# Data Preservation and Long Term Accessibility status of BaBar

@3<sup>rd</sup> DPLTA Workshop
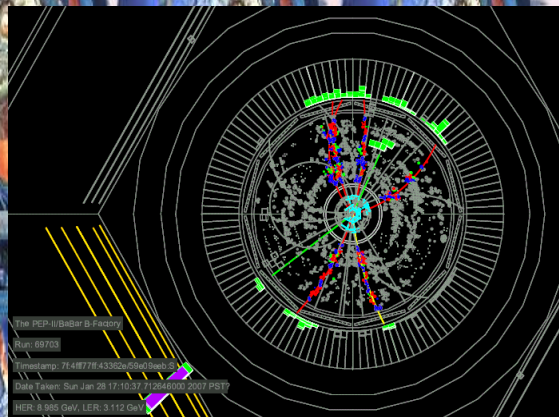
@ CERN

8 December 2009

by

Homer Neal ( **SLAC** NATIONAL ACCELERATOR LABORATORY )

BaBar Computing Coordinator on behalf of the BaBar

Long Term Data Access Group

# BaBar's Data Preservation Effort: The Long Term Data Access (LTDA) Group

- Kyle Fransham (Uvic) - migrations/virtualization

- Igor Gaponenko (SLAC) –migration/code cleanup/conditions database

- Bertrand Echenard (CalTech) –joint analysis/simple job manager

- Tina Cartaro (SLAC) –database consolidation/documentation

- David Brown (LBNL) - reconstruction

- GPDF (SLAC) –computing (platform/code design)

