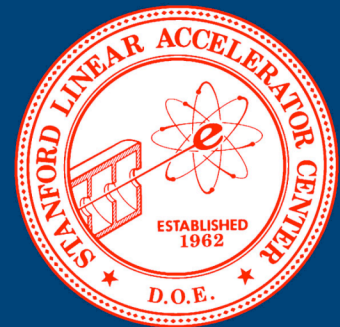


# *Extremely Large Databases for Data-intensive Computing*

Jacek Becla

Stanford Linear Accelerator Center





# *SLAC and Databases*

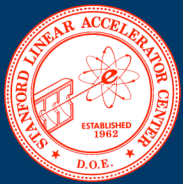
- ◆ BaBar
  - petabyte+ database
- ◆ LSST
  - SLAC leads design, O(100) PB database
- ◆ XLDB activities
- ◆ SLAC's core competency (1 of 4)
  - *Ultra-large database management for users and collaborations distributed worldwide*



# 1<sup>st</sup> XLDB Workshop



- ◆ Participation
  - data-intensive science & industries, database researchers and vendors
- ◆ Goals
  - identify trends, bridge gaps
- ◆ Very successful
  - science – db research collaboration strongly encouraged



# SciDB Mini-Workshop

## ◆ Participation

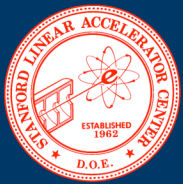
- database researchers + data-intensive science representatives (HEP, Astro, Bio, Geo, Fusion)

## ◆ Goals

- discuss common science db-requirements
- stimulate database research

## ◆ Very successful

- agreed to explore avenue of building new open-source science-oriented DBMS.  
Led by Michael Stonebraker and David DeWitt



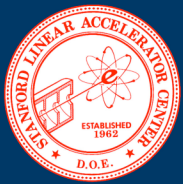
# *Features Requested*

- ◆ Scalability and fault tolerance
- ◆ Performance, extensibility and compression
- ◆ Efficient support of arrays/vectors
- ◆ Spatial and temporal support
- ◆ Provenance
- ◆ Uncertainty
- ◆ Table versioning
- ◆ Resource management



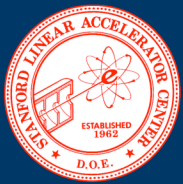
# *Scalability/Fault Tolerance*

- ◆ Scalability to tens of 1000s nodes
- ◆ Intra query fail over
- ◆ Parallel query execution
- ◆ Running on commodity hardware (cloud)



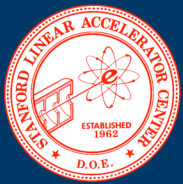
# Arrays

- ◆ Data model: multi-dimensional arrays
  - fixed or variable stride
  - chunked and partitioned
- ◆ Native array operations
  - plus user extendable operations
- ◆ Dot products on arrays 100x faster



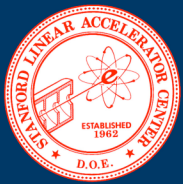
# Compression

- ◆ Can take advantage of correlations between array dimensions
  - example: delta encoding
- ◆ Will operate on compressed data



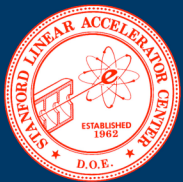
# Why “Scientifica”

- ◆ Requirements novel, unlikely to be met by existing vendors
  - arrays, spatial/temporal support, provenance, uncertainty, versioning
- ◆ Large scale and complexity prohibits roll-your-own approach
- ◆ Overlap increasing
  - significant commercial applications (R&D and non-R&D). Examples: internet, oil & gas
- ◆ Scale of “big science” requires Green DB



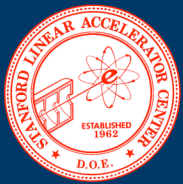
# Core Partnership

- ◆ Scientists
  - put up some resources, provide requirements, use cases, tests
- ◆ CS database brain trust
  - design, direct building of the system
  - provide some resources
- ◆ New Company
  - manage open source project
  - contribute engineering/development resources
  - provide support, services, PR



# *Partners from Science*

- ◆ LSST (astro)
- ◆ PNNL (bio, atmospheric)
- ◆ LLNL (fusion)
- ◆ FermiLab (hep, astro)
- ◆ UCSB (remote sensors)
- ◆ ...your lab/project?



