

Dealing with Large Datasets

or, “So I have 40TB of data..”

Jonathan Dursi, SciNet/CITA, University of Toronto

Data is getting bigger

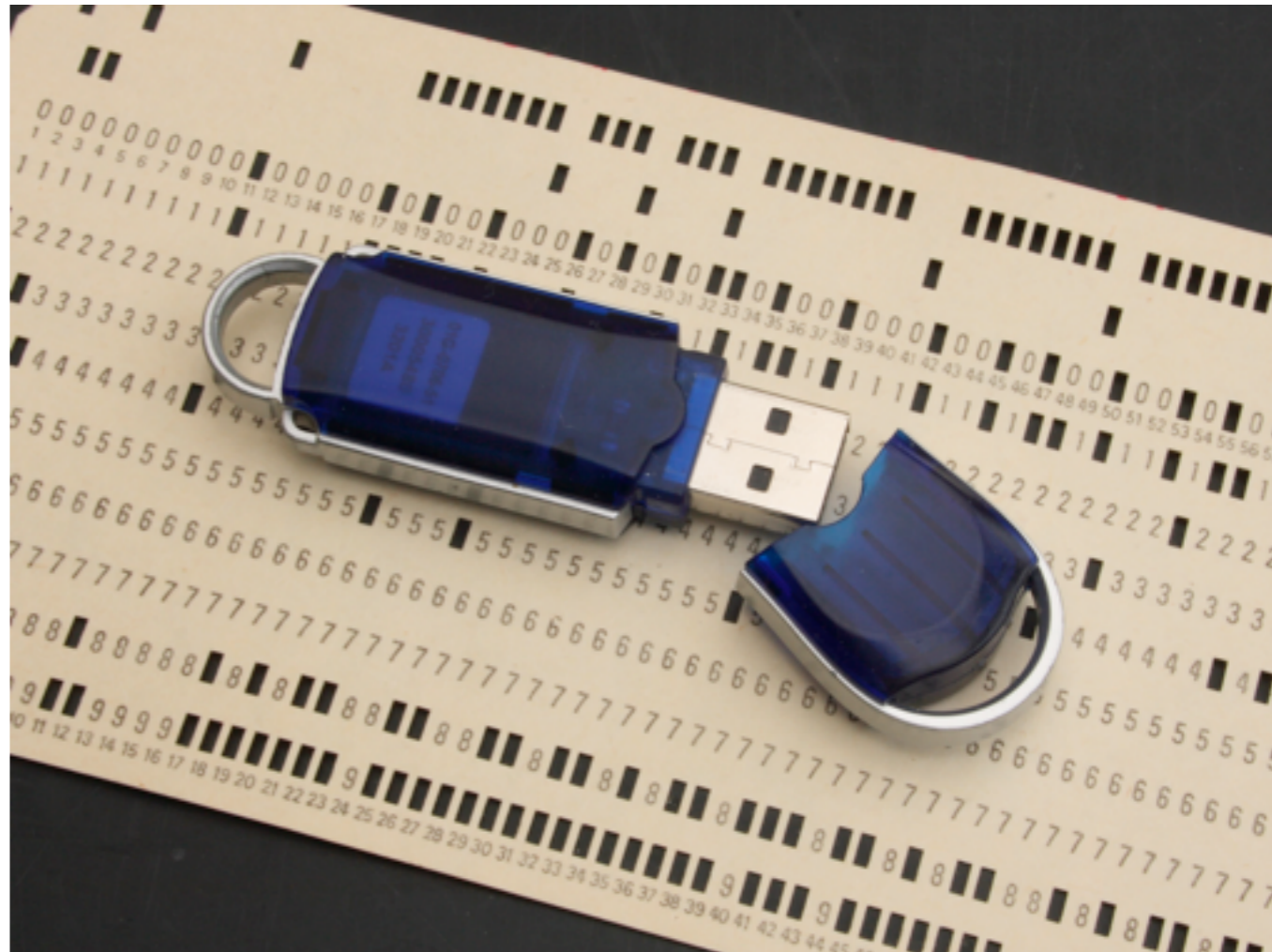
- Increase in computing power makes simulations larger/more frequent
- Increase in sensor technology makes experiments/observations larger
- Data sizes that used to be measured in MB/GB now measured in TB/PB.
- Easier to make big data than to do something useful with it!



Economist, 27 Feb 2010

What is “big data”?

- Absolute numbers don't matter
 - (and change rapidly anyway)
- Big is defined by its effects on us the scientists
 - What techniques must we use to analyze it?
- Two big milestones that define big..

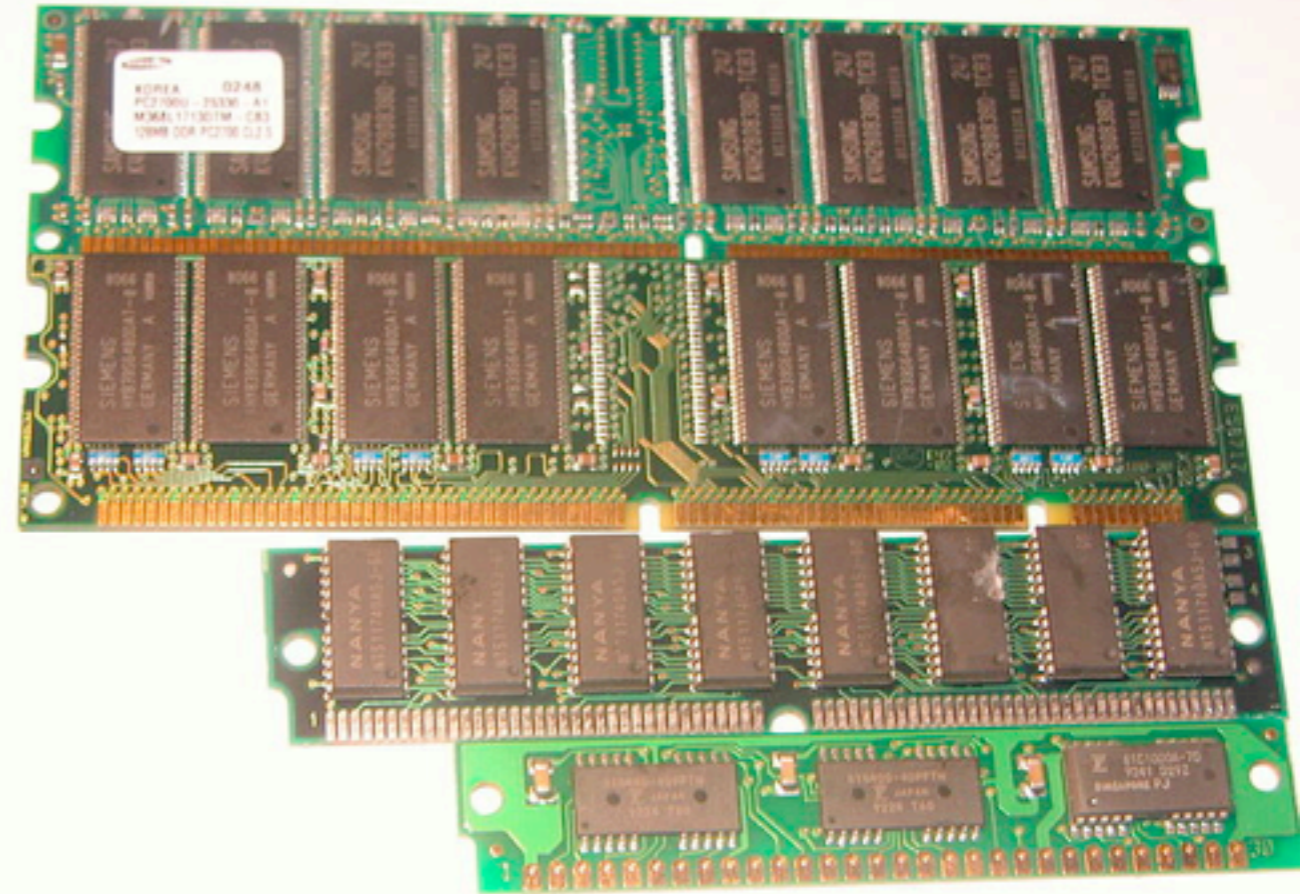


Ian-S,

<http://www.fotopedia.com/items/flickr-2152798588>

Big #1: Too Big to fit in Memory

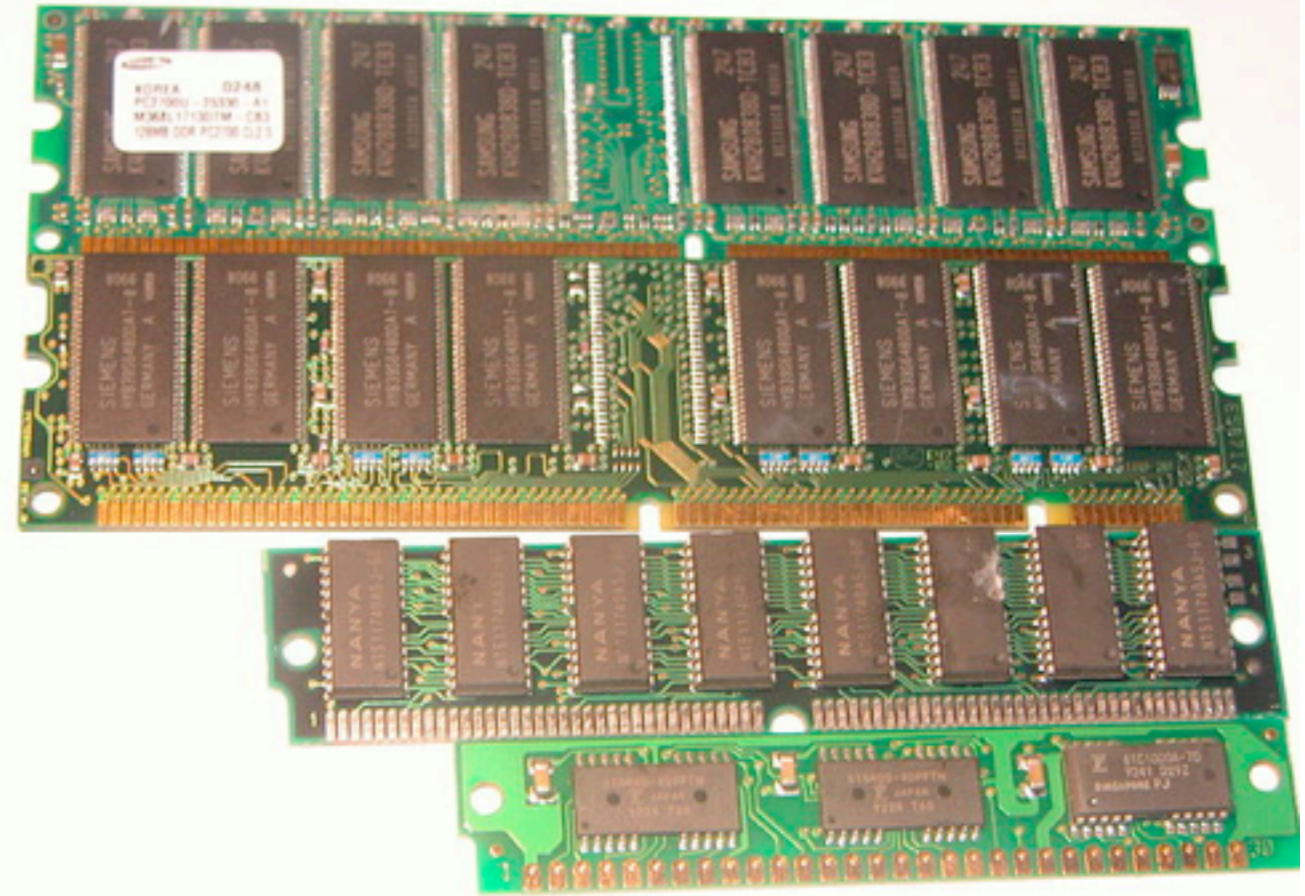
- If fits in memory,
 - easy global view of whole problem
 - simple workflow



