



Computing, STAR and the RCF

Facility Science and Technology DOE Review of the Relativistic Heavy Ion Collider

July 6-8th 2005

Guidance

- Priorities & Prospects
- Status and Performance analysis using RCF
- Accomplishments
- Productivity
- Issues

Priorities

- Our activities = Our priorities
 - Deliver quality data to the Physicists for quality science

Data mining, data production

- Evaluate, plan, and integrate new technologies, methodologies, computational techniques designed to better achieve program goals

CS R&D (calib techniques, tracking ...)

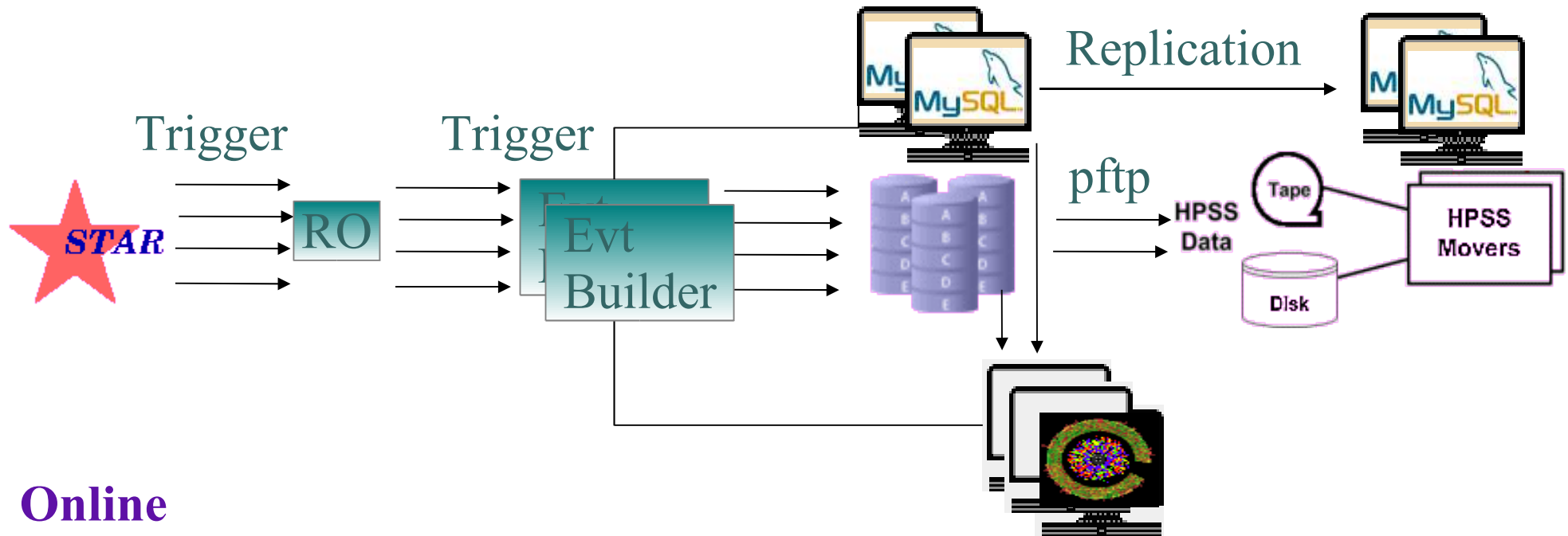
- Develop an environment fostering collaboration with others and welcoming outsourcing

Grid computing ...

- Support for our user community & analysis

Seems a bit “scholar”, order reflects priorities

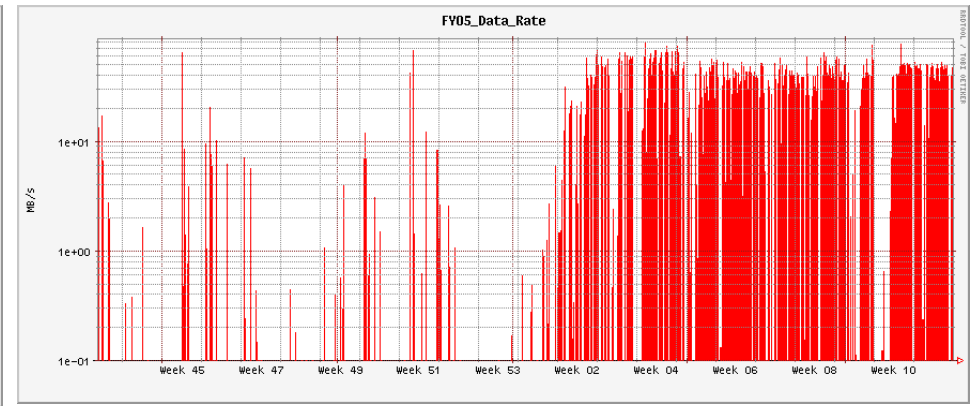
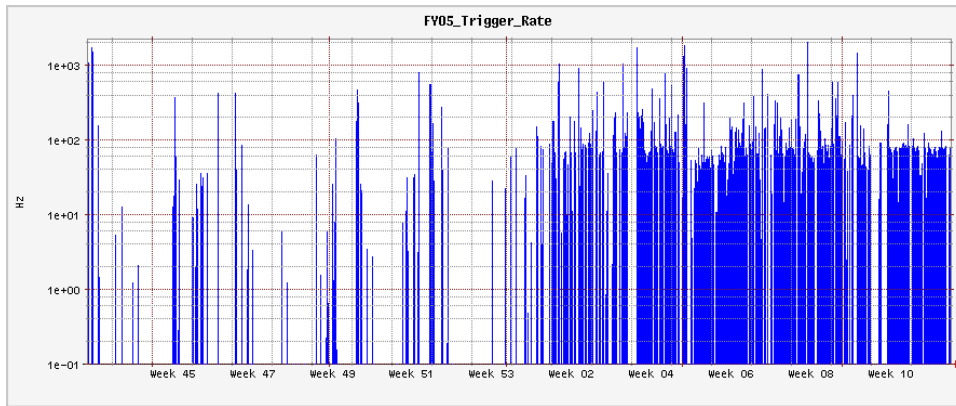
Online data flow



Online

- Event “Pool” based framework
 - “Standard” approach - used by most modern experiment
 - Designed to improve IO online (striped “cheap” disk)
- Data is pushed to HPSS (offline realm)
- Fraction used online to perform fast calibration / analysis

