

# Reading an r-file material model in SW4 using a parallel file system

USGS workshop, Menlo Park, CA

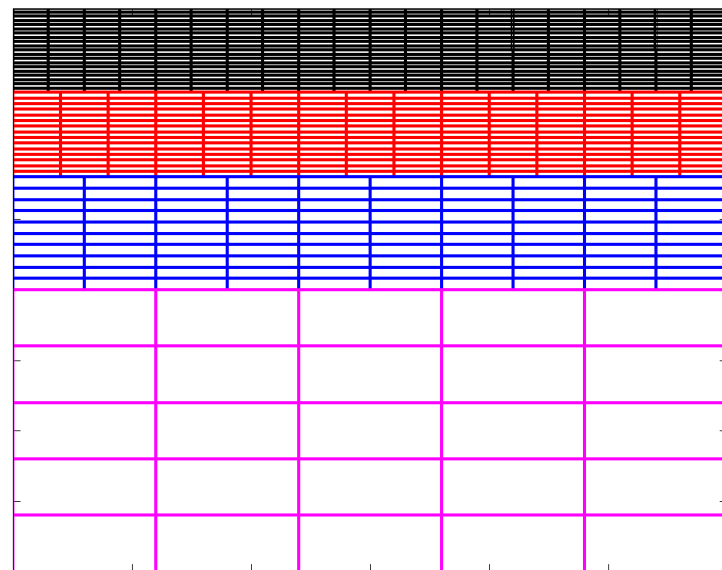
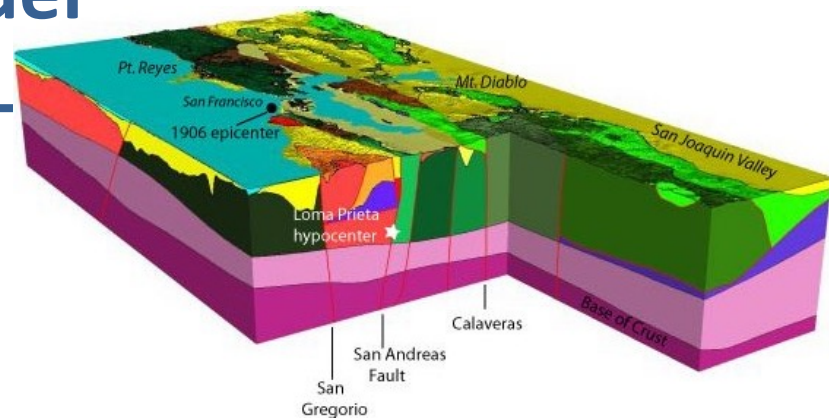
N. Anders Petersson, Bjorn Sjogreen

