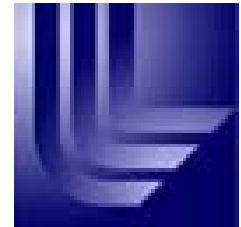


Strategies for Supporting HPC Scientists



**Presented at:
Fall Creek Falls Workshop
Terri Quinn
10/24/06
Deputy Department Head
Livermore Computing, LLNL**



Document # UCRIL-PRES-xxxxxx

**This work was performed under the auspices of the U.S. Department of Energy by the University of California
Lawrence Livermore National Laboratory under Contract No. W-7405-Eng-48.**

Lawrence Livermore National Laboratory, P.O. Box 808, Livermore, CA 94551-0808



Outline



- **Challenges with Petascale Computing Facilities**
 - **Expensive**
 - **Implications and strategies**
 - **Long time to complete one run/experiment**
 - **Implications and strategies**
 - **Complicated**
 - **Implications and strategies**
 - **Lots of Data**
 - **Implications and strategies**

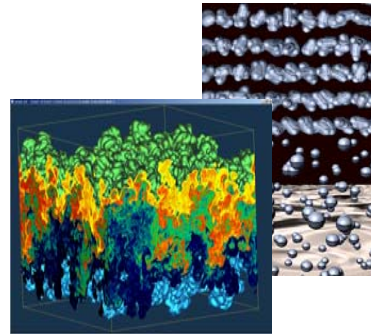
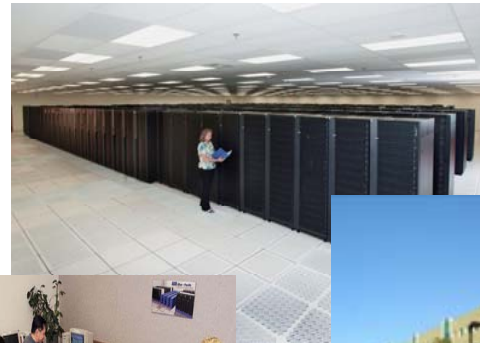




Challenges with petascale facilities



- **Expensive to build and to operate**
- **Often takes weeks or months to complete experiments**
- **Complicated to use**
- **Create extraordinary amounts of valuable data**





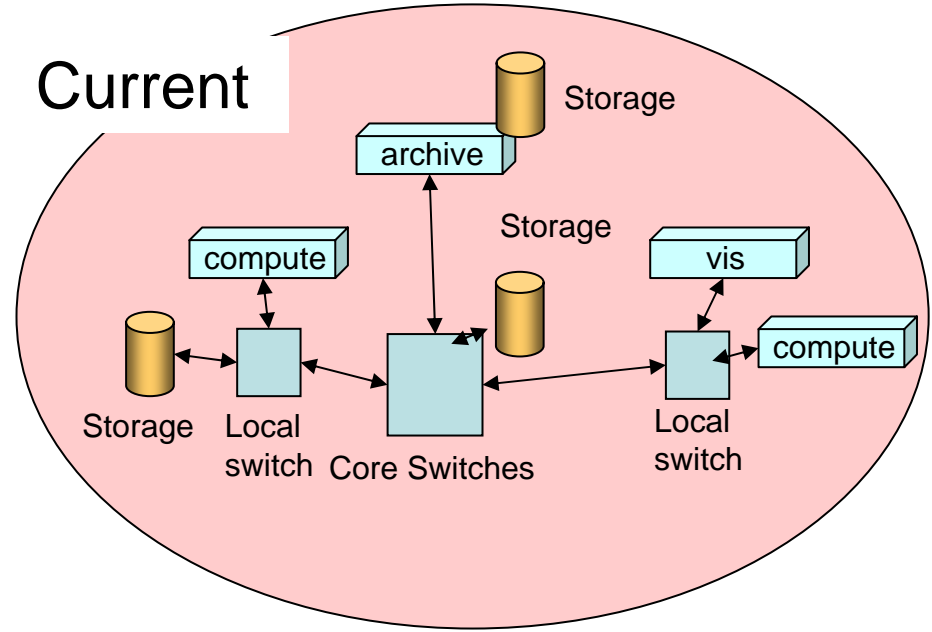
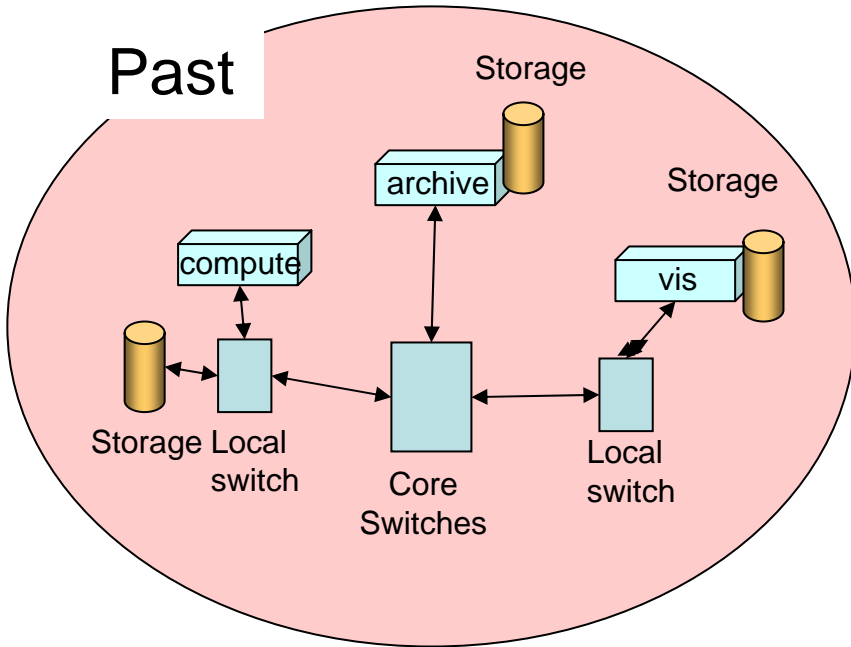
Petascale computing facilities are expensive



