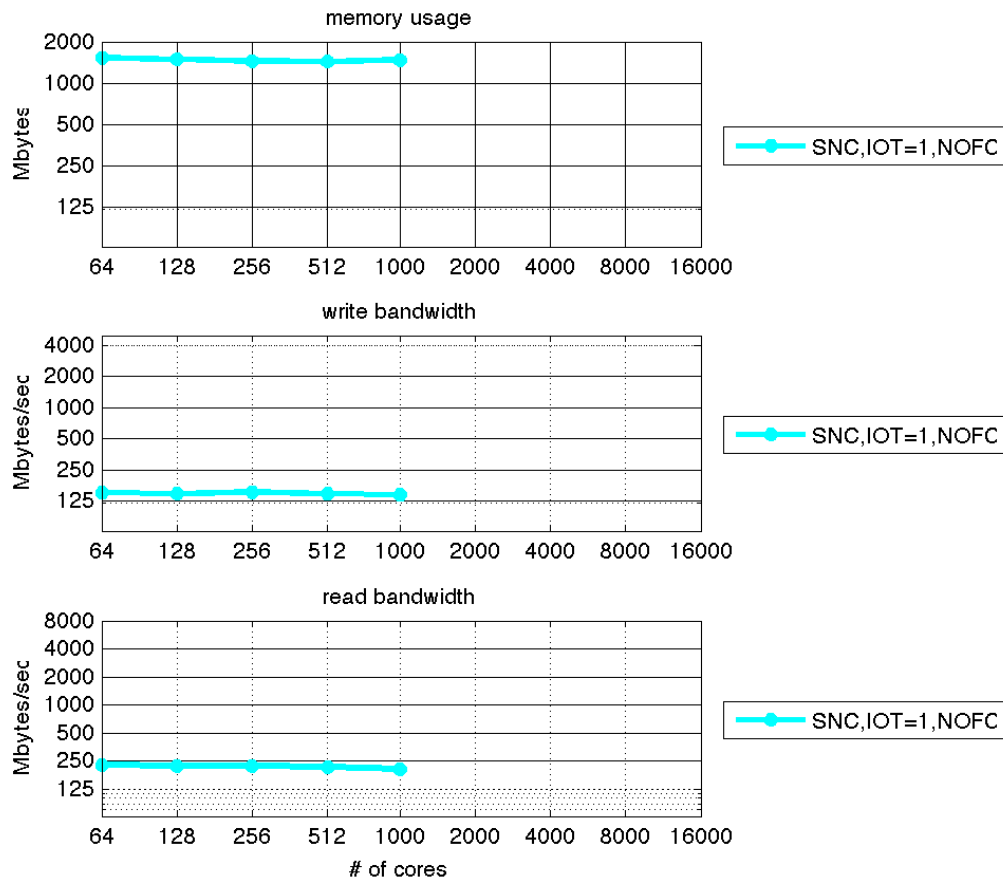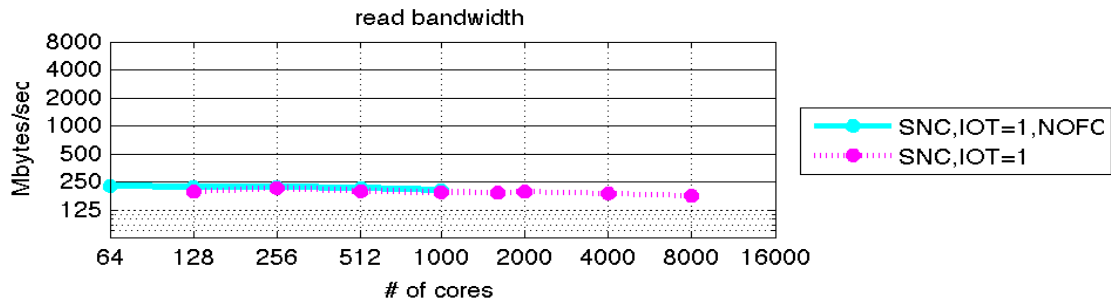# Optimizing writing data to Lustre file system
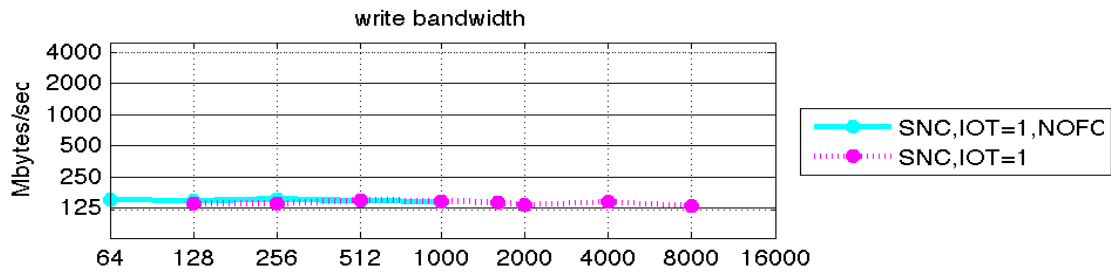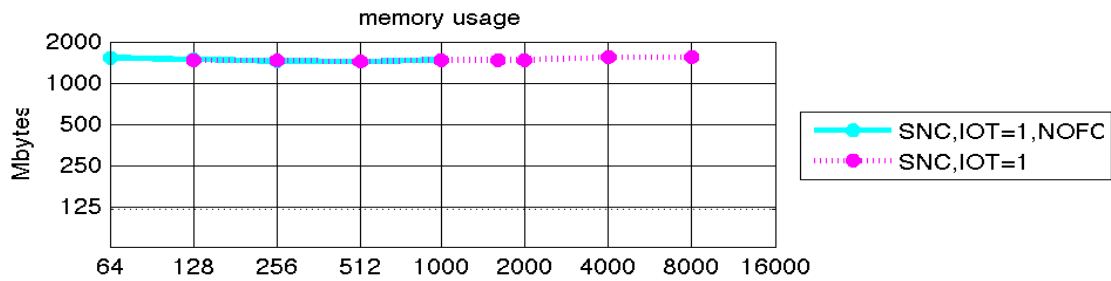
Computational decomposition