# *CS DB Brain Trust*

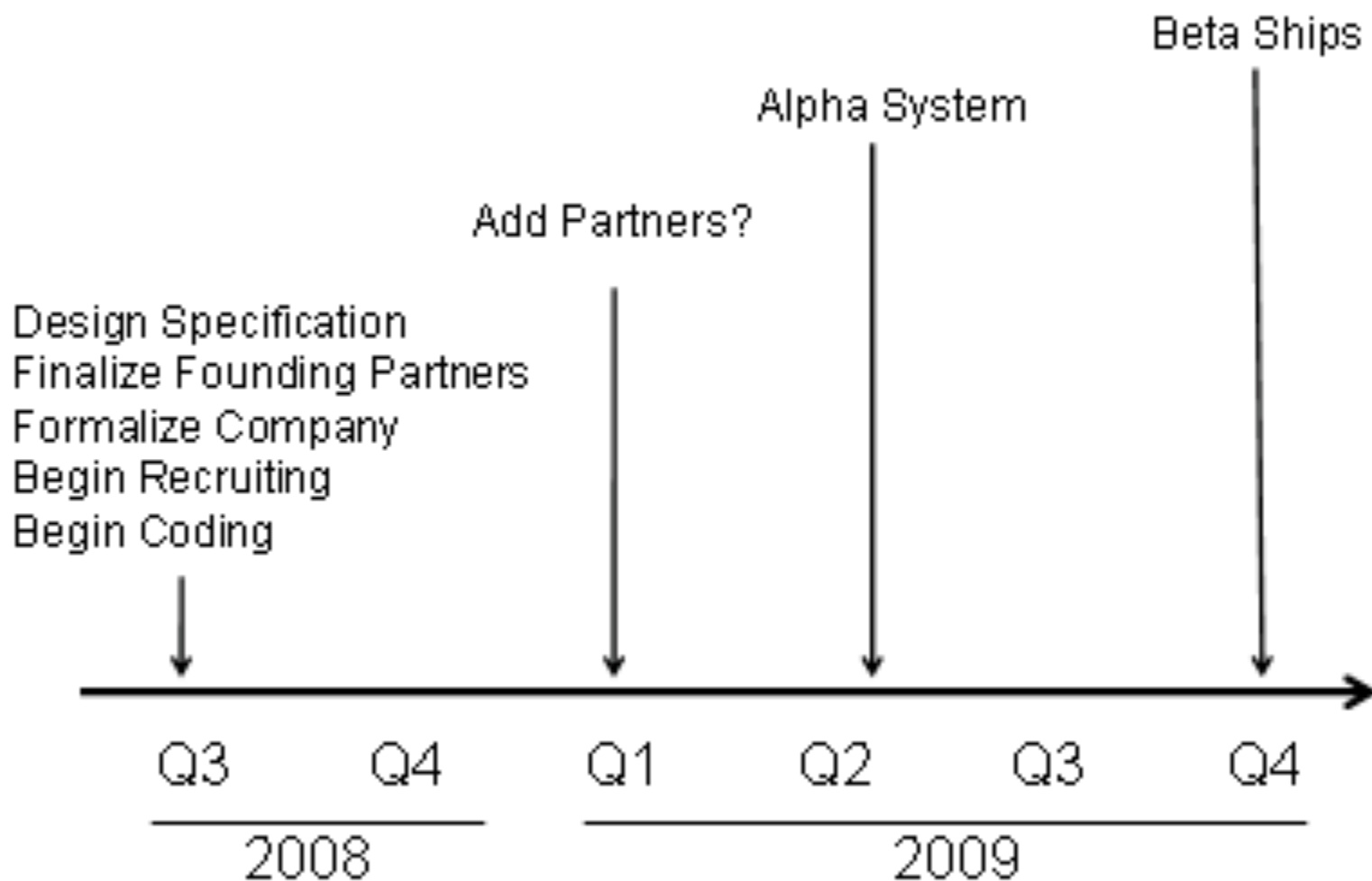
- ◆ Mike Stonebraker (MIT)
- ◆ David DeWitt (Wisconsin, Microsoft)
- ◆ Jignesh Patel (Wisconsin)
- ◆ Dave Maier (Portland State)
- ◆ Stan Zdonik (Brown)
- ◆ Sam Madden (MIT)
- ◆ Ugur Centintemel (Brown)
- ◆ Martin Kersten (Netherlands)



# *New Company*

- ◆ Funding available
- ◆ Founding happens “now”
- ◆ Office space provided by SLAC
- ◆ Details of business model etc being discussed now

# Community Development & NewCo





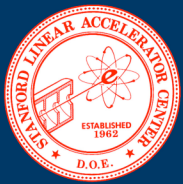
# *Partners From Industry*

- ◆ EBay: can fund it all (good and bad!)
- ◆ New Enterprise Associates: possibly \$1 million seed investment
- ◆ Facebook and Amazon interested in collaboration
- ◆ LG Electronics: possible VC funds in future
- ◆ Microsoft: discussions in progress
- ◆ Google and Yahoo!: have their own solution



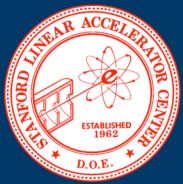
# *Design Meeting*

- ◆ Tomorrow at MIT
  - CS database brain trust + key partners
- ◆ Will discuss architecture, assignments
  - based on collected input/use-cases from scientific and industrial partners



# *Homework For Science*

- ◆ Strengthen involvement
  - Organize Scientifica Science Board, reach more labs and projects
- ◆ Strengthen “buy-in”
- ◆ Work closely with Scientifica team
  - steer, help develop, test
- ◆ Strengthen XLDB community, identify common needs
  - 2<sup>nd</sup> XLDB workshop planned for Sept 29/30 @SLAC



# Summary

- ◆ Open source, science-oriented DBMS is within reach
- ◆ Led by most influential database gurus
- ◆ Designed by most experienced database engineers
- ◆ Strong support and interest from large industrial companies

**Truly unique opportunity**