



# *BaBar Simulation Production*

BaBar simulation production – a millennium of work in under a year.

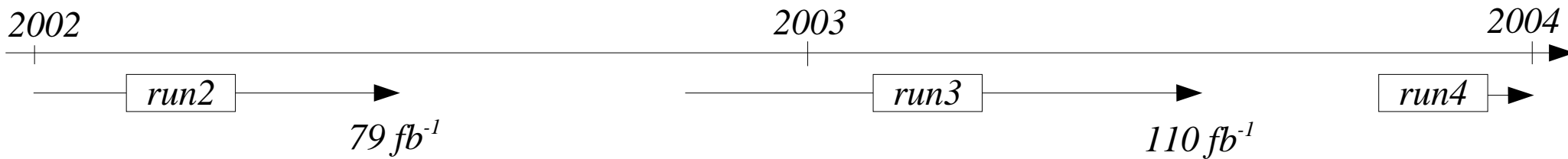
D. A. Smith, F. Blanc, C. Bozzi  
for the BaBar computing group.

CHEP 2004 - [339] - Sep. 29, 2004

# *The Effort : Needed simulation*



*BaBar data schedule:*



- ◆ Simulation production needed for physics:
  - ◆ Generic B – Bbar --> 3 times luminosity
  - ◆ Generic continuum --> same as luminosity
  - ◆ Signal decay modes --> as requested for analysis
- ◆ Simulation Cycles (each with new code):
  - ◆ SP4 --> run 1 and 2 --> 1.2 billion events
  - ◆ SP5 --> run 1-3 --> 1.6 billion events
  - ◆ SP6 --> run 4 --> 1 billion events

**This talk will focus on SP5 as a complete large scale computing effort.**

# Computing resources



- ◆ In looking at resources, there are always efficiency issues, we hoped to be able to get more than 80% use of cpus. Not all events produced were good, certain sets of events were recreated. And people always want more than stated.
- ◆ Numbers assuming a 1GHz pentium III:
  - ◆ *One event:*
    - ◆ 1 event                   --> 3-10 sec.           --> 30-45 kB
  - ◆ *Average over decay modes:*
    - ◆ 1 event                   --> 8 sec.               --> 40 kB
  - ◆ *Requested SP5 production:*
    - ◆ 1.6 billion events --> 420 years       --> 61.5 TB disk
  - ◆ *Actual SP5 production:*
    - ◆ **2.2 billion events --> 700 years       --> 84.5 TB disk**