<http://commons.wikimedia.org/wiki/File:Kinds-of-RAM.JPG>

Big #1: Too Big to fit in Memory

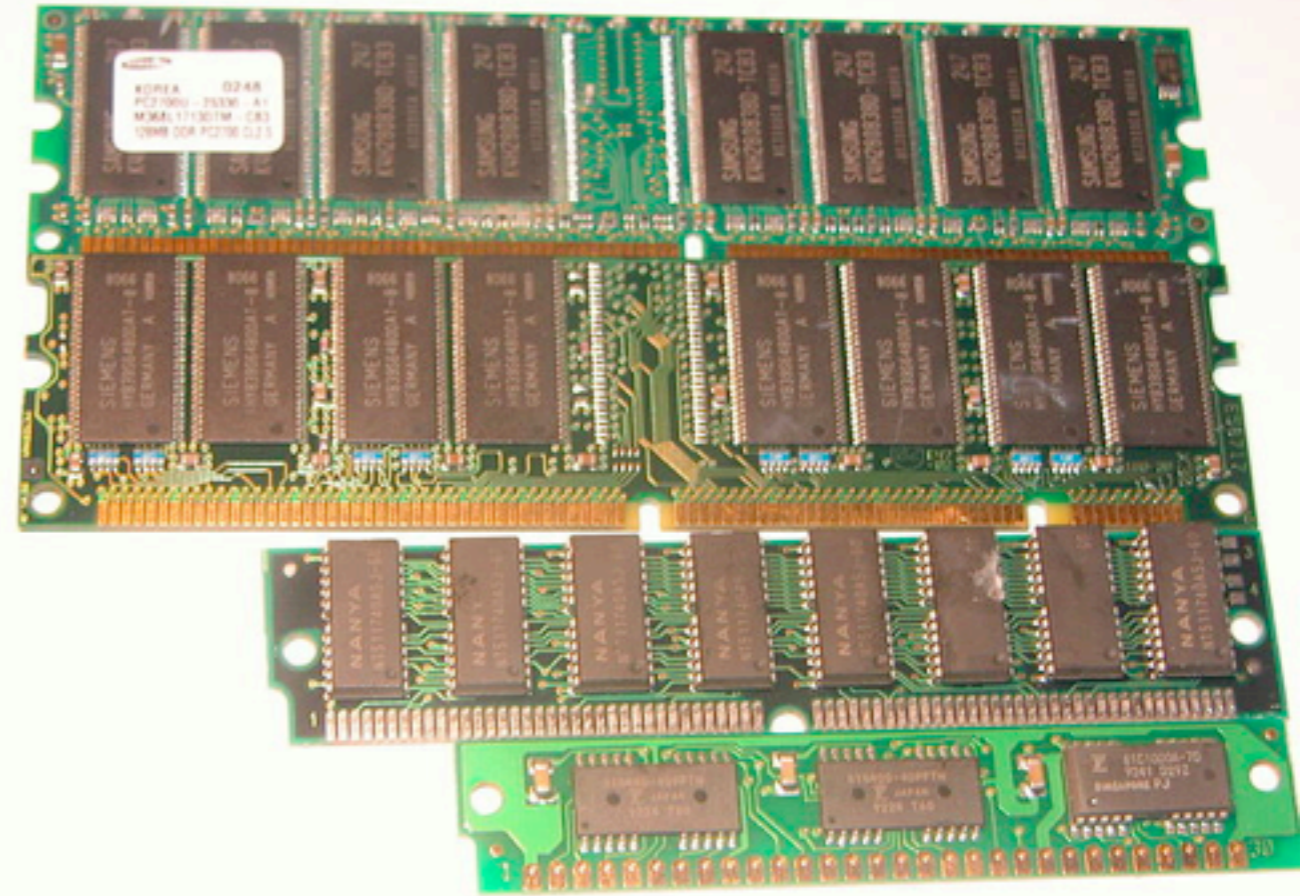
- Otherwise, must use more complicated techniques
 - Out of core
 - Multi-resolution
 - Parallel computation



<http://commons.wikimedia.org/wiki/File:Kinds-of-RAM.JPG>

Big #1: Too Big to fit in Memory

- Today:
 - ~2-16GB for workstation
 - ~128GB-256GB on specialized (shared) machine.



<http://commons.wikimedia.org/wiki/File:Kinds-of-RAM.JPG>

Big #2: Too Big to fit on one disk

- Once data size becomes comparable to typical storage medium, I/O becomes significant limitation
- Hardware has to be considered (RAID, parallel file systems)
- Almost certainly need parallel computing.



<http://commons.wikimedia.org/wiki/File:Hard-drive.jpg>

<p>2 Terabytes</p> <p>\$100 Hard drive</p>	<p>20 Terabytes</p> <p>medium-sized HPC simulation; 128 GPC nodes, 10 outputs.</p>	<p>30 Terabytes</p> <p>raw data for 30x short-read sequence of human genome</p>
<p>150 Terabytes</p> <p>Size of MERRA climate database of 1979-current obs. data</p>	<p>330 Terabytes</p> <p>Data output by LHC every week</p>	<p>1 Petabyte:</p> <p>LSST each ~month, must be searched in ~ real time for transient events</p>

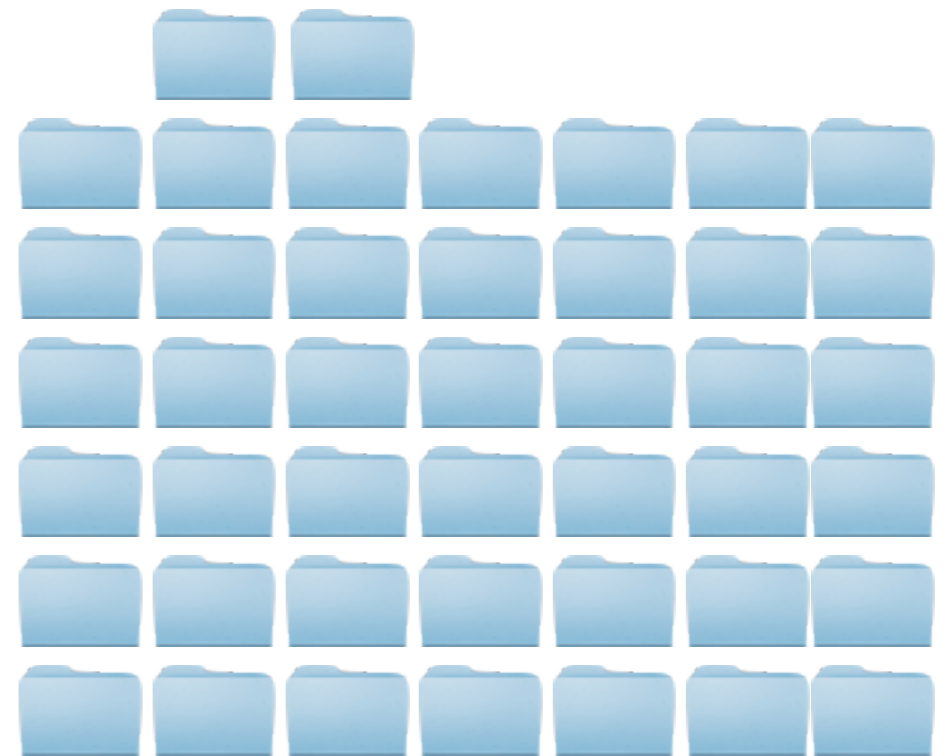
After Marian Bantjes, Wired magazine, Jun 2008

Big Data can be *big, or many,* data sets

- Few, very large, simulation outputs
- Very many, comparably small, experimental data
- Both can quickly add to TB
- Surprisingly, advice/techniques for these two cases overlap strongly.



vs



For big data, *scalability* is critical

- Many tools work great on desktop-sized datasets,
- But fail utterly at large scale:
 - on enormous files
 - when same task must be repeated on thousands of files.
- Need to use tools, techniques that handle scale well.



Not all tools will work
under arbitrarily large loads.

Scalability requires many things:

- Scalable I/O strategy
- Scalable data management
- Automatability / scriptability
- Scalable analysis flow
- Scalable tools
- Least scalable link in the workflow will bottleneck **entire process.**

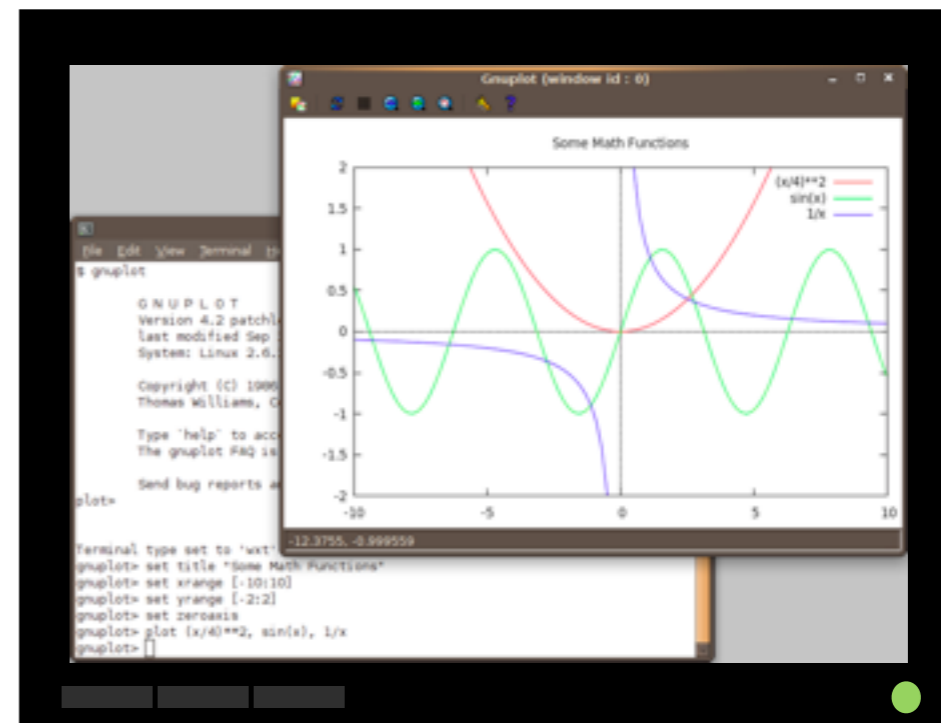
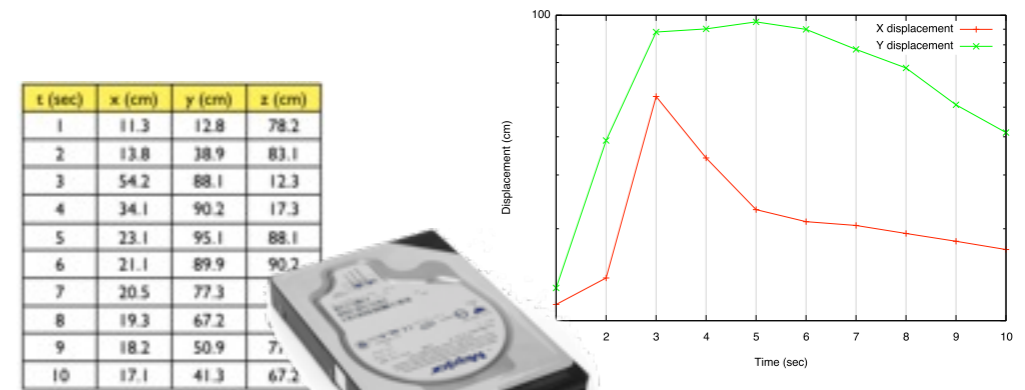


Jessie Menn

http://www.flickr.com/photos/jesse_menn/3164739580/

At scale, must plan how to do Input/Output

- Just reading and writing 40 TB from one drive at a time takes ~2weeks!
- 50MB/s read, 100MB/s write
- ..And that's **best** case
- Before planning analysis, planning *flow of data* important.



Where, how you do I/O matters.

- Binary - smaller files, *much* faster to read/write.
- You're not going to read 40TB of data yourself; don't bother trying.
- **Data:** machine-readable
- **Output:** graphs, summary tables - human-readable

Large Parallel File System	
ASCII 173s	binary 6s
Ramdisk	
ASCII 174s	binary 1s
Typical work station disk	
ASCII 260s	binary 20s

Timing data: writing 128M
double-precision numbers

Where, how you do I/O matters.

- All disk systems do best when reading/writing large, contiguous chunks
- I/O operations (IOPS) are themselves expensive
 - moving around within a file
 - opening/closing
- Seeks - 3-15ms - enough time to read 0.75 MB!

Typical work station disk

binary - one large read

14s

binary - 8k at a time

20s

binary - 8k chunks, lots of seeks

150s

binary - seeky + open and closes

205s

Timing data: reading 128M
double-precision numbers

Where, how you do I/O matters.

- RAM is much better for random accesses
- Use right storage medium for the job!
- Where possible, read in contiguous large chunks, do random access in memory
- Much better *if* you use most of data read in

Large Parallel File System

ASCII	binary
173s	6s

Ramdisk

ASCII	binary
174s	1s

Typical work station disk

ASCII	binary
260s	20s

Ramdisk

binary - one large read
1s

binary - 8k at a time
1s

binary - 8k chunks, lots of seeks
1s

binary - seeky + open and closes
1.5s

Parallel I/O and large file systems

- Large disk systems featuring many servers, disks
- Can serve files to many clients concurrently
- Parallel File Systems -
 - Lustre, Panasas, GlusterFS, Ceph, GPFS...



SciNet ~2k drives



Where, how you do I/O matters.