I/O decomposition
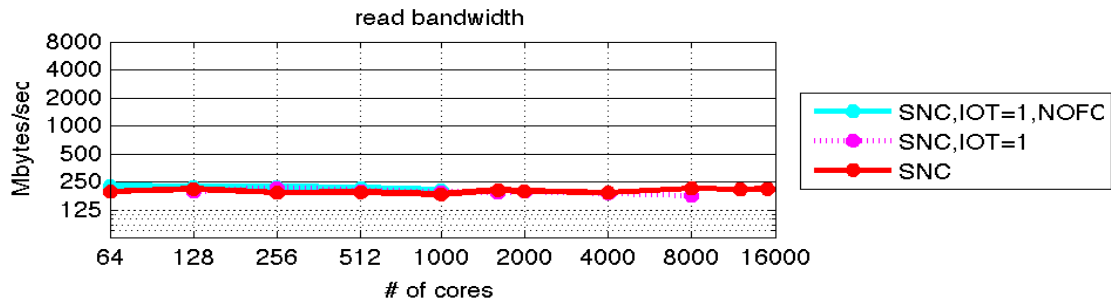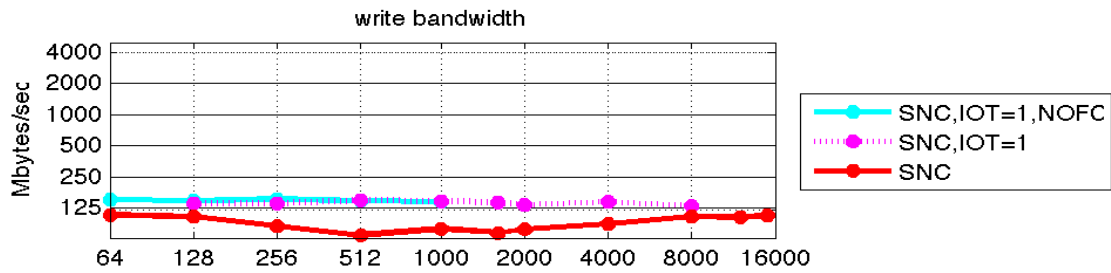
Rearrangement