# Computing Jobs

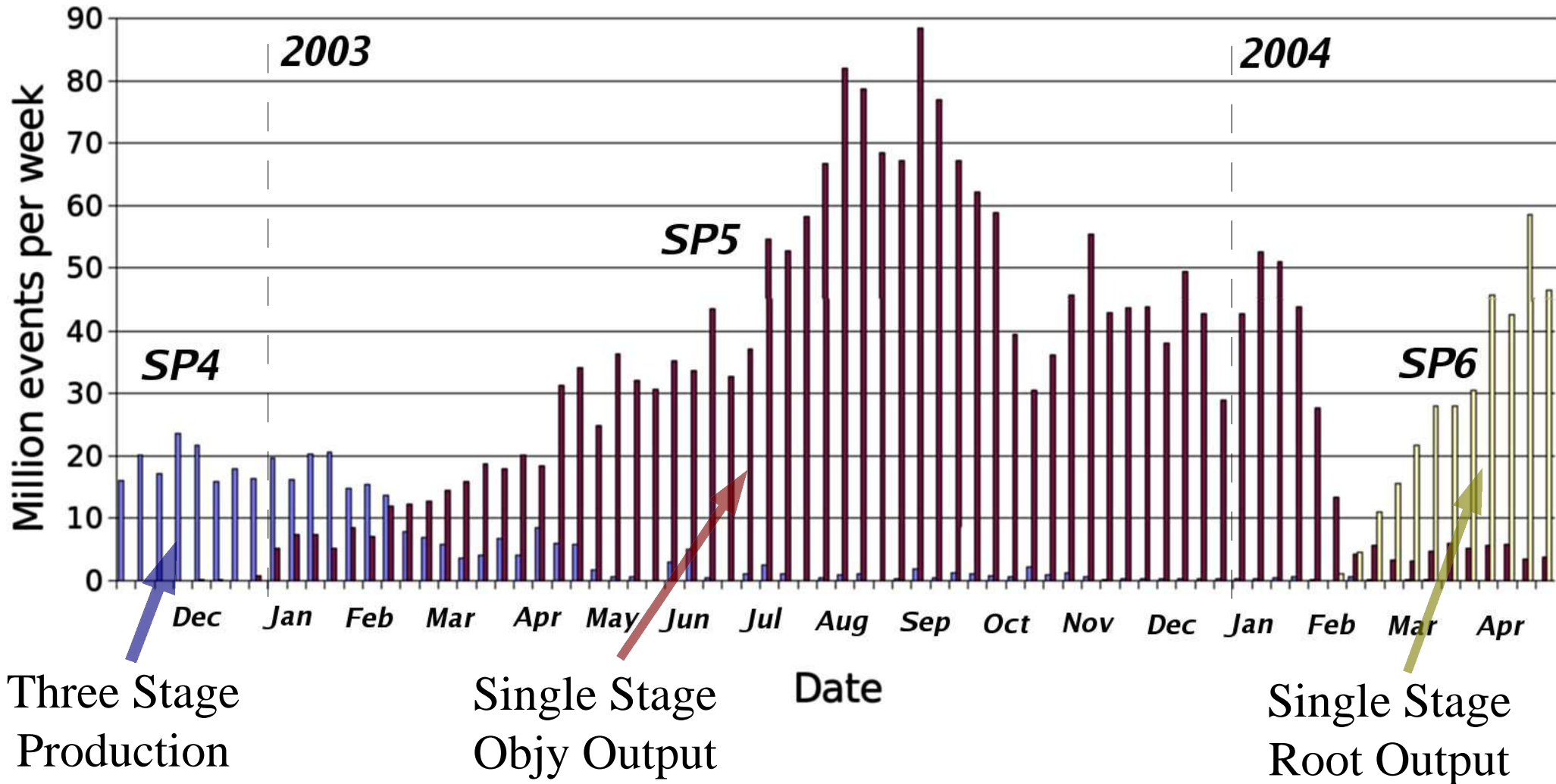


- ◆ The events are divided up into computing jobs, and each job will produce 2000 events.
- ◆ On the 1GHz pentium III machine, the computing job will take 4.5 hours on average.
- ◆ BaBar software supports Solaris and Linux systems, but 90% of SP4 and all of SP5 and SP6 were produced on Linux.
- ◆ **In SP5 The effort was managing 1.1 million jobs.**



# Production History

## Simu Production by Week



Three Stage  
Production

Single Stage  
Objy Output

Date

Single Stage  
Root Output

# *How was it to get done?*



- ◆ SLAC had 1000 cpus (1.4GHz pentium III), but they were needed for other things, like data processing and analysis.
- ◆ BaBar is a distributed collaboration, and there was a stated desire to distribute more efforts out to institutions. Simulation was a good candidate
- ◆ The effort would get done using computing farms at any institution that could participate. This was started in earlier production cycles, and the model worked, but by SP5 more farms and sites would be needed.
- ◆ Had to be careful with effort when multiplying sites, the standard became **one half time person per production site**.

