

# Computing for the LHC: The next step up

Torre Wenaus, BNL/ATLAS

October 15, 2013

CHEP 2013

Amsterdam

**BROOKHAVEN**  
NATIONAL LABORATORY

*a passion for discovery*



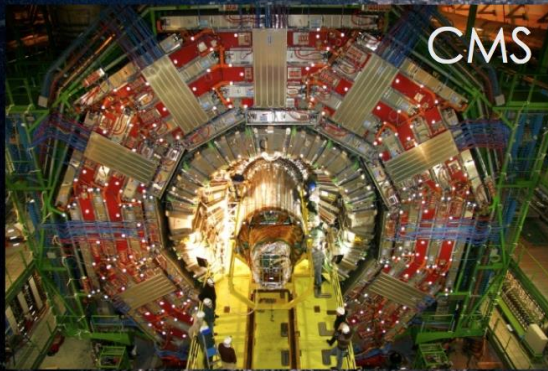
# Outline

- The LHC program and its computing demands
- Survey of (selected!) technical areas with notable lessons, developments, and implications for the future
- Conclusions

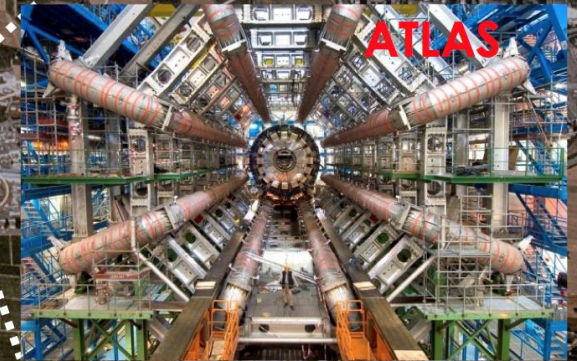
Will go lightly on important LHC computing topics covered in other plenaries – which leaves plenty to talk about!

- Concurrency, multicore, GPUs – Jim Kowalkowski, Axel Naumann
- Physics software packages – Philippe Canal
- Big Data processing in HEP – Brian Bockelman
- Intelligent, advanced networking – Inder Monga, Harvey Newman
- Data archiving and data stewardship – Pirjo-Leena Forsström

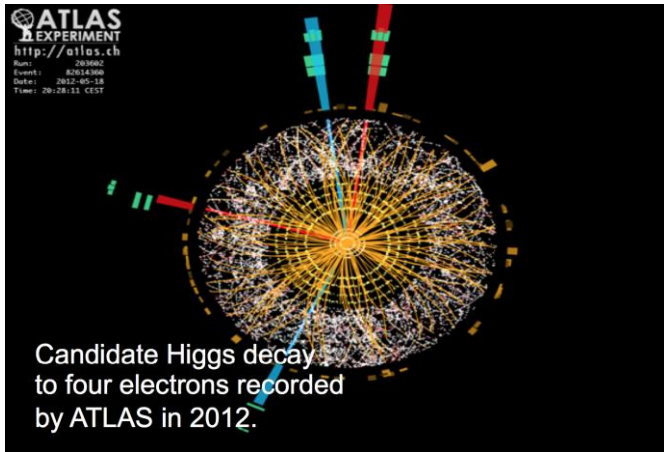
# The LHC... Needs no introduction here



Exploration of a new energy frontier  
Proton-proton and Heavy Ion collisions  
at  $E_{CM}$  up to 14 TeV



# LHC Run 1



Global Effort → Global Success July 4, 2012

Results today only possible due to extraordinary performance of accelerators – experiments – Grid computing

Observation of a new particle consistent with a Higgs Boson (but which one...?)

Historic Milestone but only the beginning

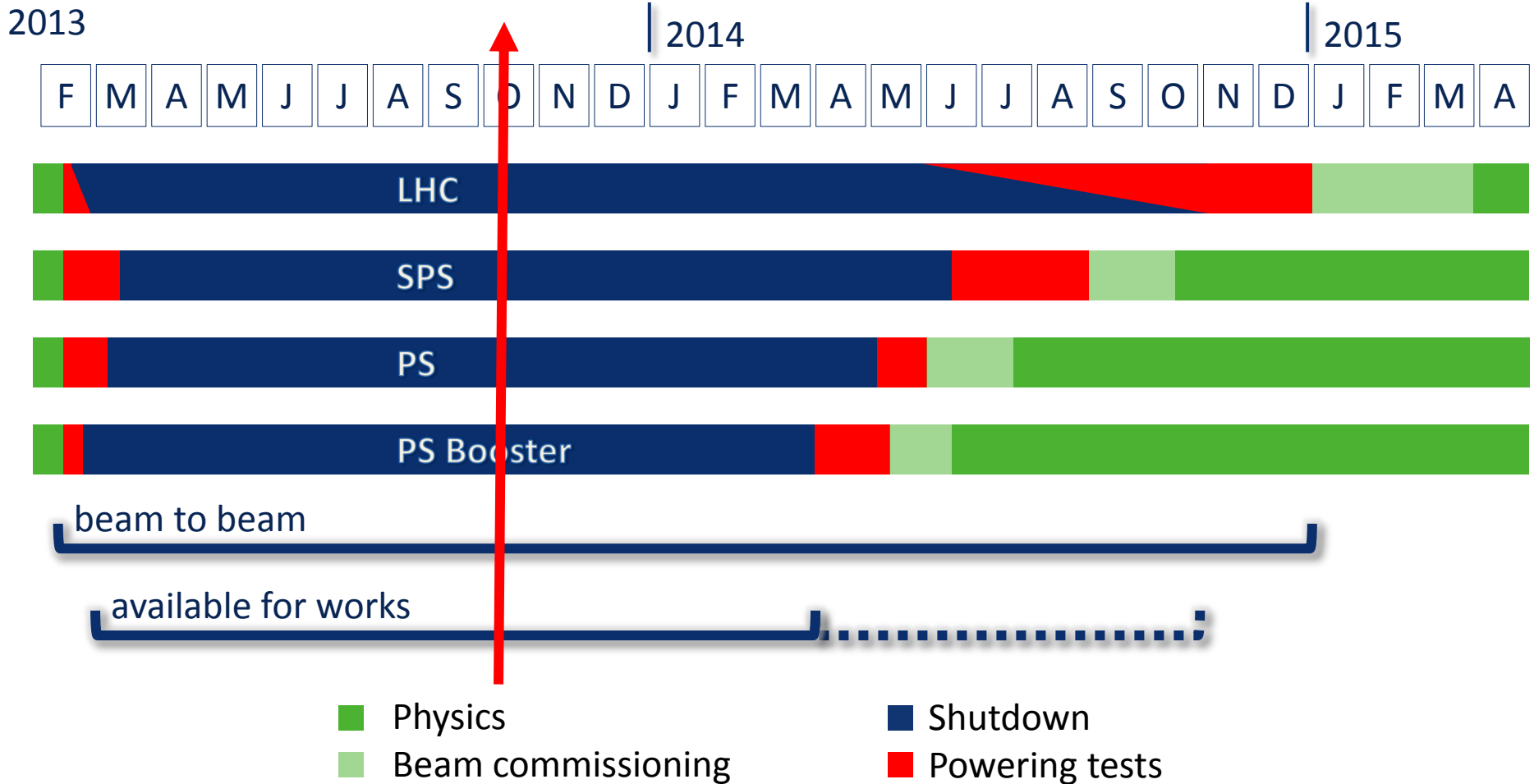
Global Implications for the future

CERN  
R-D Heuer

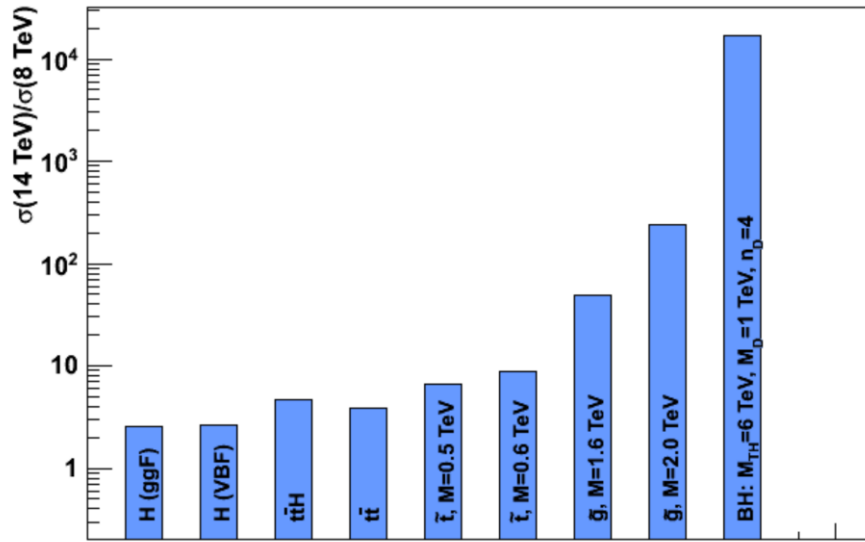
- Billions of events delivered to the experiments from proton-proton and proton-lead collisions in the Run 1 period (2009-2013)
- Collisions every 50 ns = 20 MHz crossing rate
- ~35 interactions per crossing (pile-up) at peak luminosity (LHCb has ~1.7 on average with luminosity leveling)
- ~1600 charged particles produced in every collision
- Vast and dramatic scientific output

# LHC Long Shutdown 1 (LS1)

We're a third of the way through LS1...