Match I/O decomposition to Lustre stripe size
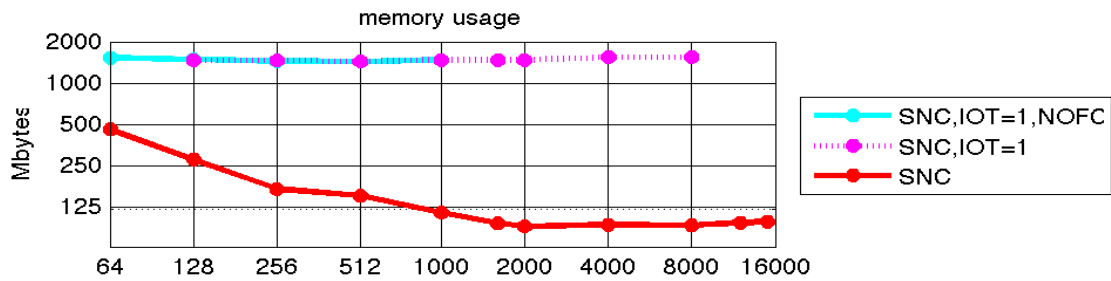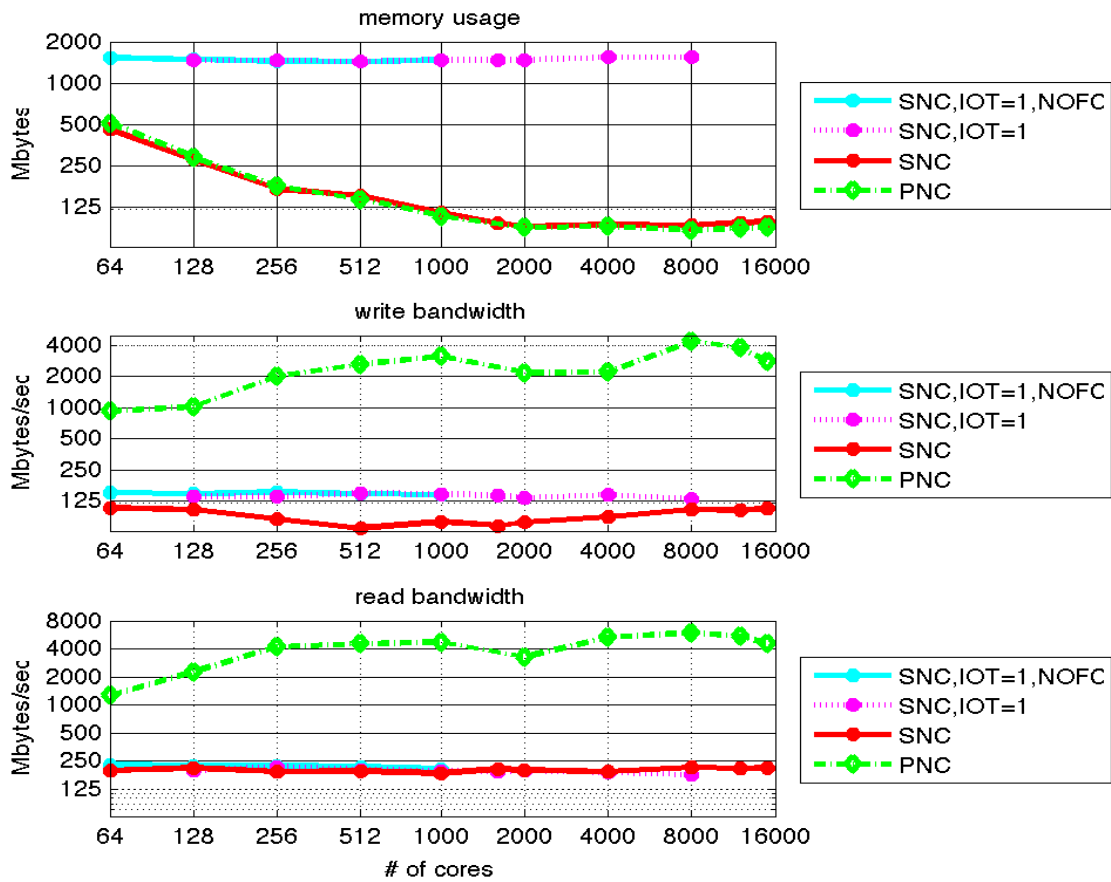
# Experimental setup

- Read/write 3D POP sized variable [3600x2400x40]
- 10 files, 10 variables per file, [max bandwidth]
- Using Kraken (Cray XT5) + Lustre filesystem
  ◦ Used 16 of 336 OST
- Impact of PIO features
  ◦ Flow-control
  ◦ Vary number of IO-tasks
  ◦ Different general I/O backends
- Did we achieve our design goals?

memory usage

write bandwidth

read bandwidth

# How do we analyze high-resolution climate data faster?

John Dennis, Matthew Woitaszek (NCAR)
Taleena Sines* (Frostburg State)