# *Management of jobs*



- ◆ A set of tools, called “ProdTools” were developed to manage the jobs. This consisted of a set of command line tools which would interact with a production database and the local batch system.
- ◆ A database was maintained at SLAC for central management, and all sites interact with this one database over the network.
- ◆ Failures always happen, and tools need to respond to reduce effort (4-6% failure rate in SP5). Also jobs failures can't hang the management tools.
- ◆ The standard was that even though there will be failures, the tools will respond correctly for all but 1 in 10,000 jobs. In SP5 there were over 1 million jobs, about 50,000 of these will fail, tools will fix the failures for all but about 100 of these jobs.
- ◆ **Lots of automation and error checking.**

# *Transfer of Data*

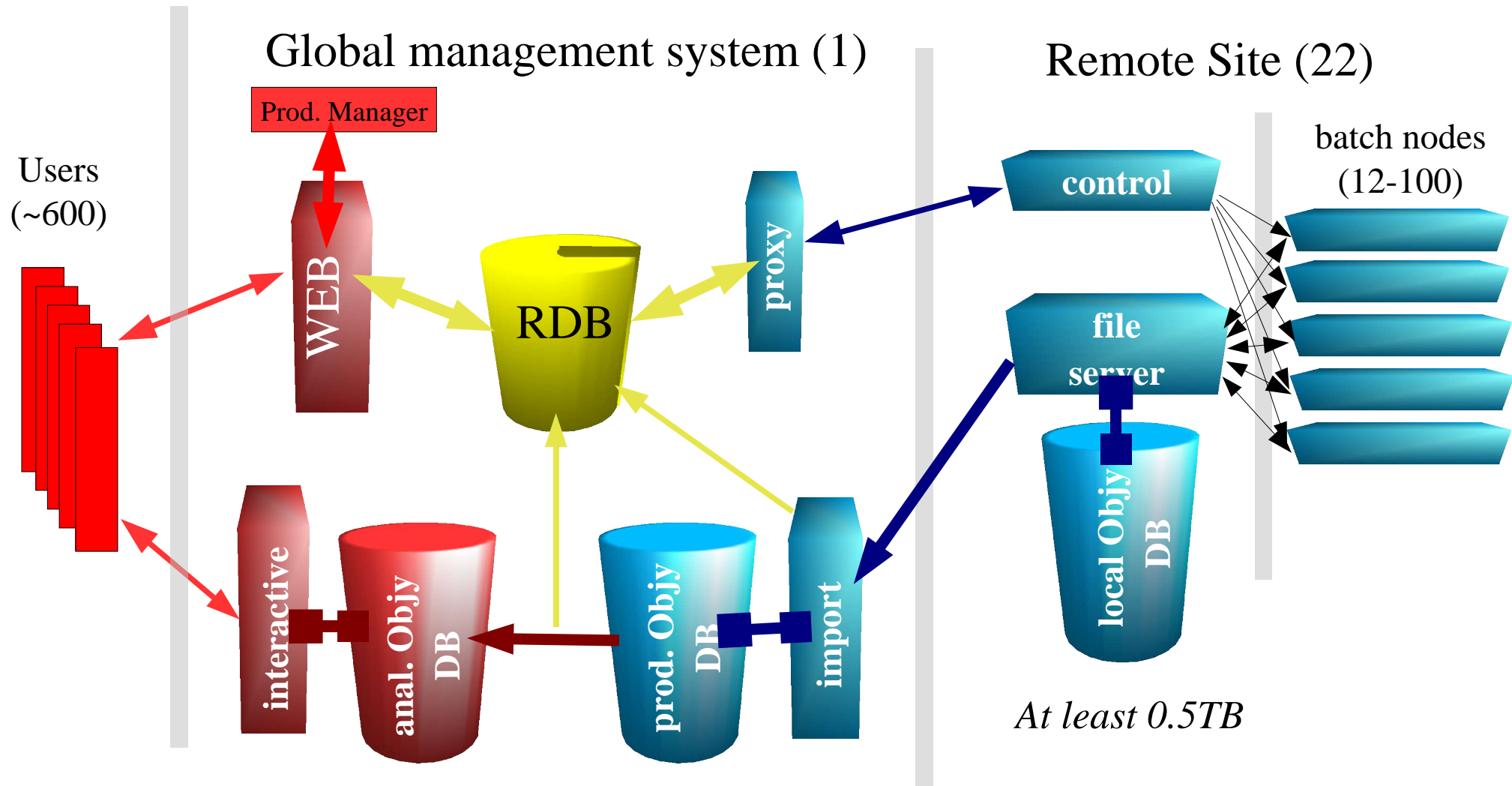


- ◆ Produced events went into Objectivity databases, and a tool “MocaEspresso” was created to manage the distribution of these back to SLAC.
- ◆ Total data created about 80TB, this required about 200GB transferred over the network each day. This varied over the year, with a max of about 500GB transferred each day.
