

# Petascale Data Management: Guided by Measurement



**SciDAC**  
Scientific Discovery through  
Advanced Computing

## petascale data storage institute

[www.pdsi-scidac.org/](http://www.pdsi-scidac.org/)

### MEMBER ORGANIZATIONS

- Los Alamos National Laboratory – [institute.lanl.gov/pdsi/](http://institute.lanl.gov/pdsi/)
- Oak Ridge National Laboratory – [www.csm.ornl.gov/](http://www.csm.ornl.gov/)
- National Energy Research Scientific Computing Center  
[pdsi.nersc.gov/](http://pdsi.nersc.gov/)
- Pacific Northwest National Laboratory – [www.pnl.gov/](http://www.pnl.gov/)

## The Computer Failure Data Repository

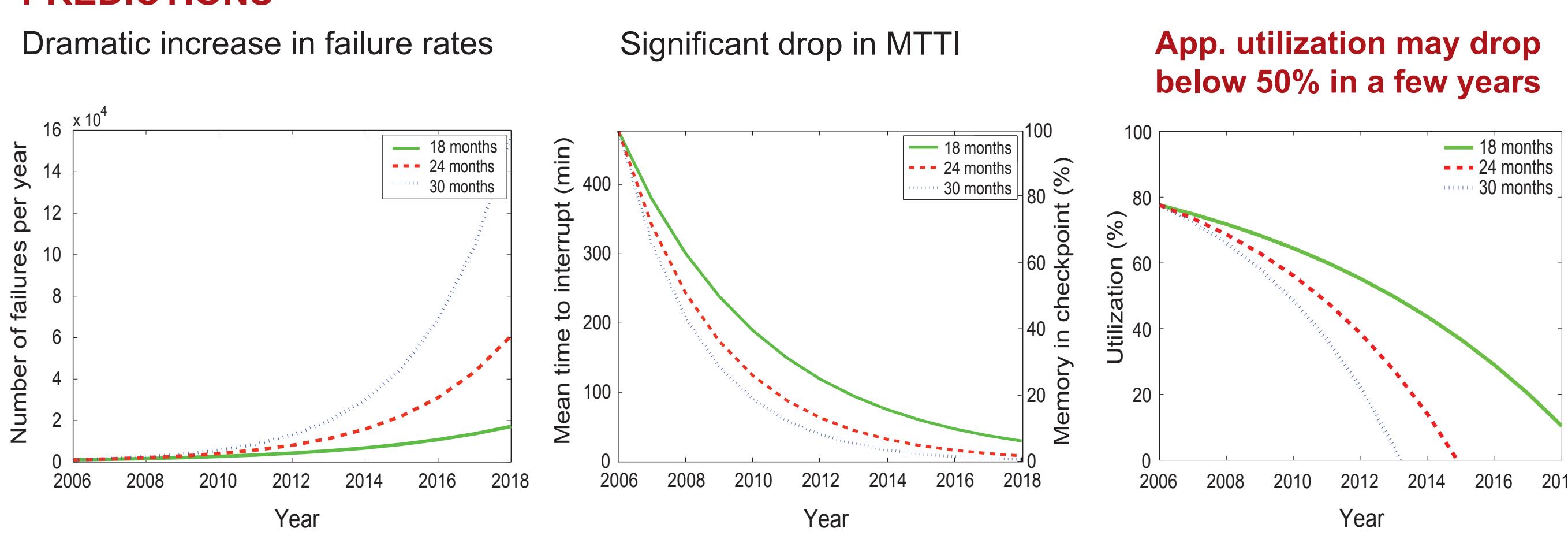
- Goal: to collect and make available failure data from a large variety of sites
  - Better understanding of the characteristics of failures in the real world
- Now maintained by USENIX at [cfdr.usenix.org/](http://cfdr.usenix.org/)

