



UPDATE ON NERSC PSched EXPERIENCES, A CONTINUING SUCCESS STORY

Tina Butler – NERSC
Brent Draney – NERSC
Mike Welcome – NERSC
Bryan Hardy – SGI
Steve Luzmoor – SGI

This work was supported by the Director, Office of Advanced Scientific Computing Research,
Division of Mathematical, Information, and Computational Sciences of the U.S. Department of
Energy under contract number DE-AC03-76SF00098.



What is NERSC?

- **National Energy Research Scientific Computing Center**
 - **Funded by DOE Office of Science**
 - **Located at Lawrence Berkeley National Lab**
 - **Provides Computational Resources to the following programs**
 - **Fusion Energy**
 - **High Energy and Nuclear Sciences**
 - **Basic Energy Sciences**
 - **Biology and Environmental Research**
 - **Computational and Environmental Research**
 - **Approximately 2500 Users from Major Universities and Government Labs**
 - **Hardware: 696 PE T3E-900, 1- J90 SE (32 CPUs) & 3 SV-1 (64 CPUs) systems**



Mcurie – The NERSC T3E

- **T3E 900 with 696 PEs running UNICOS/MK 2.0.4.67**
- **644 APP PEs**
- **256 MB per PE**
- **383 GB Swap Space – 5 partitions, each 5-way striped**
- **582 GB Checkpoint file system – 5 partitions, striped**
- **1500 GB /usr/tmp file system**
- **7 – 25 GB Home file systems, DMF managed**
- **All Large file systems “remote mounted”**



NERSC Job Mix – Application Mix

- **Applications from the fields of**
 - **Chemistry**
 - **Materials Science**
 - **Fusion Energy**
 - **Geophysics**
 - **Biology**
 - **High Energy Nuclear Physics**
 - **Climate Modeling**
 - **Astrophysics**
 - **Computational Fluid Dynamics**
- **Mostly user-written codes**



NERSC Job Mix – Diverse and Dynamic

App Size(REs)	% of all Apps	% of PE Hours
2 - 16	56	6
17 - 64	38	56
65 - 128	5	29
129 - 512	1	9

App Run Time	% of all Apps	% of PE Hours
0 – 10 min	56	1
10 – 30 min	23	10
0.5 – 3.5 hr	17	49
3.5 – 12.0 hr	4	40

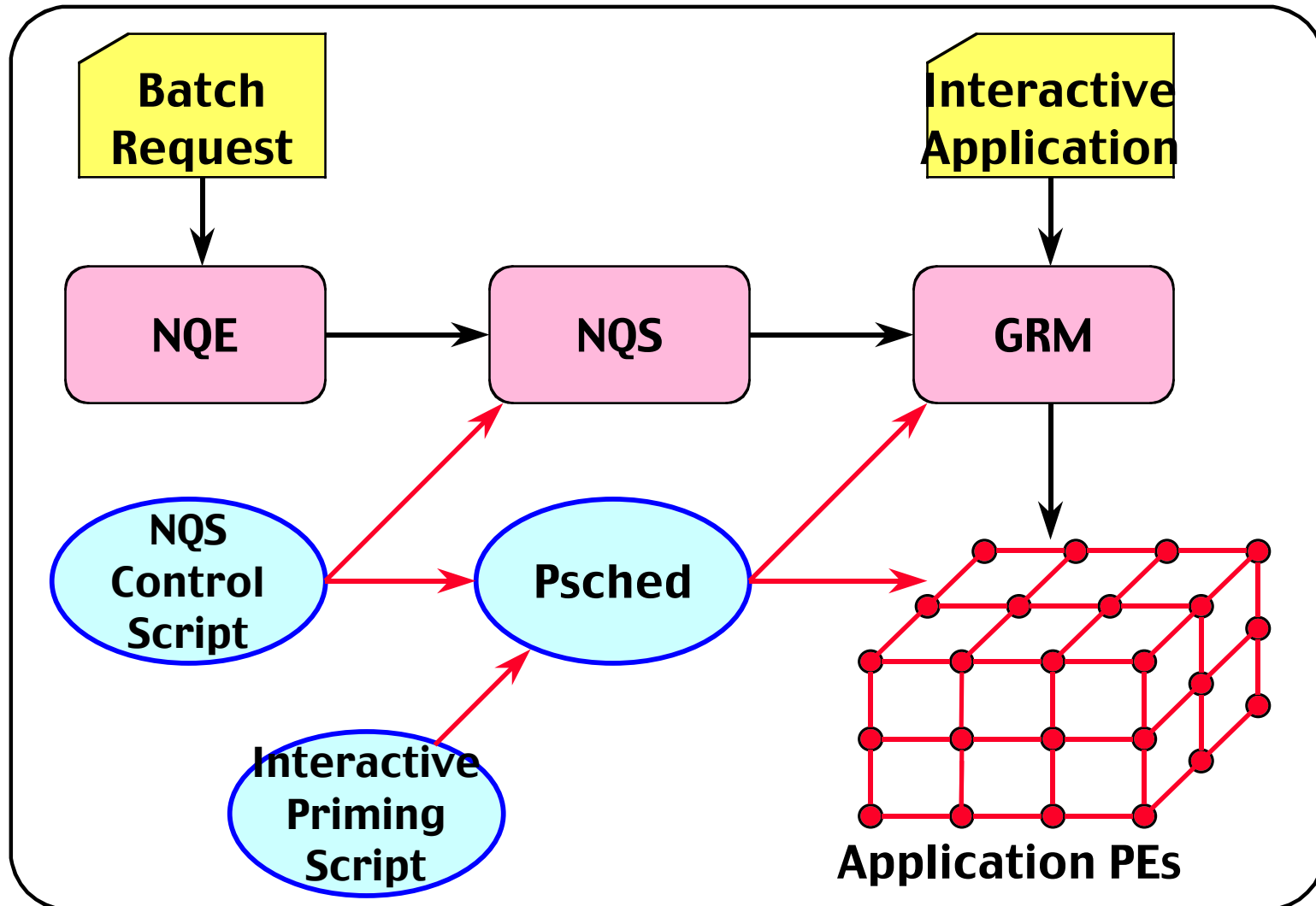
Mix of Development, Capacity and Capability computing