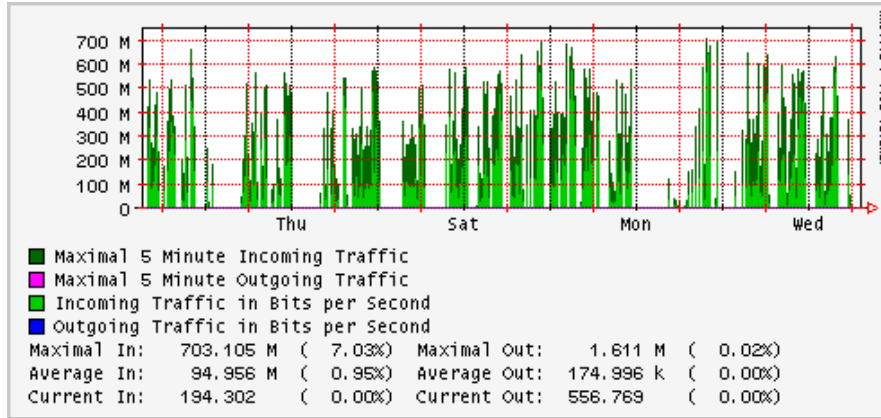
DAQ IO Rates



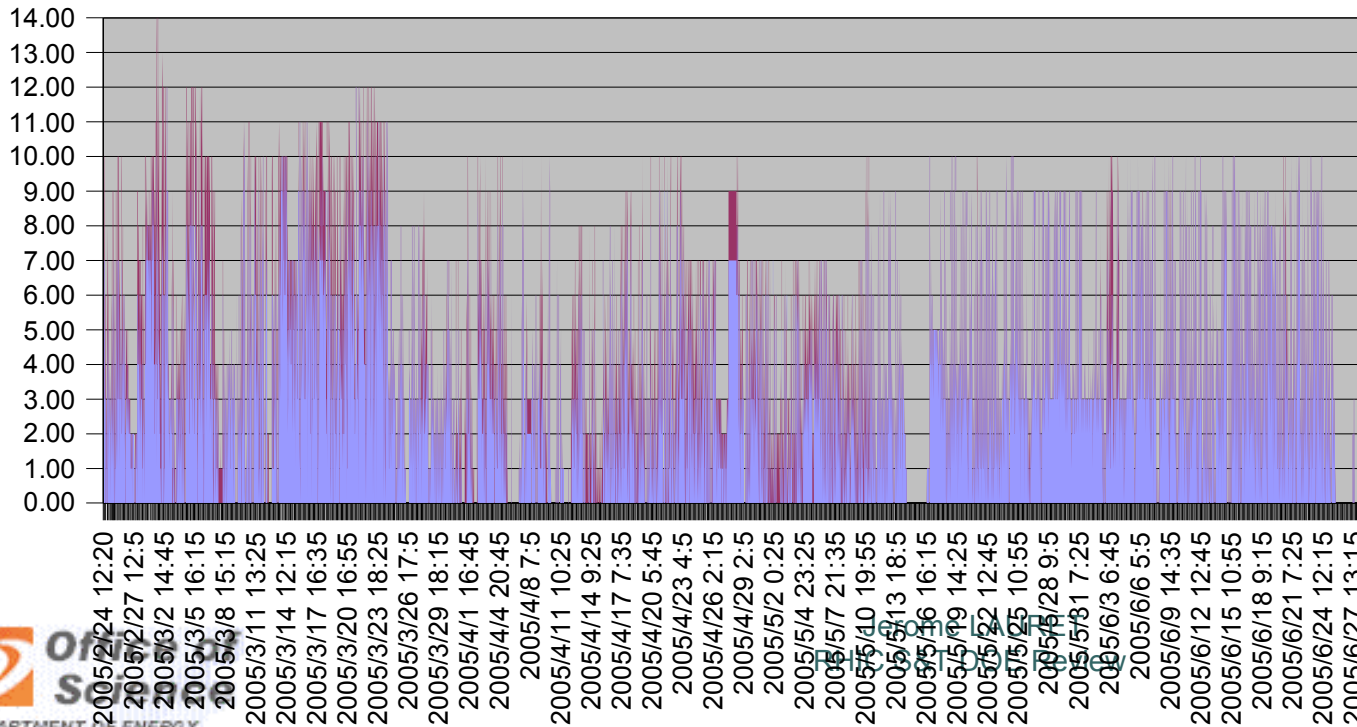
Rates ~ 100 Hz, 50-60 MB/sec sufficient to cover for the data rates and needs FOR NOW

Later program requires x10 rates (DAQ1000 – R&D)

Rates to HPSS



HPSS support for our DAQ/Raw data and network is more than adequate ...

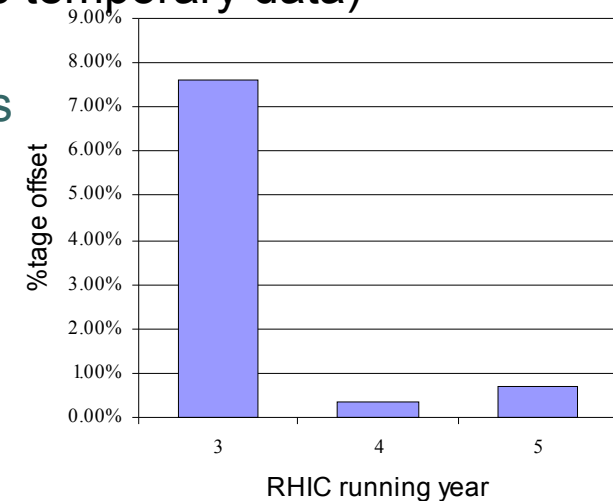


Number of tape drives maximum 10 is sufficient

Data safety

- Accounting reveals some minor losses – Two sources
 - Online
 - Year3+HPSS failures not handled and not enough buffer, miss-accounting (fraction reported missing was temporary data)
 - Corrected immediately, **now 99.3% safe**
 - Remaining losses due to hardware losses
 - Offline
 - Aging HPSS – % files regularly lost
 - Remaining un-identified losses, safety greater than 99.93%

- Some other problem with HPSS
 - Periods with reduced access concerning for later program scaled up need
 - Problems are however generally addressed and resolved with method by qualified RCF personnel



Data Sets sizes - Year4



- Raw Data Size
 - $\langle \rangle$ ~ 2-3 MB/event (HPSS)
 - Needed only for calibration, production – Not centrally (NFS) or otherwise stored
- Real Data size
 - Data Summary Tape+QA histos+Tags+run information and summary: $\langle \rangle$ ~ 2-3 MB/event
 - Micro-DST: 200-300 KB/event
- Total Year4

How long ?

Trigger	Total month	Remains	FF (* DF)
production62GeV	128.06	19.11	0.49
pp	10.33	0.84	0.04
ppMinBias	13.76	0	0.05
ProductionPP	74.91	1.41	0.28
ProductionPPnoBarrel	8.29	0.51	0.03
ProductionPPnoEndcap	1.53	0	0.01
ProductionCentral	18.42	17.09	0.07
ProductionHalfHigh	23.52	1.37	0.09
ProductionHalfLow	161.75	2.52	0.61
productionMinBiasHT	0.55	0.55	0.00
ProductionMinBias	395.24	69.63	1.50
ProductionHigh	267.50	33.53	1.02
ProductionLow	1362.49	1306.73	5.17
ProductionMid	328.13	34.77	1.25
			9.36
			10.54

Since Year4 for RHIC experiments moved out of the fast production turn around mode ...

Year scale production cycles

Before prod - Calibration notes

- Never under-estimate its importance
 - STAR is not only a large acceptance, multi-purpose detector with a TPC at its heart
 - Rule #1 Whatever can go wrong WILL go wrong
 - Rule #2 When you think you have it under control, something else comes up
- We got it all (Field distortions, twist, pile-up, ...) and survived
 - **But we “lost” our dreams for immediate data usability a while back ...**
- Typically, pre-production pass requires ~ 10% of the data pre-processed BEFORE a big production wave
 - 1/2 for TOF
 - SpaceCharge, beam line, drift velocity verification
 - dE/dx, SVT & FTPC alignment

Offline - Production model ...