- Well built parallel file systems can greatly increase bandwidth
- But typically even worse penalties for seeky/IOPSy operations.
- Parallel FS can help with big data in two ways

Large Parallel File System

ASCII **binary**

173s 6s

Ramdisk

ASCII **binary**

174s 1s

Typical work station disk

ASCII **binary**

260s 20s

Large Parallel File System

binary - one large read

7.5s

binary - 8k at a time

62 s

binary - 8k chunks, lots of seeks

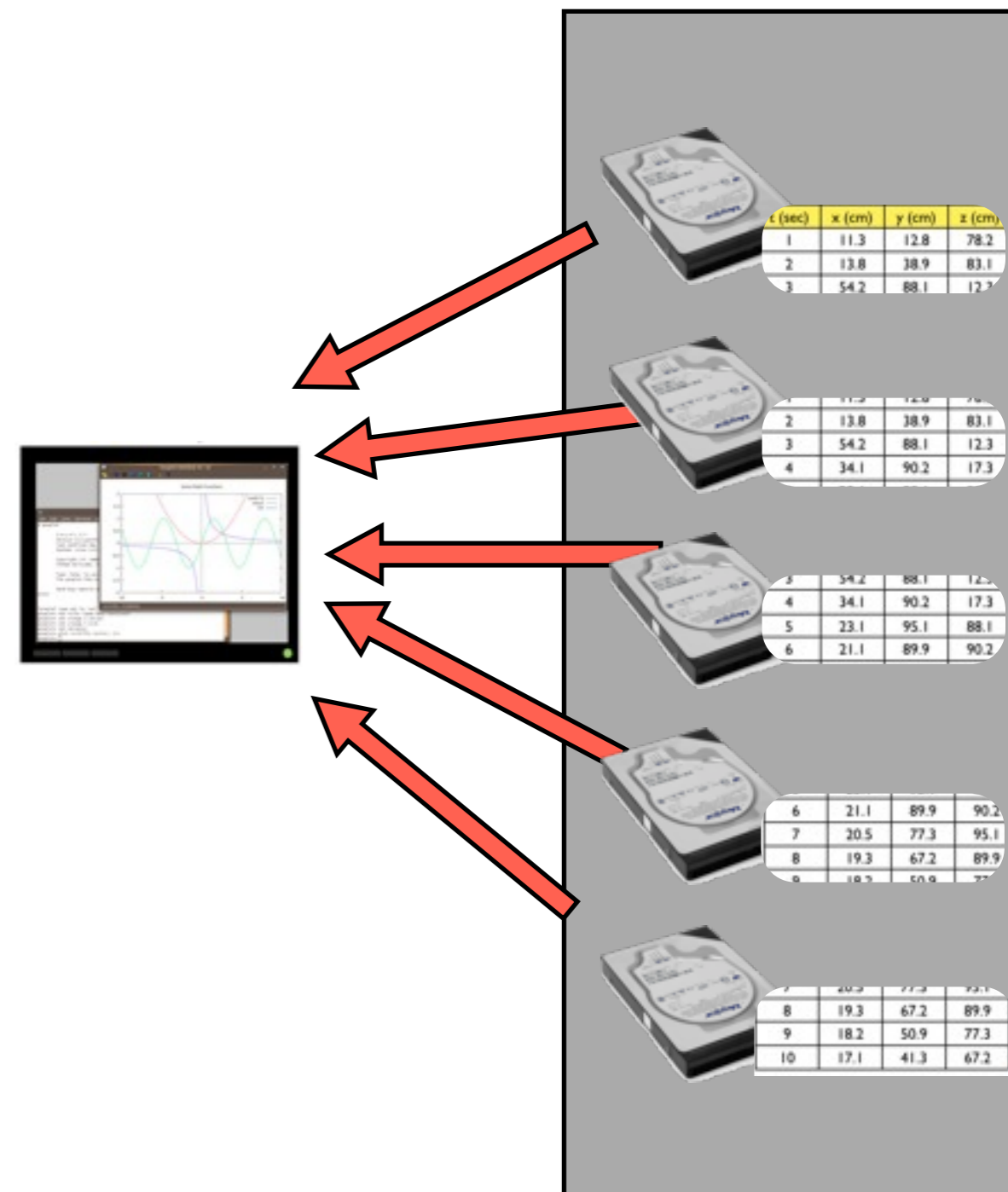
428 s

binary - seeky + open and closes

2137 s

Striping data across disks

- Single client can make use of multiple disk systems simultaneously
- “Stripe” file across many drives
- One drive can be finding next block while another is sending current block



Parallel operations on separate data

- Or can do truly parallel operations
 - multiple clients doing independent work
- Easy parallelism (good for lots of small data) - process many small files separately

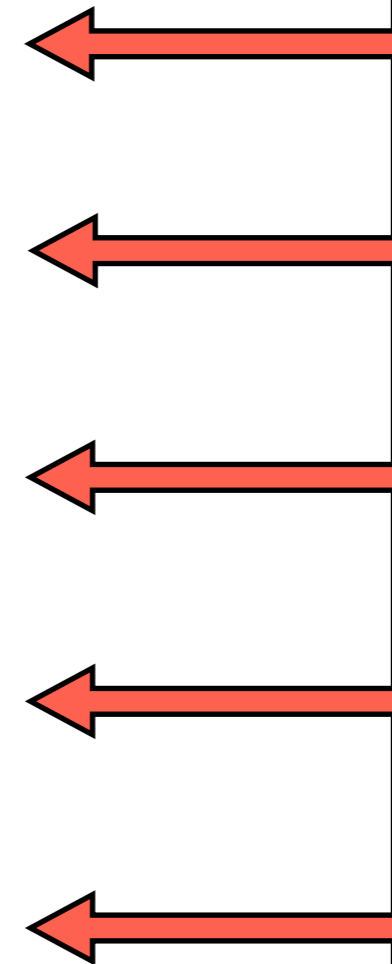


Parallel FS



Parallel operations on separate data

- Or can do truly parallel operations
 - multiple clients doing independent work
- Easy parallelism (good for lots of small data) - process many small files separately
- Harder parallelism - each does part of a larger analysis job on a big file.

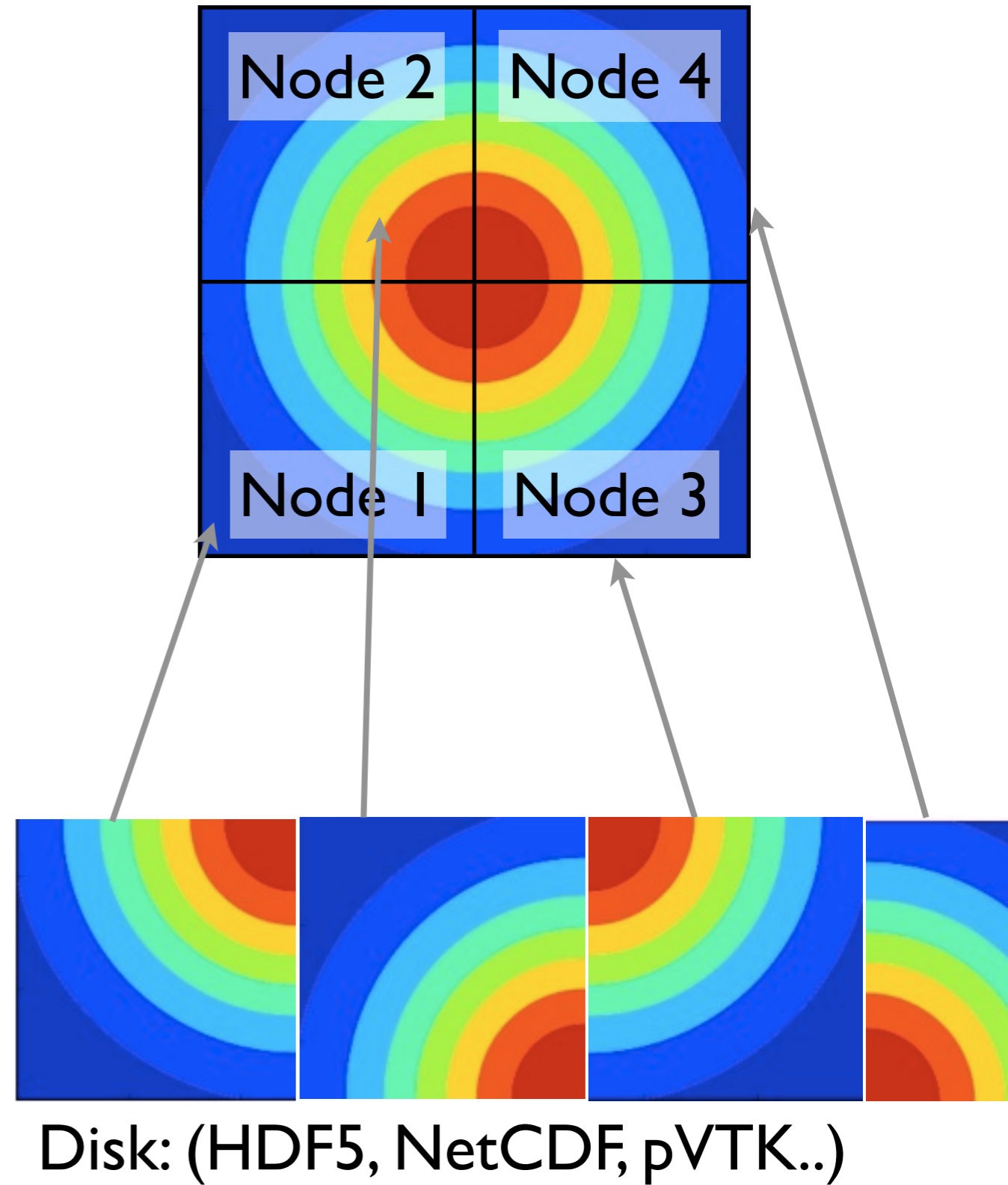


Parallel FS



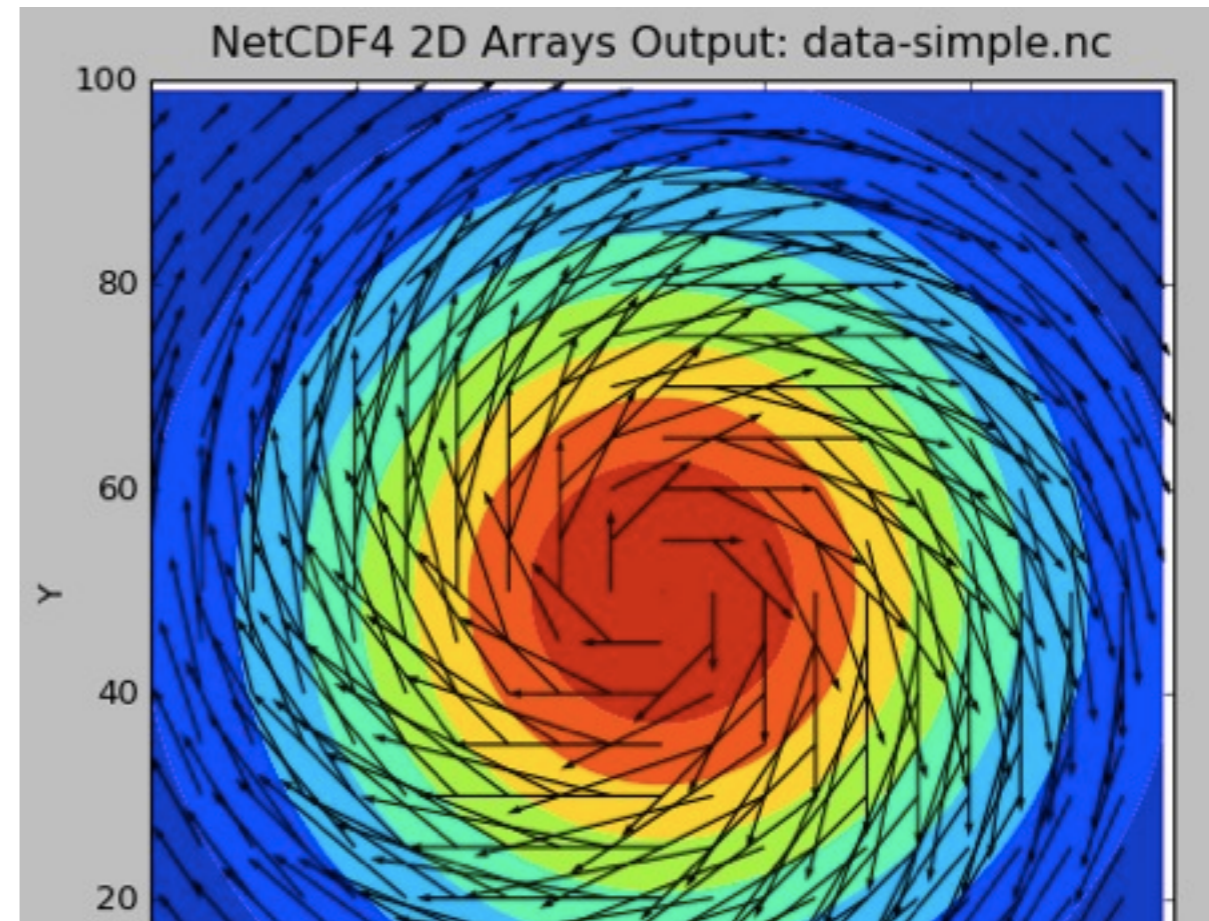
Data files must take advantage of parallel I/O

- For parallel operations on single big files, parallel filesystem isn't enough
- Data must be written in such a way that nodes can efficiently access relevant subregions
- HDF5, NetCDF formats typical examples for scientific data