NERSC T3E Scheduling Goals

- **Minimize idle time in the APP region**
- **Provide fast interactive response while managing the total interactive workload on the system**
- **Provide reasonable and even turnaround across all the batch queues**
- **Encourage users to scale applications to large number of PEs**
- **Provide “Priority Queuing” capability via NQE/NQS**

Mcurie Job Flow and Control Diagram





NERSC T3E Batch System

- **NQE – holding pen for incoming requests**
 - **Production Queues: LWS limit of 3 jobs per user**
 - **Debug Queues: LWS limit of 1 job per user**
- **NQS – Queues defined by PE size and Time Limits**

Queue	PE Lim	TimeLim	Priority
Pe512	512	4 hr	45
Pe256	256	4 hr	30
Pe128	128	4 hr	25
Pe64	64	4 hr	20
Pe32	32	4 hr	15
Pe16	16	4 hr	10
Long128	128	12hr	27
Long256	256	12m	28
Debug_md	128	10min	29
Debugsm	32	30min	23



NERSC T3E Batch System (cont.)

- **NQS Control Script (PERL 5)**
 - **Reads configuration file**
 - Contains alternate queue configurations
 - Configuration selection based on time, day of week
 - Which queues are “on”, “off”, “backfill”, etc.
 - Specifies global, complex and queue limits
 - **Gathers system state: parses output of ps, grmview, qstat, psview**
 - **Modifies NQS (via qmgr) to conform with selected configuration**
 - **Uses checkpoint/restart to switch between configurations**
 - Up to 5 checkpoints done in parallel
 - **Logs system state and all actions to time-stamped log file**



Alternate Queue Configurations

Schedule	Configuration	Queue Status
22:00 – 01:00	Full Machine	On: pe512 Backfill: pe64, pe32, pe16
01:00 -07:00	Batch Preferred	On: pe256, pe128, long128, long256, pe64, pe32, pe16, debug
07:00 – 22:00	Regular	On: pe128, long128, pe64, pe32, pe16, debug



Mcurie Configuration Prior to UNICOS/MK

2.0.4

- **GRM – two regions (manage interactive workload)**
 - 512 PE batch-only region (maximum = 512)
 - 132 PE mixed region (maximum = 64)
 - 06:00 – 18:00 weekdays: Interactive-only
 - 23:00 – 03:00 everyday: Batch-only
 - Otherwise: Both interactive and batch allowed
 - app_max = 1, abs_app_max = 1
- **Psched**
 - Two psched domains – one for each region
 - Load balancer enabled
 - No gang scheduler
 - No prime jobs



Mcurie Configuration Prior to UNICOS/MK 2.0.4

- **Problems**
 - Applications launched on region interface
 - Applications launched in “wrong” region
 - Interactive region idle if no interactive work
 - Job size “entropy”
- **Attempted Solutions**
 - Torus-Pack Script
 - De-fragment Script
 - “B-sched”