"Parallel high-resolution climate data analysis using Swift",
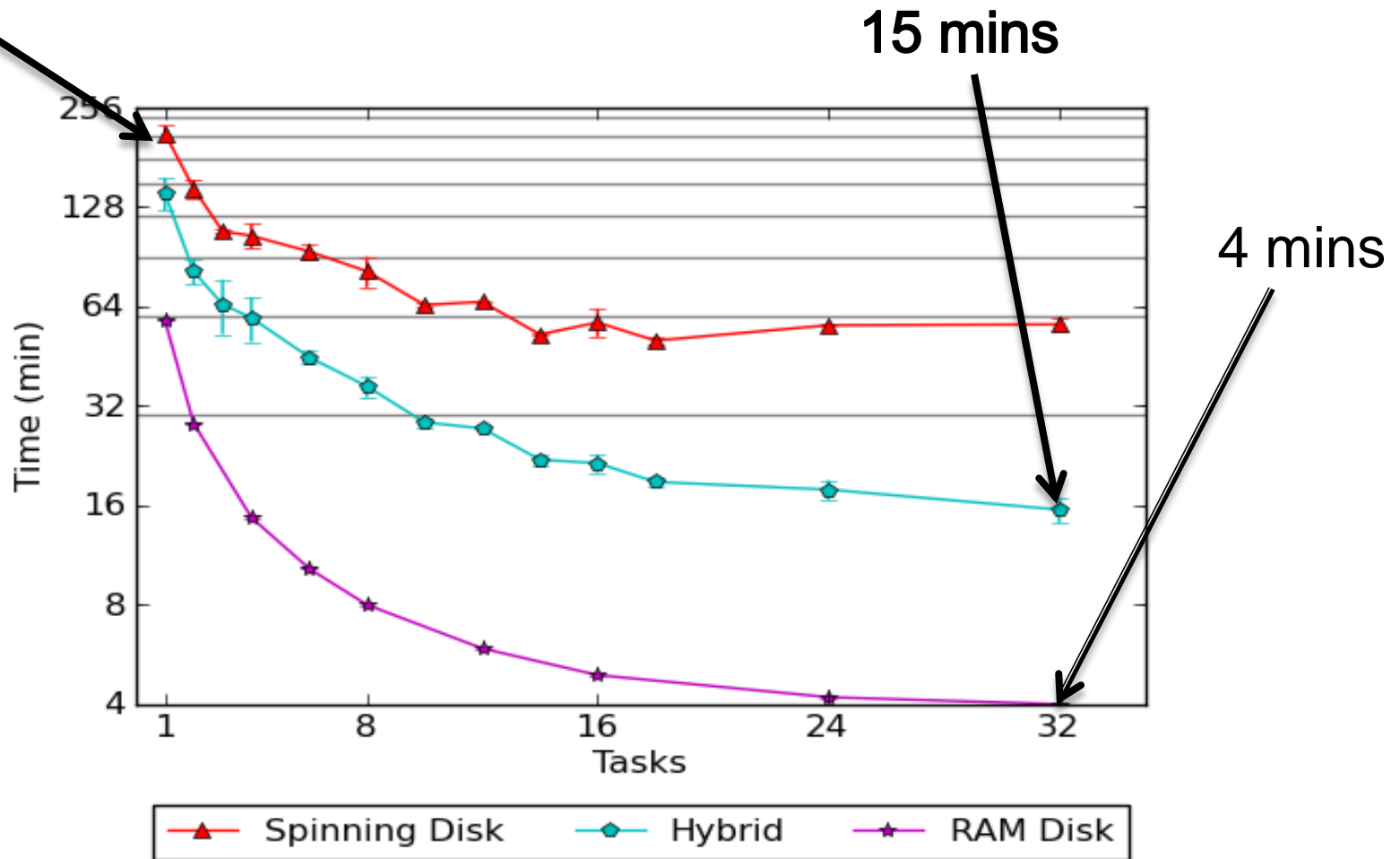Linux Clusters [under review]

\* SIParCS Intern

# Parallelizing diagnostics

- Used Swift a workflow language (UC/ANL)
  - Parallel scripting language
  - Data dependency driven
- Examine performance on several data-intensive architectures
  - Flash memory: Dash (SDSC)
  - SGI UV: Nautilus (NICS)
  - Large memory node: Polynya (NCAR)
    - 32 cores
    - 1 TB ram [ 512 GB memory/ 512 GB ramdisk]
    - 120 TB GPFS file-system (old hardware)

# Polynya: 0.5° /10yr

# Wavelet compression for Climate data

John Clyne, Yanick Polius, John Dennis (NCAR)
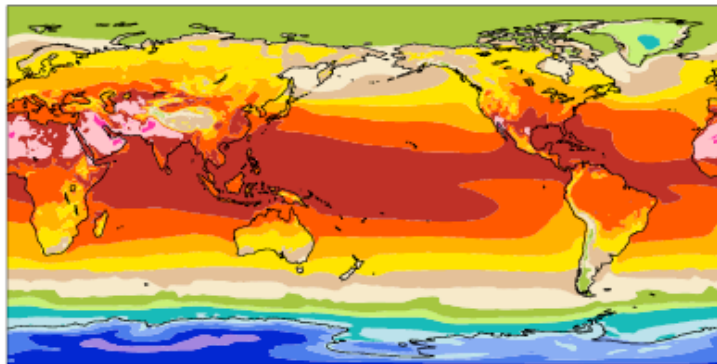
# Wavelet compression of Climate data

- Compression algorithm:
  - Apply wavelet transform to model outputs
  - Sort resulting wavelet coefficients by absolute magnitude
  - Discard smallest coefficients
  - Reconstruct an approximation of original data from remaining coefficients
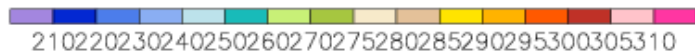- Compare original and reconstructed data using CESM AMWG Diagnostic Package