- Matt Bellis (Stanford Univ.) - outreach/data format/analysis tools

- Adam Edwards (HMC) –job manager/archival system requirements/documentation

- Rudy Resch (SLAC) –generator/simulation separation

- Ha Seong Kim (HMC) –job manager/archival system requirements

- Perry Ellis (HMC) –job manager/archival system requirements
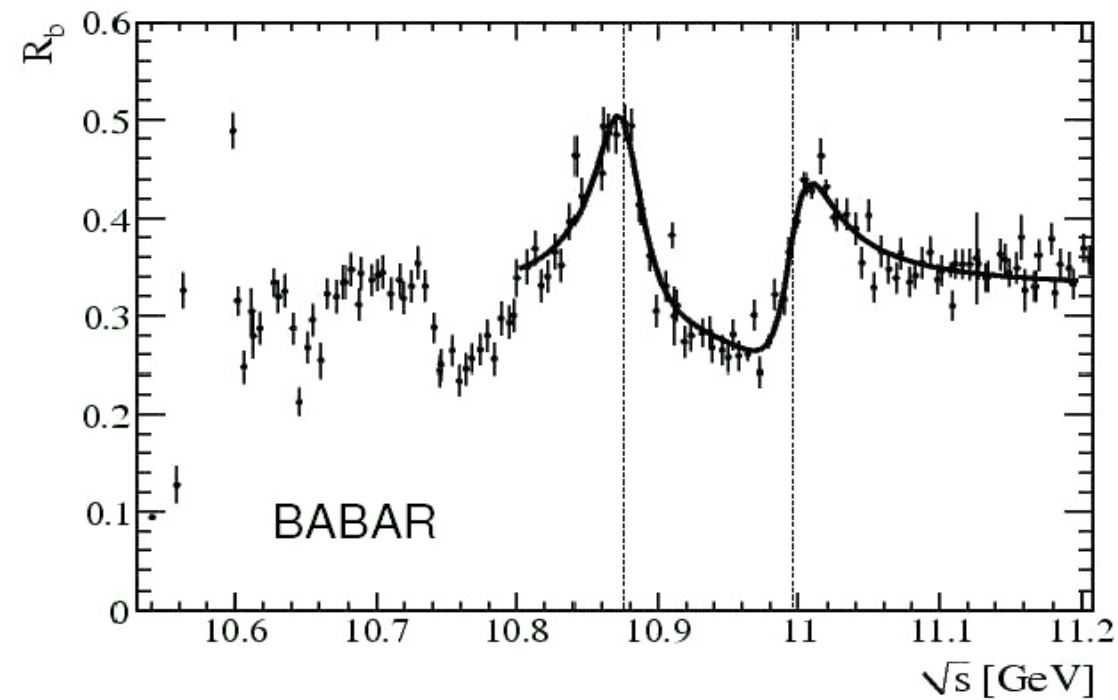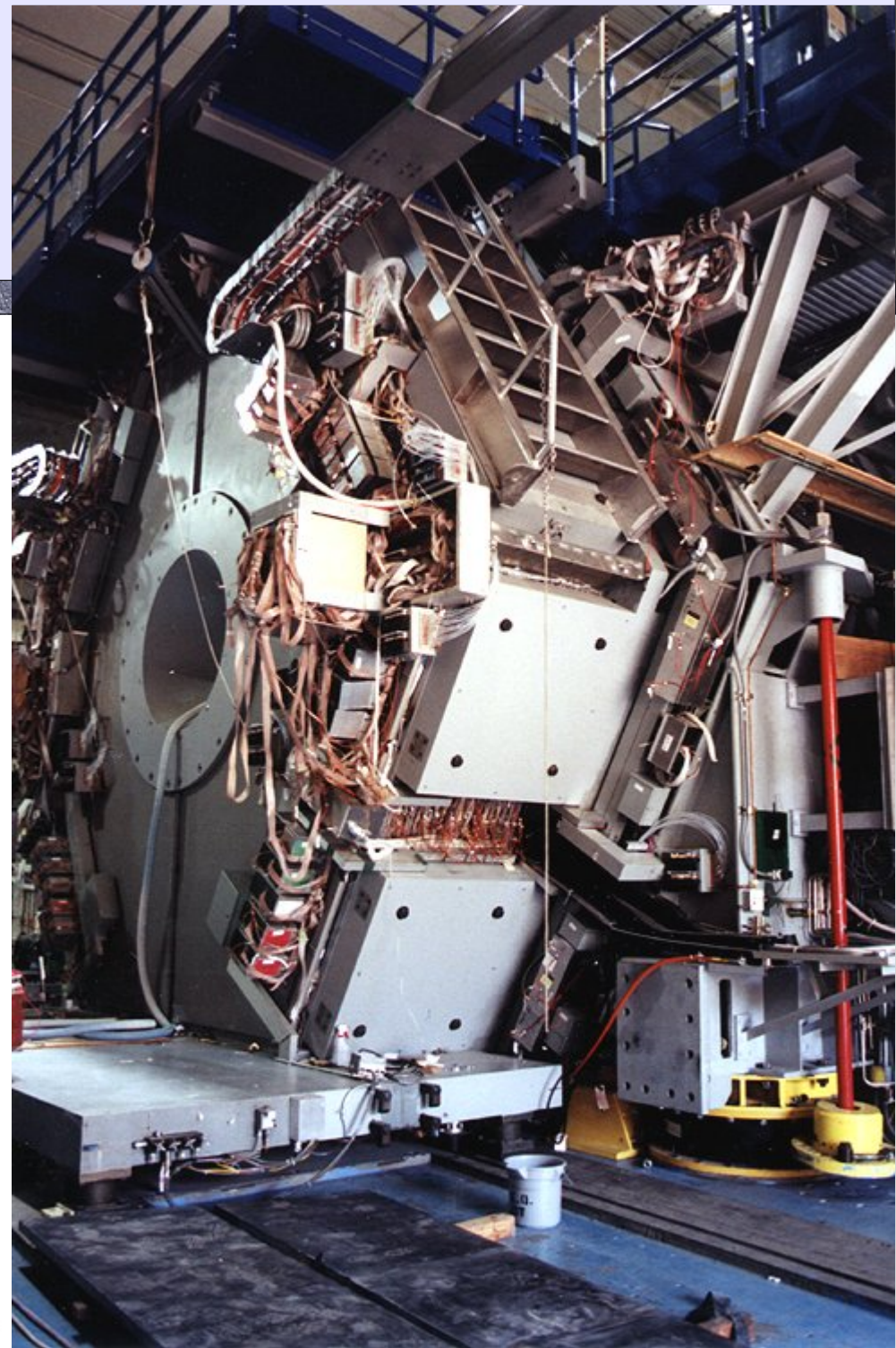
# The Story of BaBar