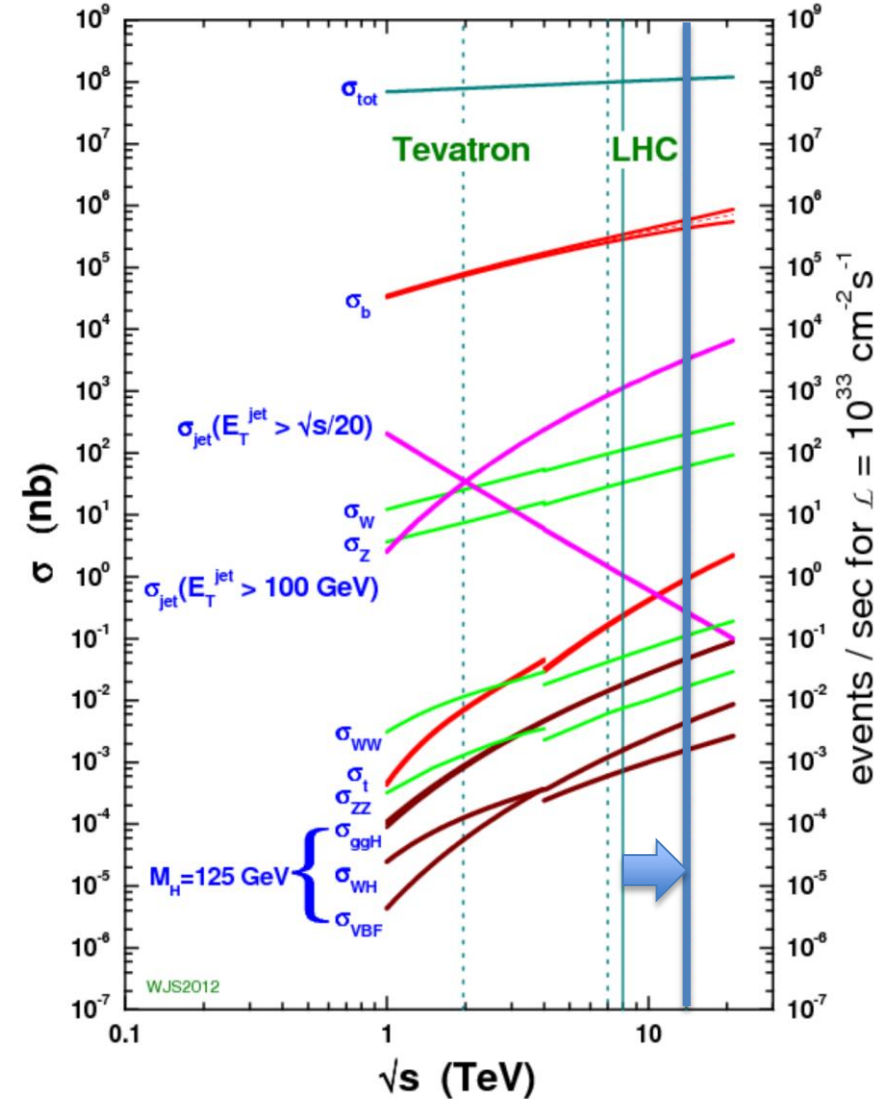
... beam still foreseen for January 2015 (physics in ~April)



- 2015 is a whole new energy domain
  - Every measurement is new
  - Every search has a new chance
- Early high energy LHC data have discovery potential
  - With a few  $\text{fb}^{-1}$  many searches reach or surpass current sensitivity
  - Discovery of TeV scale particles possible
- Center of mass energy  $\sim 2x$
- ATLAS, CMS event rate to storage  $\sim 2x$ , pile-up above 30, 25 ns bunch spacing
- Experiments asking for much less of a computing increment than extrapolation from 2013 practice would suggest (e.g. for CMS,  $\sim 2x$  instead of  $\sim 6x$ )

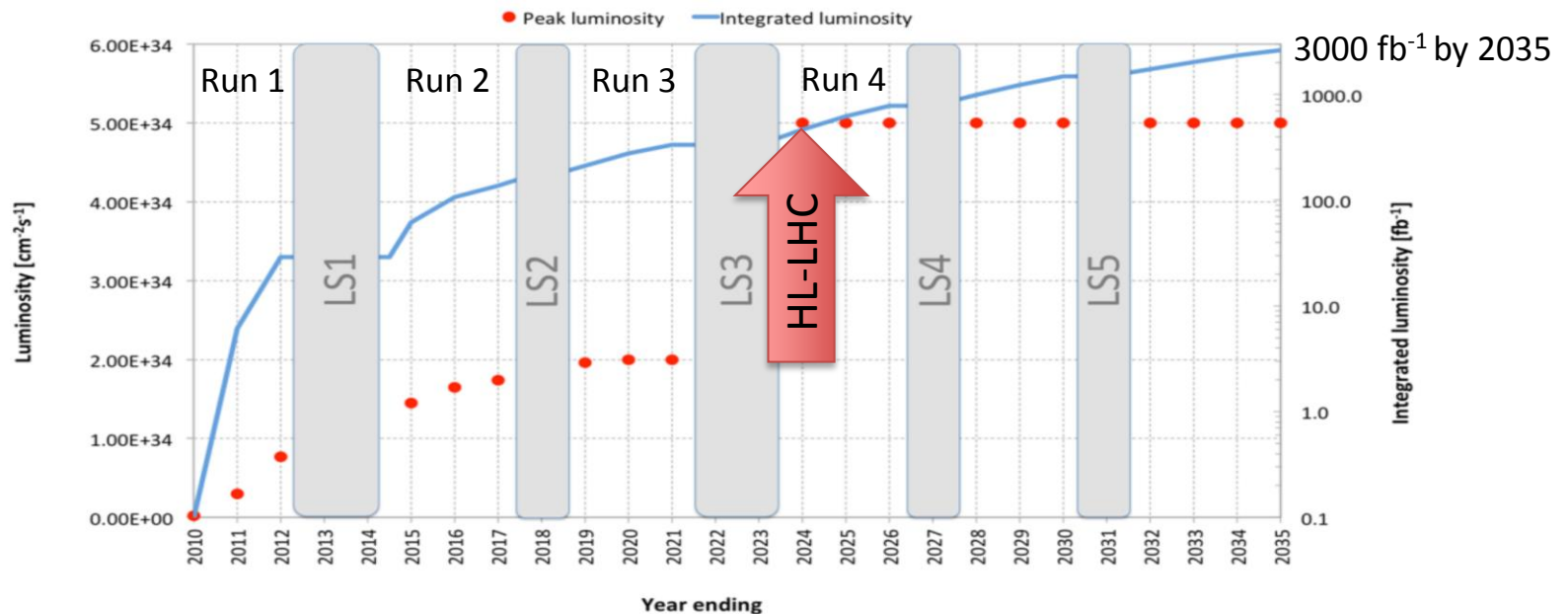
# LHC Run 2

proton - (anti)proton cross sections



# LHC Beyond 2015

## to the High Luminosity LHC (HL-LHC)



- Trigger rates, event complexity increase steadily through machine and detector upgrades
  - **~15 PB/year LHC raw data now; ~130 PB/year in 2021**
  - Very rough estimate for **new raw data per year in Run 4: 400 PB**
  - **Raw data only the beginning**, e.g. ATLAS dataset is ~140 PB, ~70 PB on disk
  - Pile-up reaches ~150 at HL-LHC, multiplicity up 8x
  - Very rough CPU estimates lie within Moore's law limits but **presume performance gains such that we are still able to leverage Moore's law**

# LHC Computing in Run 2 and Beyond

- **Storage and processing extrapolations lead to unacceptable costs (flat budget assumption)** – we must work on performance and efficiency
- **Storage is largest cost**, e.g. ATLAS spends ~60% more money on disk than on CPU
- As to processing cost, **must track Moore's law as effectively as possible**
  - *Adapting to new processors is much more challenging than in the past*
- One approach is of course to do less
  - Cost constrains the data rate already, driving choices on triggers and analyses
  - Write some data to tape only, pending a physics case to analyze?
    - A penny puts 1000 CMS events on tape
- Before taking steps like these, first make the most of what we have
- Therein lies the LS1 (and beyond) computing upgrade program



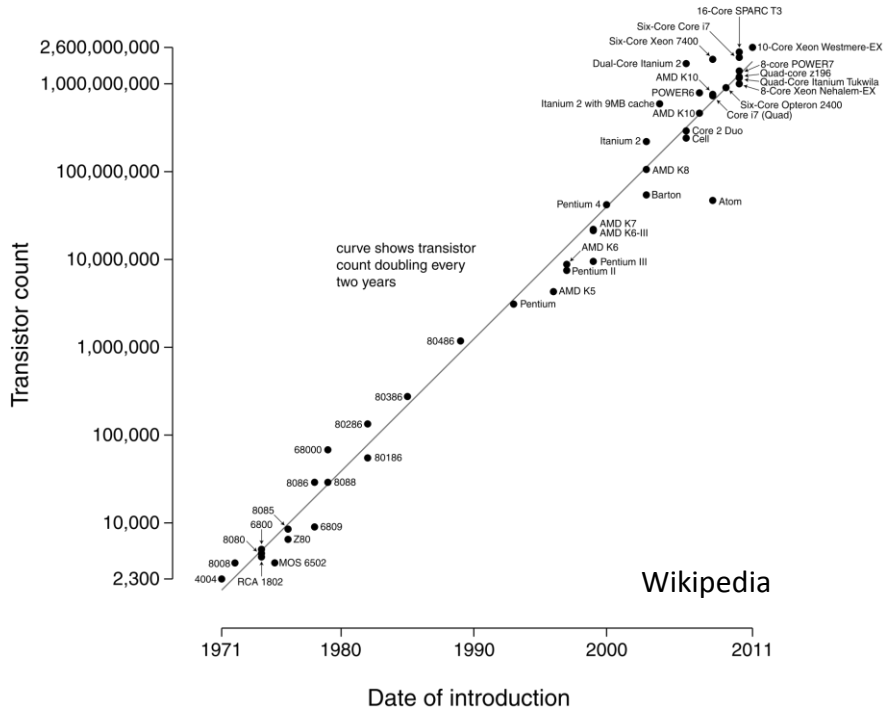
# Upgrading LHC Computing in LS1

- LS1 is a valuable opportunity to assess
  - the lessons and experiences of Run 1
  - the computing demands of Run 2
  - the technical and cost evolution of the computing milieu within which we will operate
- And so informed, undertake an intensive planning and development cycle to ready LHC computing for 2015 and beyond
  - While sustaining steady state full scale operations
  - With an assumption of flat (at best) funding
- This has been happening internally to the experiments and collaboratively with CERN IT, WLCG, common software and computing projects