- ▸ Preliminary experiment
- ▸ 10 years of 0.5 deg CAM
- ▸ Compression
  - ◦ 2D variables – 2:1
  - ◦ 3D variables – 8:1
- ▸ Evaluate using student T-test
- ▸ Outcomes/Issues
  - ◦ Looks great!
  - ◦ Limited temporal variability
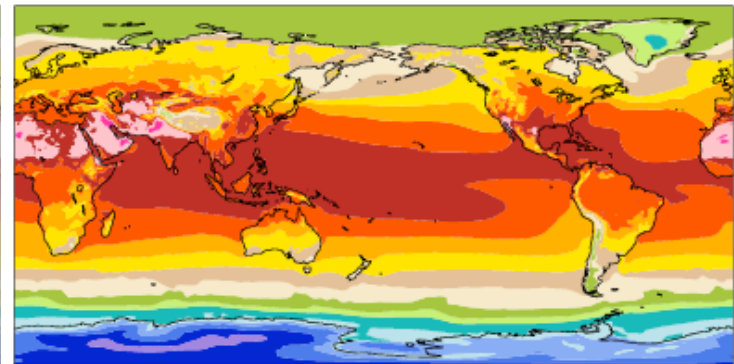  - ◦ Does not preserve zeros

# TS (Surface Temperature)



Surf Temp (radiative) mean= 290.09          K    Surf Temp (radiative) mean= 290.09          K

Min = 206.35 Max = 315.44          Min = 206.35 Max = 315.50

21022023024025026027027528028528529029530030510          21022023024025026027027528028528529029530030510

LRC01 - LRC01          T-test of the two means at each grid point

mean = −0.00    rmse = 0.07          K    Colored cells are significant at the 0.05 level

Min = −0.64 Max = 0.69

−12−10−8 −6 −4 −2 −1 0 1 2 4 6 8 10 12

# CLDTOT (Total Cloud)



Total cloud          mean=   44.17          percent

Min =    1.89 Max =   94.85

5  10 15 20 25 30 40 50 60 70 75 80 85 90 95

Total cloud          mean=   44.17          percent

Min =    1.89 Max =   94.84

5  10 15 20 25 30 40 50 60 70 75 80 85 90 95

LRC01 - LRC01

mean =    0.00     rmse =    0.26          percent

Min =   −2.11 Max =    2.46

−50−40−30−20−15−10−5  0  5  10 15 20 30 40 50
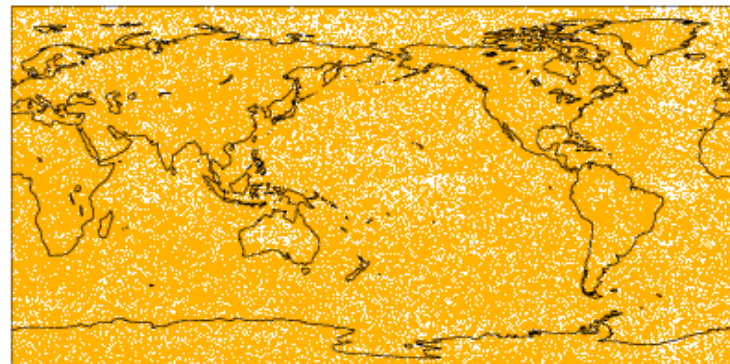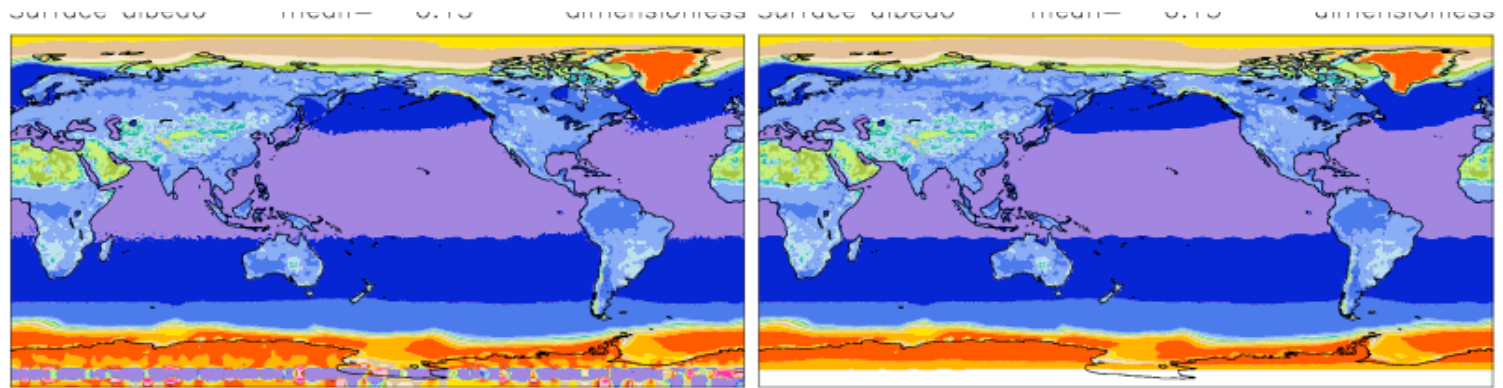
T-test of the two means at each grid point