As fast as possible

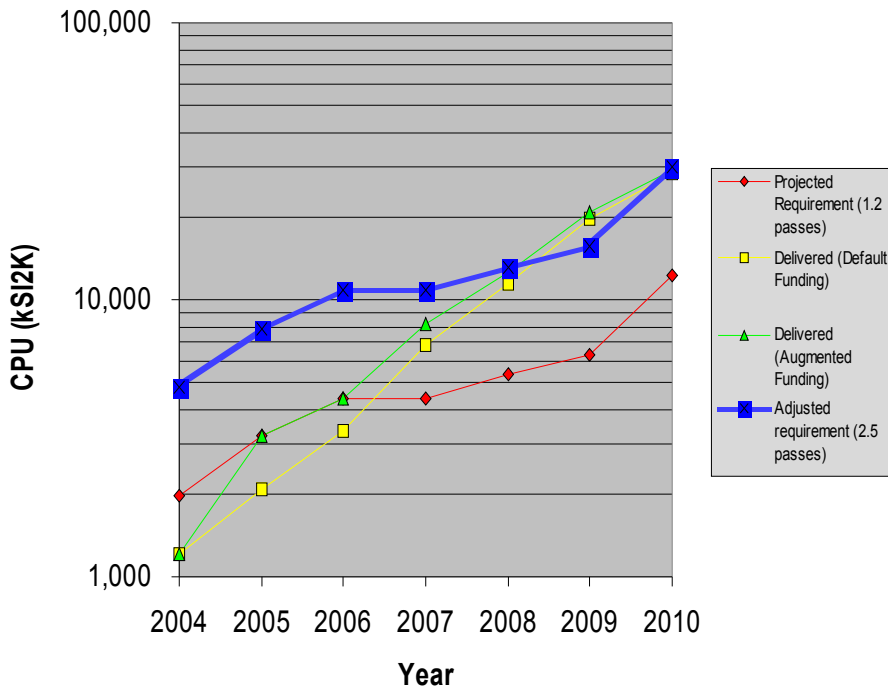
- Centralized – Tier0 BNL data production
 - Note: User Analysis balanced on Tier0/Tier1
 - ALL REAL DATA produced at BNL. EVENT files get copied on HPSS at the end of a production job
- Achievable during the run
 - Online: QA & Fast Online Calibration
 - Offline: Fast-Offline production
 - Fraction of the data ; up to 5-10 % processed
 - TPC Laser runs identifiable (naming convention) all processed
 - Automatic calibration of TPC drift velocity, offline QA, calorimeter, TOF, FTPC, ...
- Re-distribution
 - When production done, system is automated
 - If “sanity” checks (integrity and checksum), files become immediately available to the end-user
 - **30 seconds after the file is produced at Tier0**
 - **30 mnts to Tier1 (PDSF) – Strategy implies dataset IMMEDIATE replication**

Past resource estimates

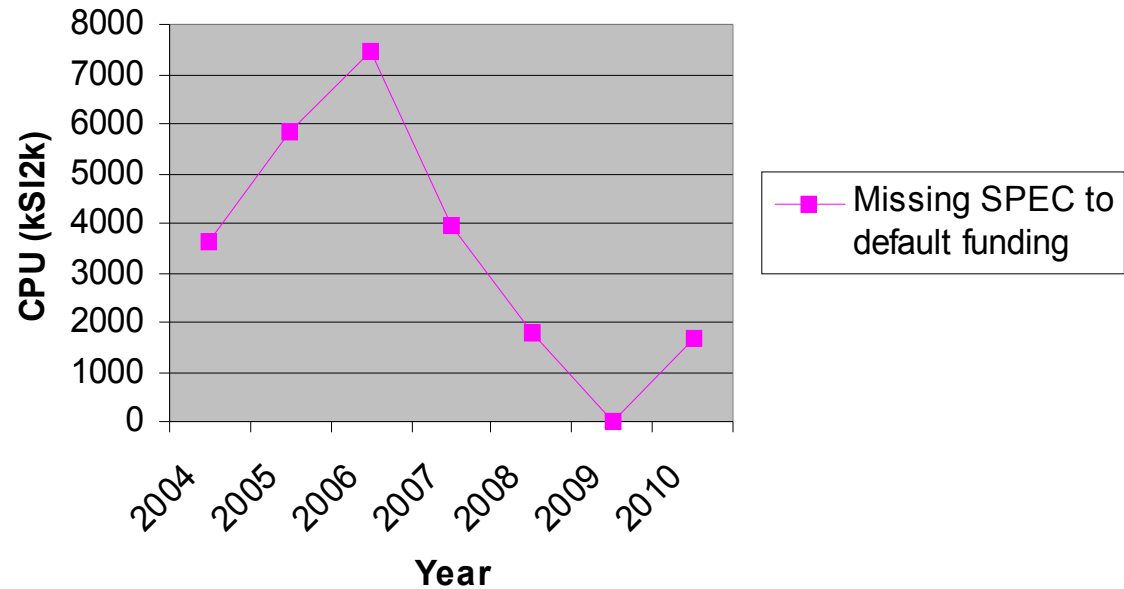
- 1.2 pass based model
 - We know we need 20% needed for calibration (0.2)
 - 1 pass = 1 time chance, **if something goes wrong, the data set CANNOT be reproduce, the science CANNOT be delivered**
- We STRONGLY believe in a minimum 2.5 passes
 - 10% fast calibration (as data arrives)
 - 10% slow calibration
 - 10% R&D
 - **2 passes, each pass twice as fast = better & faster quality data and science**
 - **A breathing margin for an already over-worked team**