Mcurie Configuration after UNICOS/2.0.4 Upgrade

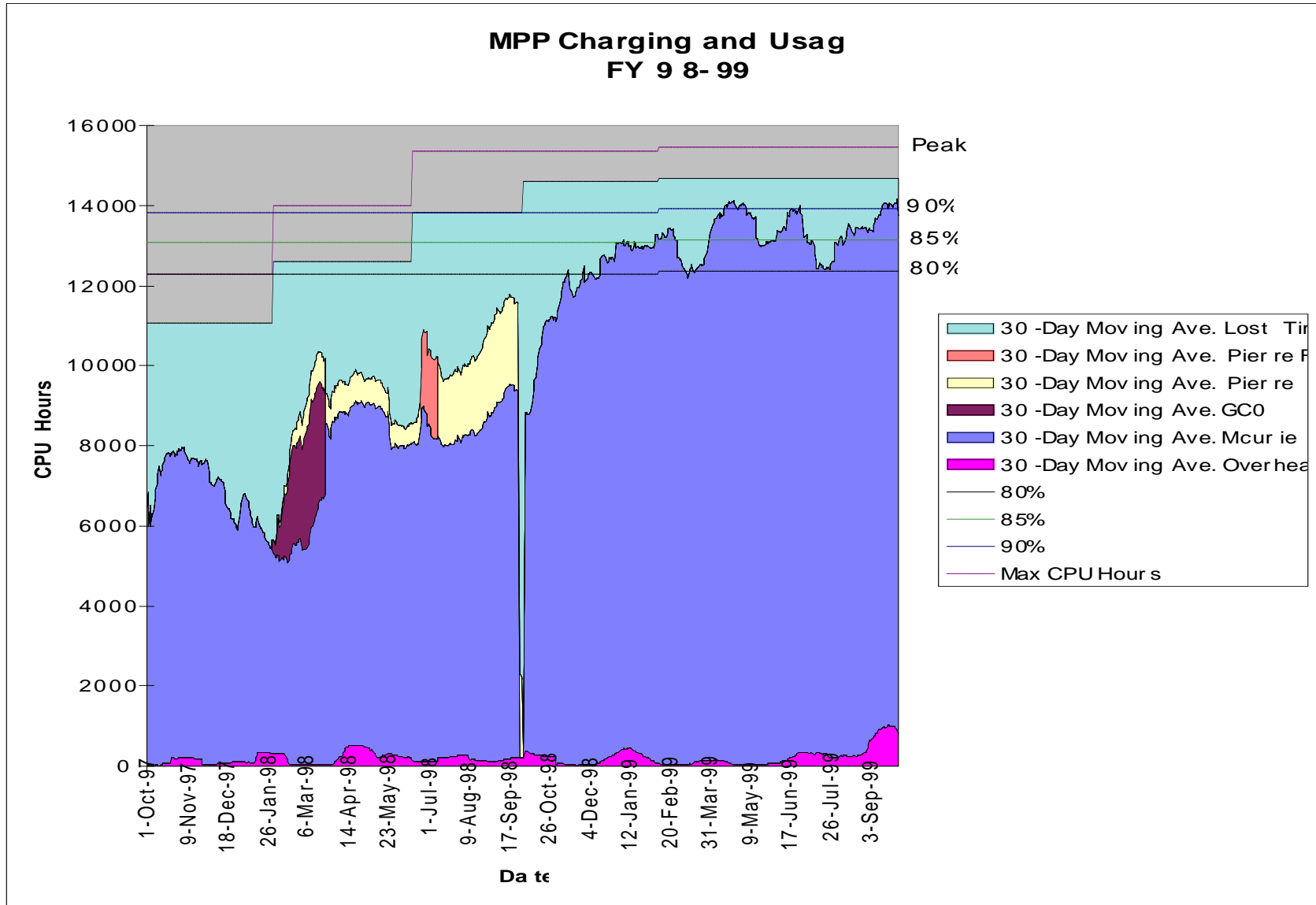
- **GRM**
 - single uniform 644 PE APP region
 - service limits to control interactive workload (132-day/4-night)
 - app_max = 1, abs_app_max = 2
- **Psched**
 - Load balancer – 5 sec heartbeat
 - Gang scheduler – 1 hr time-slice
 - Resource manager – prime jobs
- **Interactive Priming Script**
 - All interactive work is “prime” from 05:30 – 22:00
- **NQS Control Script**
 - Large Jobs run “prime”
 - 30% over-subscription (global MPP_limit=960)