- Implications
 - Every \$1 spent on the computing facility is \$1 less for the scientists
- Mitigation Strategies
 - Include total cost of ownership in calculus for selecting the computer – in addition to initial cost consider power, maintenance, and “uniqueness” of the computer
 - Open source vs. proprietary software
 - Take a hard look at how you are doing things: take some risks
 - Look for disruptive technologies to dramatically reduce these costs



Our model for high performance file systems is a single global scalable parallel file system



- **With little or no loss of service our infrastructure costs have decreased**
 - **File system bandwidth need only meet the bandwidth requirement of the most capable system and not the aggregate bandwidth requirement**
 - **Less capacity since users do not need to copy data sets**



Adopting disruptive technologies can dramatically reduce costs



Purple

Peak Speed = 93 Teraflops

Footprint = 12,000 sq. ft.

Power = 4.5MW

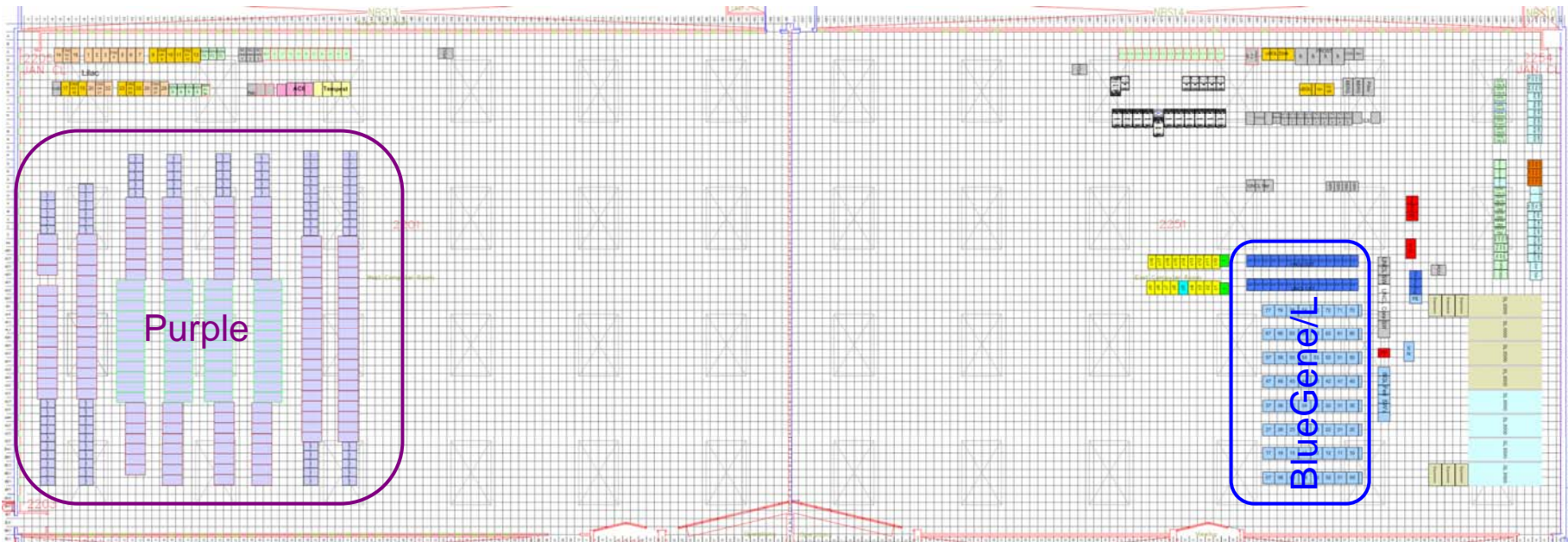
BlueGene/L

Peak Speed = 180/360 Teraflops

Footprint = 2,500 sq. ft.