Required

Comparison of CPU Delivered to Projected



Missing SPEC to default funding



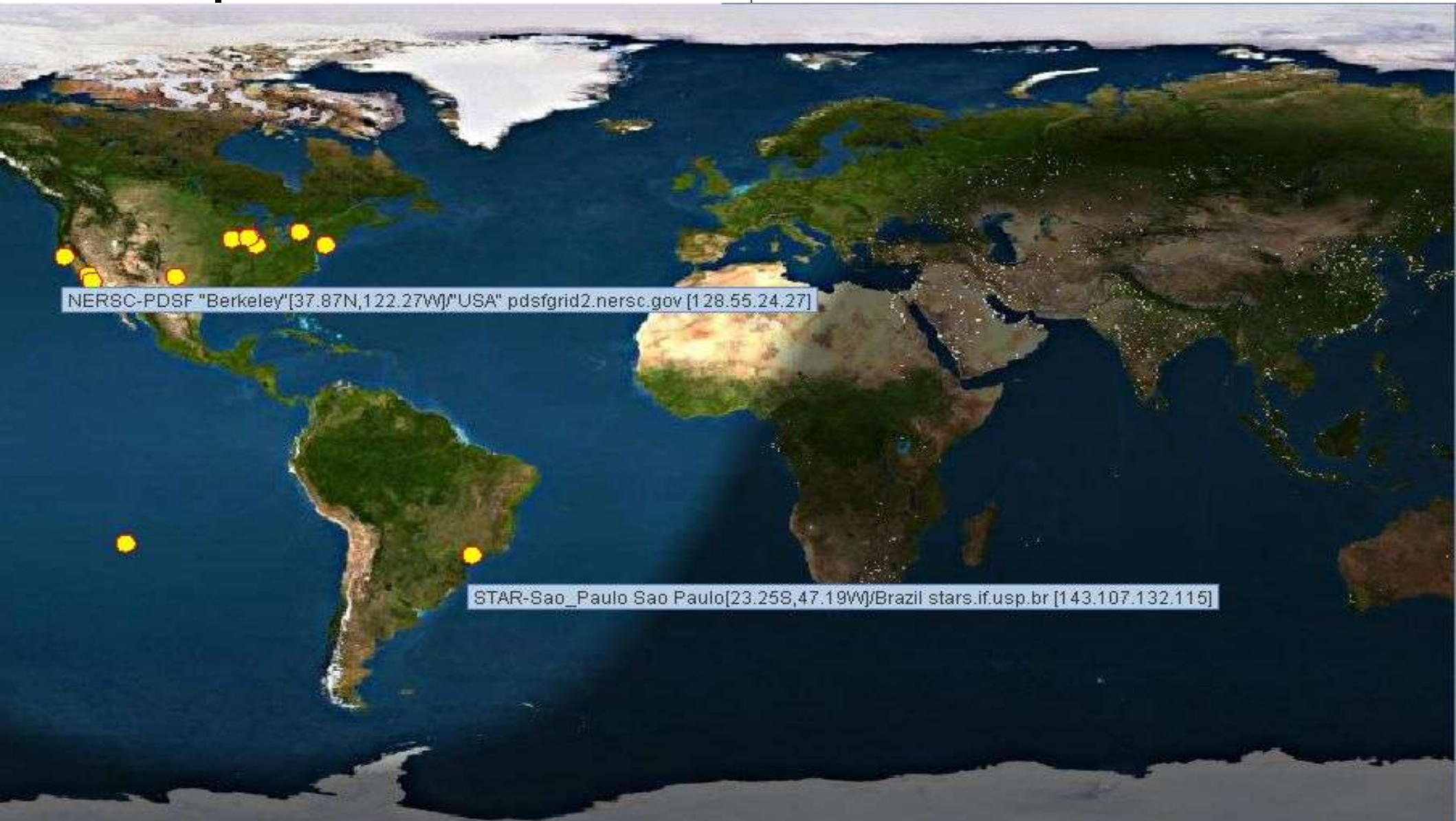
Several solutions to recover resources

- More money flows to RHIC computing
- Resources are borrowed from other providers (NSF / TerraGrid, ...)
- Resource are gathered from other collaborators

Opportunities & prospects

- Allocation to NSF TeraGrid
 - Provides a fraction of the missing spec on first year, less later
 - Also provide a superb opportunity for DOE/NSF resource exchange
 - Grid interoperability exercise and building the future in this area
 - Remote institutions
 - Several coming with 100ds of CPU within the next 2 years
 - May be at reach if Grid activities continue to be viably supported (funding for PPDG ends soon)

Current Status - OSG





Efficiencies and productivity with current RCF resources

- Every time one speaks of “efficiency”, “performance” or “productivity”, the word “business” crosses my mind ...

Speed:

Windows crashes:

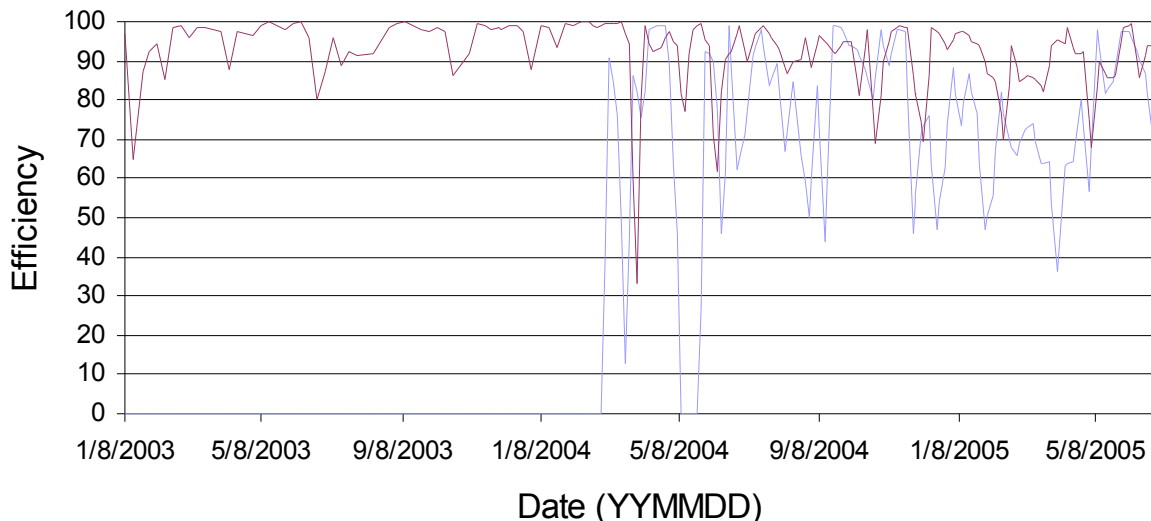
Memory Leaks:

Efficiencies and productivity with current RCF resources

- In general good
- Our code itself is very robust
 - Losses (crash) are below 0.1% at worst
 - Main reason for low rate: *FastOffline* or automated calibration catches problems
 - Problems found are fixed on a weekly schedule
- Technology factor & limitations
 - Using the RCF job submission software “black box”
 - New system designed to be scalable (good)
 - Efficiency purely based on success / failure trapped by the system (HPSS staging, miscellaneous) shows some concerning trends [later period may show improvements]

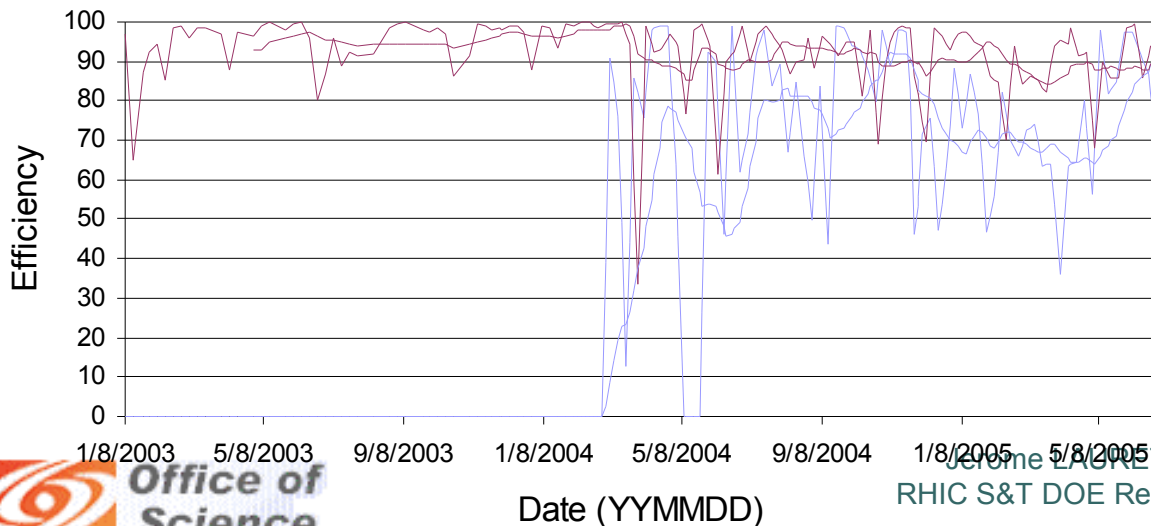
Efficiencies current resources

New and Old Reco systems



Old system efficiency slightly dropped mainly due to HPSS tape drive allocation
Trend is similar with other experiment (in fact, the graph represents an average)