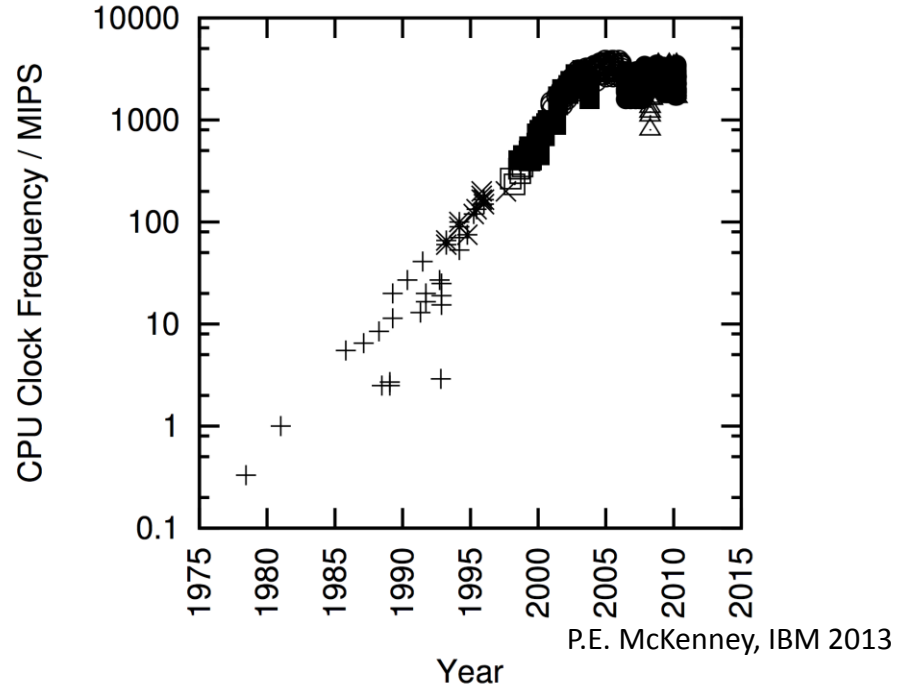
# Processing

Transistor count growth is holding up...

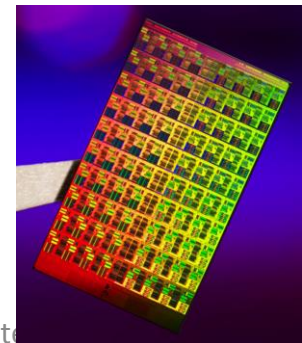
Microprocessor Transistor Counts 1971-2011 & Moore's Law



...but clock speed growth died a heat death...



... replacing the free lunch of ever faster processors with the necessity of sustaining throughput growth by leveraging growth in core count, co-processors, concurrency features



Intel 80-core chip

# Adapting Software to New Processors

- High concurrency, modest memory/core, GPUs, ... is the new environment
  - Multi-core now → many-core soon → finer grained parallelism needed
  - GPUs present challenges in programmability and data bottlenecks
  - *Many or most of our codes require extensive overhauls*
- *The whole world faces it, which is good news* – tools, libraries, compilers are emerging that can help the migration from serial software
- Our common tools are being adapted, to the benefit of all
  - Geant4, ROOT, math libraries, reconstruction tools, ...
- But it's a large effort on the part of scarce software experts
  - The software investment is necessary – *living with inefficient software is much more expensive*
  - Objective is products widely usable in the HEP community
- CERN-led Concurrency Forum is providing a collaborative context for this

# Simulation

- Agreement between simulation and data in Run 1 was spectacular
- Challenge now is to preserve the physics quality while maximizing performance, particularly on new architectures
- Most LHC CPU cycles go to simulation (60-70%) – a lot to gain
- Multithreaded Geant4 release imminent, R&D to utilize GPUs, ...
- Tailor simulation detail detector by detector appropriately for the physics
  - Optimally balance precision and resource consumption
  - With fast digitization, reconstruction to avoid bottlenecks
  - E.g. ATLAS' Integrated Simulation Framework (ISF) and Very Fast Simulation (VFS)
- ALICE migrating from Geant3 to Geant4, working with G4 experts on improving performance
  - Also work on fast parameterized simulation, more use of embedding, event mixing, ...
- Be ready and able to use opportunistic resources

# Reconstruction

- CPU needs grow rapidly with track multiplicity (pile-up). Tracking a major focus for improvements
- Substantial code rewrites to optimize quality and performance, introduce concurrency
- Studying and selecting new linear algebra and other libraries optimized for concurrency on new architectures
- Leverage tracking reconstruction already performed in HLT
- Leverage GPUs (e.g. 1 GPU = 3 CPUs for ALICE online tracking use case)
- Optimize data organization/access further – disk access an ongoing constraint

# Analysis Software & Systems

- ROOT is the universal basis for analysis, analysis level event data storage
  - ROOT role still evolving... e.g. new analysis format/tools (ATLAS, LHCb) will enable ROOT and Gaudi based analysis of both AOD and NTUPLE level data
  - Major new release ROOT 6 being prepared
- Experiments moving more analysis activity into organized workflows – aim to shrink ‘chaotic analysis’ to only what really is user specific
- Remove redundancies in processing and storage, reducing operational workloads while improving turnaround for users
- ALICE train model drawing interest from other experiments as basis for more managed analysis production

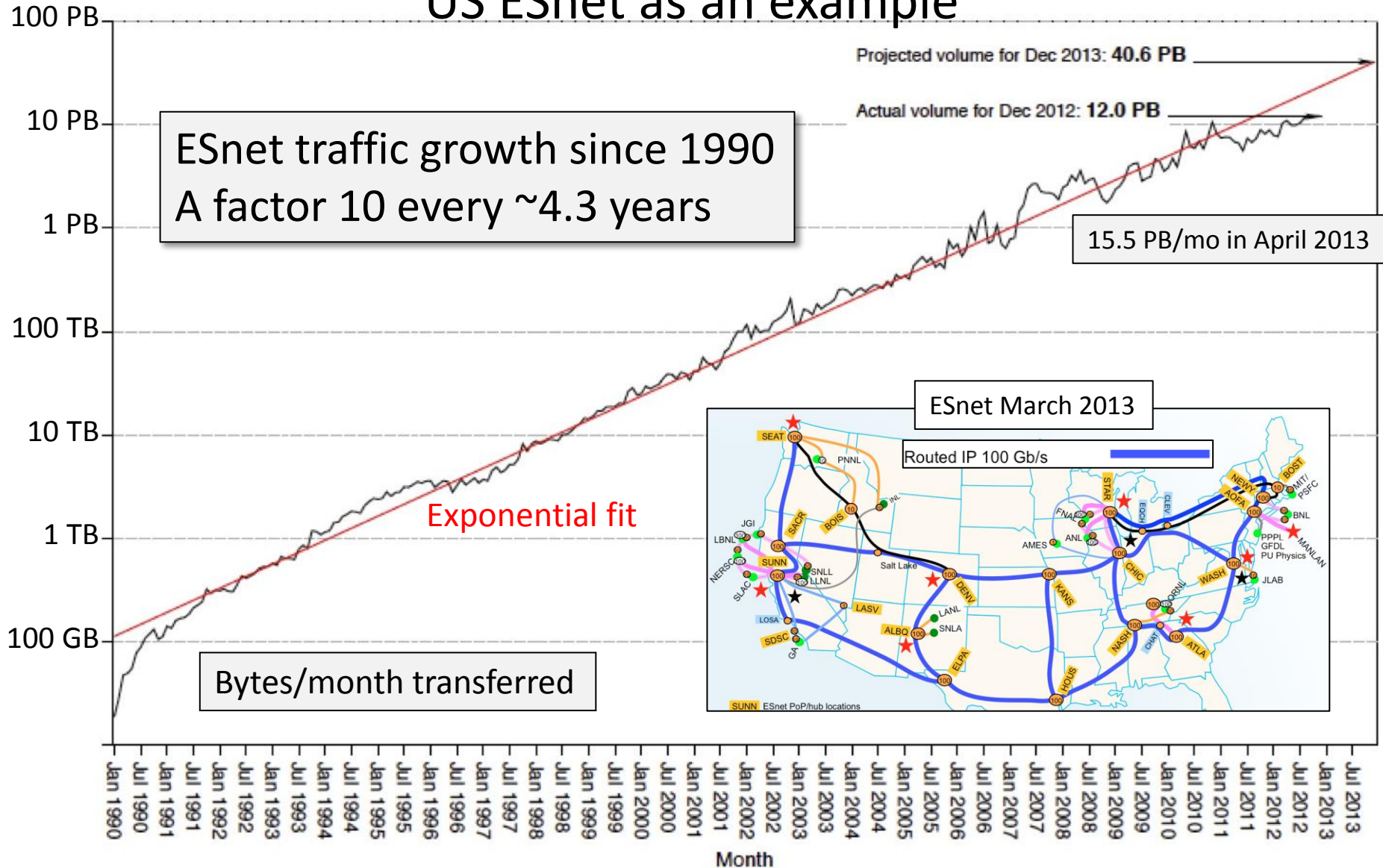
# Networking

- Enabling backbone of LHC computing is reliable, high-bandwidth, feature-rich networks
- HEP was a pioneer in network-intensive science and international research networks, and continues to lead
  - Networks optimized for massive data flows, e.g. now testing the first 100Gb transatlantic production link
- **Making the most of the network translates to more science at lower computing cost**
  - Important that we design our workflows around this fact
- Next generation networks allow applications to interact with the network, reacting to conditions and proactively controlling it
  - e.g. work underway to integrate network awareness in job brokerage data distribution (PheDex) and job brokerage (PanDA)