**conceived 15 years ago (LOI 1994)**

**collected frontier HEP data over almost a decade**

**through the efforts of 10 countries and**

**~600 collaborators**

**April 7th, 2008 marked the end of BaBar data taking but the analysis activity has been increasing**

# BaBar: Some notable events since the end of data collection

July 9, 2008 - Physicists Discover New Particle: The Bottom-most "Bottomonium"



The Nobel Prize in Physics 2008 and the B FACTORIES



Congratulations 10 years
August 13, 1999
August 13, 2009

**BaBar Collaboration Caps Meeting Week with 400th Scientific Publication**

*by Lauren Knoche*

The BaBar Collaboration reached another milestone Tuesday—just in time for celebration during the group's meeting, which ends today at SLAC. The collaboration published its 400th paper Tuesday, less than nine years after publishing its first in 2001. That's an average of one publication per week, every week, for nearly nine years straight.

"I do not know of any other collaboration that has achieved such a production rate of outstanding quality science in particle physics, it is really something rare," said BaBar spokesperson Francois Le Diberder.

The milestone paper was published online Tuesday and appears in the November 1, 2009 issue of *Physical Review D* (Volume 80, Number 9). The study examines differences in the rates at which subatomic particles called B+ mesons and their antiparticle partners, B- mesons, decay to related particles called "charm" and "strange"

(Photo by Brad Plummer.)

5 Nov. 2009 – SLAC Today

- We foresee >100 papers being published over the next several years and *about 35 papers to be published after the loss of the existing large BaBar computing infrastucture provided by SLAC, other TierA sites and ~20 universities*.

There are 35 (20 of these are "MUST DO") analyses already foreseen to be users of the archival system. Among these are:

- Initial State Radiation Physics These are very delicate analyses aiming at the determination of the cross-sections : $e^+e^- \to 3\pi^o$ , $K_S K_L$, $K_S K_L \pi$, $K_S K_L \pi\pi$, 7(8) pions with the goal of providing a complete set of cross-sections for $e^+e^- \to hadrons$, at low energies, where the contribution to g-2 is the most sensitive.

- Charm Physics : rare decays, many body final states, etc.

- Y(nS) Physics like $Y(mS) \to \pi\pi$ (or $\gamma\gamma$)$Y(nS)$

- sin2beta analyses (sic: a few could be done, although they are not high profile) like $Y(4S) \to J/\psi\rho$

In addition, there is a strong likelihood that new models will need to be tested in the archival period and that checks against the BaBar data will need to be done by new projects such as SuperB.

# Current Computing Resources

- 5700 cores at SLAC accessible to BaBar for general (typically analysis) work

- 1100 dedicated cores principally for BaBar only production

- Old xrootd cluster servers recently replaced
  (620 TBytes currently, normally 450 TBytes)

- All old batch systems (352:4 core machines and
  256:2 core machines) being
  replaced with 160 Nehalem dual-quad core machines

- XFER data distribution machines recently replaced

- NFS and AFS servers being replaced

# Remote Resources

- Much of our simulation production capacity currently comes from CCIN2P3, RAL, and various INFN sites

  - ~20 sites total (Uvic, GridKa, and universities)