NAME	SYSTEM TYPE	SYSTEM SIZE	TIME PERIOD	TYPE OF DATA
Los Alamos National Laboratory	22 HPC clusters	5000 nodes	9 years	Any node outage
Pittsburgh Supercomputing Center	1 HPC cluster	765 nodes 3,400 disks	5 years	Hardware/disk drive replacements
1 Internet service, Various HPC sites	3 storage, many HPC clusters	>10,000 nodes >100,000 disks	1 mth - 5 yrs	
NERSC	HPC cluster	A number of production systems	5 years	I/O specific failures
COM 1	Internet services cluster	Multiple distributed sites	1 mth	Hardware failures
COM 2	Internet services cluster	Multiple distributed sites	20 mths	Warranty service log of hardware failures
COM 3	Internet services cluster	Large external storage system	1 yr	Aggregate quarterly stats of disk failures

### BASE ASSUMPTIONS

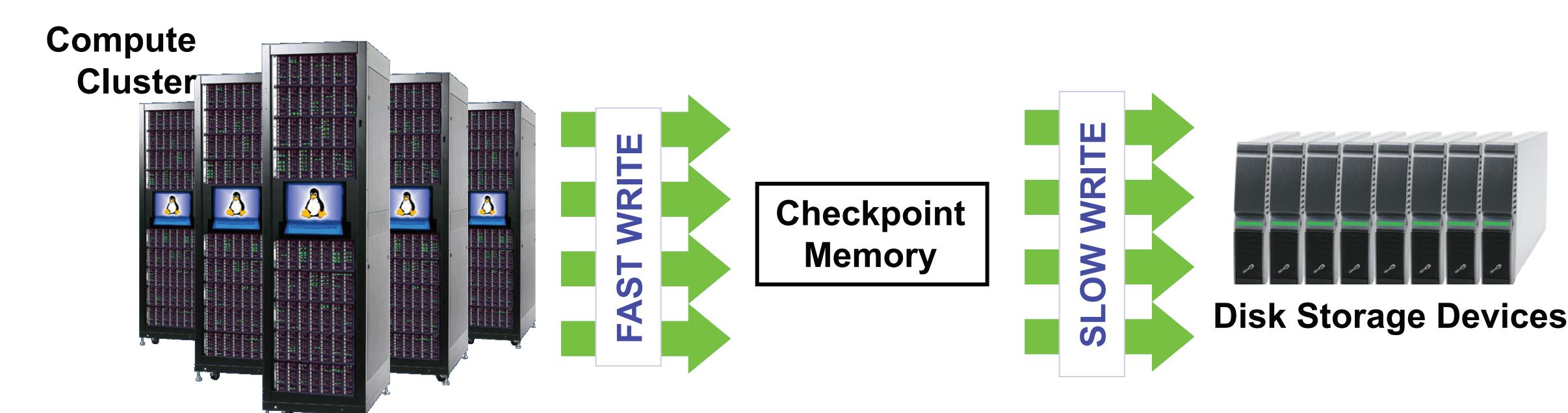
- Top500.org: Performance of top installations doubles every year
- Technology trends:
- Little/no increase in clock speed → cycles coming from more sockets and more cores per socket
- # cores per socket doubles every 18-30 months
- Failure data:
- Hardware is not getting more reliable over time
- Failure rate proportional to # sockets (optimistically 0.1 failures/socket/year)

### PREDICTIONS



### POSSIBLE SOLUTIONS

- Make more reliable chips (not likely given history data we have seen)
- Use non-growing portion of bigger machines (not an answer for demanding applications)
- Checkpoint smaller fraction of memory (compression)
- Increase storage bandwidth (significantly more drives)
- Specialized checkpoint devices – use memory to stage checkpoint (maybe flash if organization for write bandwidth is improved)



- Investigate alternative methods, e.g. process pairs
- Expensive, but will pay off once application utilization drops below 50%



[www.pdsi-scidac.org/](http://www.pdsi-scidac.org/)

- Parallel Data Lab, Carnegie Mellon University – [www.pdl.cmu.edu/](http://www.pdl.cmu.edu/)
- Sandia National Laboratories – [www.sandia.gov/](http://www.sandia.gov/)
- Center for Information Technology Integration, U. of Michigan  
[www.citi.umich.edu/projects/pdsi/](http://www.citi.umich.edu/projects/pdsi/)
- University of California at Santa Cruz – [www.pdsi.ucsc.edu/](http://www.pdsi.ucsc.edu/)