**In general it's much cheaper to transport data than to store it**

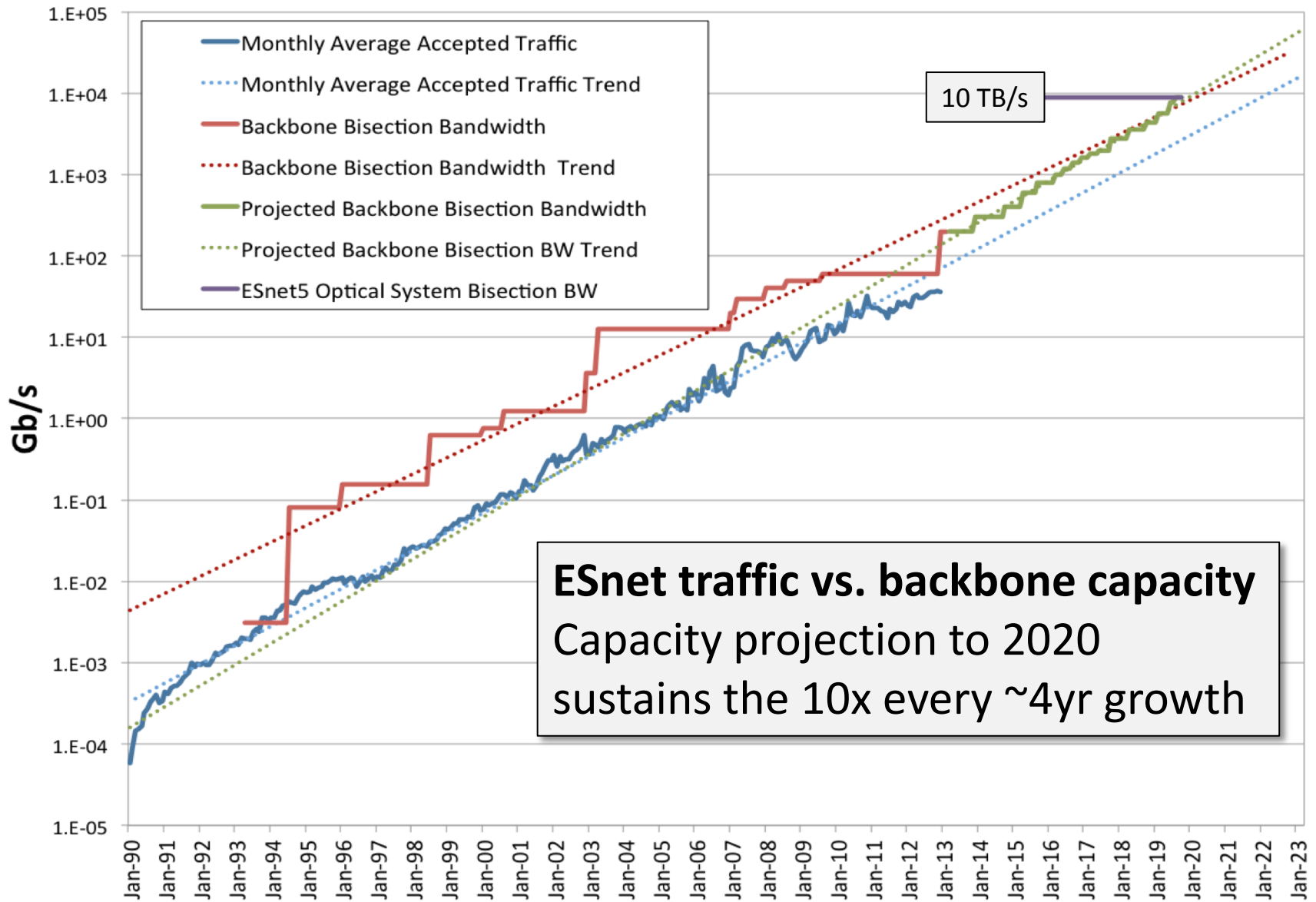
# Networking growth has been dramatic

## US ESnet as an example





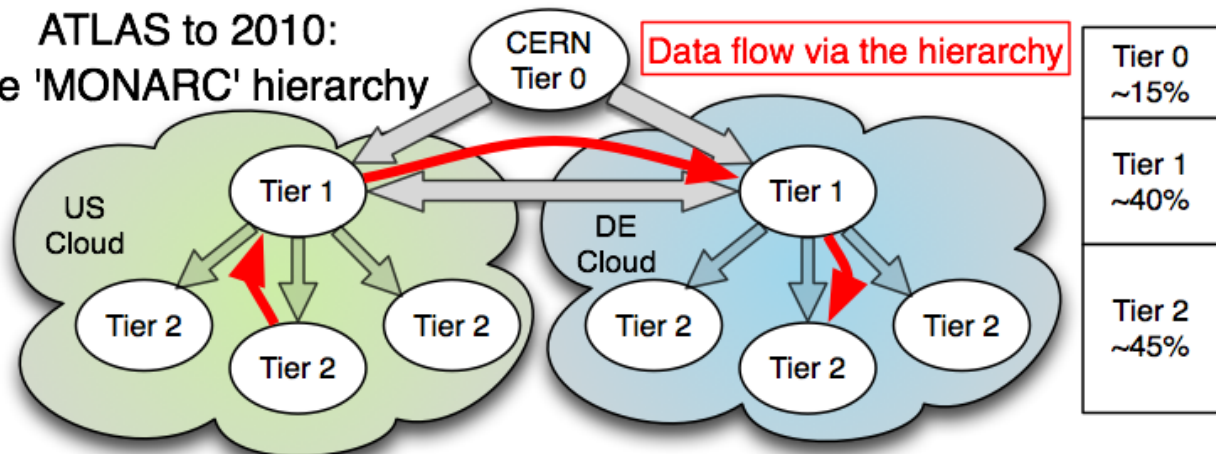
# Planned capacity growth sustains the trend



# Networking has been a critical enabler for evolving LHC computing models – ATLAS as example

ATLAS to 2010:

The 'MONARC' hierarchy



... 10 clouds/Tier 1s, ~70 Tier 2 sites

**Original model:**

Static strict hierarchy  
Multi-hop data flows  
Lesser demands on  
Tier 2 networking  
Virtue of simplicity

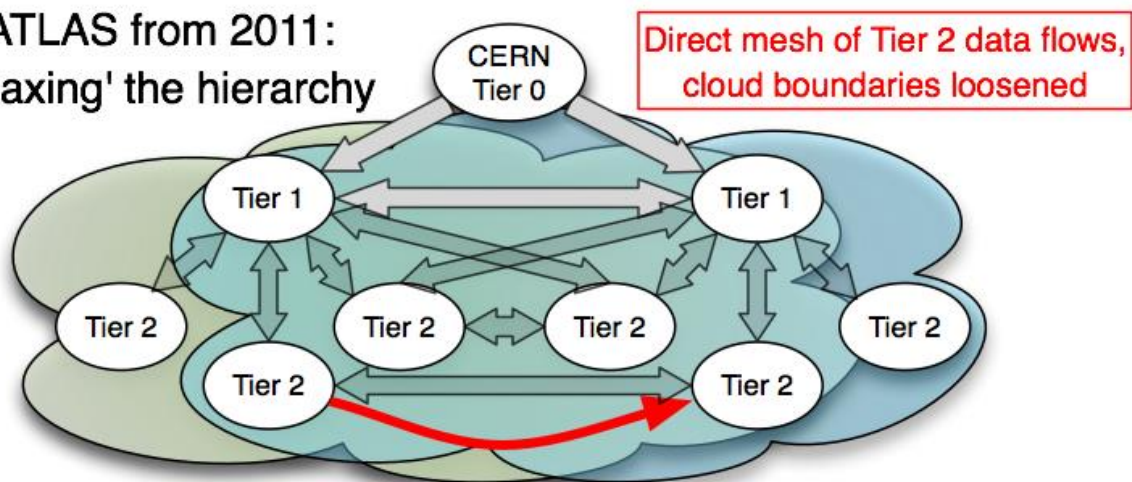
**Designed for <~2.5 Gb/s  
within the hierarchy**