Psched Success at NERSC

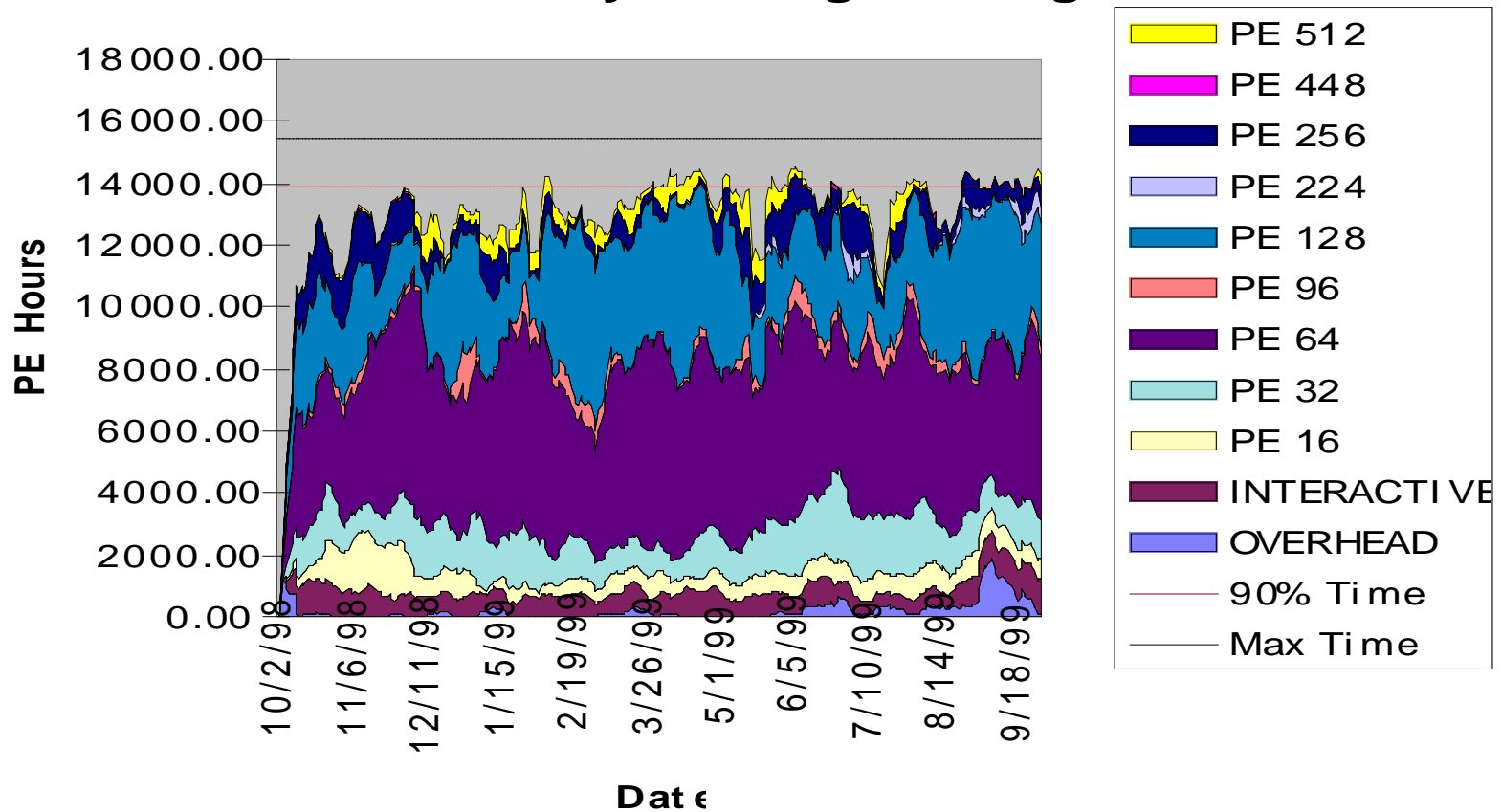
- **Average System Utilization (Connect Time)**

Dates	Utilization	Comments
10/01/98 – 03/04/	79.4%	Prior to 2.0.4
03/05/99 – 03/24/	85.6%	Post 2.0.4
03/25/99 – 05/08/	90.2%	Current Configuration
05/09/99 – 09/30/	87.3 %	Allocation Problems