## SP10 - KiloEvents per week

### The date is the time the simu job completed

| Week Beginning | TOTAL | caltech | ccin2p3 | cu-boulder | fzk | infngrid | infnta | osu | ral | slac | slac2 | tud | udo | uofl | utd | uvic2 | westgrid |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2-8-2009 | 303663 | 0 | 91590 | 16 | 0 | 6884 | 0 | 9523 | 109404 | 65804 | 0 | 2327 | 0 | 0 | 0 | 17828 | 287 |
| 5-10-2008 | 613867 | 6348 | 135569 | 0 | 0 | 21151 | 0 | 12808 | 98767 | 231212 | 60411 | 2300 | 0 | 0 | 0 | 35426 | 9875 |

  - RAL will be lost soon

  - Commitment from CCIN2P3  to maintain current level of resources for BaBar for as long as analysis activity remains high

8

# Analysis resources

- Analysis Working Groups continue to be hosted at
  - SLAC
  - CCIN2P3
  - CNAF
  - GRIDKA
  - UVIC

# Platforms

- Platforms for batch have switched at most sites to SL5

- A few stragglers still using old releases of data/code (*long term analyses that are mature and need more time to complete*)

- Only very latest releases buildable under SL5

  or virtual machines for another ~1 year

# Current Media Migration Status

- Nov. 12$^{th}$ – after much waiting, the new silos and HPSS system became ready for use

- Preparing to start a ~1 year migration from ancient (some are 9 years old) 9940a/b tapes to T10000 tapes

- The raw data plus two most recent processings of the data will be migrated (2 Pbytes to be migrated … ~1/3 of total)

- Tapes from some old processings were ejected and boxed to allow an old silo to be returned. Ejected tapes have 1 year expiration date.

# How soon could BaBar be superceded

It will take ~8 years for SuperB to be approved by the governments and funding agencies, constructed, commissioned and obtain a dataset more significant than BaBar's.

- **2012 to 2018 is the <u>minimum</u> required existance of the BaBar archival system**

_____

There may also be a need to validate initial results against those of BaBar and Belle

# BaBar's approach

- Preserve all capabilities of doing full analyses from inception to publication including some **new simulation, access to raw data and latest reprocessings of the detector and simulated data, all databases, documents, simplified/improved interfaces for job submissions and introduction of new models, frozen trusted reconstruction in an** <span style="color:red">**archival system**</span>

| Preservation Model | Use case |
| --- | --- |
| 1. Provide additional documentation | Publication-related information search |
| 2. Preserve the data in a simplified format | Outreach, simple training analyses |
| 3. Preserve the analysis level software and data format | Full scientific analysis based on existing reconstruction |
| 4. Preserve the reconstruction and simulation software and basic level data | Full potential of the experimental data |