**Today:**

**Bandwidths 10-100 Gb/s, not limited  
to the hierarchy**

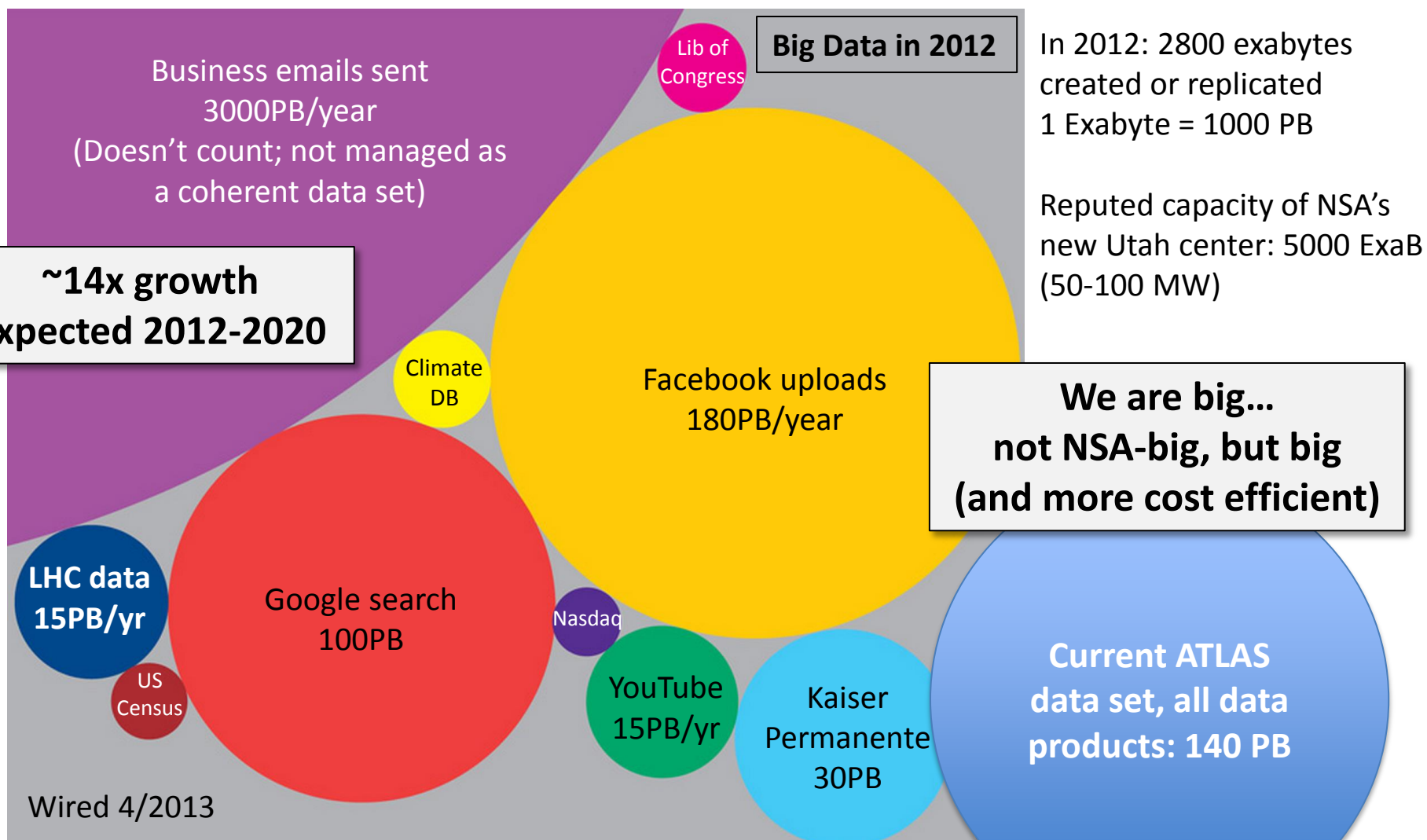
Flatter, mostly a mesh  
Sites contribute based on capability  
**Greater flexibility and efficiency**  
**More fully utilize available resources**

ATLAS from 2011:  
'relaxing' the hierarchy



# Data Management

## Where is LHC in Big Data Terms?



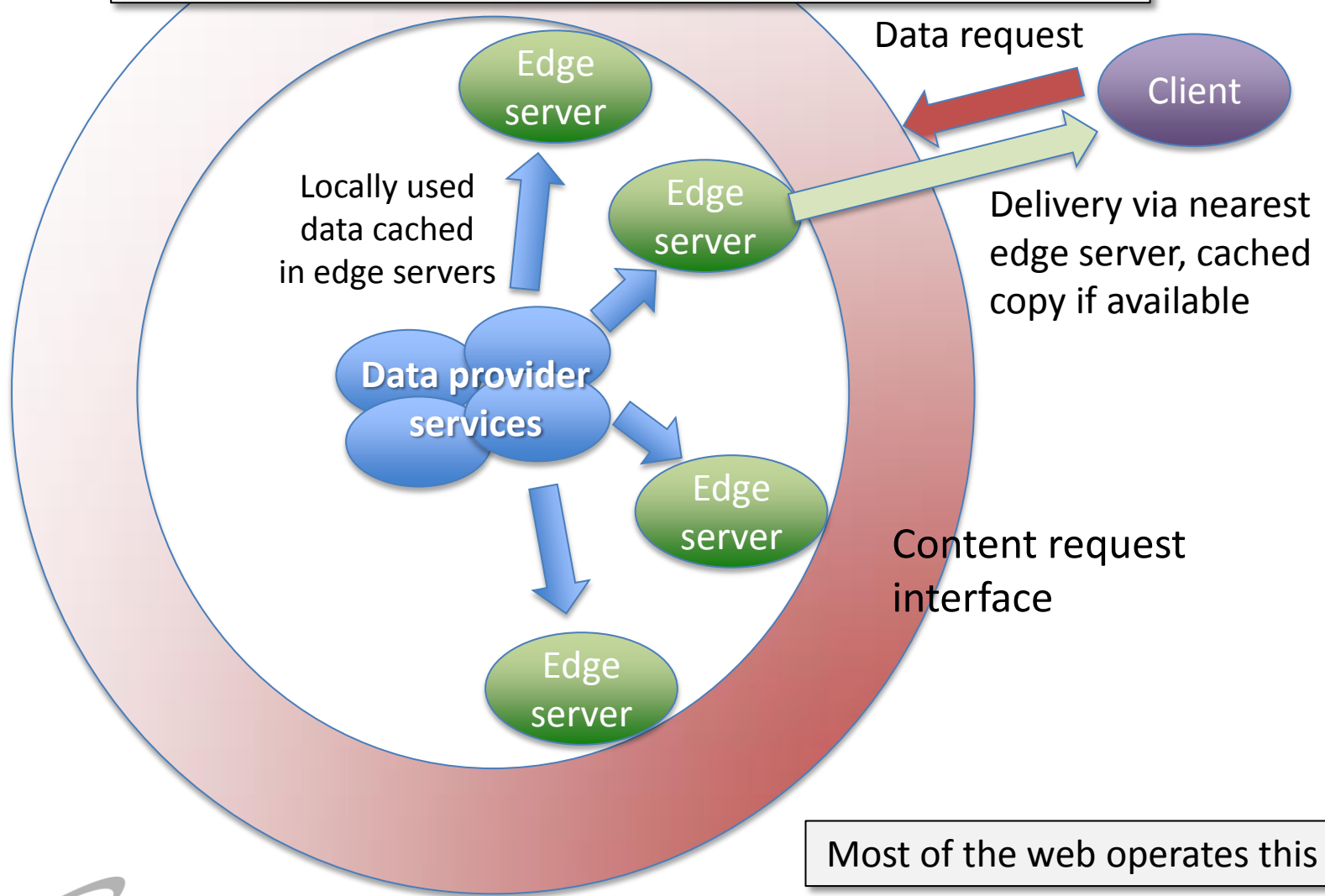
<http://www.wired.com/magazine/2013/04/bigdata/>

# Evolution of LHC Data Management

- Surfing the Big Data wave helps – e.g. we have cases where new technologies such as **Hadoop** help us scale – but we can't afford to simply scale up
- We need **more efficient distributed data handling, lower disk storage demands, lower operational load** (storage is highly labor intensive for operations)
- Storage savings can be turned into CPU capacity
- (Aspire to) *send only the data you need, only where you need it* (and cache it when it arrives)
- One successful approach in Run 1 is being pursued more widely: building **intelligent dynamic data placement** into workflows
- Another now reaching production: **Federated data storage**
- Further steps underway or planned could transform our traditional approach
  - Dispensing with file-based management and delivering events – **event service**
  - Once you're delivering events you can trade off retrieval against on-demand generation – **virtual data**
- Industry has been at this approach for years, in **content delivery networks**

# The Content Delivery Network Model

Content delivery network: deliver data quickly and efficiently by placing data of interest close to its clients



# The Content Delivery Network Model

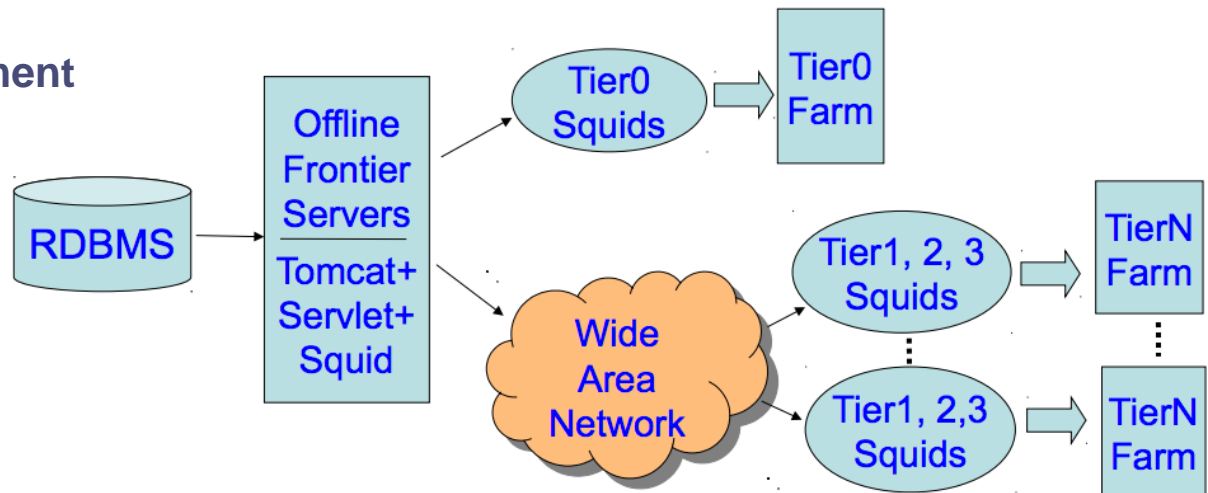
A growing number of HEP services are designed to operate broadly on the CDN model

Service	Implementation	In production
Frontier conditions DB	Central DB + web service cached by http proxies	~10 years (CDF, CMS, ATLAS, ...)
CERNVM File System (CVMFS)	Central file repo + web service cached by http proxies and accessible as local file system	Few years (LHC expts, OSG, ...)
Xrootd based federated distributed storage	Global namespace with local xrootd acting much like an edge service for the federated store	Xrootd 10+ years Federations ~now (CMS AAA, ATLAS FAX, ...) <i>See Brian's talk</i>
Event service	Requested events delivered to a client agnostic as to event origin (cache, remote file, on-demand generation)	ATLAS implementation coming in 2014
Virtual data service	The ultimate event service backed by data provenance, regeneration infrastructure	Few years?