Power = 1.5MW

Over a 4 year period Purple will cost \$12M more for power than BG/L
Cost per Teraflop for BG/L was much less than Purple's





Petascale computing “experiments” need weeks to months to complete



- **Implications**
 - **Scientists demand that the facility stay up for the duration and that help is available throughout**
 - **Special scheduling is required – facility needs to be capable of being reserved for the “experiment”**
- **Mitigation Strategies**
 - **Implement a reservation process similar to what is done at other more traditional experimental facilities such as observatories and high-energy facilities**
 - **Use highly reliable software and hardware**
 - **Provide 24x7 user support**



BlueGene/L has exhibited outstanding reliability



Reliability Features

- N+1 redundancy on power and cooling components
- Error detection and correction in memory, caches, networks and registers
- Data integrity checks
- Very nice RAS data base for identifying FRUs for preventive maintenance



Reliability Results

- HW MTBF is 6+ days; MTB Application Failure (errors that kill jobs) is 5+ days
- For August and September MTBF was 12.2 days, ~2x the contractual goal.
- This same failure rate on MCR (losing 1 of 64K nodes every 6 days) means a node failure every 8.75 months



Petascale compute systems are complex to use



- **Implications**
 - Scientists need to become compiler experts, file system experts, and architecture experts to make the most out of their effort
 - Complex systems require much care which leads to many “fixes”
- **Mitigation Strategies**
 - Hide the complexity by developing smart tools that can be the experts for the scientists
 - Coordinate the” fixes” to minimize the noise in the system
 - **Minimize the varieties hw and sw configurations**
 - **Invest heavily in user support services**



Transitioning to a new configuration for all our computing needs (serial, parallel, visualization)



- We are rolling in AMD/IB systems to replace Intel/Quadrics, Tru46 and some old SP systems
- Leverages staff expertise and gives the scientist a familiar environment

System	Program	Manufacturer/ Model	OS	Interconnect	Nodes	CPUs	Memory (GB)	Peak TFLOP/s
Classified Network (SCF)								
BlueGene/L	ASC	IBM	Linux	IBM	65,536	131,072	32,768	547.69
Purple	ASC	IBM SP	AIX	Federation	1,532	12,288	49,152	93.39
UM (pEDTV)	ASC	IBM p655	AIX	Federation	128	1,024	2,048	6.14
UV (pEDTV)	ASC	IBM p655	AIX	Federation	128	1,024	2,048	6.14
Ice	ASC	IBM SP	AIX	Colony	80	1,280	1,280	1.92
Tempest	ASC	IBM Power5	AIX	N/A	12	84	480	0.55
<i>Minos (Peloton)</i>	ASC	<i>Appro</i>	<i>Linux</i>	<i>IB</i>	864	6,912	13824	33.18
<i>Rhea (Peloton)</i>	ASC	<i>Appro</i>	<i>Linux</i>	<i>IB</i>	576	4,608	9,216	22.12
<i>Hopi</i>	ASC	<i>Appro</i>	<i>Linux</i>	<i>N/A</i>	76	608	1,408	2.92
Gauss	ASC	GraphStream	Linux	IB	256	512	2,048	2.46
Lilac (xEDTV)	ASC	IBM xSeries	Linux	Elan3	768	1,536	3,072	9.19
Ace	ASC	Rackable Systems	Linux	N/A	176	352	704	1.97
Queen	ASC	Rackable Systems	Linux	N/A	63	126	252	0.71
Klein	ASC	GraphStream	Linux	Elan4	10	20	40	0.14



We invest heavily in customer support to help meet the increasing demands for support



LC customer data and statistics	2000	2001	2002	2003	2005
Number of active users	1959	2219	2304	2450	3189
Classified	964	1080	1231	1360	1655
Unclassified	1609	1826	1847	1904	2474
Remote Users	576	676	648	709	731
Sandia	85	128	133	163	156
LANL	70	88	108	119	139
ASC Alliances	116	122	297	307	93
Other	305	338	297	307	343
Average Number of hotline contacts per day	100	110	116	109	135
Average number of web page hits per day	1080	7145	7308	8675	16721
WWW documentation					
Number of pages of documentation	5131	3798	4315	4331	3878
Total number of Remedy tickets					114703

Customer satisfaction is gauged with bi-annual formal surveys and more recently with quick surveys that we just re-instated



Petascale computations create extraordinary amounts of valuable data



- **Implications**
 - Infrastructure will be stressed
 - Existing data analysis tools and processes won't cut it
- **Mitigation Strategies**
 - We have many approaches that I think have not yet proven themselves
 - Tools that hide complexity
 - Rethink what data needs to be stored and for how long
 - Global parallel file system
 - More ...



DISCLAIMER

This document was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor the University of California nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial products, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or the University of California, and shall not be used for advertising or product endorsement purposes.

This work was performed under the auspices of the U.S. Department of Energy by University of California Lawrence Livermore National Laboratory under contract No. W-7405-Eng-48.