These formats are *self-* *describing*

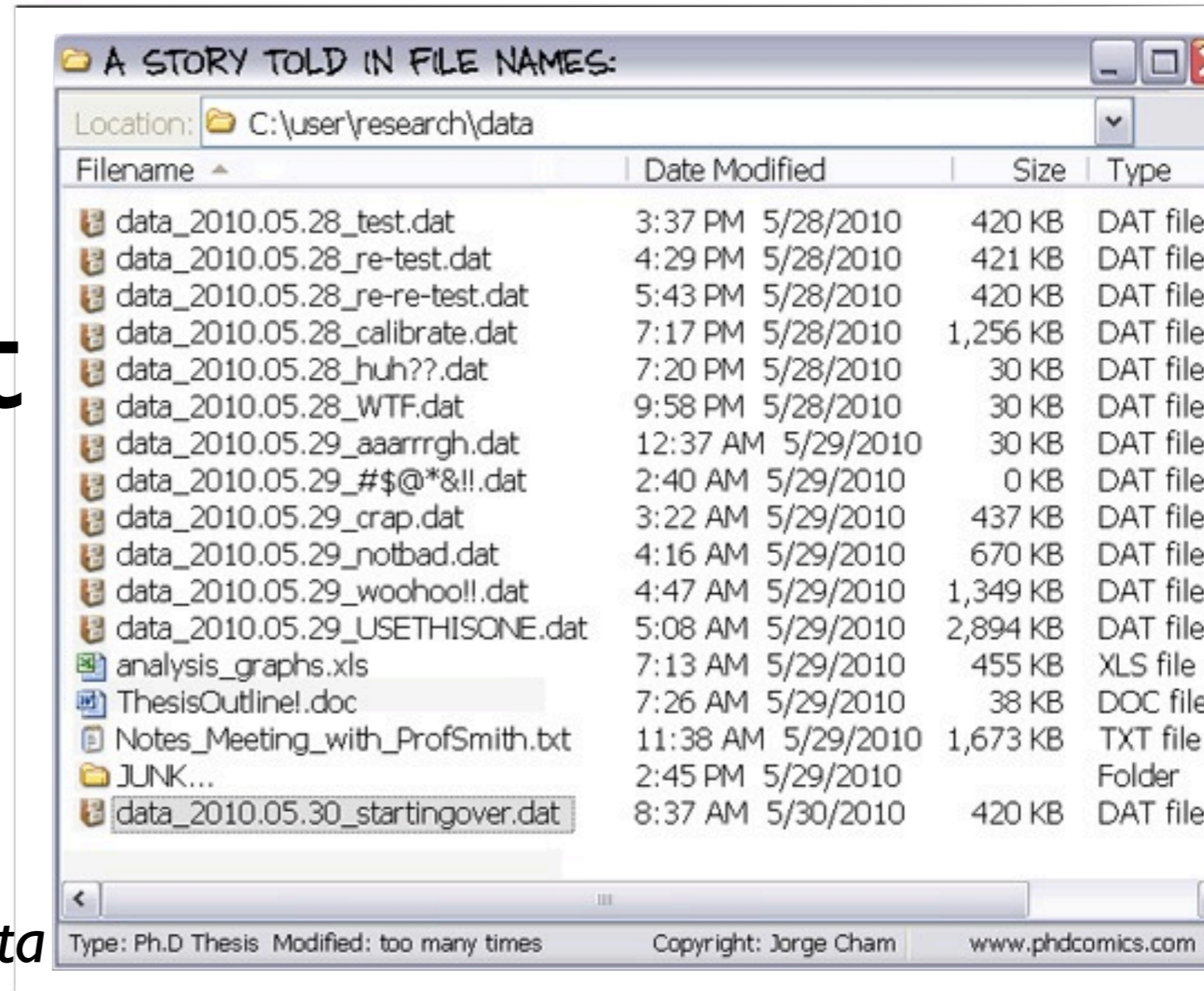
- HDF5, NetCDF have other advantages anyway
 - Binary
 - Self describing - contains not only data but names, descriptions of arrays, etc
 - Many tools can read these formats
- Big data - formats matter



```
$ ncdump -h data-simple-fort.nc
netcdf data-simple-fort {
  dimensions:
    X = 100 ;
    Y = 100 ;
    velocity components = 2 ;
  variables:
    double Density(Y, X) ;
    double Velocity(Y, X, velocity components) ;
}
```

Data Management

- Human-interpretable file-names lose their charm after few dozen files
- (or even after a few months pass)...
- Rigorously maintained *metadata* becomes essential.
- Also, need to avoid zillions of files in same directory

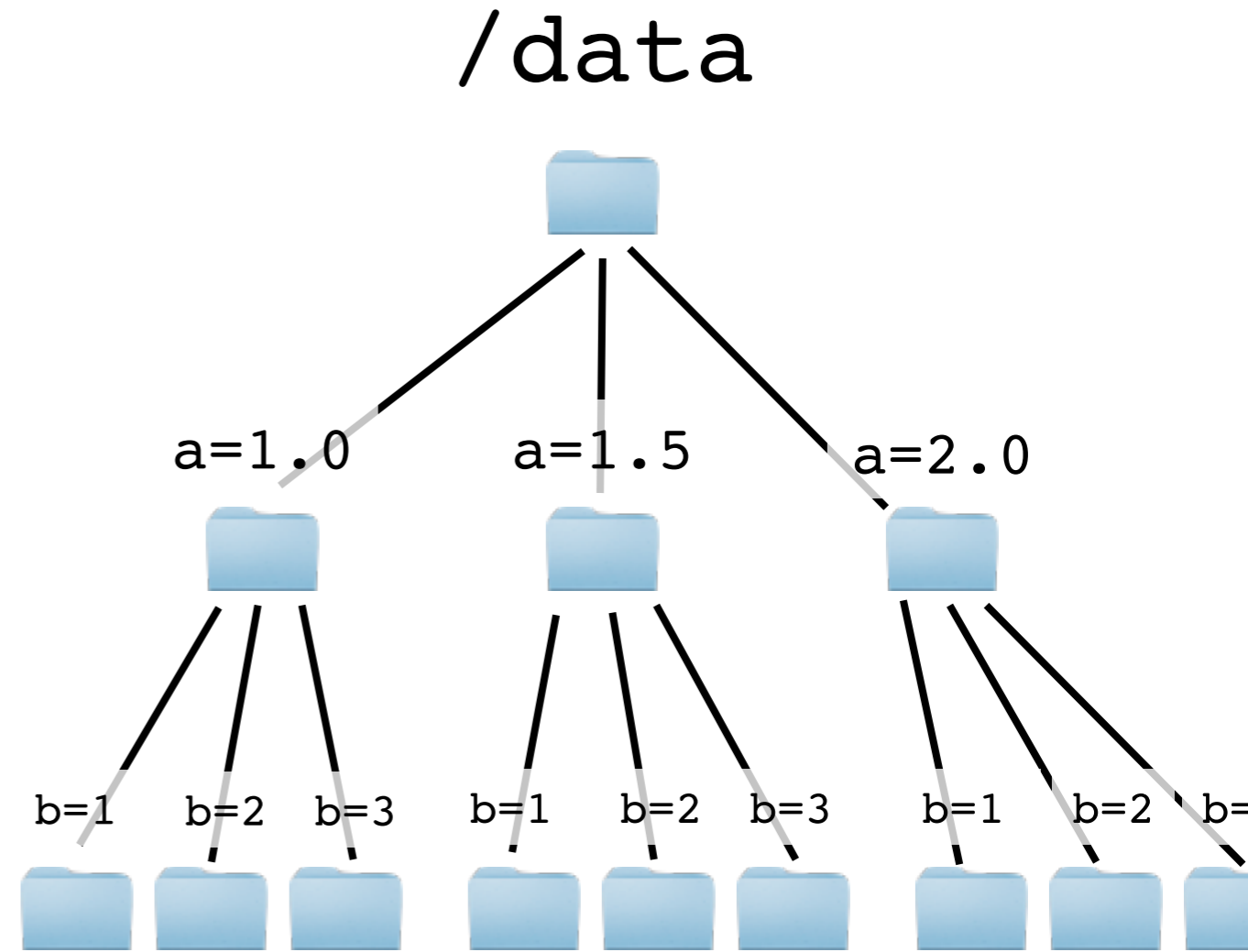


PhD Comics

<http://www.phdcomics.com/comics/archive.php?comicid=1323>

Data Management

- Hierarchy of directories work better,
- As long as layout will meet all analysis needs
- Some metadata is included in the directory structure itself
- Avoids tonnes of files in single directory



Databases for science

- Databases?
- Overhead (seeky), but not so bad if chunks in database very large
- Very handy if will be often regrouping the data
 - Don't yet know what relations you will highlight, or
 - Will highlight several overlapping relations

run#	success	size	transport	...
93	no	12k	eth	
1	yes	512	eth	
87	yes	64	ib	
13	no	32	eth	

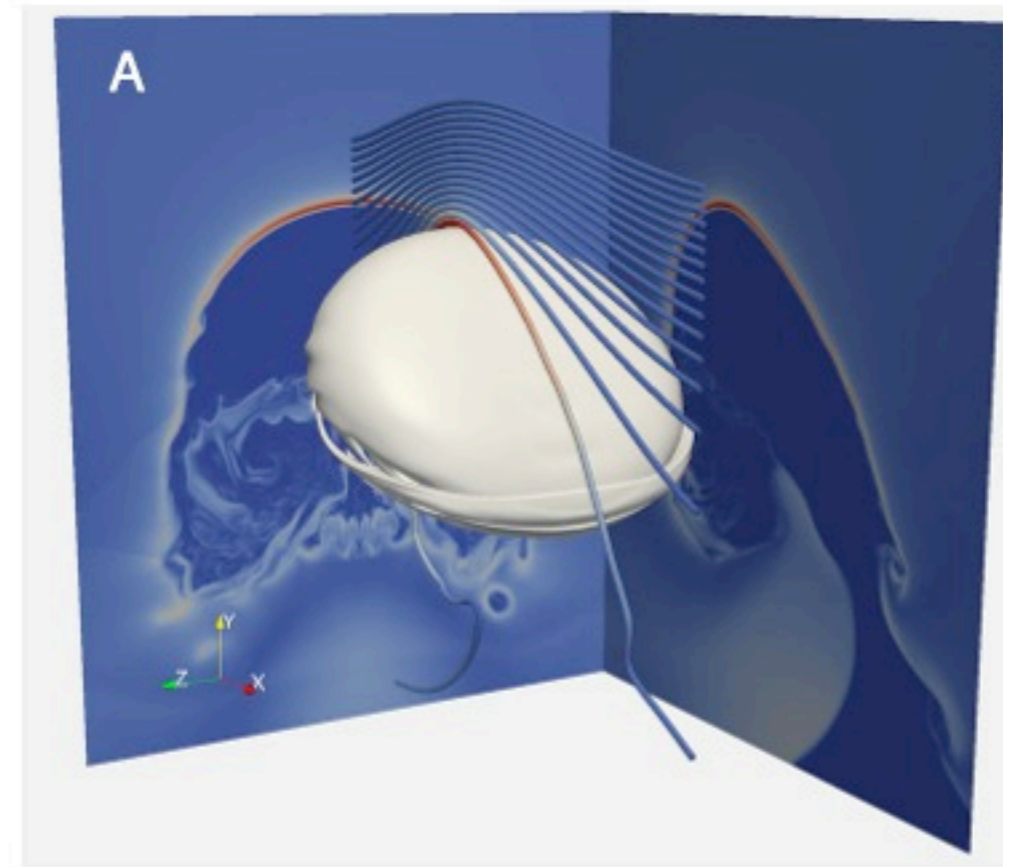
```
insert into benchmarkruns  
values (newrunnum, datestr,  
timestr, juliannum)
```

...

```
SELECT  
nprocs, test, size, transport, mpi  
ype, runtime, mopsperproc, run  
FROM mpirundata WHERE  
(success=1)
```

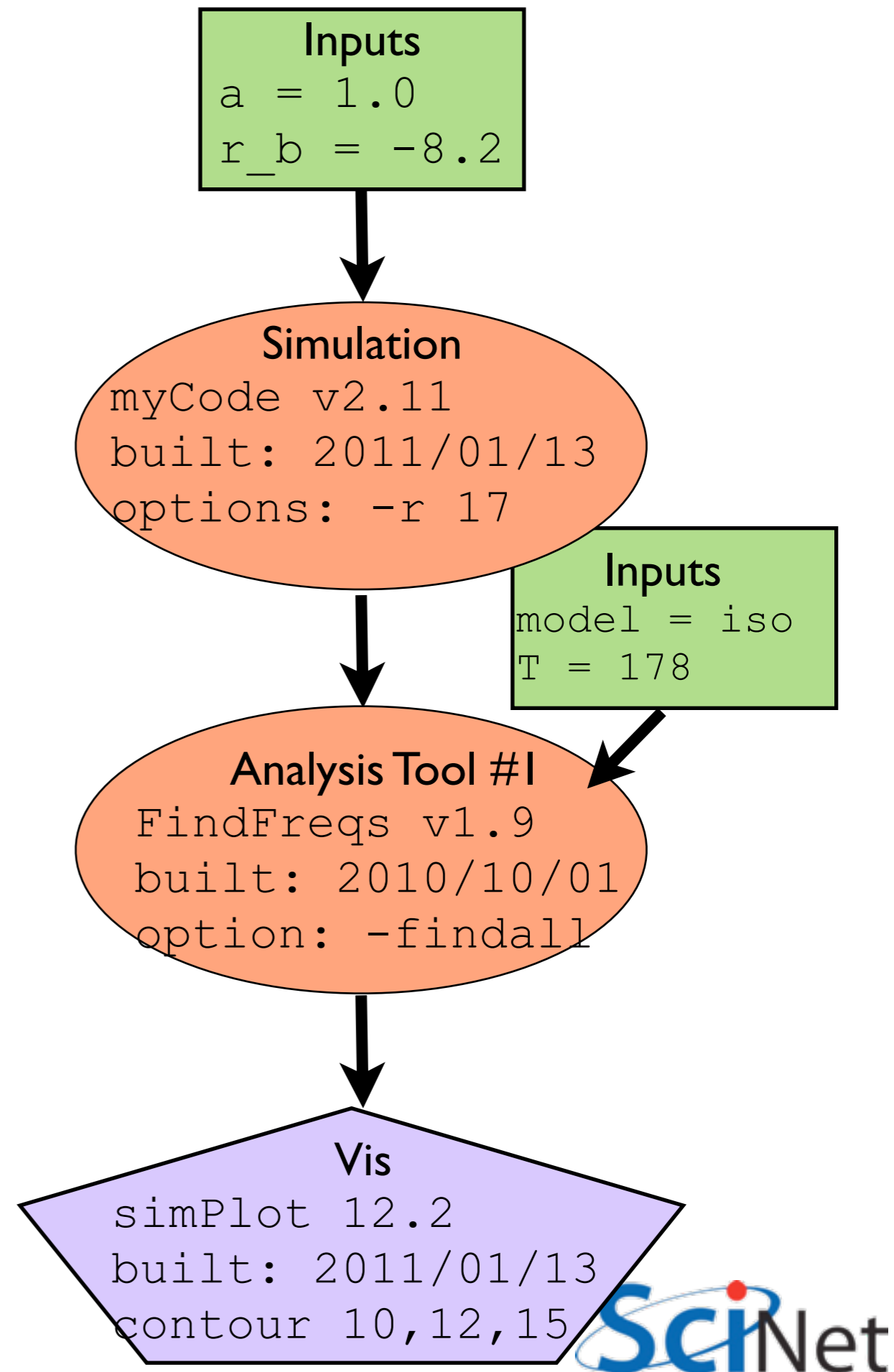
At scale, *crucial* to track data provenance