#4) BaBar's choice

13

# The BaBar Archival System

- A modest many core box with enough storage for the data, databases, documents, releases, and new generators and simulated signals.
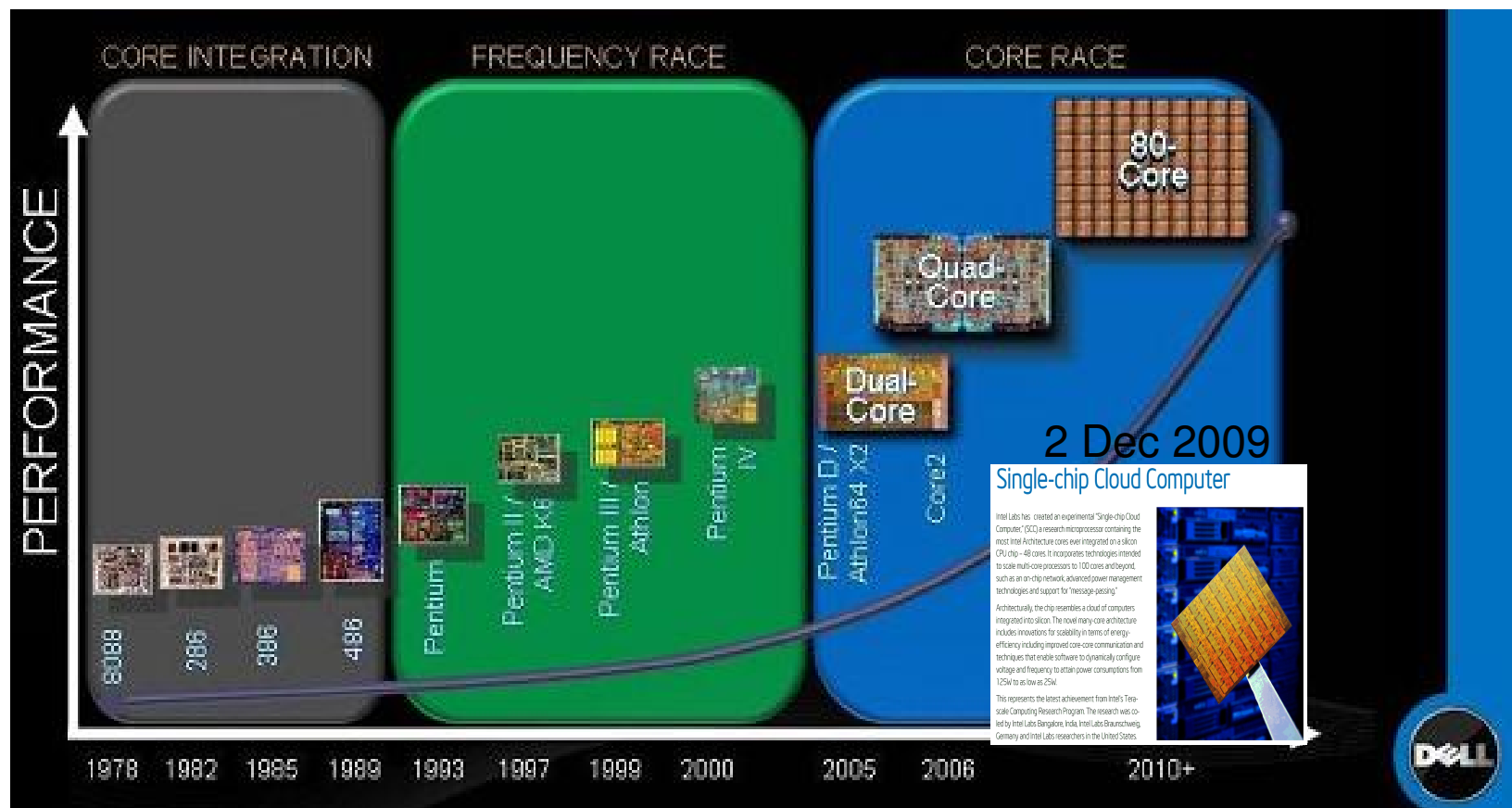


- Running a current platform that will be protected by living in virtualization containers

14

# Many core systems are the future

- ## From SC08:

http://news.cnet.com/8301-13924_3-10101987-64.html

**How well will they handle high data output jobs???**

# The BaBar Archival System

# The BaBar Archival System

- Currently working on selecting prototype archival system hardware

  32-cores, enough storage for micro data needed by most analyses

- Set it up and get testers to identify faults

# Status of BaBar Effort

- Migration to SL5 and ROOT 5.22 ~ completed

  - Full validation in progress and currently attacking one reconstruction problem

  - Progression: 5.20-00/ 5.21-06/ 5.22-00/ 5.22-00d/ 5.24-00b/(installed and ready to use in releases and validate)

- Simplified job managers ~ completed

- Tests using virtualization and virtual machines in batch on a single node successful