March 2018



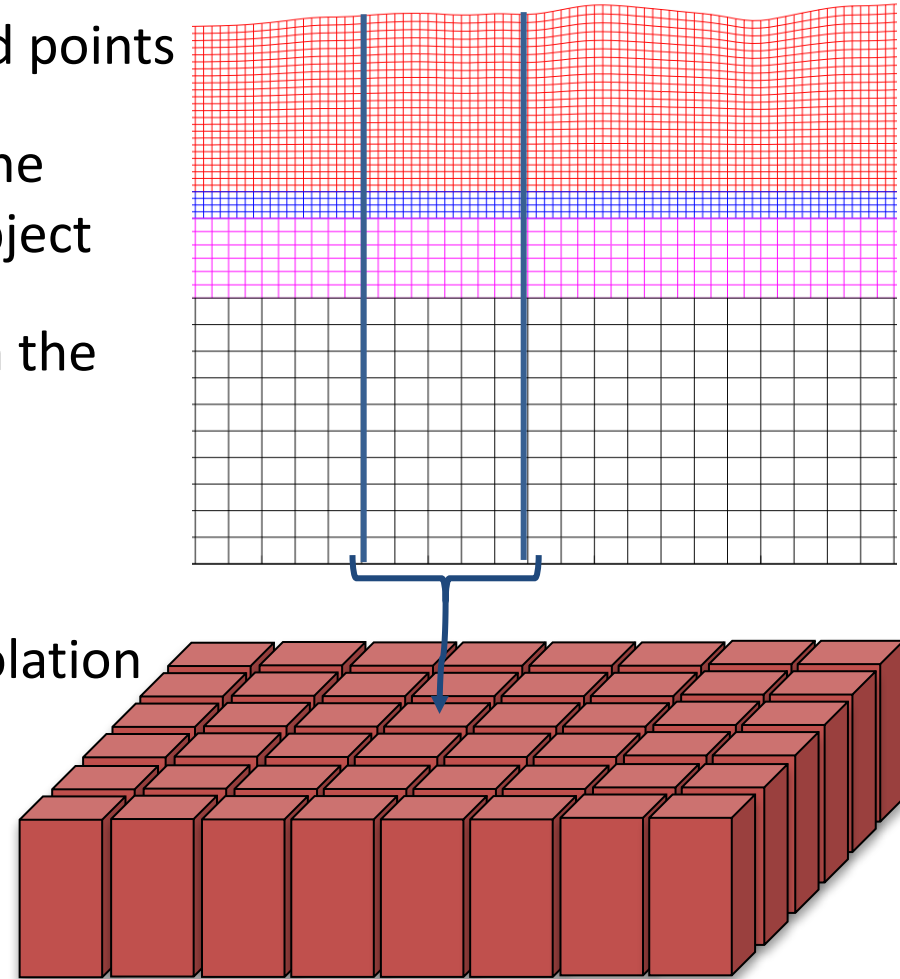
# USGS Bay Area Material Model

- Oct-tree variable resolution cell-centered data structure
- 8 Gbyte file (regional)
- Proj4 cartographic mapping (lon,lat,depth)-> (x,y,z)
- cencalvm: Point-wise query functions in C++
- R-file binary format developed in 2014
  - Byte-swapping for IBM BG/Q system at LLNL
  - Faster access; parallel file-system
- Header + rectangular blocks of vertex-centered data
- Same data as in the USGS e-tree model
- Read ( $N_c \times N_i \times N_j \times N_k$ ) blocks of data



# In SW4, only a subset of MPI-tasks read the r-file from the parallel file system

- Each MPI-task owns a “pencil” of grid points
- Each MPI-task requests a subset of the material model from a Parallel\_IO object
- Only a subset of MPI-tasks read from the parallel system (8-128 readers)
- Material data is distributed by MPI
- Grid point values by tri-linear interpolation



# Example in 2D, 6x6 matrix on 9 MPI-tasks

(1, 1) (1, 2) 1 (2, 1) (2, 2)	(1, 3) o 2 o o	o (1, 6) 3 o o
(3, 1) o 4 (4, 1) o	o o 5 o o	o o 6 o o
(5, 1) o 7 (6, 1) (6, 2)	o o 8 o o	o o 9 o (6, 6)

On disk in column order: (1,1), (2,1), (3,1)... (6,1), (1,2), (2,2),  
... (6,6).