- ◆ This required a farm of file servers specify for transfer to keep up with the rate, and specialized tools created at SLAC.





# Resources at sites



# *Resources for each site*



- ♦ Sites setup were very diverse. Many were just as needed : one server, one disk array, and ~32 dedicated dual cpu machines. But others were large shared batch farms, jobs running as background to other production efforts, multi-use disk servers, and so on.
- ♦ Many different batch systems, tools needed to be flexible. In most cases, the system could not define what resources and systems would get used at a site (had to accept LSF, PBS, BQS, SGI, Codine, and more)
- ♦ If you were in BaBar and you could dedicate some number of cpu, we would try to find a way to use you, and tools needed to be flexible enough to allow this.

# *Production Sites*



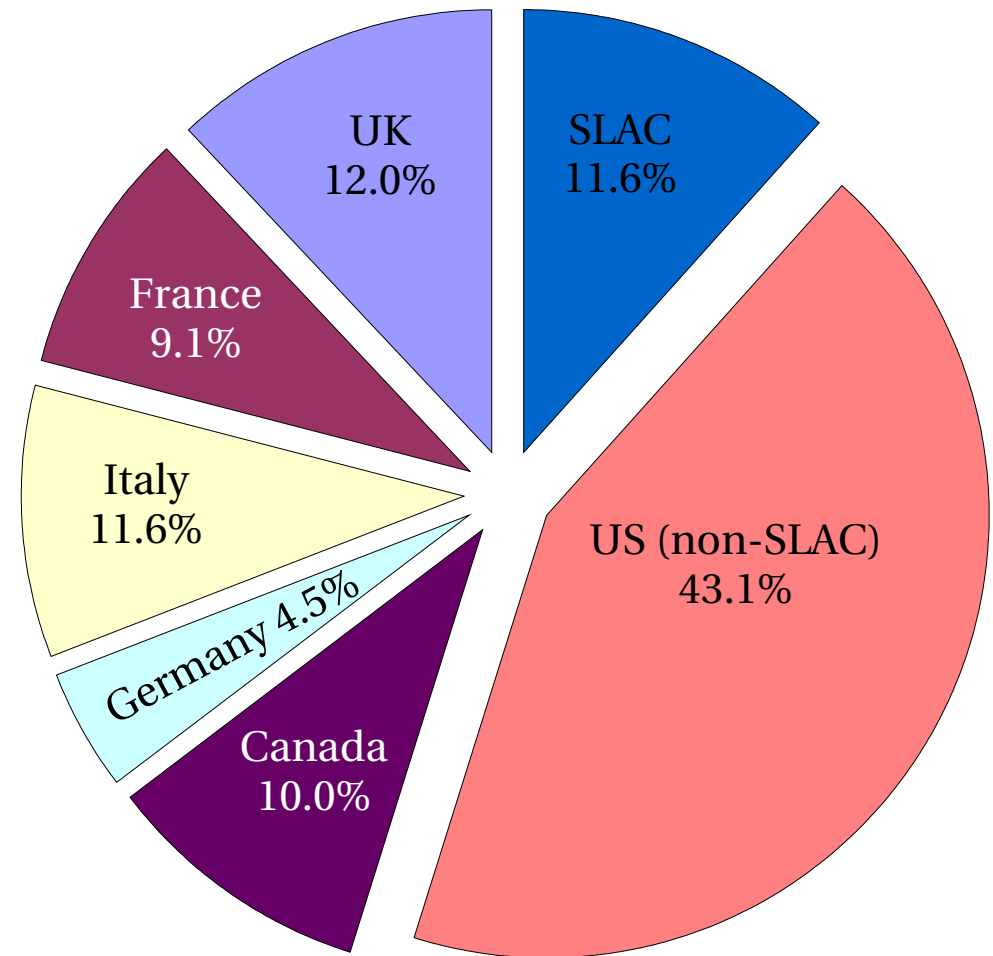
- ◆ Total of 22 sites setup, in 6 countries, on 2 continents:
- ◆ **North America:** Cal. Tech, CO State Univ., CO Univ. at Boulder, Iowa State Univ, Ohio State Univ, SLAC, SUNY Albany, Univ. Tech. Dallas, Univ. Tenn., Univ. Victoria, Vanderbilt
- ◆ **Europe:** Birmingham, Bristol, IN2P3, FZK, INFN, Liverpool, Queen Mary, RAL, Royal Holloway, Scot Grid, Tech. Univ. Dresden