- Progress on migrating databases to or current test system

- Starting documentation effort

- New BaBar MediaWiki site in use and working on improving server setup

- Working on acquiring a prototype archival system

- Working on finding analysts to test it and provide feedback

18

# BaBar Data Preservation

BaBar & Belle collaborating



In real life: B± → K±π+π± decay

Same exercise with the master at Caltech (Los Angeles), one worker at SLAC and the other worker at ccin2p3 (France) with secured connections.
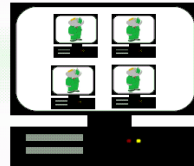
Los Angeles

SLAC                France

Minuit

Fit performed in a bit less than 20 minutes. Note that we had slow 32-bits machines, a fit SLAC-SLAC-SLAC took almost 4 minutes
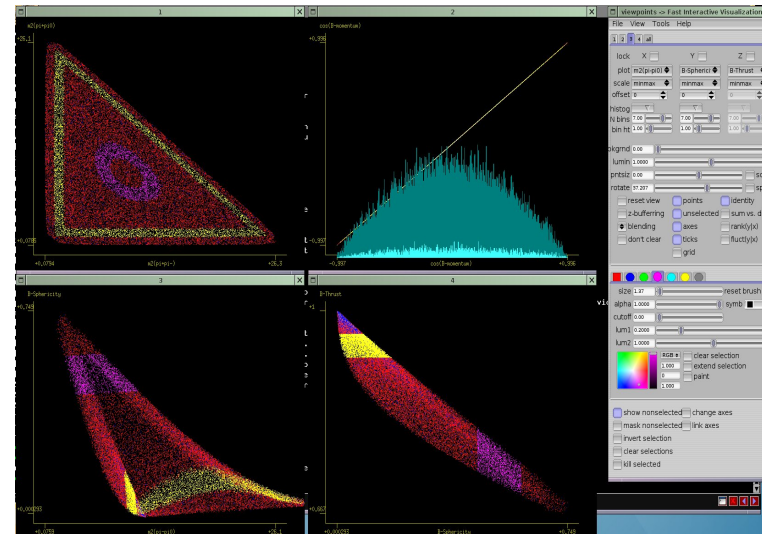
It worked very well

B. Echenard / E. Ben-Haim        BaBar Collaboration Meeting / November 2009        p. 11

## Virtualization

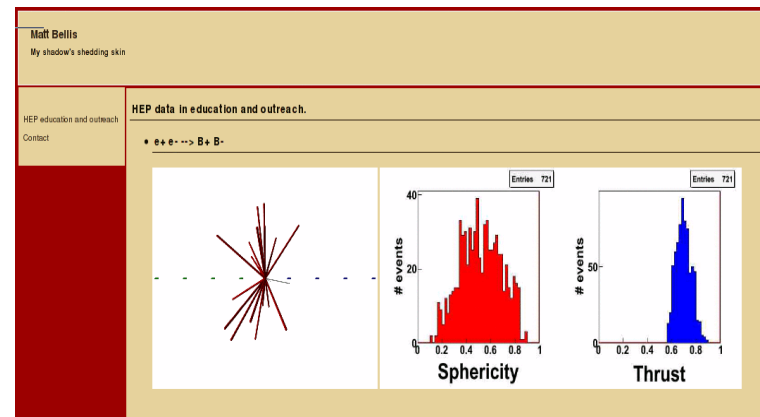- The status at SLAC: 4 SL5.3 VMs installed on yakut13.
- VMs were added to a special batch queue.
- SL5 migration checks to be done on virtual machines.
- Simultaneously validates the SL5 build and the VM technology.