## Filesystems Statistics Survey

### GOALS

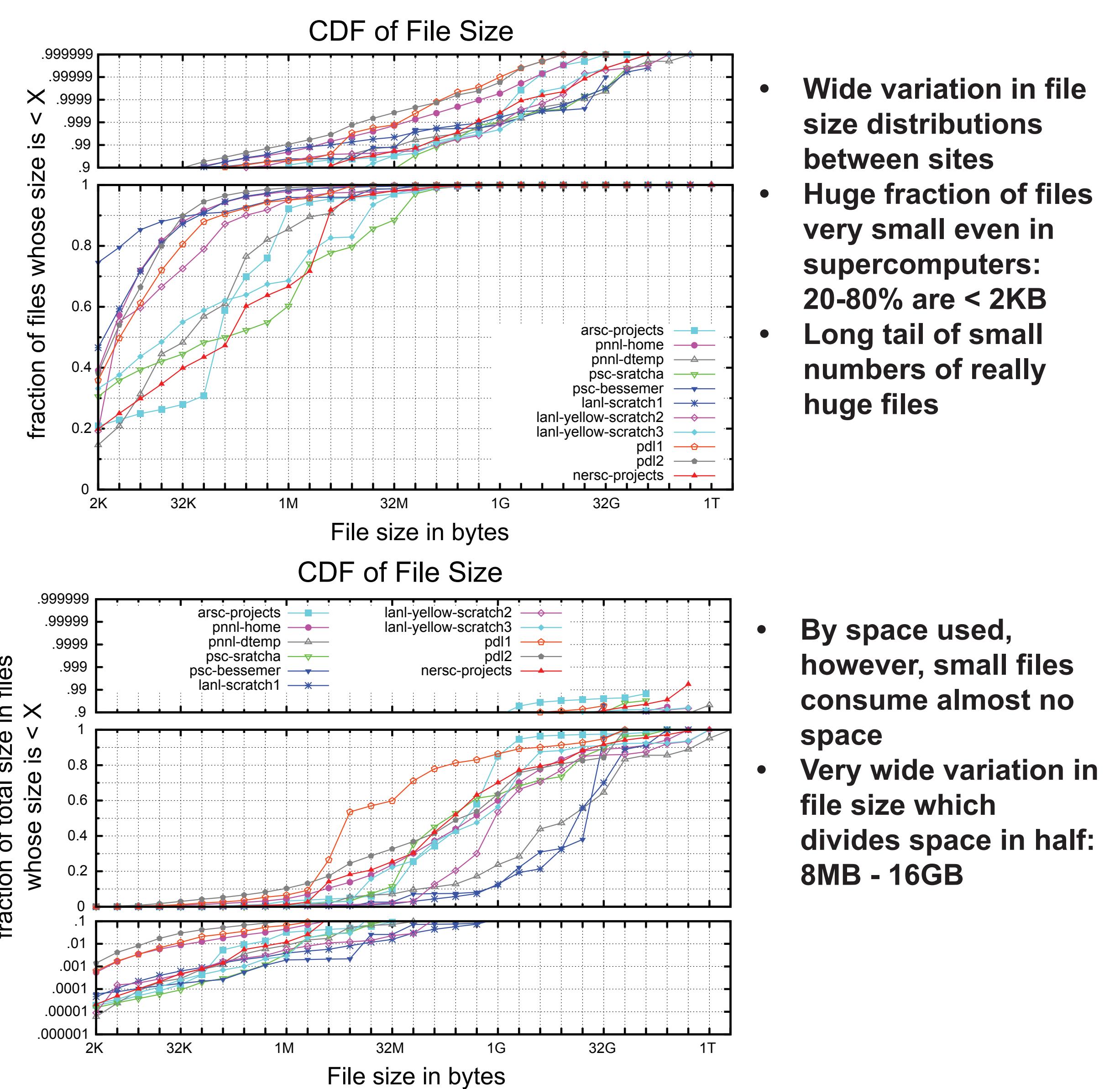
- Gather & build large DB of static filetree summary
  - Build small, non-invasive, anonymizing stats gather tool
  - Distribute fsstats tool via easily used web site
  - Encourage contributions (output of tool) from many FSs
  - Offer uploaded statistics & summaries to public