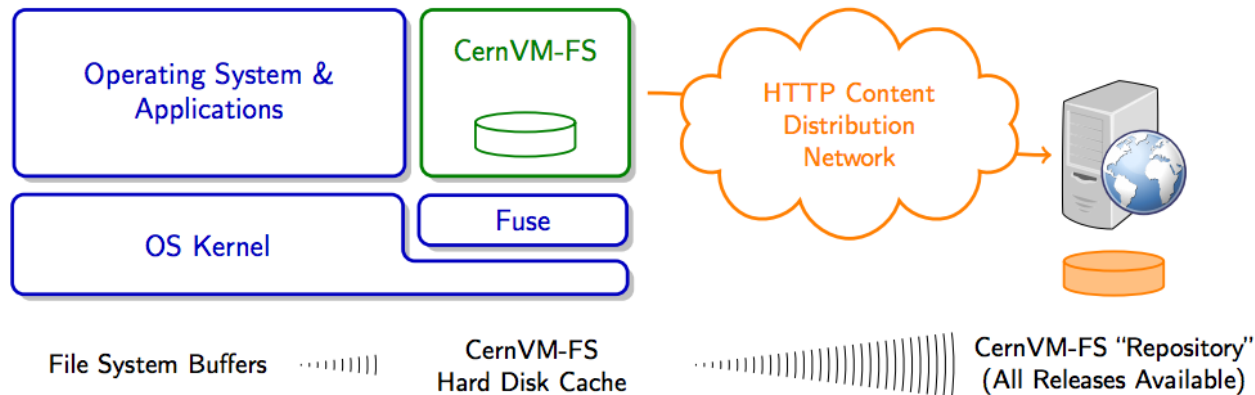
# Frontier for Scalable Distributed DB Access

- How to provide scalable access from tens of thousands of global grid jobs to centrally Oracle-resident detector conditions data?
- The answer, first developed by CDF/CMS at FNAL, and adopted for LHC: Frontier
- Frontier is a web service that translates DB queries into HTTP
  - Far fewer round trips between client and DB than using Oracle directly; fast over WAN
- Because it is HTTP based, caching web proxies (squid) can provide hierarchical, highly scalable cache based data access
- Very successful

## CMS Frontier deployment



# CERNVM File System (CVMFS)



**Caching HTTP file system optimized for software delivery**

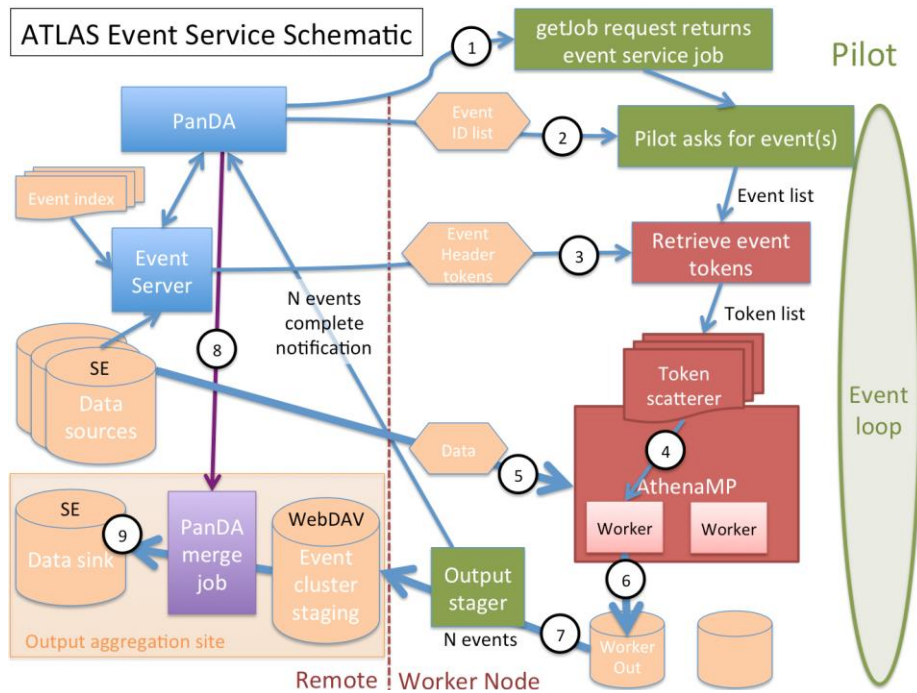
**Efficient:** compression over the wire, duplicate detection

**Scalable:** works hand in hand with proxy caches

- Originally developed as lightweight distributed file system for CERNVM virtual machines
  - ◆ Keep the VM footprint small and provision software transparently via HTTP and FUSE, caching exactly what you need
- Adopted by all LHC experiments and beyond, very successful for general software distribution (and other file data)
- Other new services leveraging it, e.g. OASIS from OSG uses CVMFS to provide easy provisioning of software on grid worker nodes via a centrally managed repository



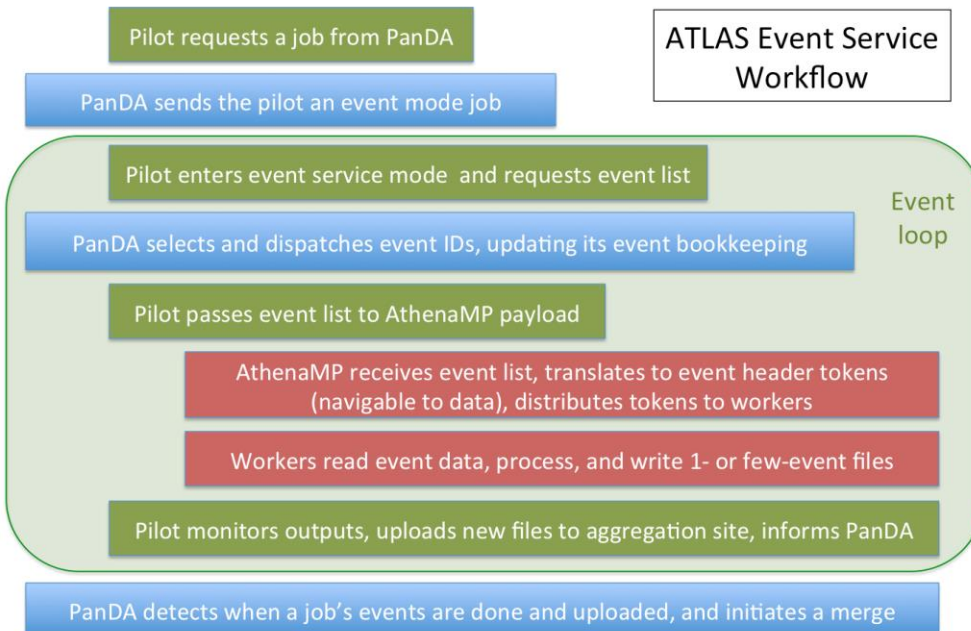
## ATLAS Event Service Schematic



# ATLAS Event Service

- Ask for exactly what you need, have it delivered by a service that knows how to get it to you efficiently
- Return the outputs in a ~steady stream, such that a WN can be lost with little lost processing
- Well suited to transient opportunistic resources, hole filling, volunteer computing
- Draws on developments in PanDA, Prodsys2, AthenaMP, event I/O
- Probably ready to try in first half 2014, on simulation