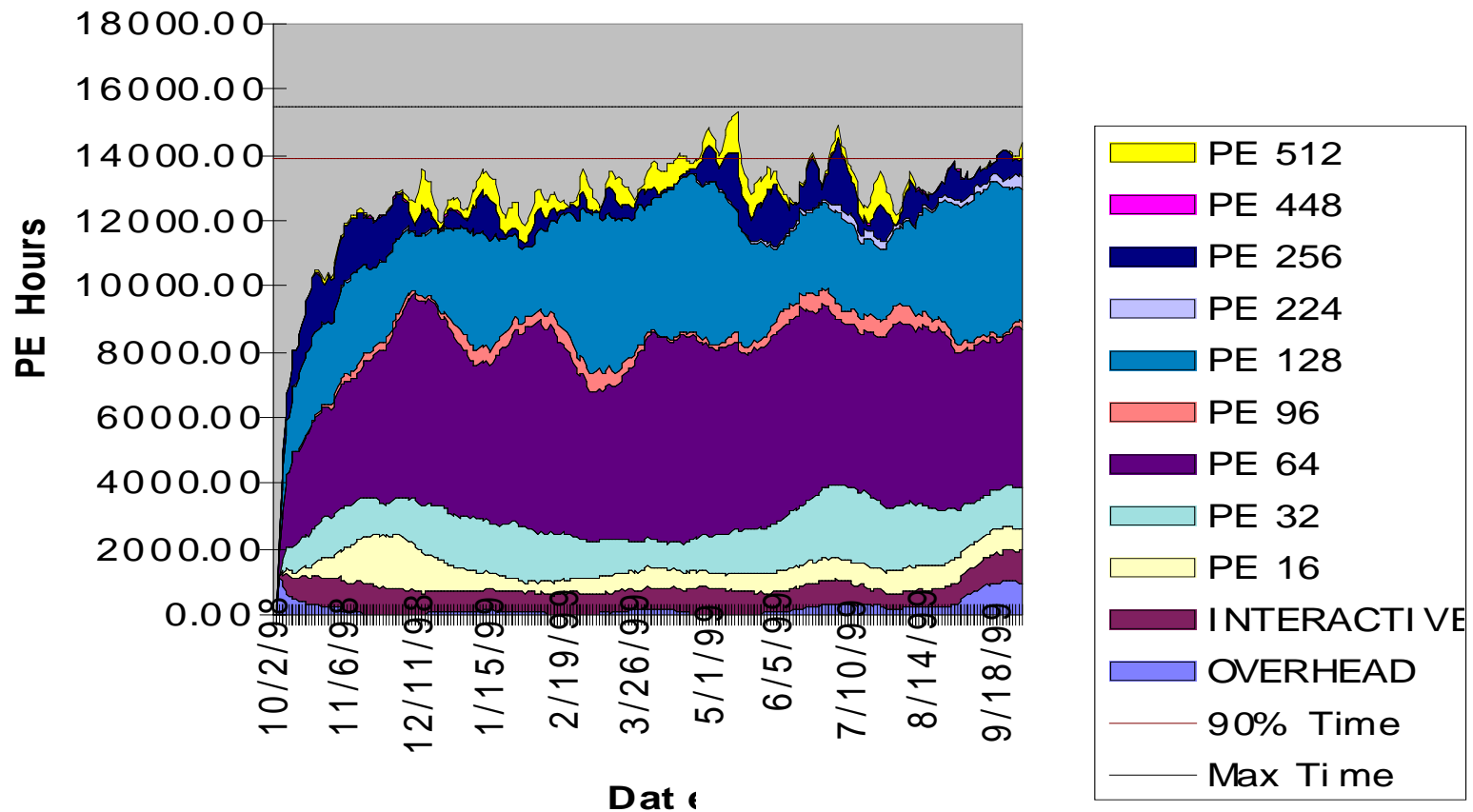
- **Average queue wait time**
 - reduced
 - decreased for large queues
- **Interactive workload**
 - restricted but given priority