- Two of my inputs and one analysis routine changed; do I need to redo this figure?
- How - What steps were involved in making this figure?
- 2 years later, someone questions the result - can I reproduce this key figure with new code, etc?
- Different at scale - more plots, longer to generate



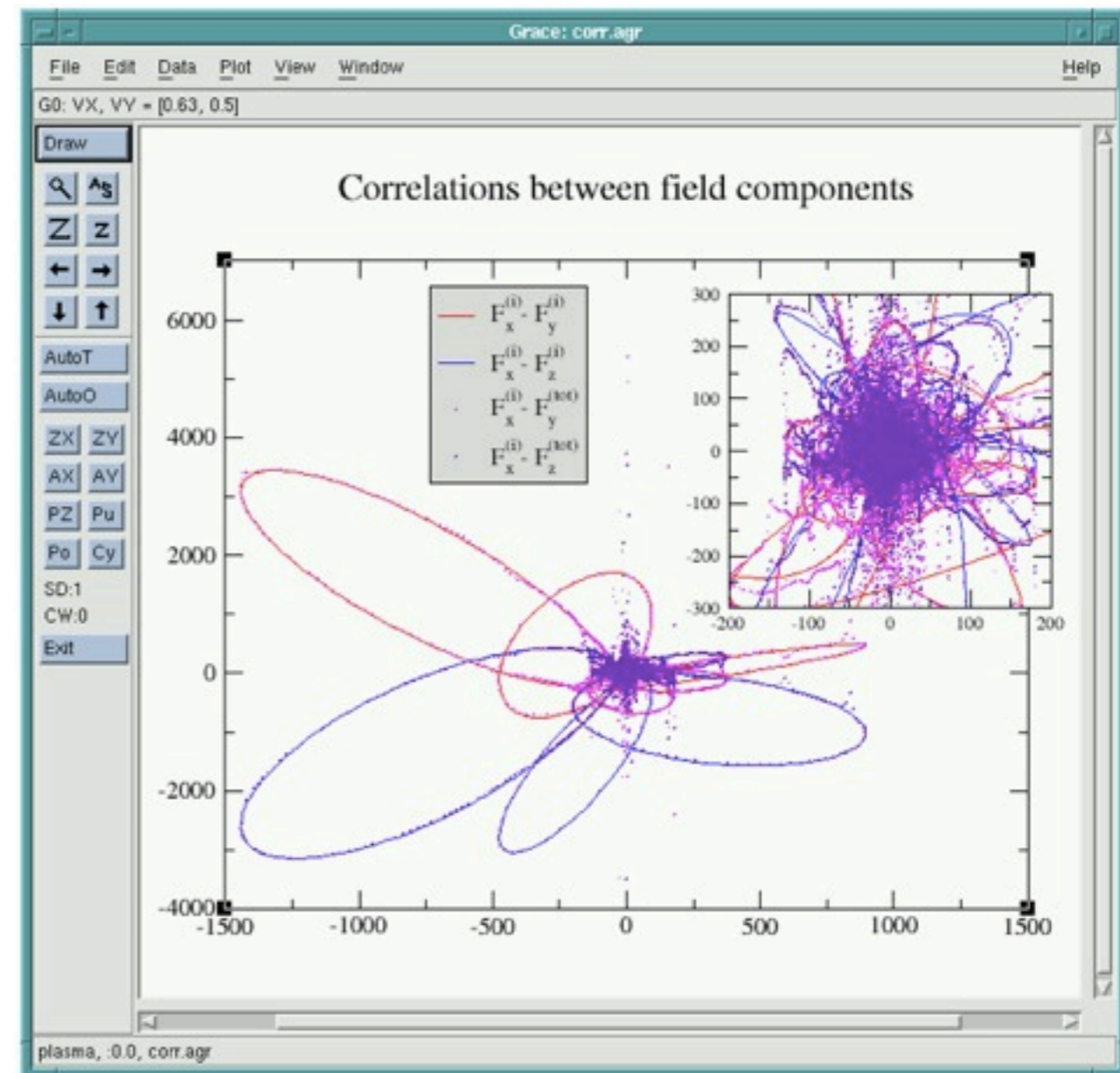
At scale, *crucial* to track data provenance

- Provenance - documentation of origin of a work
- Inputs, tools, steps in process
- Versions, dates, options
- Version control can greatly help you with this
- Then propagate the data along steps, keep it in (say) comment fields
- http://software-carpentry.org/4_0/essays/provenance



Scalability requires Automation

- Need for automation clear when dealing with thousands of small datasets..
- But large sets, too. (Sitting at a GUI for hours while waiting data to load not an option)

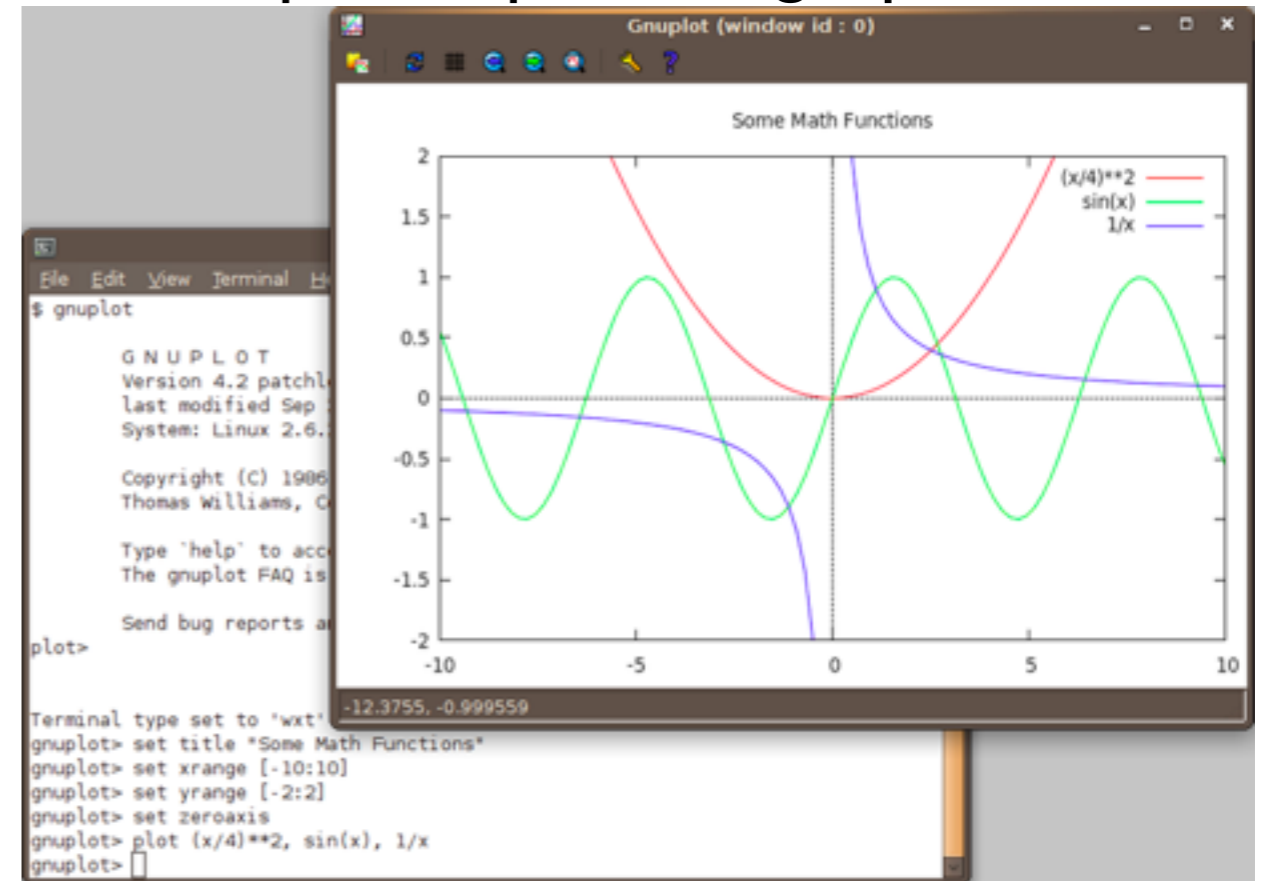


XMGrace,

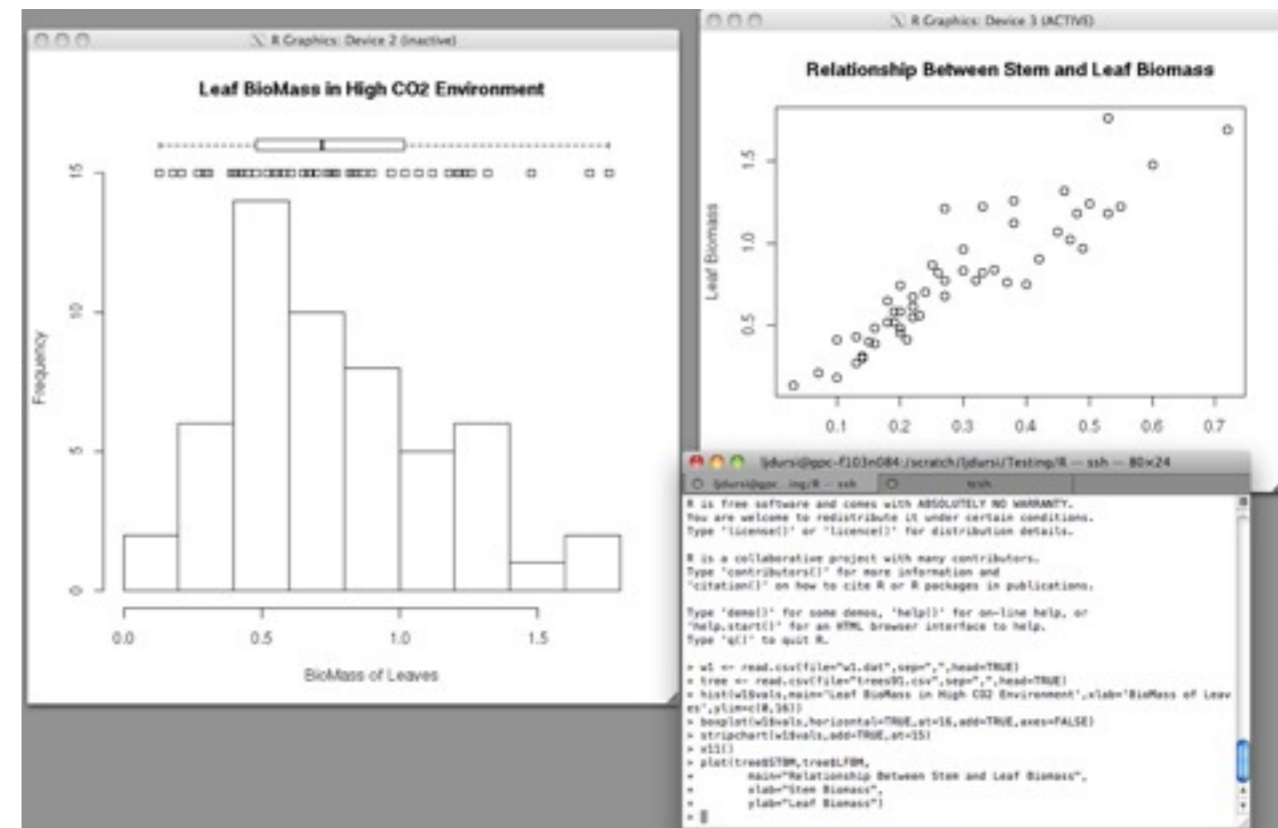
<http://plasma-gate.weizmann.ac.il/Grace/>

Scalability requires Automation

- Scripting based packages like gnuplot, matplotlib, R...
- Implicitly automatable
- Harder learning curve
- Learning basic Unix shell scripting *priceless* for automation
- http://software-carpentry.org/4_0/shell/