New and Old Reco systems

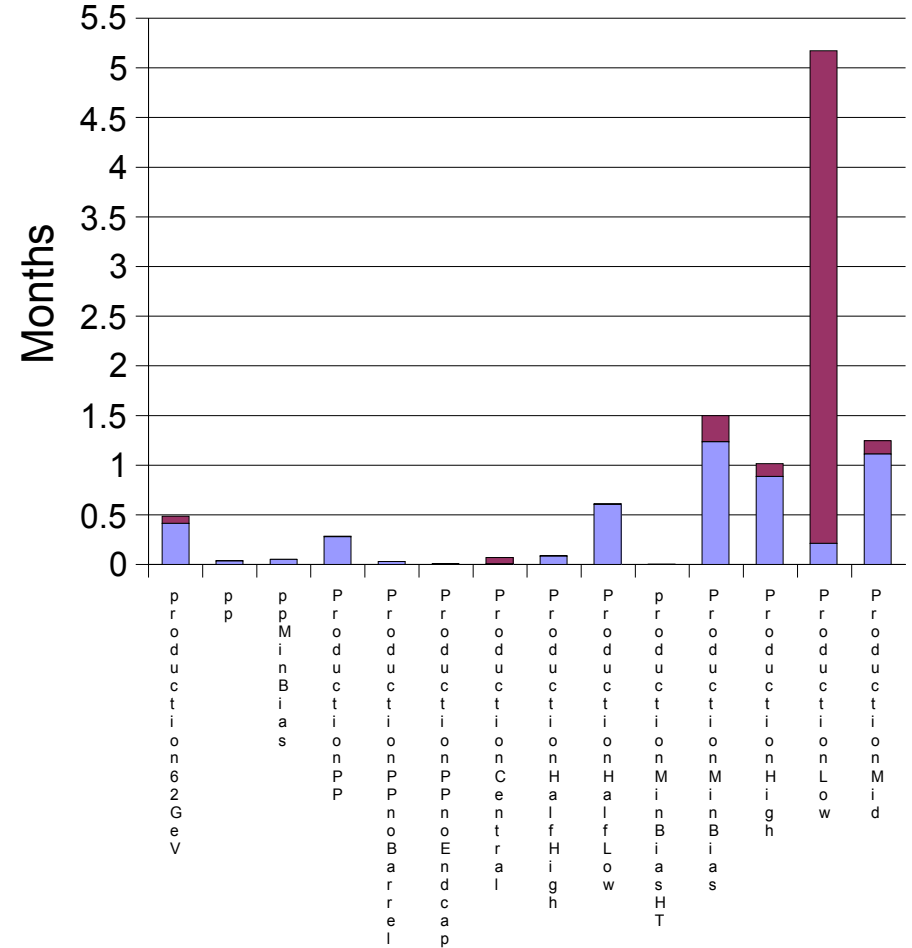
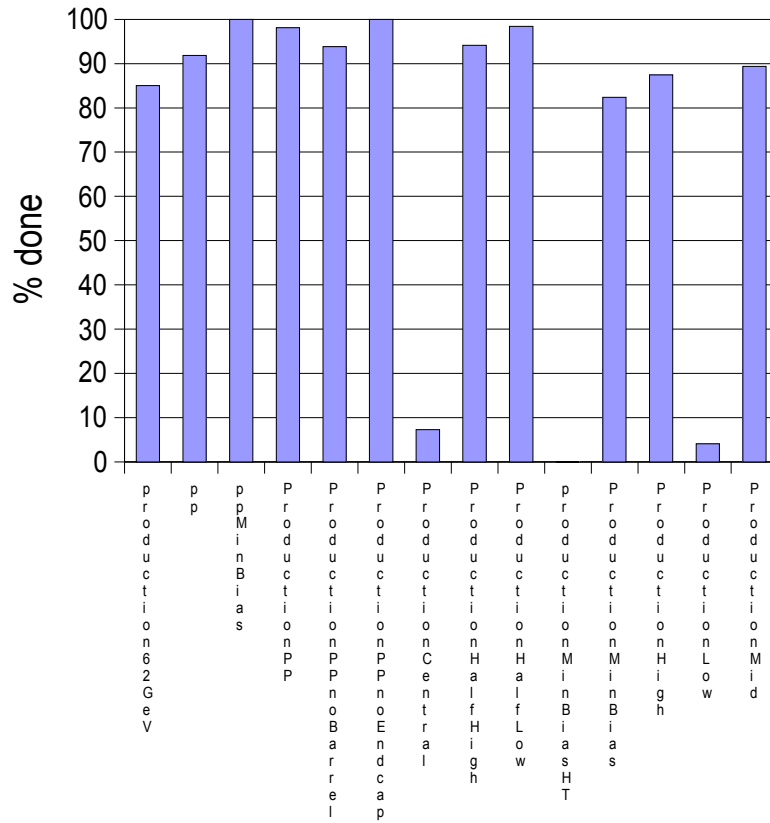


Efficiency is 20% lower average over a year
Trend lately to the increase
(hopefully will remain)

Efficiencies current resources & Communication with the RCF in other area

- CPU/Linux team much improved with a (new) set of motivated, qualified and friendly people
- Outstanding communication with the Storage / disk team
 - Helped in the evaluation of several storage system
 - In fact, IO stressed tested all of them
 - **Best tuned to real life science, lead to better scalable solution - Currently invested in PANASAS**
 - Has resolved some of our most important IO bottlenecks
- Concerns
 - Grid support seemed slow (was slow?)
 - Tickets response time of months scale (operation downtime)
 - Discussed and hopefully corrected
 - Good support for Virtual Organization related tasks
 - HPSS scalability
 - Lots of mysterious features and behavior
 - System will change in future, the knowledge is there though ...

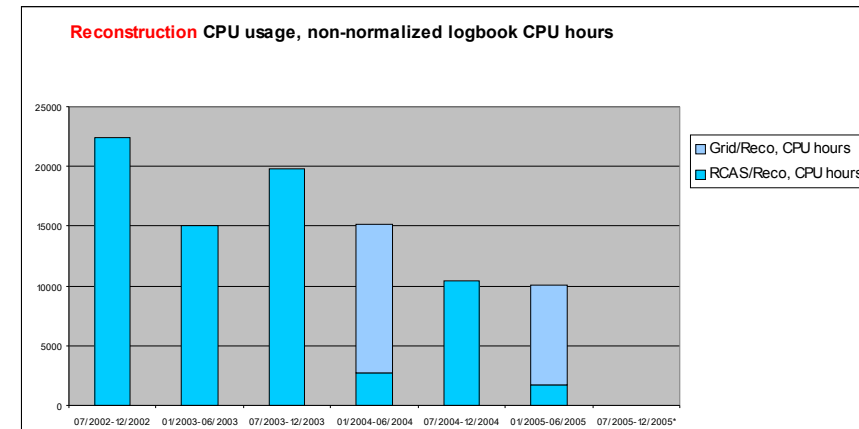
Year 4 data produced to date



Nonetheless, all data planned for production is now produced

Production status

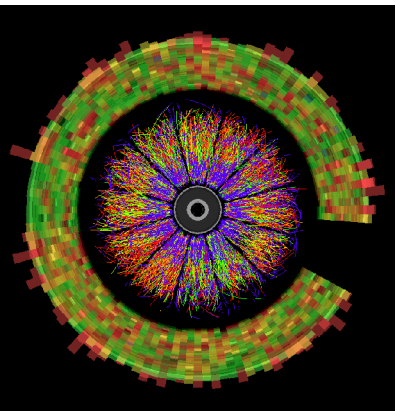
- Initial prediction based on 85% duty factor are -10% off.
 - New model –Merging analysis and reconstruction resources
 - Moved event generator simulation to Grid-based production
 - success rate reported in PPDG DOE quarterly report Jan-Mar-05 was 100% over 500 jobs
 - Average success rate [85-90]%
- Similar reached target for year5 data
 - Still however ~ 6 months worth of year 4 data
 - Year 5 lags behind as a side effect
- Resource OFFLOAD surely helps whenever resources are reachable / viable



The facility is providing resources and support to get the job done within the 1.2 passes expectations

Other Accomplishments

- Regardless of the current resource situation, several projects were carried to success or ongoing
 - DAQ100 new cluster finder
 - New Integrated Tracker for STAR (ITTF)
 - Strong validation procedure compares significant dataset for High Pt bias and several other Physics
 - **Project delayed due to lack of resources (human & CPU)**
 - Pileup-proof Vertex finder development (depends on previous)
 - Drop of the legacy gstar framework, now starsim
 - Allows for transition to integrated simulation / reconstruction a-la Alice Virtual Monte-Carlo
 - Project will be reviewed later this year
 - Three new sub-systems code in production (reconstruction, simulation and embedding)
 - ROOT / Qt development made at BNL
 - ...