# Let 2 MPI-tasks be designated IO processors

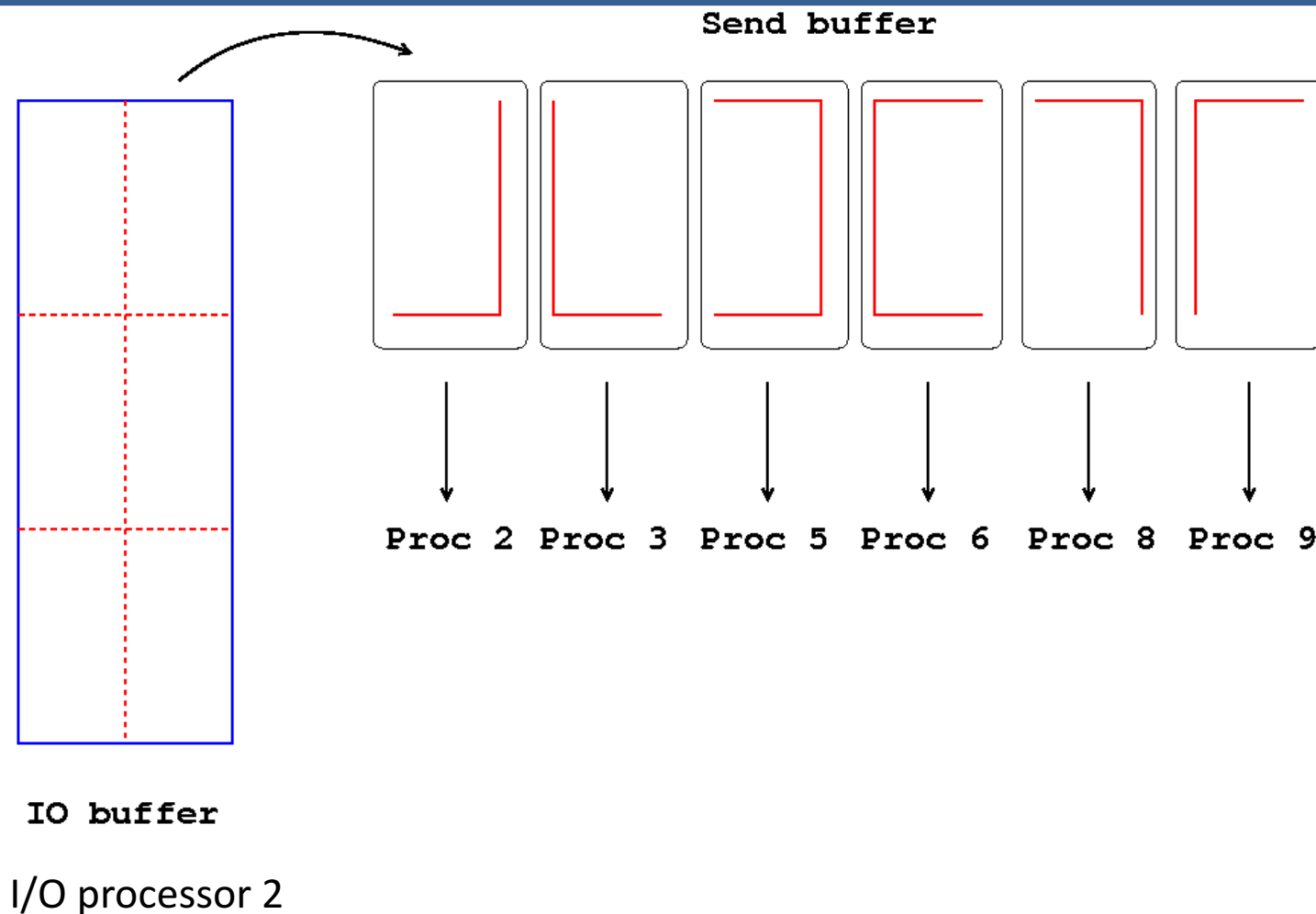
(1, 1) (1, 2)	(1, 3)	○	○ (1, 6)
(2, 1) (2, 2)	○	○	○ ○
(3, 1) ○	○	○	○ ○
(4, 1) ○	○	○	○ ○
(5, 1) ○	○	○	○ ○
(6, 1) (6, 2)	○	○	○ (6, 6)

1

2

Designated I/O processors

# Each IO-processor reads data into a buffer and sends the requested subsets with MPI



# Summary

- R-file format holds the same data as the USGS e-tree format
- Each MPI-task requests a subset of the material model
- 8-128 I/O processors read data from the parallel file system
- ~5min to read and initialize 60 billion grid point model on 2,048 nodes (131,072 MPI-tasks, 64 readers)
- Parallel IO routines originally developed for CFD applications:
  - B.Sjogreen, H.C.Yee, M.J.Djomehri, A.Lazanoff, and W.D.Henshaw (2010) "Parallel performance of ADPDIS3D - A high order multiblock overlapping grid solver for hypersonic turbulence", in "Parallel Computational Fluid Dynamics", R.Biswas ed., DEStech Publications.