# *Distribution of Production*

- ◆ Work nicely distributed among sites.
- ◆ Sites becoming better balanced at each cycle. SLAC was 85% in SP3, and 35% in SP4.

## SP5 Production by Country





## *SP5 – the continuing story*

- ◆ There was an epilogue to the SP5 story this year. Over the past year BaBar changed the event store from Objectivity databases to Root files. This required a conversion of all SP5 data to the new event store format.
- ◆ This conversion would take about 1 sec. per event, and only the requested 1.6 billion events would get converted.
- ◆ A farm of ~1400 cpu was put to the conversion task, and shared with other conversion needs. This was done only at SLAC, to save the over-head of distribution and training a team of people. (I did almost all the conversion myself.)
- ◆ Took about 6 weeks in the end, with a max of a 1200 jobs running at the same time, converting on ave. 40 million events per day.



# *SP6 improvements*

- ◆ SP6 almost finished (already matched minimum request).
- ◆ Include BaBar's Computing Model 2 changes, from an objectivity database output to ROOT file output.
- ◆ This provided a lot more control over production, lower server load, less overhead in job startup and event writing.
- ◆ We were able to reduce the failure rate from 4-6% to now 0.2-0.5%.
- ◆ With lower failure rate and no objy. server management, the effort for production was also reduced, from 50% for each site to now 10%.
- ◆ Report from production managers, things are now stable enough to usually leave running for a week.





# *Comments on GRID use*

- ◆ Large scale distributed computing effort, perhaps should be a GRID project, but mostly it is not.
- ◆ GRID tools are being used as part of production: Two projects to do this exist, one in the U.K., and one in Italy, using different GRID software and methods.
- ◆ There is a poster on the Italy GRID project – C. Bozzi – id 129, 10:00 Wed.
- ◆ These projects can only handle a small amount of needed production. GRID tools do not have a large enough installed base, are not stable enough at this time, and require more effort than current tools in use.
- ◆ GRID production has been proven to work, with heroic efforts. But for now, not the whole answer for BaBar production.



# *Conclusions*

- ◆ Simulation of events for BaBar analysis is a huge computing effort, requiring over 1000 cpu, distributed throughout the world.
- ◆ Even though this is a large and difficult effort, it was done, and with a reasonable number of people, and provided needed data on time for analysis.
- ◆ Good tools are needed to reduce effort, be robust to handle failures, and stable enough to not hang or fail over nights and weekends so production managers can also have a life.
- ◆ Things continue to improve, sites continue to come on-line and add more cpu, things still scale, and we look forward now to SP7 starting this fall.