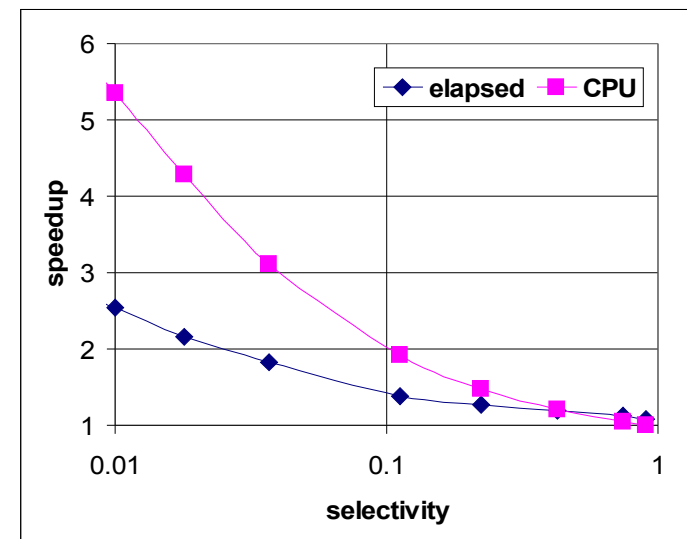
Other accomplishment made in collaboration with others

- Production-mode data transfer using the SDM DataMover
 - Strong and long standing partnership with the SDM centre at LBNL
 - Successful development of production Grid-aware tools
 - “Data Grid” architecture
- Replica Registration Service (RRS)
 - Allows on arrival file registration / availability to analysis
- GridCollector
 - Serves sub-events from distributed files to users – Speed x4
 - An interactive Grid analysis framework
 - Relies in SRM technology, second generation of Data Grid development

GridCollector

- “tags” (*bitmap index*) based
 - need to be define a-priori [production]
 - Current version mix production tags AND FileCatalog information (derived from event tags)

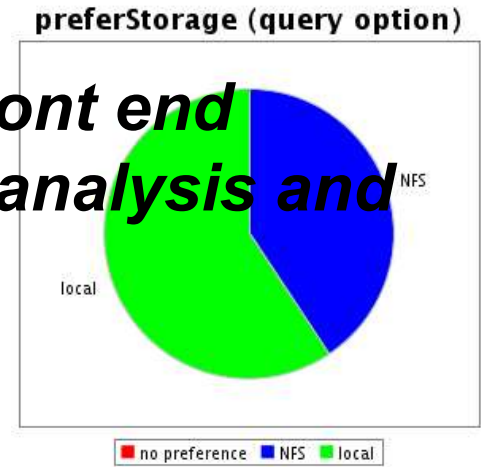
- **Usage in STAR**
 - Rest on now well tested, deployed and robust SRM (DRM+HRM)
 - Next generation of SRM based tools - “caching out” on past R&D
 - Immediate Access and managed storage space
 - Files moved transparently by delegation to SRM service
 - Easier to maintain, prospects are enormous
 - “Smart” IO-related improvements and home-made formats no faster than using GridCollector (a priori)
 - **Physicists could get back to physics**
 - **And STAR technical personnel better off supporting GC**



It is a WORKING prototype of Grid interactive analysis framework
Generalized, user analysis may gain ins peed

SUMS

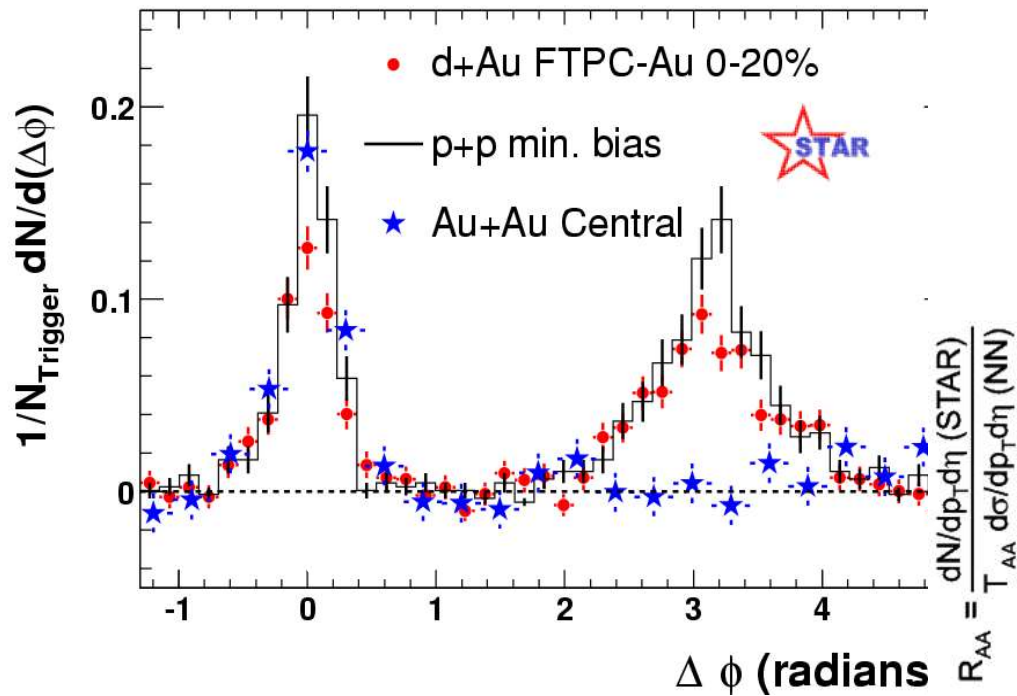
The STAR Unified Meta-Scheduler, A front end around evolving technologies for user analysis and data production