## ATLAS Event Service Workflow



# Workload Management

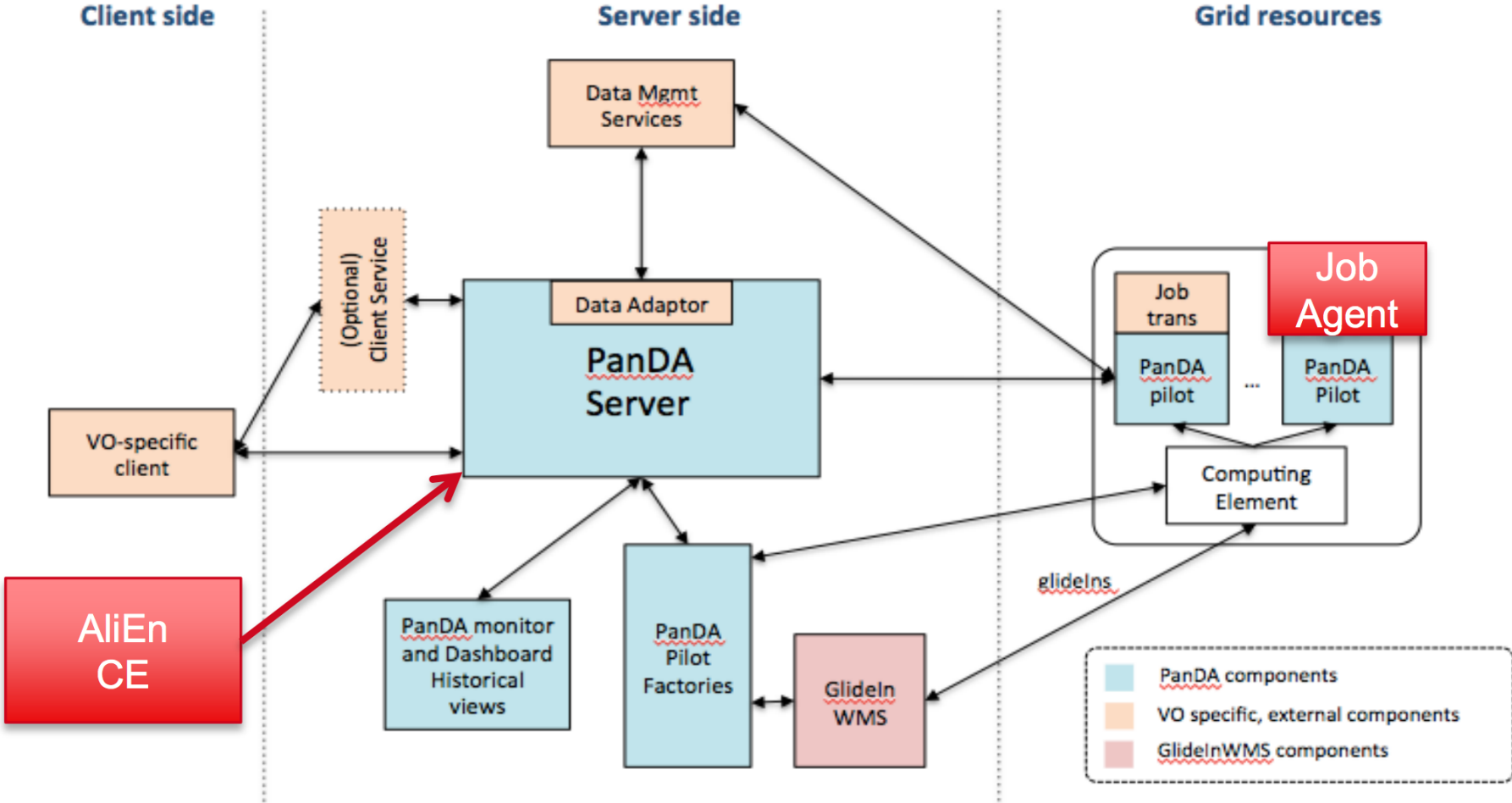
- Pilots, pilots everywhere
  - Approach spread from ALICE and LHCb to all experiments
- Experiment-level workload management systems successful
  - No further need for WMS as middleware component
- Scales very well – PanDA does 1.5M jobs/day without breathing hard
- Is WMS commonality still possible then? Yes...
  - DIRAC (LHCb) long established as a tool common to multiple experiments, as is AliEn (ALICE)
  - PanDA (ATLAS) is spreading to others: Common Analysis Framework (CAF) project with CMS and CERN IT; ALICE now evaluating; in use beyond LHC (e.g. AMS)
    - PanDA project has support to generalize it for the wider community, extend it to supercomputers, integrate intelligent networking
- Pilot submission and management has a powerful enabler in HTCondor, and its glideinWMS layer

# CMS, ALICE integration with PanDA

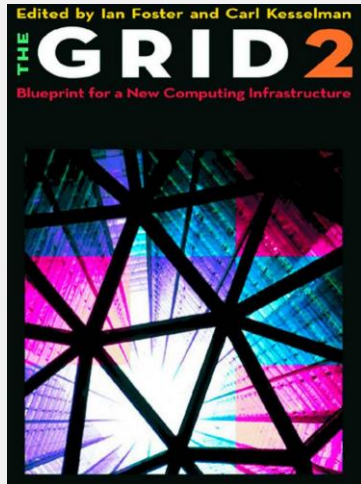
PanDA core

Refactoring for CMS (et al)

AliEn integration



# Utility Computing: Grids to Clouds



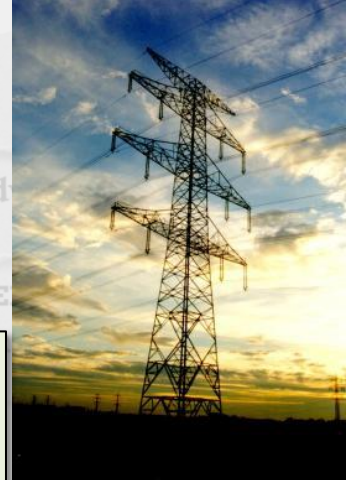
**The Grid: 1998 and 2003 (2<sup>nd</sup> Ed.)**

*Grid is used by analogy with the electric power grid... has had a dramatic impact on human capabilities...*

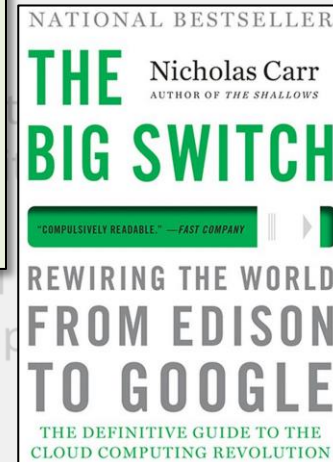
- Clouds complement and extend the grid
  - Decrease heterogeneity seen by the user
  - Shrink middleware, interface software accordingly
- Flexible, dynamic resource sharing, eg Tier 2/Tier 3
- New resources (commercial, research clouds)
- VMs provide a uniform user interface to resources
  - Integrate diverse resources manageably
  - Isolate software from physical hardware
- Basis for volunteer computing
- Long term data preservation solution
- Draws on a growing tool suite, in-house and wider
  - CernVM, CVMFS, OpenStack, ...

**The Big Switch [to the Cloud]: 2009**

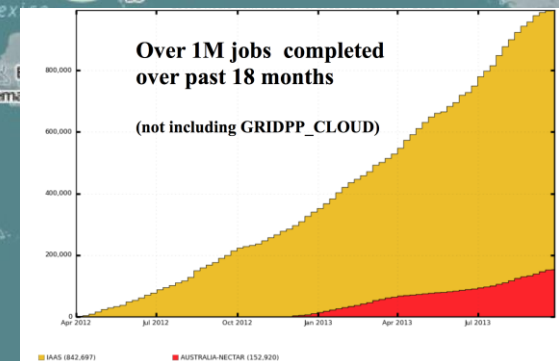
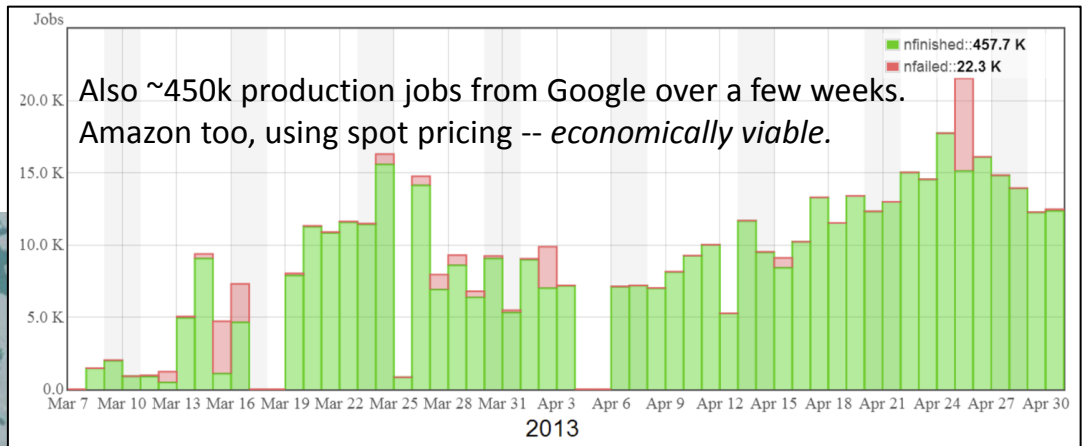
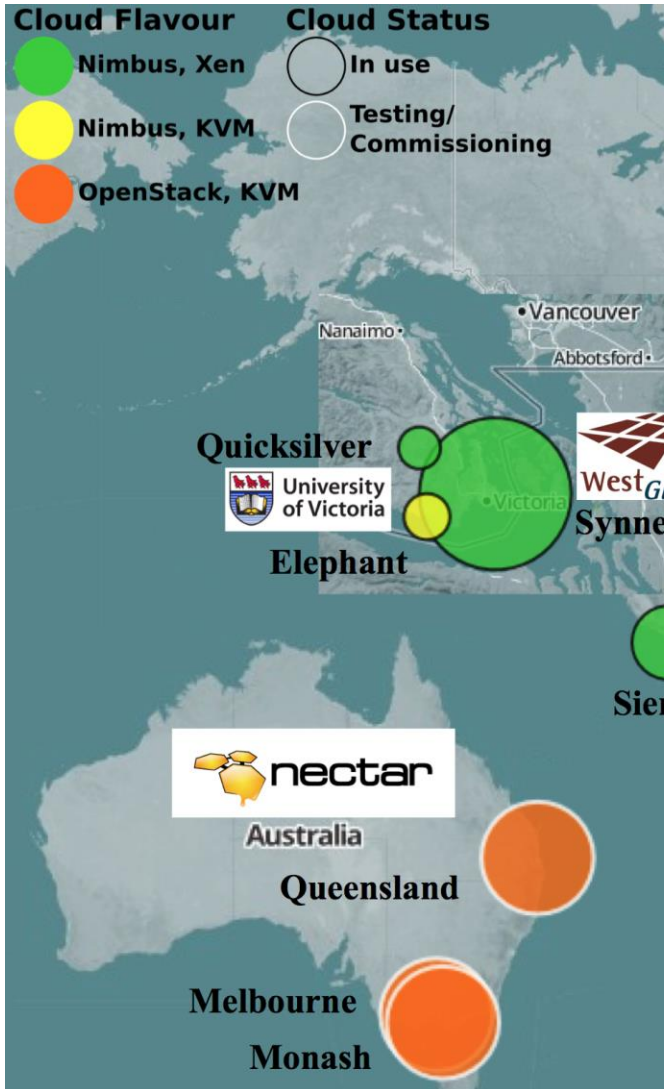
*Computing is turning into a utility... will ultimately change society as completely as cheap electricity did...*