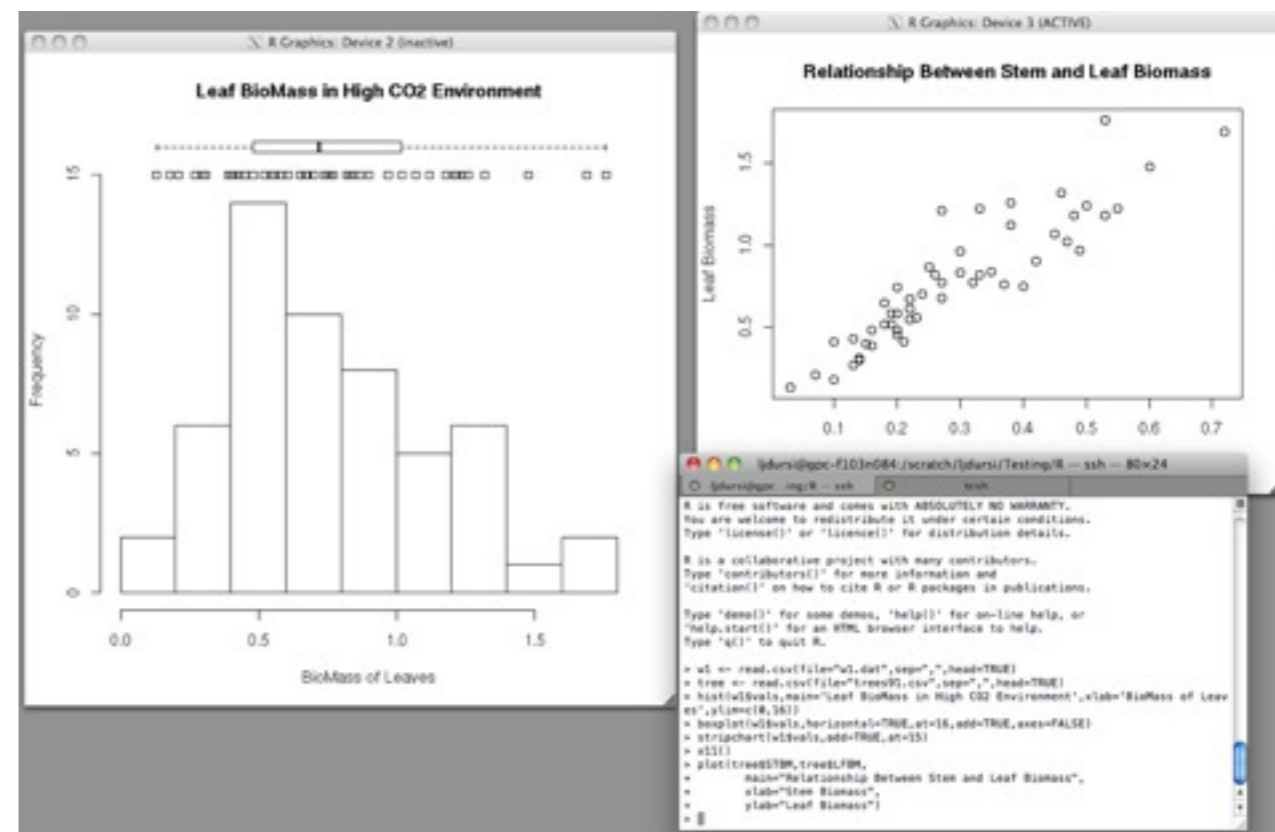
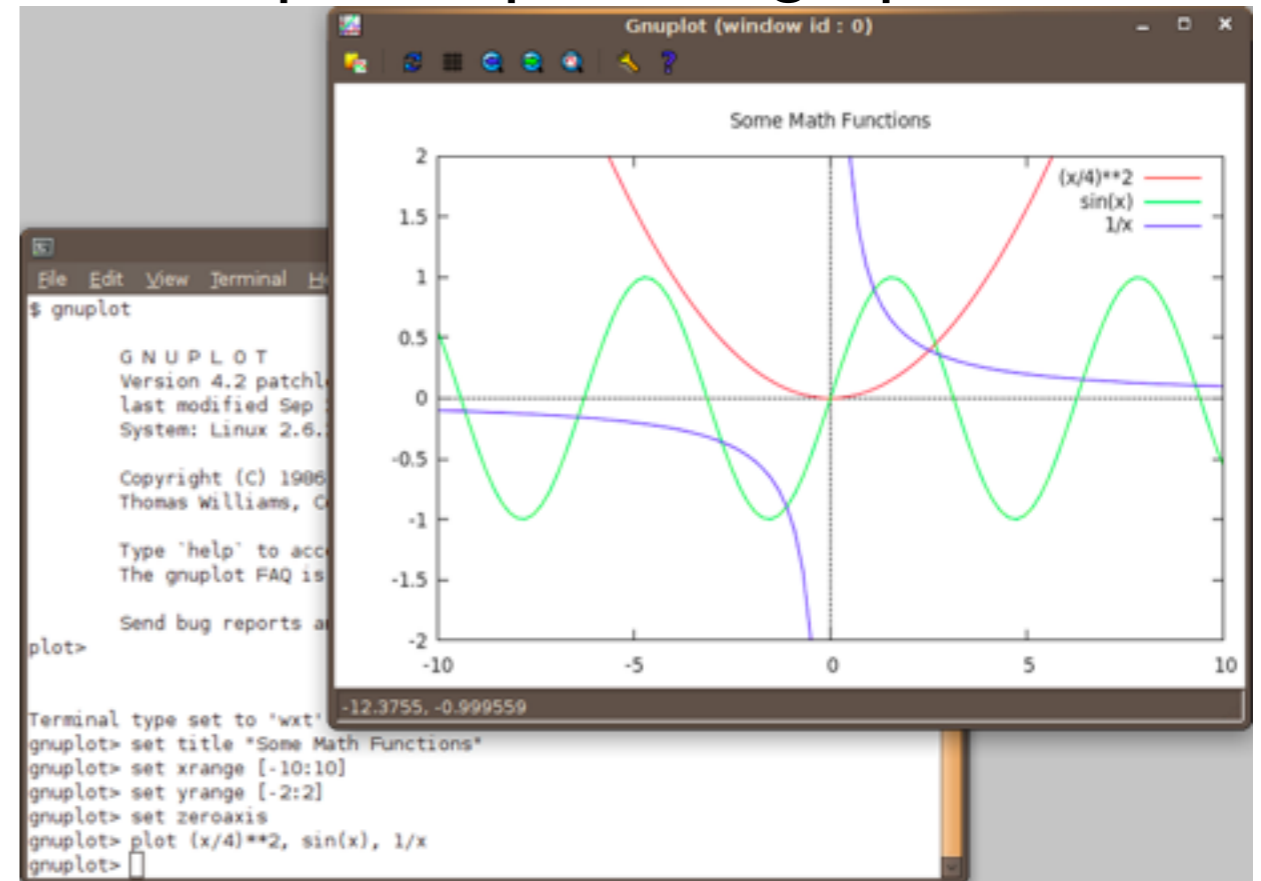


R <http://www.r-project.org/>



Scalability requires Automation

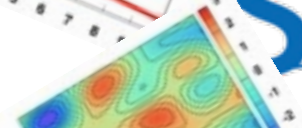
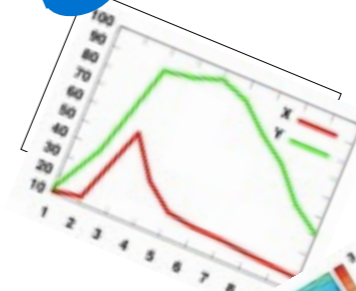
- Scripting makes processing thousands of files feasible, redoing huge visualizations less tedious
- Provides reproducibility, some form of documentation of process
- Scripts can be kept in version control



Planning your analysis pipeline

- Start w/ 40 TB, will presumably end with much less
- Do as much of that reduction as early as possible in the process
 - Average, bin, integrate, combine, contour
- Automate everything
- Easier if you know exactly what you'll be doing

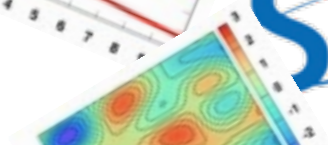
```
29 14 42 78 63 41 8 7 9 48 2 98 3 43 74 3 9 29 58 76
62 29 58 44 3 3 46 2 56 25 53 38 79 38 75 89 91 24 33 5
14 28 89 36 8 16 66 8 76 58 11 17 14 58 4 4 8 13 57 18
17 17 12 48 65 1 19 67 77 57 34 29 23 25 58 23 64 86 69
18 89 72 18 96 22 55 18 1 74 26 15 87 41 44 98 35 36 85 13
18 88 25 53 16 36 31 46 53 188 32 88 27 19 32 53 39 45 87 14
24 61 38 38 74 14 21 78 27 39 5 74 63 63 7 36 7 45 14 81
3 83 57 19 86 11 6 75 64 88 65 23 7 58 37 67 29 3 85 68
26 53 46 21 2 81 3 45 72 12 95 21 1 38 61 31 58 28 67 61
45 73 41 39 92 54 86 6 76 89 41 98 47 14 71 33 66 52 28 68
33 63 18 48 33 87 43 38 81 65 38 96 65 91 48 42 23 51 27 12
65 35 52 92 15 58 47 54 35 79 48 61 61 78 37 22 38 93 46 25
61 36 8 24 65 94 35 48 98 18 58 84 53 27 32 34 31 61 71 98
42 19 33 92 29 87 55 43 68 27 12 21 17 76 88 67 78 68 94 68
85 61 99 84 12 63 3 82 92 14 56 48 71 5 92 73 56 57 53 24
96 71 98 45 98 25 58 68 25 66 12 67 88 88 26 42 26 78 23 94
89 39 61 95 36 15 94 84 33 66 18 92 91 49 2 51 76 39 188 91
61 75 71 35 89 43 78 39 1 37 25 34 93 4 58 18 96 68 72 87
55 93 17 19 2 48 4 94 48 14 98 37 76 61 43 65 56 94 38 18
41 46 75 19 84 48 79 91 27 45 51 23 33 59 65 93 64 27 5 3
34 17 29 3 62 64 45 19 68 82 53 54 88 67 28 23 55 25 13 38
51 35 73 45 49 67 23 59 52 15 22 89 27 76 81 63 16 43 15 75
48 13 37 58 22 47 188 28 32 9 24 98 38 188 59 16 36 58 82 79
11 42 71 91 28 88 77 68 29 94 19 49 29 57 98 39 32 57 1 68
91 98 66 92 19 41 89 188 81 58 78 45 68 98 75 58 19 48 13 66
51 18 51 58 41 48 48 13 6 98 98 67 86 17 78 85 7 28 75
54 68 93 11 21 4 15 43 24 72 58 88 78 7 42 87 71 67 74 86
76 47 34 36 58 58 63 75 8 86 27 84 78 57 34 58 38 91 7 16
67 45 78 41 53 21 33 11 73 36 77 33 81 35 37 55 79 33 45 5
74 41 26 8 41 54 61 54 73 75 98 57 93 32 36 75 188 98 8 83
22 98 87 59 71 8 23 62 58 23 39 11 84 36 79 92 64 74 85 72
43 188 87 86 83 34 59 65 88 3 57 46 98 78 99 17 1 57 7 96
14 17 18 78 78 48 15 58 65 71 55 28 92 79 14 58 85 22 78 27
68 25 47 28 72 75 38 88 63 65 49 51 31 68 3 54 68 56 22 59
43 89 76 96 89 18 86 58 17 4 55 28 58 23 32 3 84 11 25 21
75 66 16 59 188 21 86 15 83 45 61 53 2 11 3 36 78 79 68 58
99 17 26 47 7 87 75 76 27 54 68 88 53 34 4 63 32 56 28 77
88 56 24 188 55 63 25 38 29 18 5 53 69 48 37 19 82 77 48 97
31 32 46 81 27 47 19 78 46 55 52 29 28 54 63 57 48 37 59 58
6 38 5 94 74 75 84 48 85 19 11 95 84 6 35 25 33 66 98 87
3 79 98 35 98 9 56 21 19 42 66 77 38 91 18 62 59 73 1 89
27 73 7 4 52 39 87 91 22 95 92 66 83 3 99 82 52 13 61 44
47 18 1 21 95 51 27 21 7 28 83 41 46 41 27 22 59 47 93 48
82 38 29 98 96 5 86 99 51 71 38 43 53 98 65 98 18 62 36 65
```



```
29 14 42 78 63 41 8 7 9 48 2 98 3 43 74 3 9 29 58 76
62 29 58 44 3 3 46 2 56 25 53 38 79 38 75 89 91 24 33 5
14 28 89 36 8 16 66 8 76 58 11 17 14 58 4 4 8 13 57 18
17 17 12 48 65 1 19 67 77 57 34 29 23 25 58 23 64 86 86 69
18 89 72 18 96 22 55 18 1 74 26 15 87 41 44 98 35 36 85 13
18 88 25 53 16 36 31 46 53 188 32 88 27 19 32 53 39 45 87 14
24 61 38 38 74 14 21 78 27 39 5 74 63 63 7 36 7 45 14 81
3 83 57 19 86 11 6 75 64 88 65 23 7 58 37 67 29 3 85 68
26 53 46 21 2 81 3 45 72 12 95 21 1 38 61 31 58 28 67 61
45 73 41 39 92 54 86 6 76 89 41 98 47 14 71 33 66 52 28 68
33 63 18 48 33 87 43 38 81 65 38 96 65 91 48 42 23 51 27 12
65 35 52 92 15 58 47 54 35 79 48 61 61 78 37 22 38 93 46 25
61 36 8 24 65 94 35 48 98 18 58 84 53 27 32 34 31 61 71 98
42 19 33 92 29 87 55 43 68 27 12 21 17 76 88 67 78 68 94 68
85 61 99 84 12 63 3 82 92 14 56 48 71 5 92 73 56 57 53 24
96 71 98 45 98 25 58 68 25 66 12 67 88 88 26 42 26 78 23 94
89 39 61 95 36 15 94 84 33 66 18 92 91 49 2 51 76 39 188 91
61 75 71 35 89 43 78 39 1 37 25 34 93 4 58 18 96 68 72 87
55 93 17 19 2 48 4 94 48 14 98 37 76 61 43 65 56 94 38 18
41 46 75 19 84 48 79 91 27 45 51 23 33 59 65 93 64 27 5 3
34 17 29 3 62 64 45 19 68 82 53 54 88 67 28 23 55 25 13 38
51 35 73 45 49 67 23 59 52 15 22 89 27 76 81 63 16 43 15 75
48 13 37 58 22 47 188 28 32 9 24 98 38 188 59 16 36 58 82 79
11 42 71 91 28 88 77 68 29 94 19 49 29 57 98 39 32 57 1 68
91 98 66 92 19 41 89 188 81 58 78 45 68 98 75 58 19 48 13 66
51 18 51 58 41 48 48 48 13 6 98 98 67 86 17 78 85 7 28 75
54 68 93 11 21 4 15 43 24 72 58 88 78 7 42 87 71 67 74 86
76 47 34 36 58 58 63 75 8 86 27 84 78 57 34 58 38 91 7 16
67 45 78 41 53 21 33 11 73 36 77 33 81 35 37 55 79 33 45 5
74 41 26 8 41 54 61 54 73 75 98 57 93 32 36 75 188 98 8 83
22 98 87 59 71 8 23 62 58 23 39 11 84 36 79 92 64 74 85 72
43 188 97 86 83 34 59 65 88 3 57 46 98 78 99 17 1 57 7 96
14 17 18 78 78 48 15 58 65 71 55 28 92 79 14 58 85 22 78 27
68 25 47 28 72 75 38 88 63 65 49 51 31 68 3 54 68 56 22 59
43 89 76 96 89 18 86 58 17 4 55 28 58 23 32 3 84 11 25 21
75 66 16 59 188 21 86 15 83 45 61 53 2 11 3 36 78 79 68 58
99 17 26 47 7 87 75 76 27 54 68 88 53 34 4 63 32 56 28 77
88 56 24 188 55 63 25 38 29 18 5 53 69 48 37 19 82 77 48 97
31 32 46 81 27 47 19 78 46 55 52 29 28 54 63 57 48 37 59 58
6 38 5 94 74 75 84 48 85 19 11 95 84 6 35 25 33 66 98 87
3 79 98 35 98 9 56 21 19 42 66 77 38 91 18 62 59 73 1 89
27 73 7 4 52 39 87 91 22 95 92 66 83 3 99 82 52 13 61 44
47 18 1 21 95 51 27 21 7 28 83 41 46 41 27 22 59 47 93 48
82 38 29 98 96 5 86 99 51 71 38 43 53 98 65 98 18 62 36 65
```