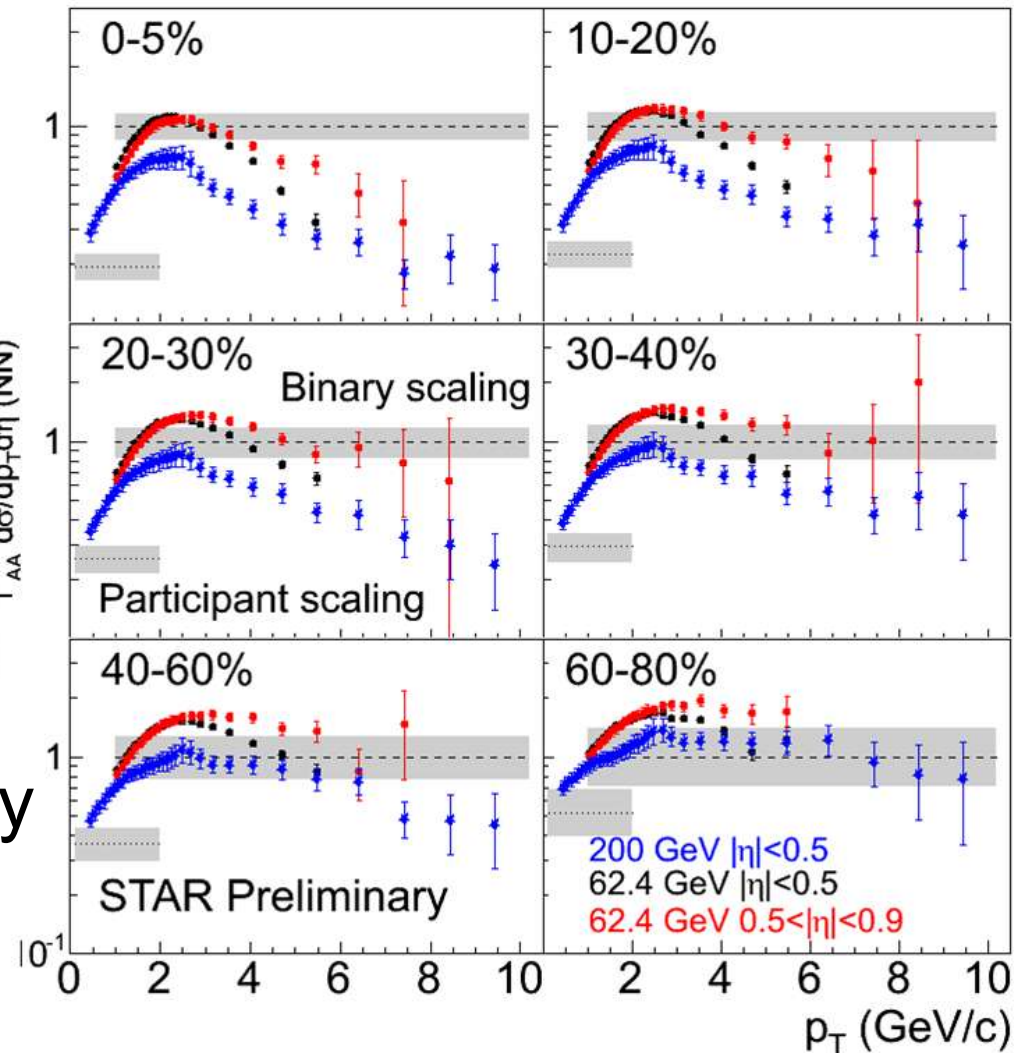
- **STAR Unified Meta-Scheduler**
 - Gateway to user batch-mode analysis
 - Flexible policy-based grid-ware, Collects usage statistics
 - User DO NOT need to know about the popular "batch" flavour of the day (adaptable technology - plug-in)
- **Has allowed to optimize resource usage**
 - IO throttling is automated
 - Best queue / resource found – Grid AWARE
 - Integrated to the file catalog, a distributed disks approach (locally attached to sparse nodes) is possible
- Scavenger hunt for resources is in place
 - Non negligible actually, 100 TB of distributed disk comparing to 130 TB central (NFS, PANFS)
- Still a lot to be done
 - Some spiky features in resource utilization needs better understanding and development of enhanced meta-scheduling policies