June 22, 2009        Long Term Data Access        6



Outreach tools/data already being used in classrooms



Also major advancements in the use of cloud computing

19

# BaBar's Plan

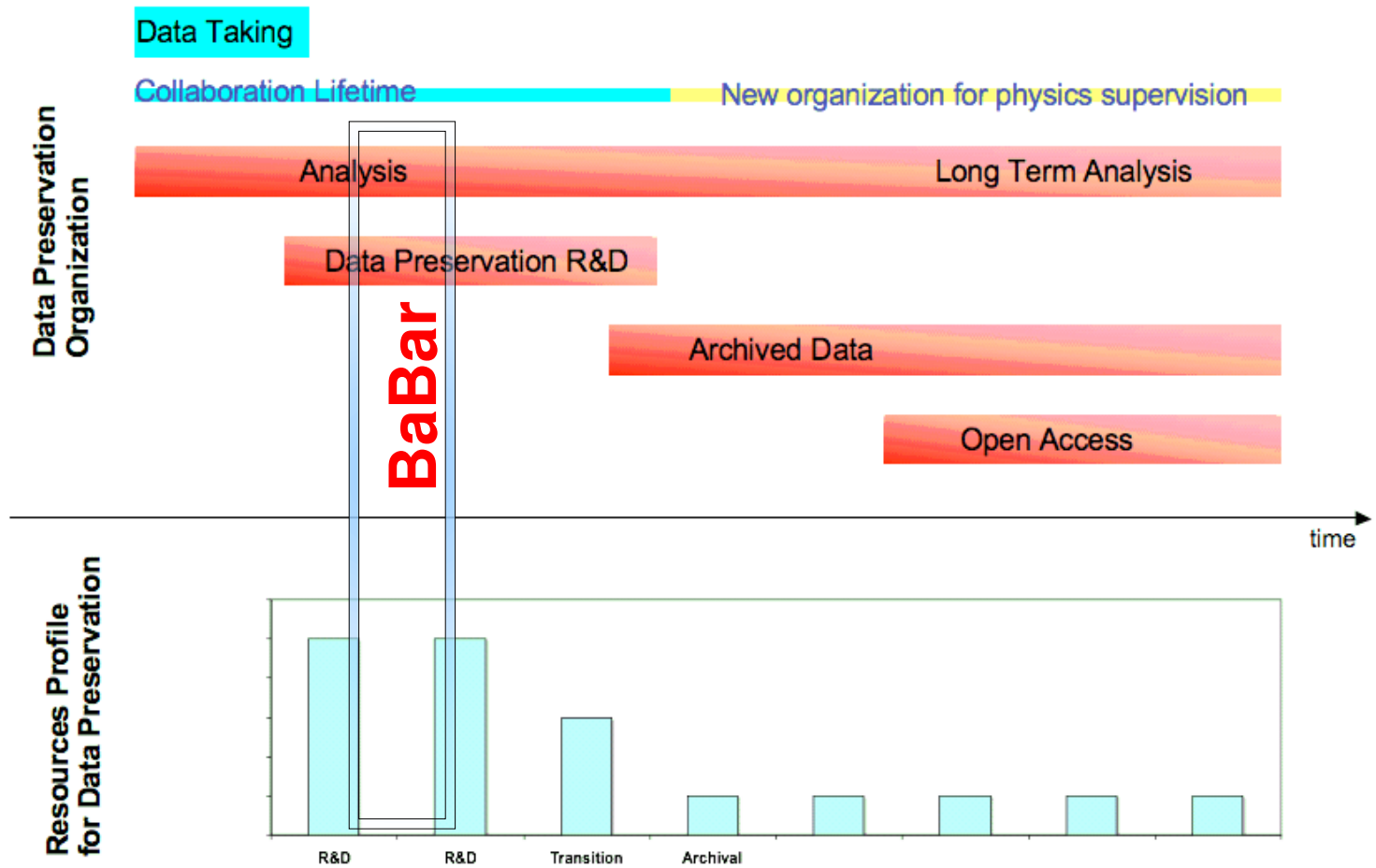# DPHEP: data preservation timeline



Figure 1: A possible model for data preservation organisation and resources.

21

# Concerns

- Virtualization technology lifetime (Xen, VMware,...)

- One box as a database, batch, interactive analysis, web server … will we have a melt down

- Data flow from a many core system – sufficient bandwidth?