Planning your analysis pipeline

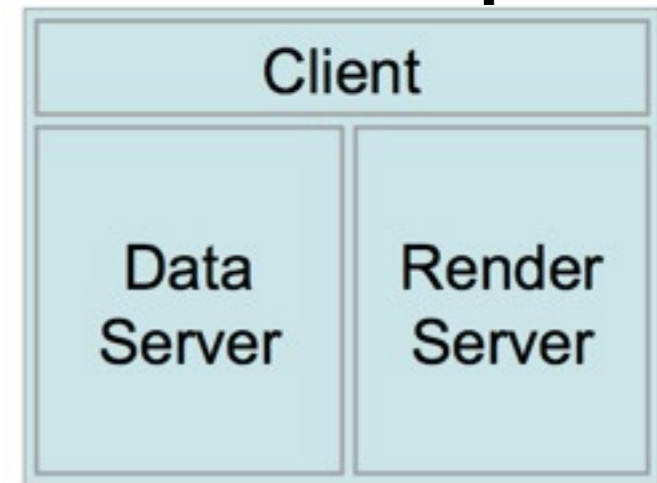
- Don't generate dozens of TB of data without knowing what you're going to do with it!
- Start small
- Explore on smaller data, or subsets



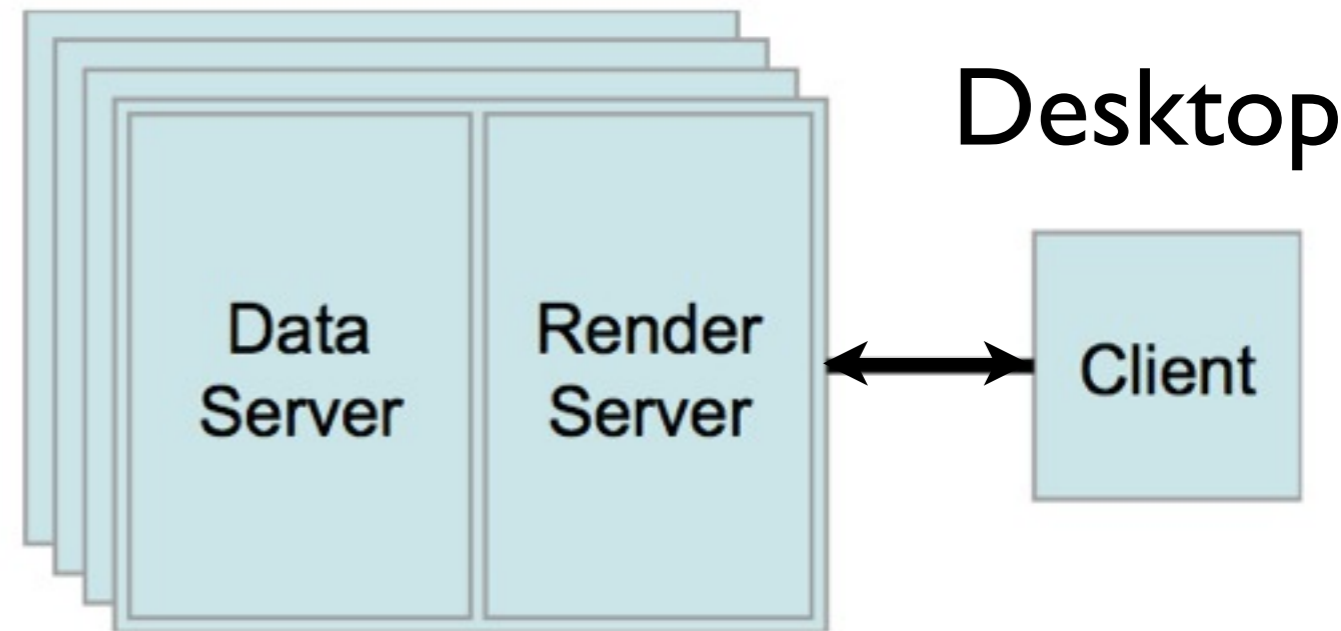
Scalability may require parallel

- Scalable visualization packages
 - eg, ParaView, VisIT
- Client/Server model
- On desktop, client/server coexist
- Or, servers can be **many** nodes on cluster!

Desktop



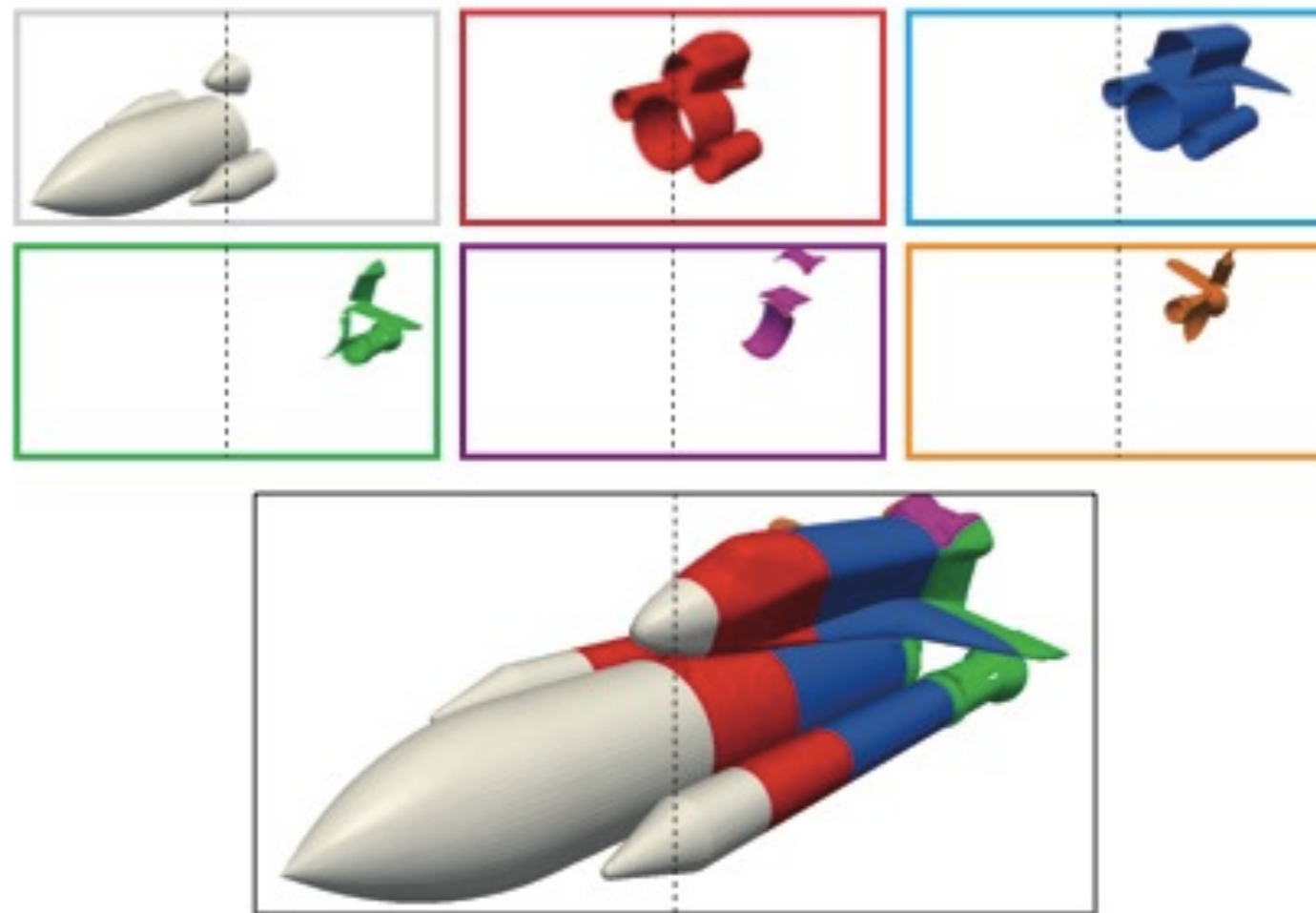
Parallel Cluster



Paraview Tutorial,
<http://paraview.org/>

Parallel Visualization

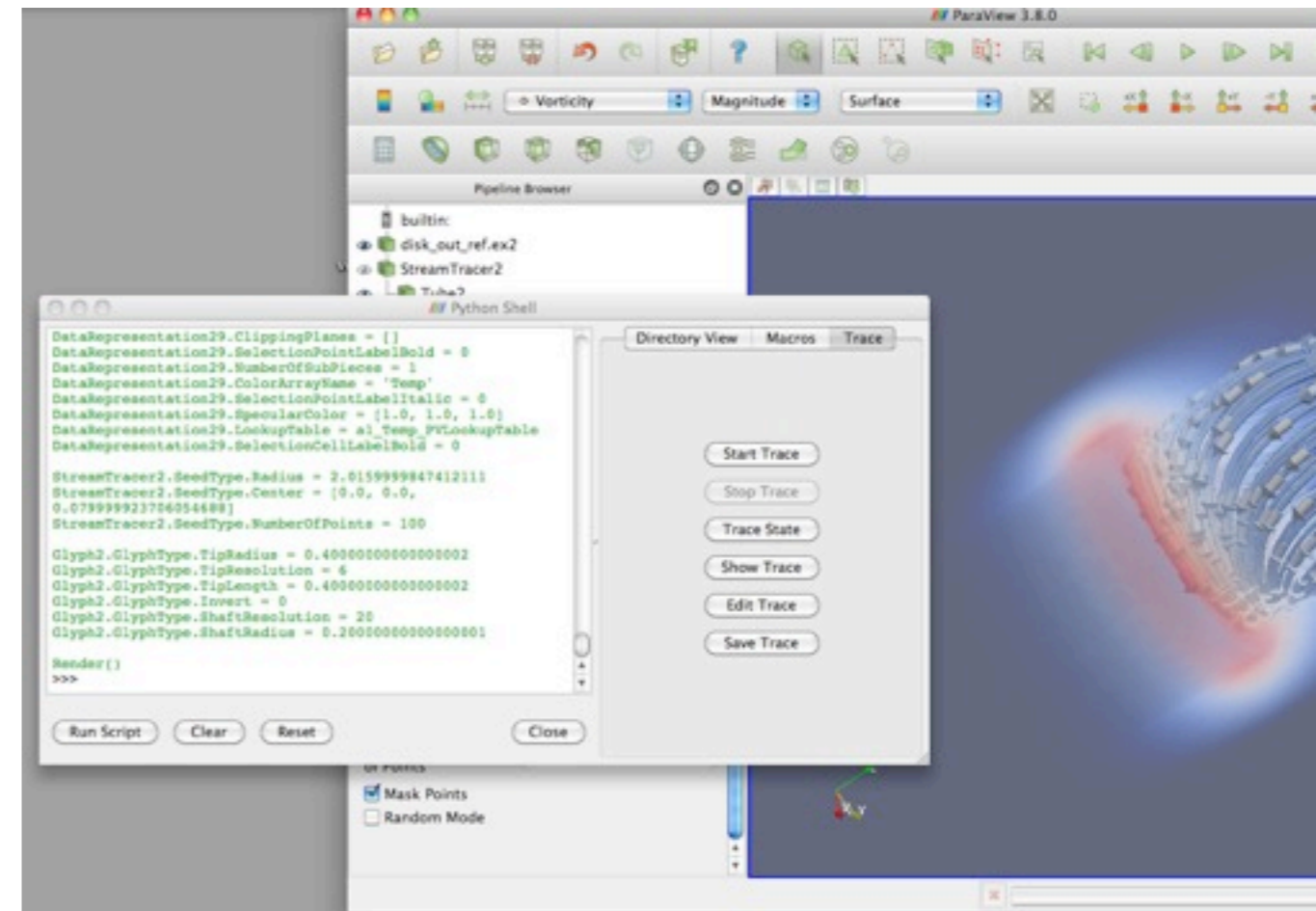
- Decompose data onto many processors (need parallel file systems, format!)
- Each processor generates its portion of the geometry, or image
- Composited en route to client
- Can control visualization interactively from desktop



IceT users guide,
<http://www.cs.unm.edu/~kmorel/IceT/>

Remains scriptable

- Even in client/server mode, these tools are scriptable and automatable
- Client end can be run as a script
- Client coordinates all communication with server nodes
- Can also be run without GUI, pure batch mode on cluster.