The real measure of accomplishments and productivity

Phys Rev. Letter 91 (2003) 072304



Beautiful results internationally known

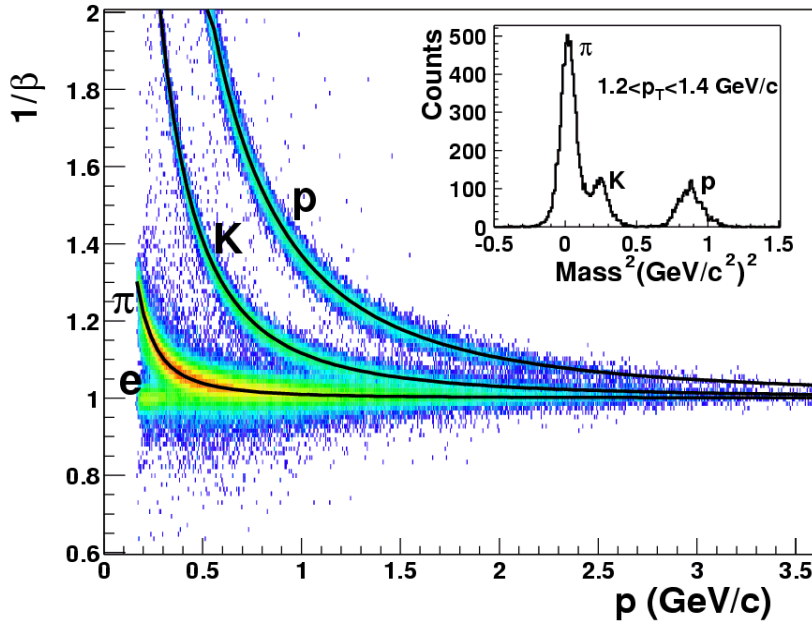


The real measure of

Phys. Lett. 94 (2005) 062301

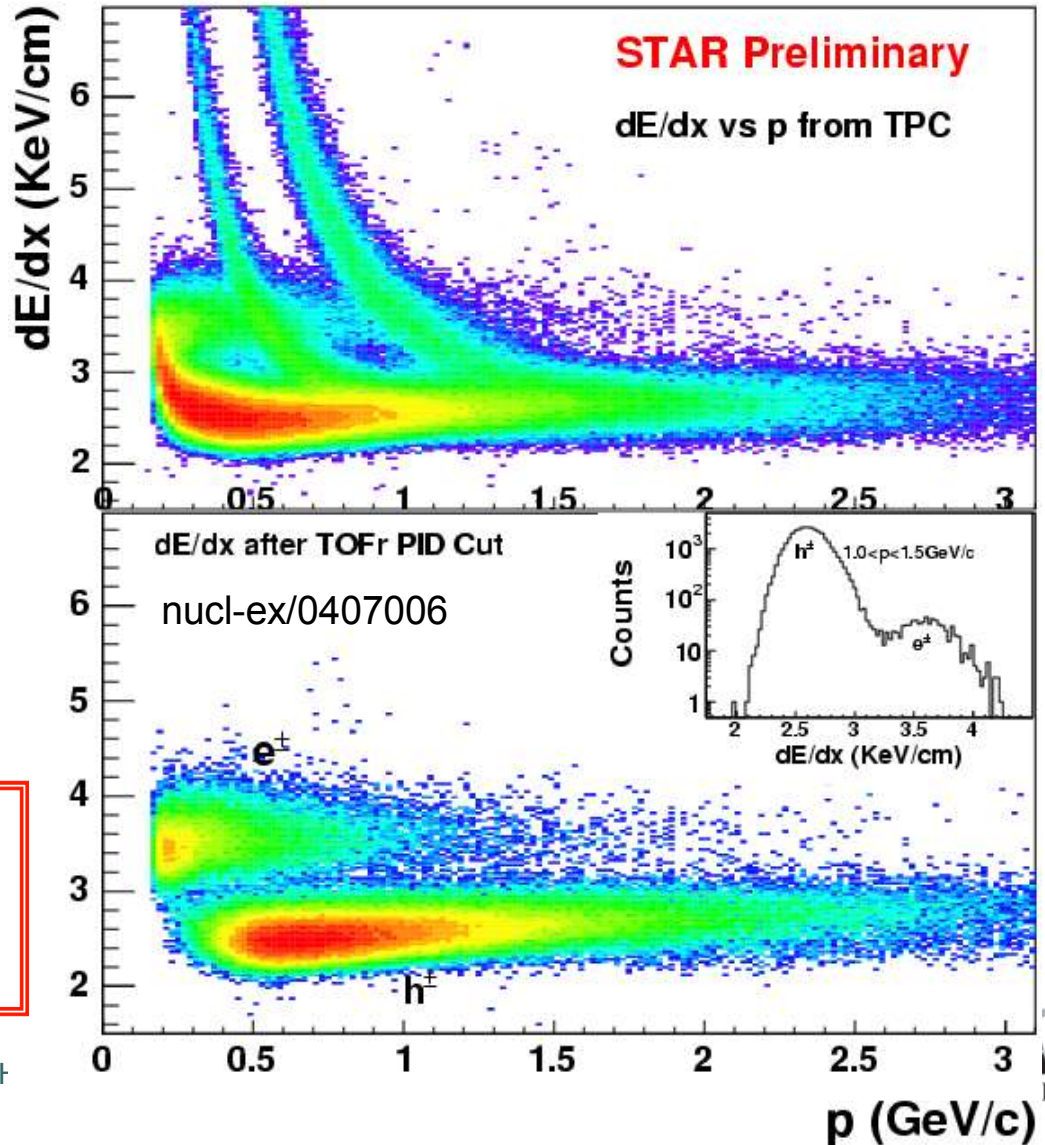
Phys. Lett. B 616 (2005) 8

accomplishments and productivity



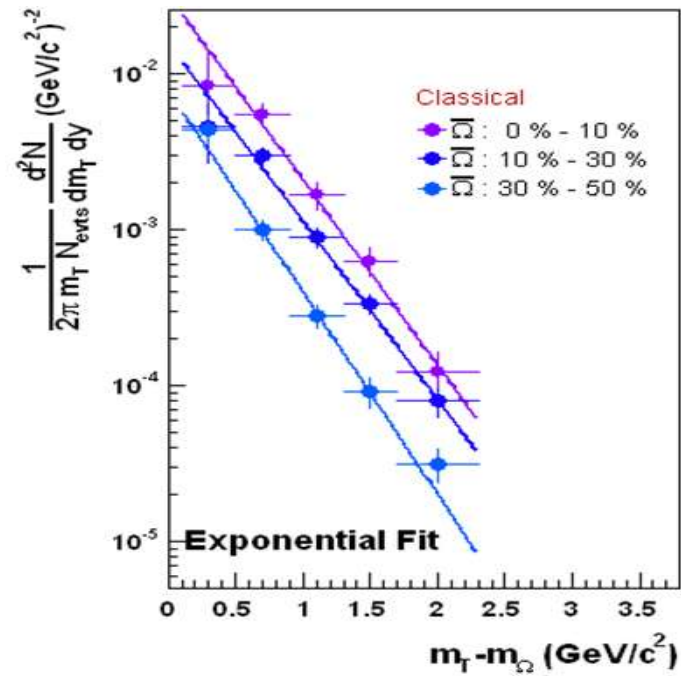
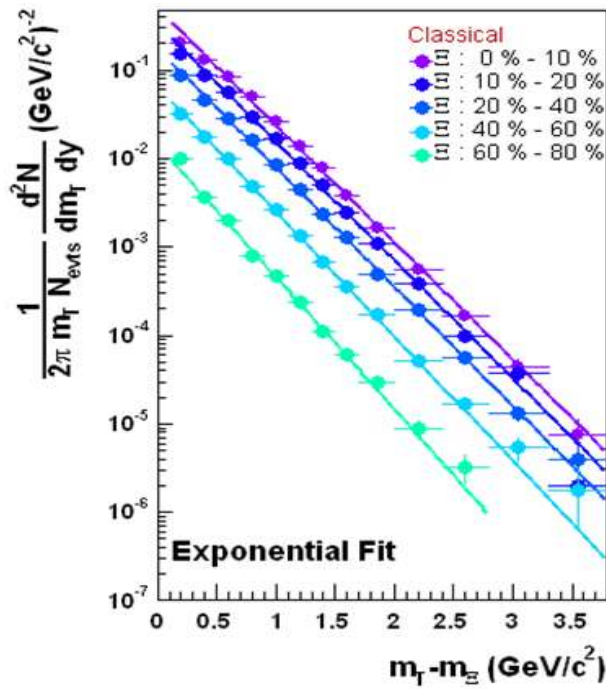
Hadron identification:
STAR Collaboration, *nucl-ex/0309012*

Electron identification:
TOFr $|1/\beta - 1| < 0.03$
TPC dE/dx electrons!!!



Rf

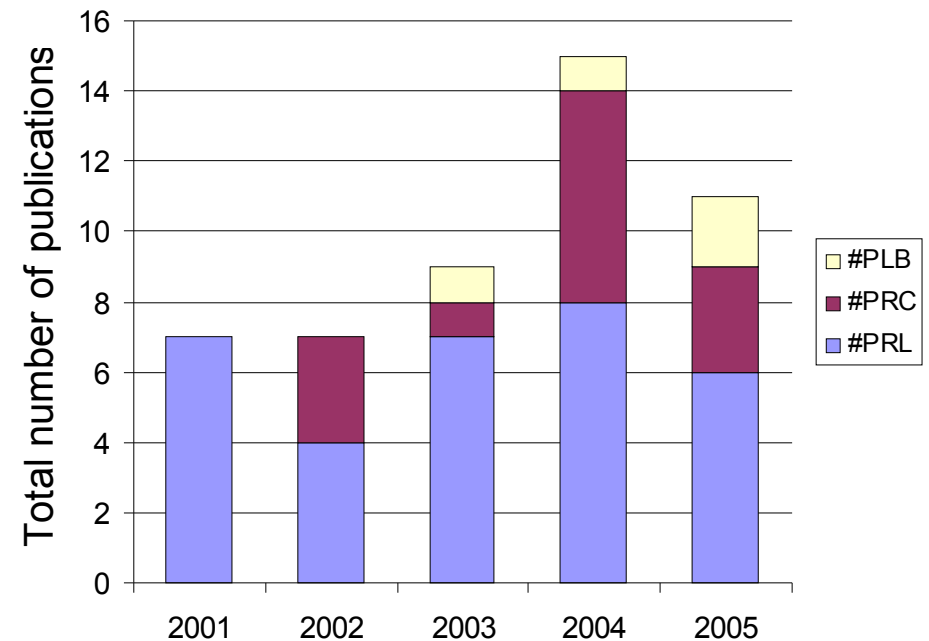
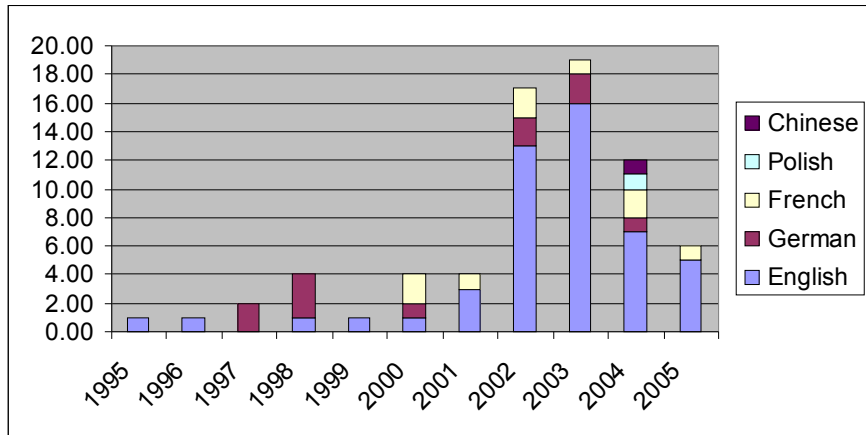
The real measure of accomplishments and productivity



Productivity in STAR

- Scientific Achievement in STAR

STAR Publication trend



2692 citations

28 PRL, 15 PRC, 4 PLB published to date

Summary & conclusions

- RCF related
 - Provides resources adequately within 1.2 pass model
 - New technology integration and support time scale (i.e. production scheduling system, Grid support) has side effects and needs emphasis
 - Mass Storage reliability and scalability concerns
 - Communication and knowledge much-improved = right direction
 - Several collaborative activities leading to better tuned solutions
- STAR related
 - 2.5 passes needed in STAR
 - Implies search for additional resources
 - We have a clear plan toward success but
 - Grid computing a strategic choice already at the heart of our production
 - Difficulty to perform computing R&D within resource figure: both (in)human and hardware
 - Even harder to make long term support for projects (ITTF,...)
 - Lots of activities however needed for the future
- Scientific productivity, the best measure of success is outstanding