- Integrity of results produced from publicly accessible HEP data

This all-sky view from the Fermi telescope reveals bright emission in the plane of the Milky Way (center), bright pulsars and super-massive black holes. (Image: NASA/DOE/International LAT Team.)

**Ever since the Large Area Telescope launched aboard the Fermi Gamma-ray Space Telescope in June 2008, the LAT team has been analyzing data, searching for answers to some of the most pressing questions in astrophysics. <u>Now everyone else can join in.</u>**

Today, the collaboration and the Fermi mission makes the first year of LAT gamma-ray data publicly available.

"This is a way of maximizing the scientific return from the mission," said Fermi Project Scientist Julie McEnery. "There is a very large number of scientists in the community with very good ideas of what to do with this data. By sharing it among a large group of people, we really get a lot more."

To ensure that others in the astrophysics community can take full advantage of the data, the LAT collaboration, working with the Fermi Science Support Center at NASA Goddard Space Flight Center, has spent a considerable amount of time preparing for the release.

"It took significant effort both on our side and the Goddard side to both get the data out and to get it out in a form that's usable by the whole community," said Astrophysicist Jim Chiang, LAT Collaboration member who works on the analysis software for FGST.

The data set released today includes more than 150 million detected gamma rays. In contrast, in the more than nine years that the LAT's predecessor, EGRET, operated, it collected 1.4 million gamma rays. In all, the LAT has collected more than 100 times as many photons in about one-tenth the time.

…

—Kelen Tuttle

SLAC Today, August 25, 2009

**A new concept for HEP**

# Large Area Telescope First Year Data Released

…

As in all particle physics experiments, Chiang said, LAT data are unique to the instrument and require unique software. With this in mind, the collaboration will also make available high-level software that other researchers will need in order to analyze the data. In addition, NASA is offering further resources and funds to guest investigators who successfully submit proposals.

"We can see both from the large number proposals submitted to the guest investigator program and the large number of references in papers that the community is excited about the data," McEnery said.

LAT Principal Investigator **Peter Michelson added: "The LAT team has made significant discoveries and significant progress in many areas. I expect that the collaboration will continue to come out with the most results, but I also expect others to make discoveries. Releasing this data is good for the project, good for the collaboration, and good for science."**

—Kelen Tuttle

SLAC Today, August 25, 2009

# Words from the Archivists:
# The government demands the preservation of data

## Scientific Data:

- Raw data (all levels)
    - 10 year retention (N1-434-07-01, item 4c(12)

- Evaluated or Summarized data
    - Level 1: permanent retention (N1-434-96-9, item1B13a)
    - Level 2: 25-year retention (N1-434-96-9, item1B13b
    - Level 3: 10-year retention (N1-434-96-9, item1B13c)

# On Demand Data Analysis

- Use virtual machines on clouds like Amazon EC2 to simulate/analyze archived data as needed

- Currently being investigated for BaBar

- Used when extra processing power is needed for Belle

# Advisory Committee

- Report on status and plans of the BaBar Long Term Data Access group to be scrutinized by an LTDA Advisory Committee made of external members

# Funding of BaBar Data Preservation

- Summer 09 International Finance Committee meeting revealed continuing strong support for the LTDA effort from all TierA sites (including SLAC)

- Request to be voted on in January '10 includes funding of FTEs for:

  - the current data preservation effort (FTEs)

  - an archival system manager in the archival period and

  - funding for the next migration of the data to new media

# Summary

- BaBar's data preservation effort continues to go strongly

- Looking forward to approval of funding of the data preservation effort for 2010 and beyond this January.

- Working on acquiring a modest many core server to be the 1$^{st}$ prototype archival system.

# BaBar and DPHEP

Looking forward to continuing strong ties with the ICFA DPHEP group for guidance in our choices towards a successful preservation of the BaBar data and analysis capabilities and to provide feedback to DPHEP