Paraview, <http://paraview.org/>

Highly Scalable Parallel Visualization

- Can work to extremely large scales
- VisIt - 4 trillion zone simulations
- ParaView - billions of polygons/sec



Red RoSE visualization cluster
(credit: Sandia Nat'l Lab)

Paraview on SciNet

- Launch job on cluster
- Module load paraview
- Mpirun paraview server
- Interactive? Connect your desktop client
- Run visualization

Your Desktop

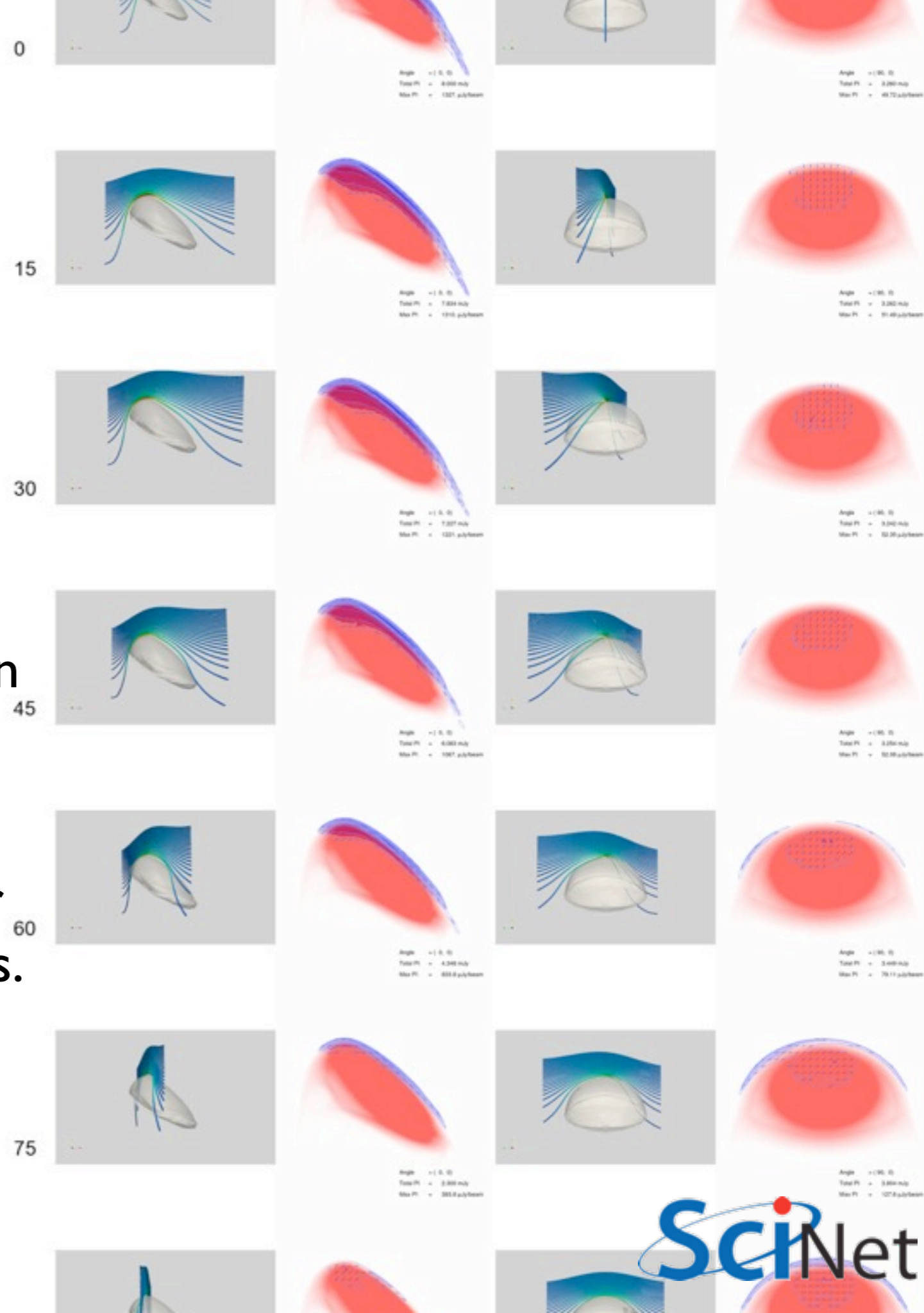


SciNet



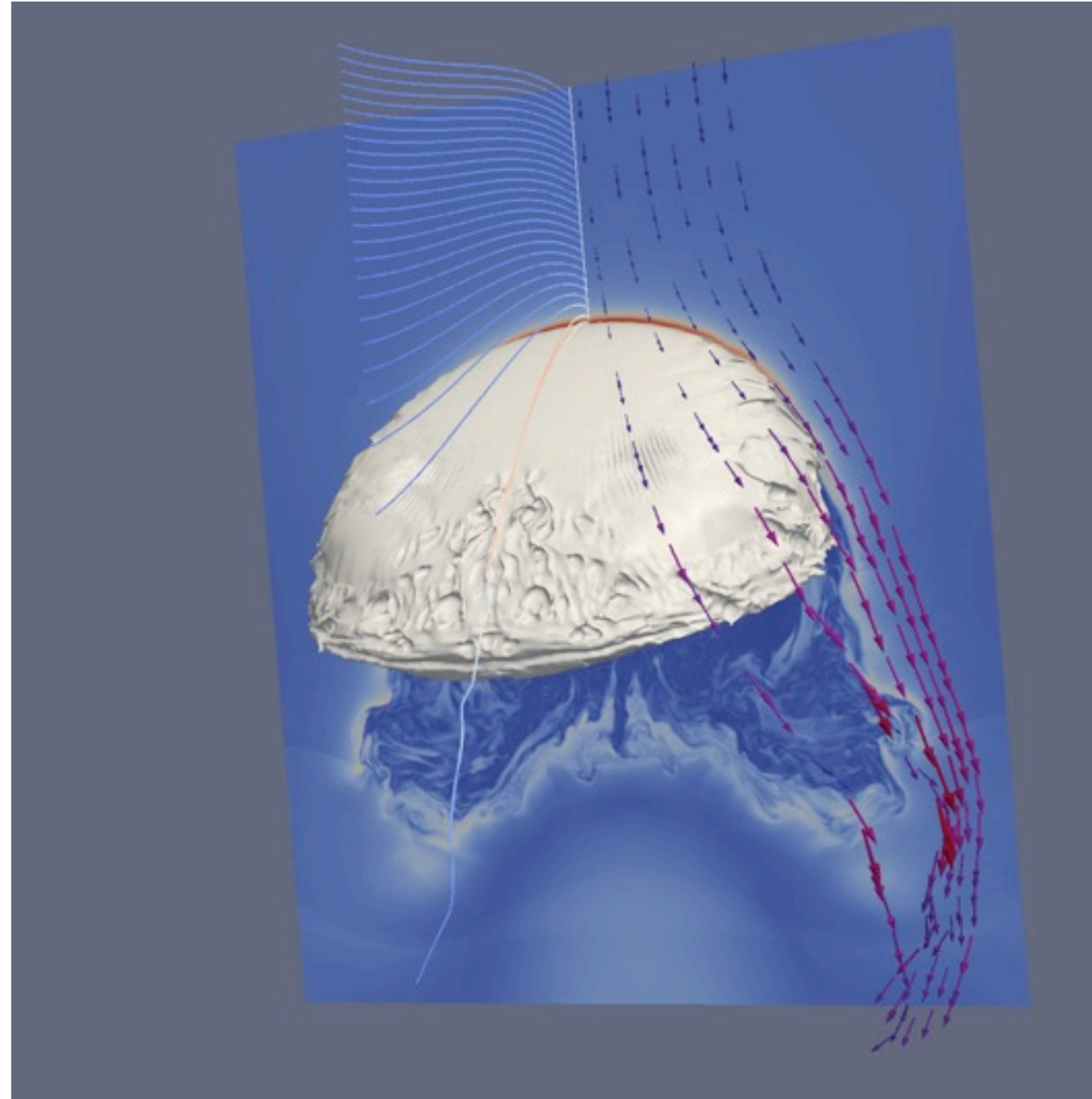
Paraview on SciNet

- One of first simulations run on SciNet's GPC
- Full data set was about 40 TB
- Included collections of smaller simulations and big simulations.



Paraview on SciNet

- High resolution simulation
- 2.5 billion zones
- 330,000 cpu hours
- 256 processors just to load, visualize data



Plan I/O carefully

- Binary files
- As much as possible, large files
- File formats that can be read in parallel, subregions extracted
- Avoid zillions of files
- Especially avoid zillions of files in single directory

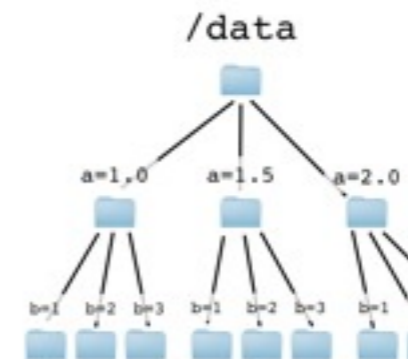
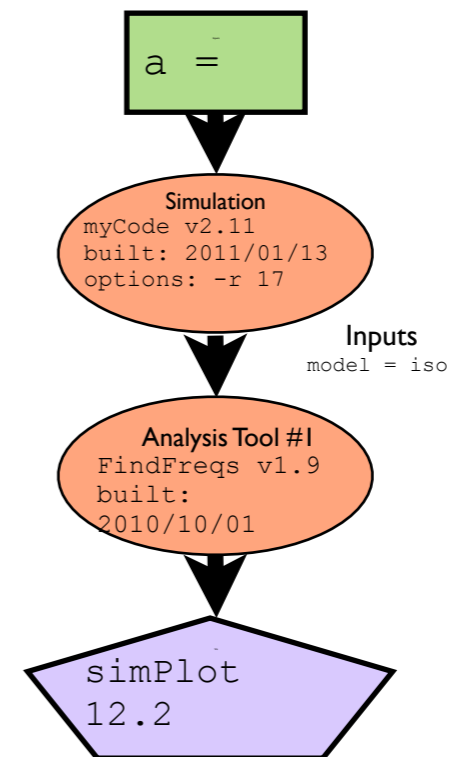


010101001111



Data management must scale

- Include metadata for provenance
 - Reduce need to re-do
- Sensible data management
 - Hierarchy of data directories only if that will always work
 - Data bases, formats that allow metadata



Use scalable, automatable tools.

- For large data sets, parallel tools
- ParaView, VisIt, etc
- Scriptable tools - gnuplot,
python, R, ParaView/Visit...
- Scripts provide reproducibility!



```
gnuplot> reset
gnuplot> set xrange [-5:5]
gnuplot> set yrange [-5:5]
gnuplot> unset key
gnuplot> set palette rgbformulae 33,13,10
gnuplot> p 'test.dat' with image, 'cont.dat' w l lt -1 lw 1.5
gnuplot> show term

terminal type is aqua 0 title "Figure 0" size 846,594 font "T
```



Carefully plan your workflow

- Reduce data as early as possible in the process
 - Average, integrate, combine, contour
- Automate everything



```
gnuplot> reset
gnuplot> set xrange [-5:5]
gnuplot> set yrange [-5:5]
gnuplot> unset key
gnuplot> set palette rgbformulae 33,13,10
gnuplot> p 'test.dat' with image, 'cont.dat' w l lt -1 lw 1.5
gnuplot> show term

terminal type is aqua 0 title "Figure 0" size 846,594 font "T
```



Thanks to

- Mubdi Rahman, for putting this all together
- Ramses van Zon, Scott Northrup, Danny Gruner -
importance of I/O