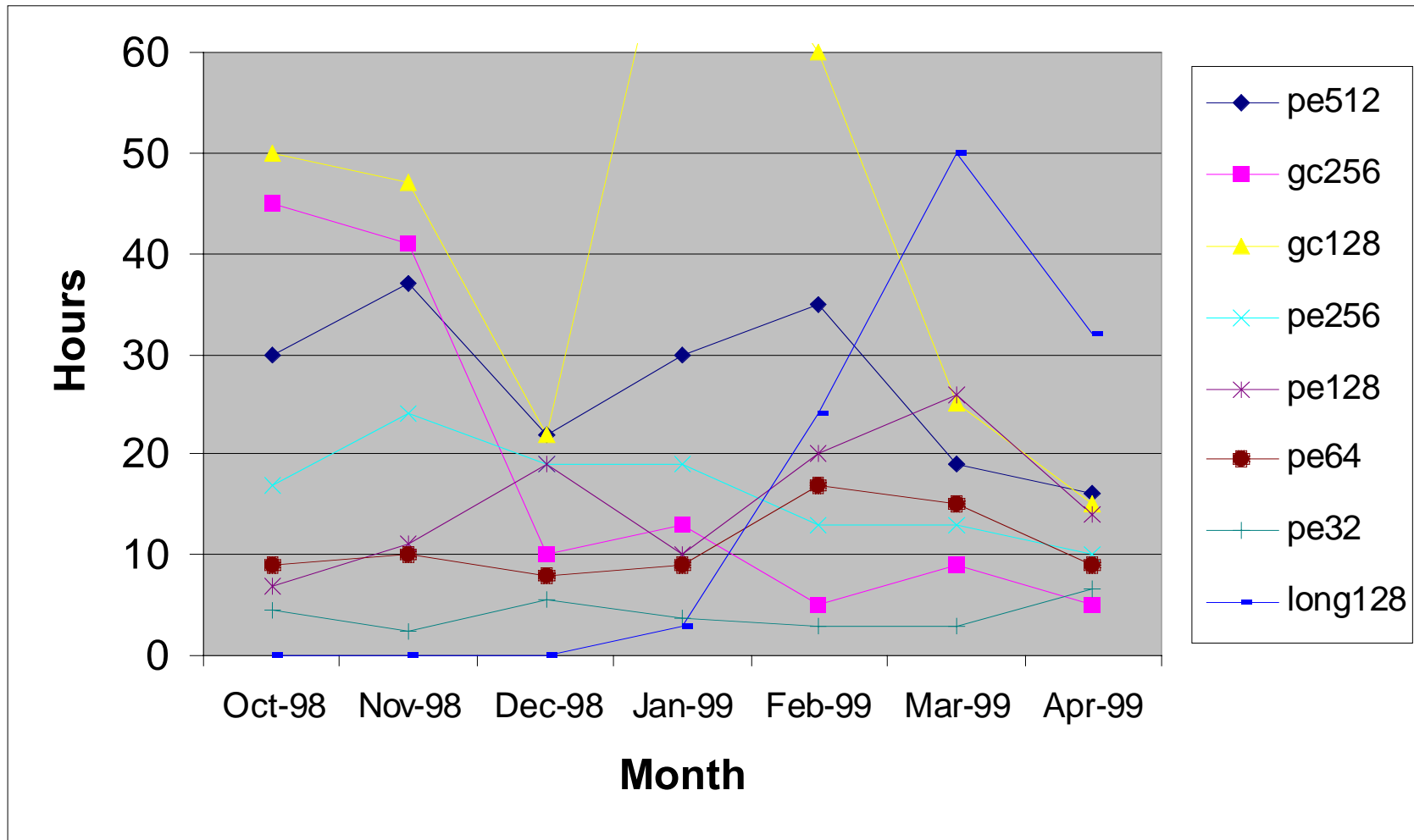
Mcurie Connect Time by Applicatic Size 7-Day Moving Average



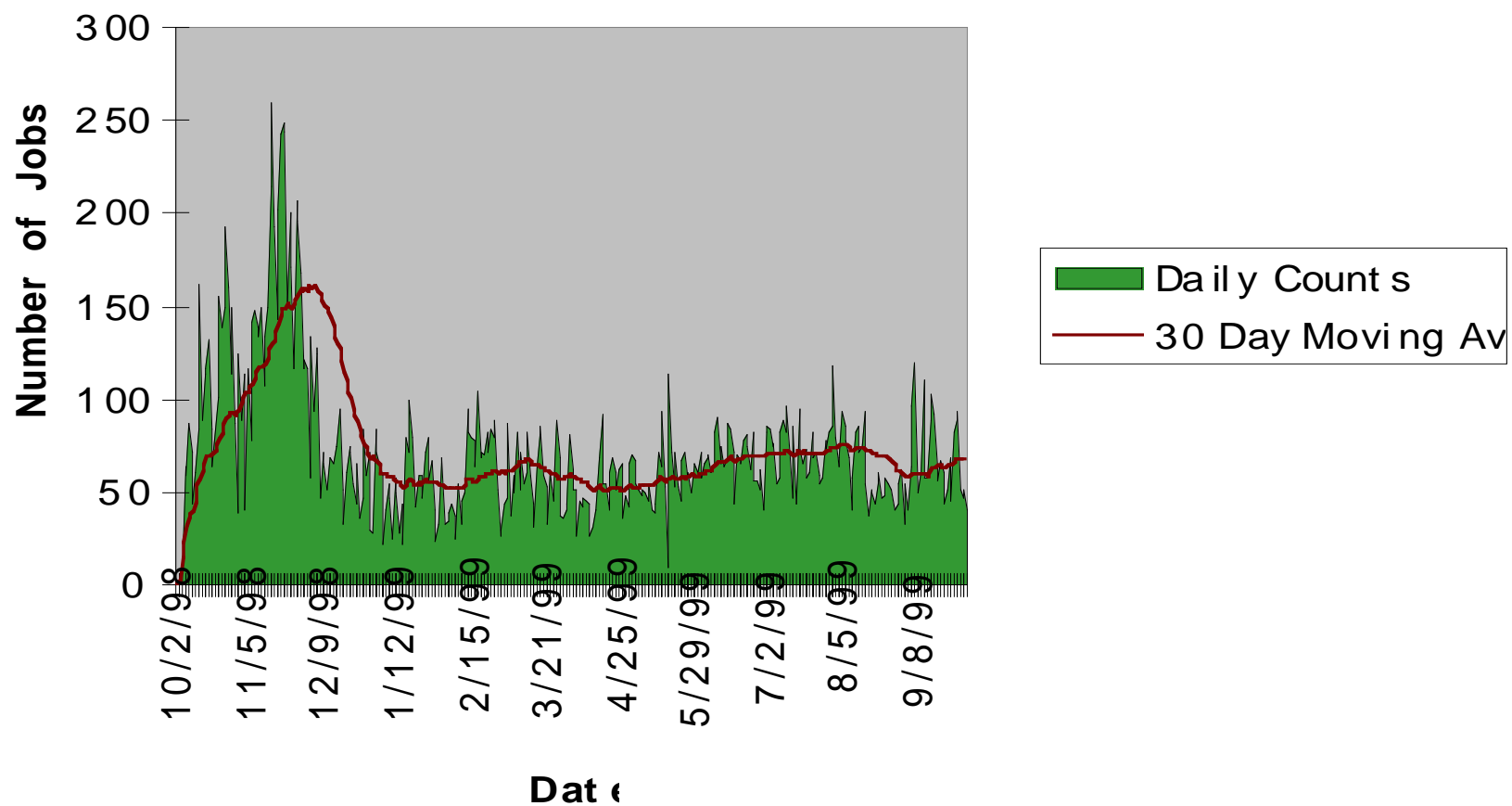
Mcurie Connect Time by Applicatic Size 3 0 Day Moving Average