Colored cells are significant at the 0.05 level

# ALBSURF (Surface Albedo)



Surface albedo    mean = 0.13    dimensionless    Surface albedo    mean = 0.13    dimensionless

Min = −419.87 Max = 0.99

0.050.10.150.20.250.3 0.4 0.5 0.6 0.70.750.80.850.90.95

Min = 0.04 Max = 0.85

0.050.10.150.20.250.3 0.4 0.5 0.6 0.70.750.80.850.90.95

LRC01 - LRC01

mean = −0.00    rmse = 0.42    dimensionless

T-test of the two means at each grid point
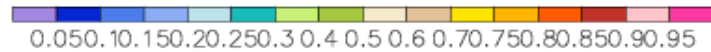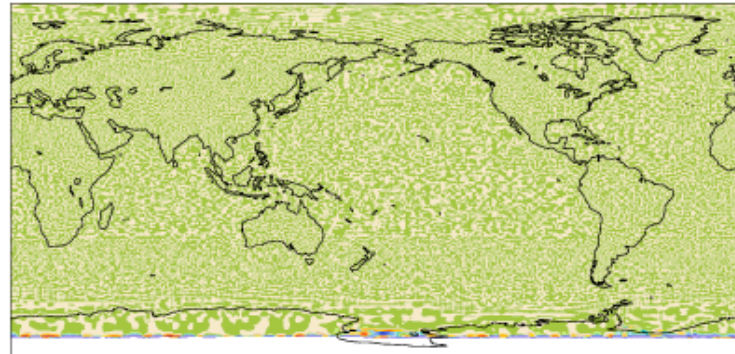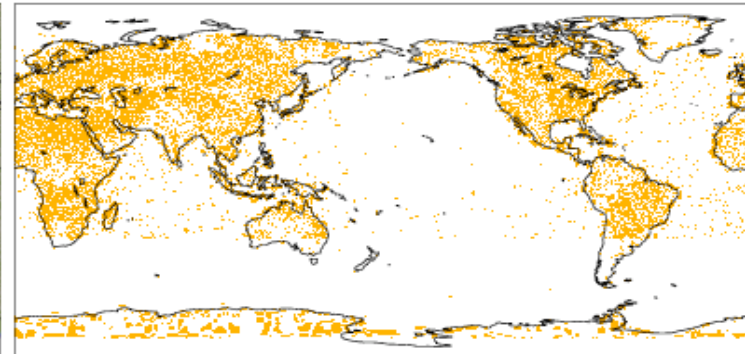
Colored cells are significant at the 0.05 level

Min = −272.14 Max = 0.22

−0.250.20.150.10.070.050.030.00.030.050.070.10.150.20.25

# Acknowledgements

- NCAR:
  - D. Bailey
  - F. Bryan
  - T. Craig
  - B. Eaton
  - J. Edwards [IBM]
  - N. Hearn
  - K. Lindsay
  - N. Norton
  - M. Vertenstein
- COLA:
  - J. Kinter
  - C. Stan
- U. Miami
  - B. Kirtman
- U.C. Berkeley
  - W. Collins
  - K. Yelick (NERSC)
- U. Washington
  - C. Bitz

- NICS:
  - M. Fahey
  - P. Kovatch
- ANL:
  - R. Jacob
  - R. Loy
- LANL:
  - E. Hunke
  - P. Jones
  - M. Maltrud
- LLNL
  - D. Bader
  - D. Ivanova
  - J. McClean (Scripps)
  - A. Mirin
- ORNL:
  - P. Worley