Label	Date	Type	File System	Total Size TB	Total Space M	# files K	# dirs	max size GB	max space GB	max dir ents	max name bytes	avg file MB	avg dir ents	
Satyendarayanan81	<1981	Home	TOPS10	<0.016	.086							.012		
Irlam03	Nov 1993			.259	12							.022		
SFS97	<1997		NFS									.027	30	
Douceur99	Sept 98	Desktops	NTFS	10.5		141						.079		
VU2005	2005	Home	UNIX			1.7		2				.327		
SFS2008	<2008		NFS					.32		30		.531	30	
CMU gg1	4/10	OS	HFS+	.044	.046	1.0	258	2.1	1.00344	252	.046	.5		
CMU gg2	4/10	Home	HFS+	.0098	.028	3.2	328	328	448	123	.37	10		
CMU gg3	4/10	Media	HFS+	.065	.066	042	2.6	2.2	536	129	.16	.17		
CMU pd1	4/9	Project	WAFL	3.93	11.3	821	37.7	23.4	56960	255	.37	.15		
CMU pd2	4/9	Project	WAFL	1.28	1.09	8.11	694	37.7	23.4	89517	255	.17	.14	
NERSC	4/8		GPFS	107	20.5	917	616	523	143365	152	.5.3	.23		
PNNL nwfs	3/17	Archival	Lustre	265	264	1824	1074	1074	57114	232	19.3	9		
PNNL home	3/17	AdvFS	4.7	4.3	10.1	682	268	35	23556	255	.46	.16		
PNNL dtemp	3/17	Scratch	Lustre	22.5	19.2	2.2	51	1074	8004	89	10.3	.44		
PCS scratch	3/27	Scratch	Lustre	32	32	2.07	451	173	64010	160	15.6	.6		
PCS besemer	3/27	Project	Lustre	3.7	3.7	0.38	15	51	8226	89	9.6	.26		
LANL scratch1	4/1	Scratch	PanFS	9.2	10.7	1.52	120	134	154	14420	90	6.0	.14	
LANL scratch2	4/10	Scratch	PanFS	25	26	3.30	241	978	1076	50000	73	8.2	.15	
LANL scratch3	4/10	Scratch	PanFS	26	29	2.58	374	998	1099	45002	65	10.9	.8	
ARSC sea1	3/13	Archival	SAM-QFS	305	4.3	10.5	326	386	13.7	62803	245	.29	.34	
ARSC sea2	3/14	Archival	SAM-QFS	115	4.6	5.3	116	366	7.0	25008	144	21.7	.47	
ARSC nanu1	3/12	Archival	SAM-QFS	69	4.5	6.7	338	601	13.6	56648	234	10.4	.21	
ARSC projects	3/13	Archival	SAM-QFS	32	.93	6.2	898	171	3.7	24153	81	5.2	.8	

### THE FSSTATS TOOL

- Authored by Shobhit Dayal, CMU and Marc Unangst, Panasas
  - Portable perl code, GPL license
  - One file, built-in doc (perldoc)
- Gathers size, links and age histograms
  - User data, blocks used, directory entries, name length, symlink lengths
- Outputs ASCII histogram, CSV, checkpoints
  - Can restart from checkpoint so long captures can be killed
- Upload the CSV output
- Cumulative distribution graphs show % of “metric” below given size
- Graphs special case “five 9s” at the bottom or top end of scale as appropriate

### RESULTS



- Wide variation in file size distributions between sites
- Huge fraction of files very small even in supercomputers: 20-80% are < 2KB
- Long tail of small numbers of really huge files
- By space used, however, small files consume almost no space
- Very wide variation in file size which divides space in half: 8MB - 16GB

**Better Data Management Comes From Better Understanding of Data Monitor Your Data Systems and Publish their Statistics!**



Sandia National Laboratories