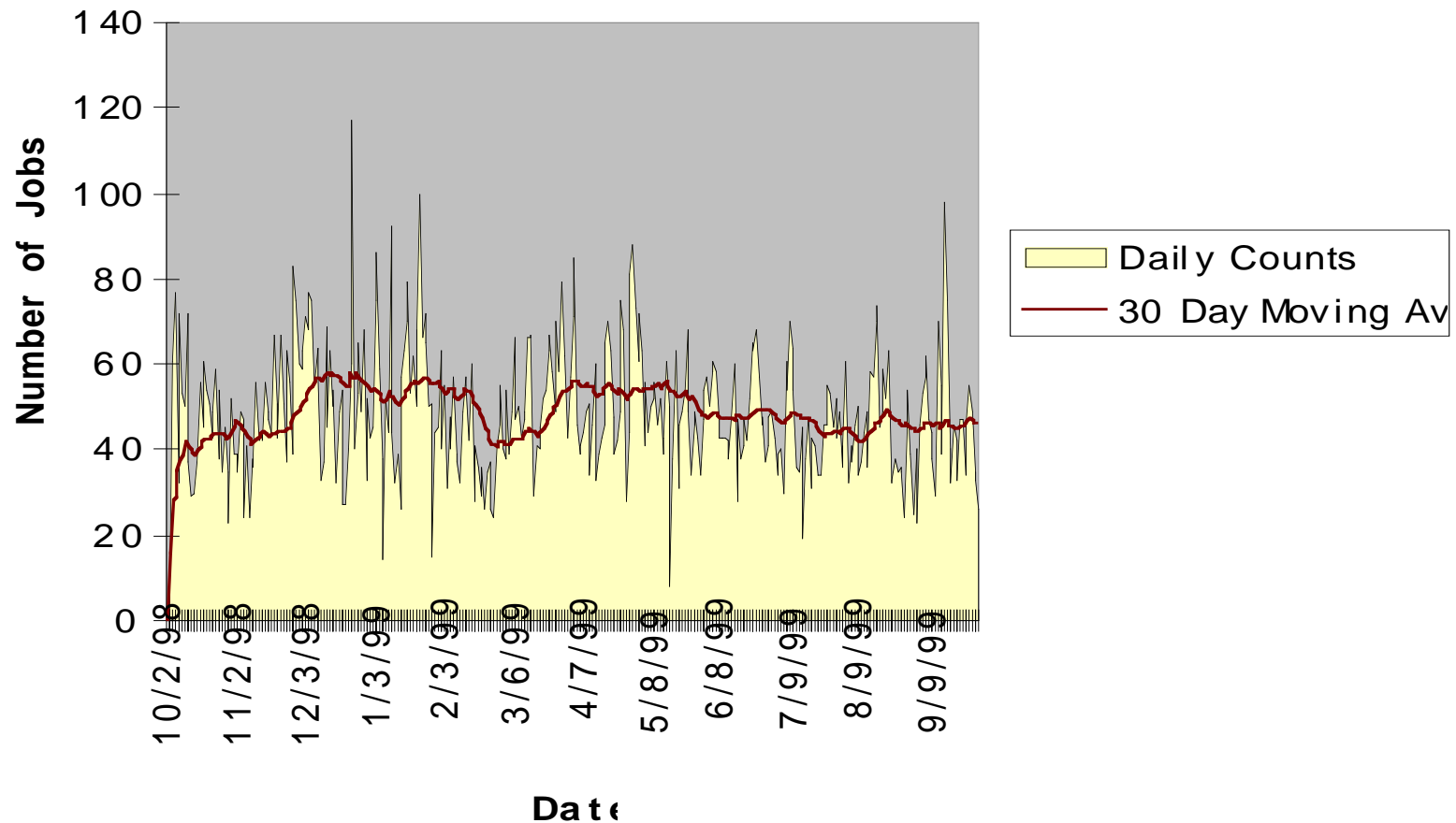
Mcurie: Average Wait Time per Queue



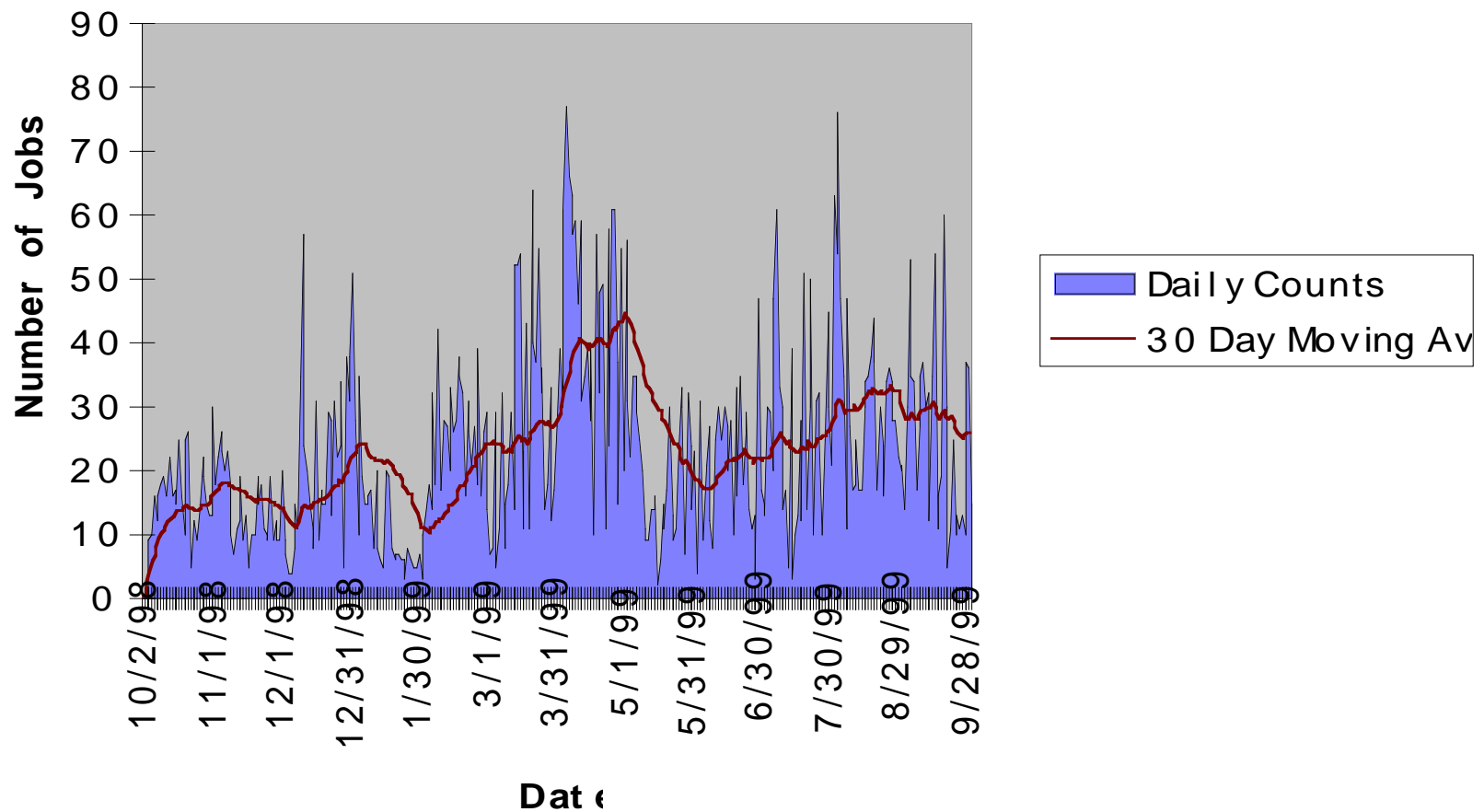
Mcurie Production Jobs less than 33 PE's



Mcurie Production Jobs 33 to 96 PE's



Mcurie Production Jobs Greater than 96 PE's





Conclusions

- **Psched has been very stable**
- **GRM Service Limits are an effective means of managing the interactive workload**
- **Prime job feature is an effective tool for**
 - **providing quick interactive response**
 - **scheduling large jobs**
- **System management is simplified**
- **Utilization is high**
- **Too early to declare victory with “priority queuing”**