and many more…
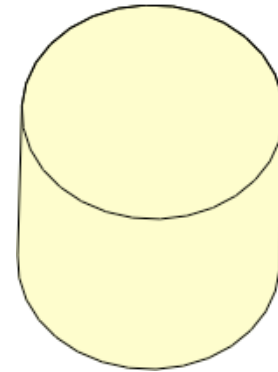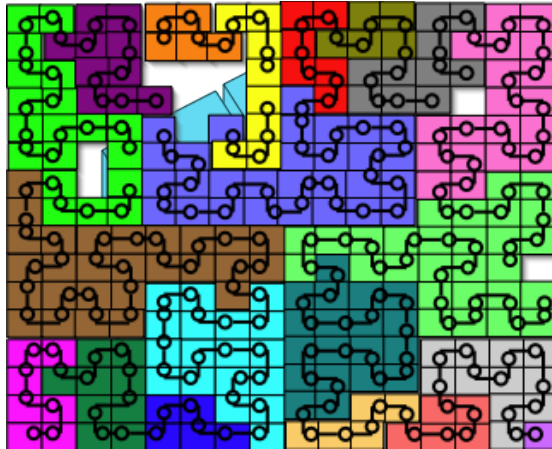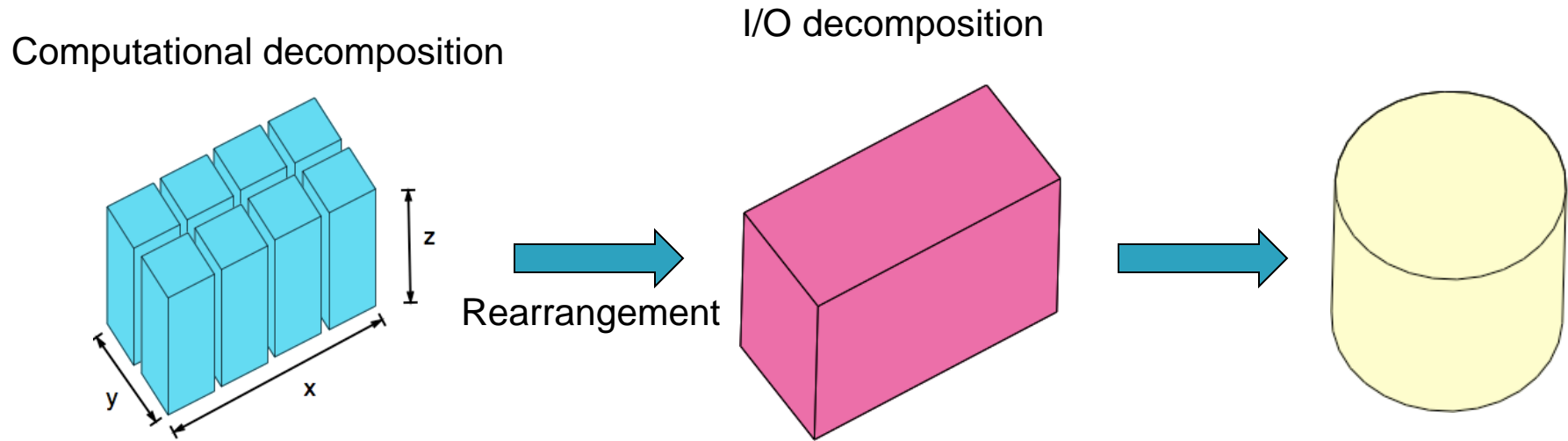
# Questions?

dennis@ucar.edu

# PIO:
# Writing distributed data (I)

Computational decomposition



+ Simple

+ Most versions of MPI-IO will do aggregation

- Computational decomposition may not be optimal for disk access

- pNetCDF requires block cyclic decompositions
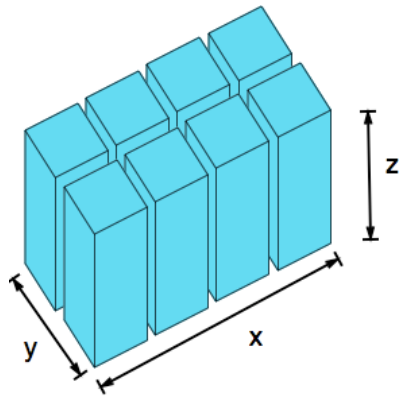
# PIO: Writing distributed data (II)

Computational decomposition

I/O decomposition



Rearrangement

+ Maximize size of individual io-op's to disk
- Non-scalable user space buffering
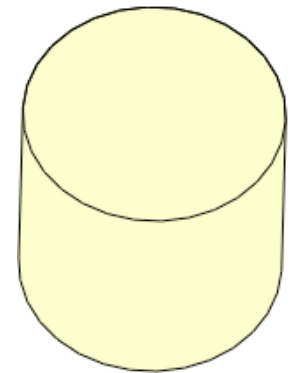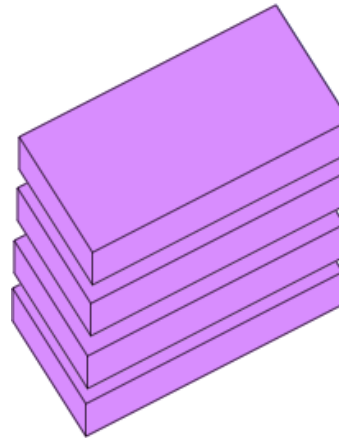- Very large fan-in → large MPI buffer allocations

# PIO: Writing distributed data (III)
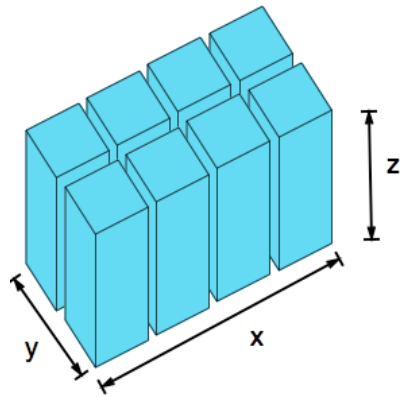
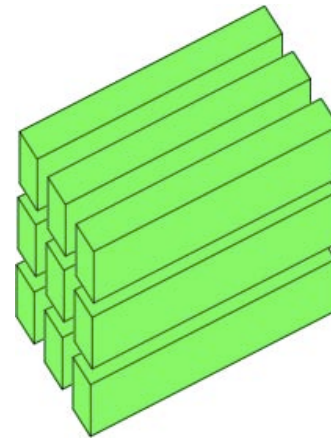Computational decomposition

I/O decomposition



Rearrangement

+ Scalable user space memory
+ Relatively large individual io-op's to disk
- Very large fan-in → large MPI buffer allocations
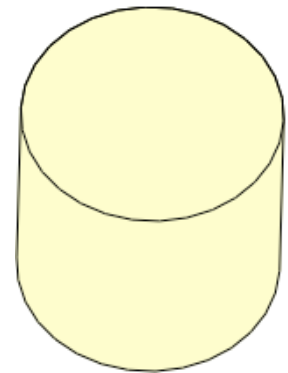
# PIO: Writing distributed data (IV)
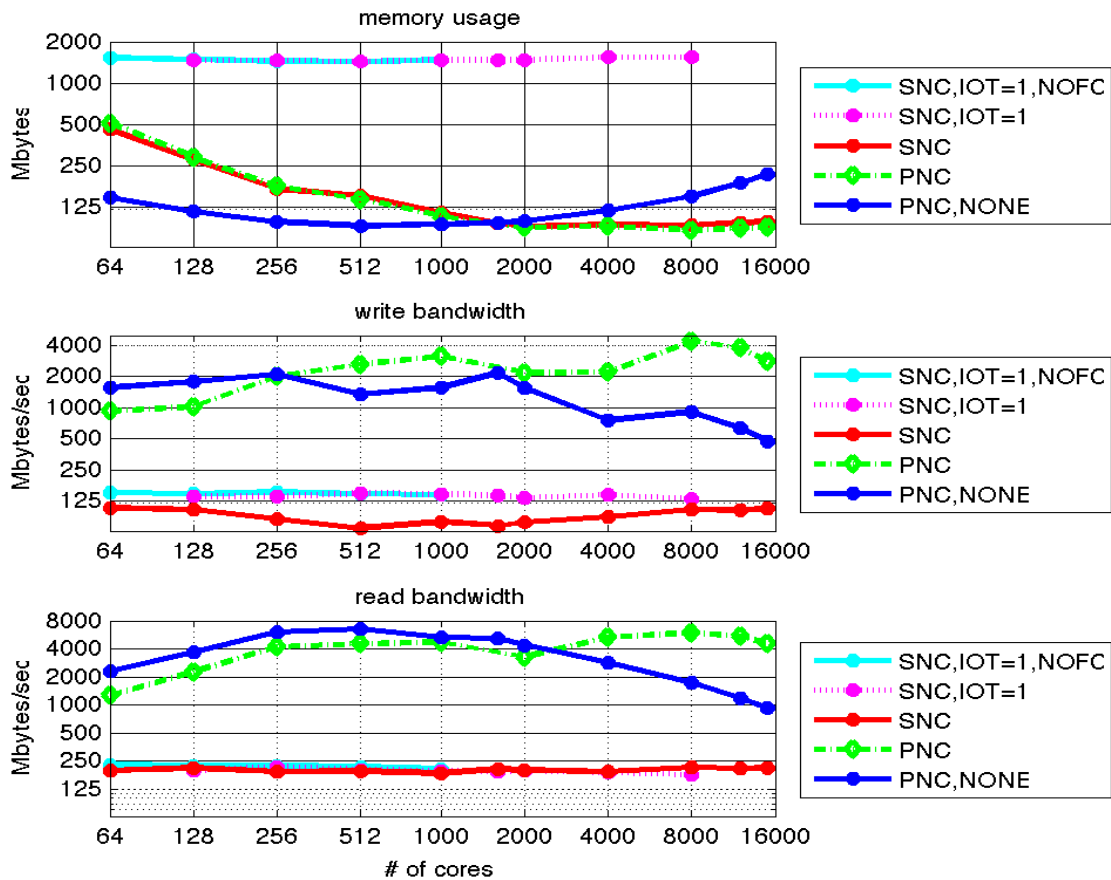
Computational decomposition

I/O decomposition

Rearrangement

+ Scalable user space memory
+ Smaller fan-in -> modest MPI buffer allocations
- Smaller individual io-op's to disk

# PIO Status

- Supported parallel I/O library in CCSM4 & CESM1 release
- Addition of Flow-control algorithms (Worley)
- Initial documentation using Doxygen
- Small but growing user base
  - ESMF
  - VAPOR + wavelet compression