omputing the  
ig Thing?"  
ve become  
ctioning of  
r stations."  
conomist

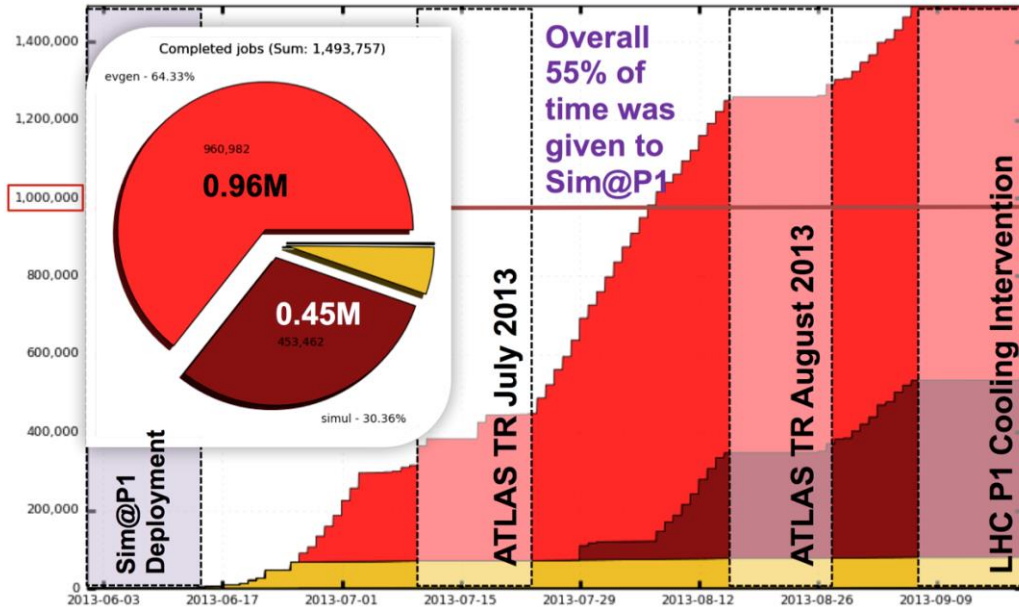


# “Grid of Clouds” used by ATLAS



# Opportunistic Resources – HLT

## ATLAS Production at P1 HLT Farm



June 1 – September 18, 2013 Total 1.5M jobs

The largest ATLAS grid site when running

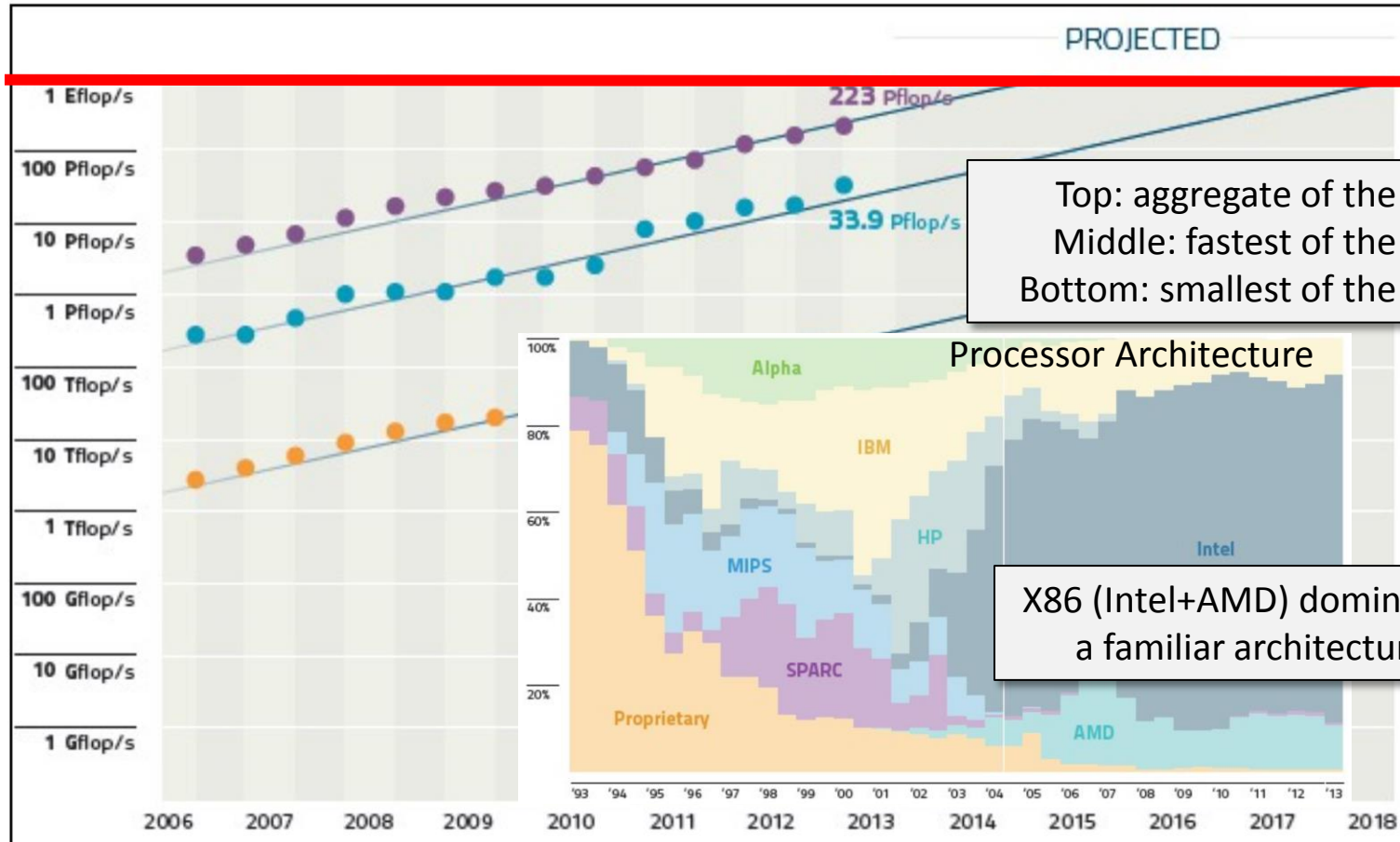
- Experiments' own HLT farms can be a major resource for offline processing, opportunistically
- LHCb has been using theirs since early 2013, ~20% of their resources in 2013
- CMS, ATLAS came into production this year, building OpenStack based cloud platforms

# Opportunistic Resources – HPCs

- HPC (supercomputing) resources can be valuable to HEP computing
- They have cycles open to us – many – even though we wouldn't build the machines that way if we were paying for them (the point is, we aren't)
  - They have holes we can fill: cycles instead of sitting idle would be going to high profile science
  - The *current* US national HPC allocation for HEP is comparable to global CMS+ATLAS computing in 2012, ~1.5B hours
- Also there is increasing convergence, making our apps more appropriate
  - HPC has a growing number of data intensive use cases, future architectures will have to take this into account
  - **More concurrency, leveraging architectures used in HPCs make our applications more suited to HPC**
- We're porting appropriate applications (generators, simulation) and extending workflow and data management systems to support them
- We've begun to put HPC facilities into production

# HPC growth remains rapid... Exaflop systems by 2020?

1 Exaflop



Top: aggregate of the Top 500  
Middle: fastest of the Top 500  
Bottom: smallest of the Top 500

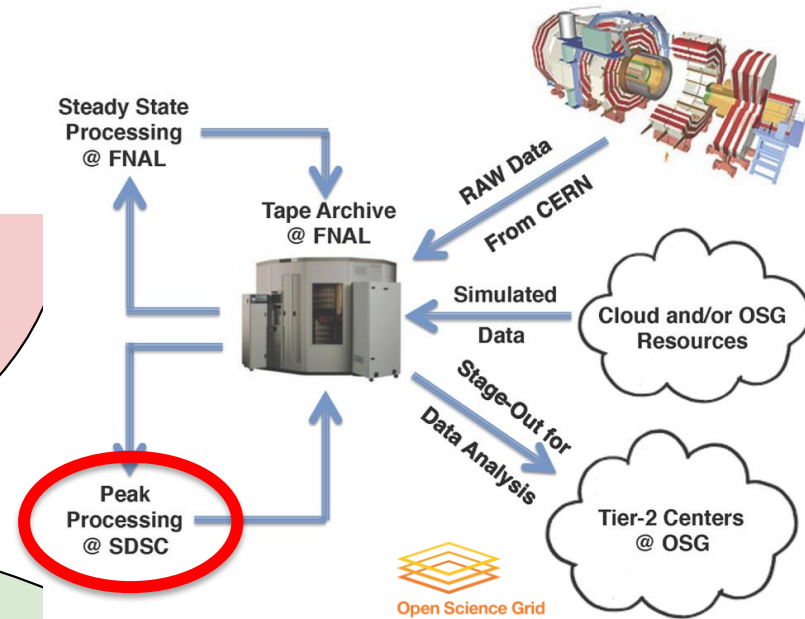
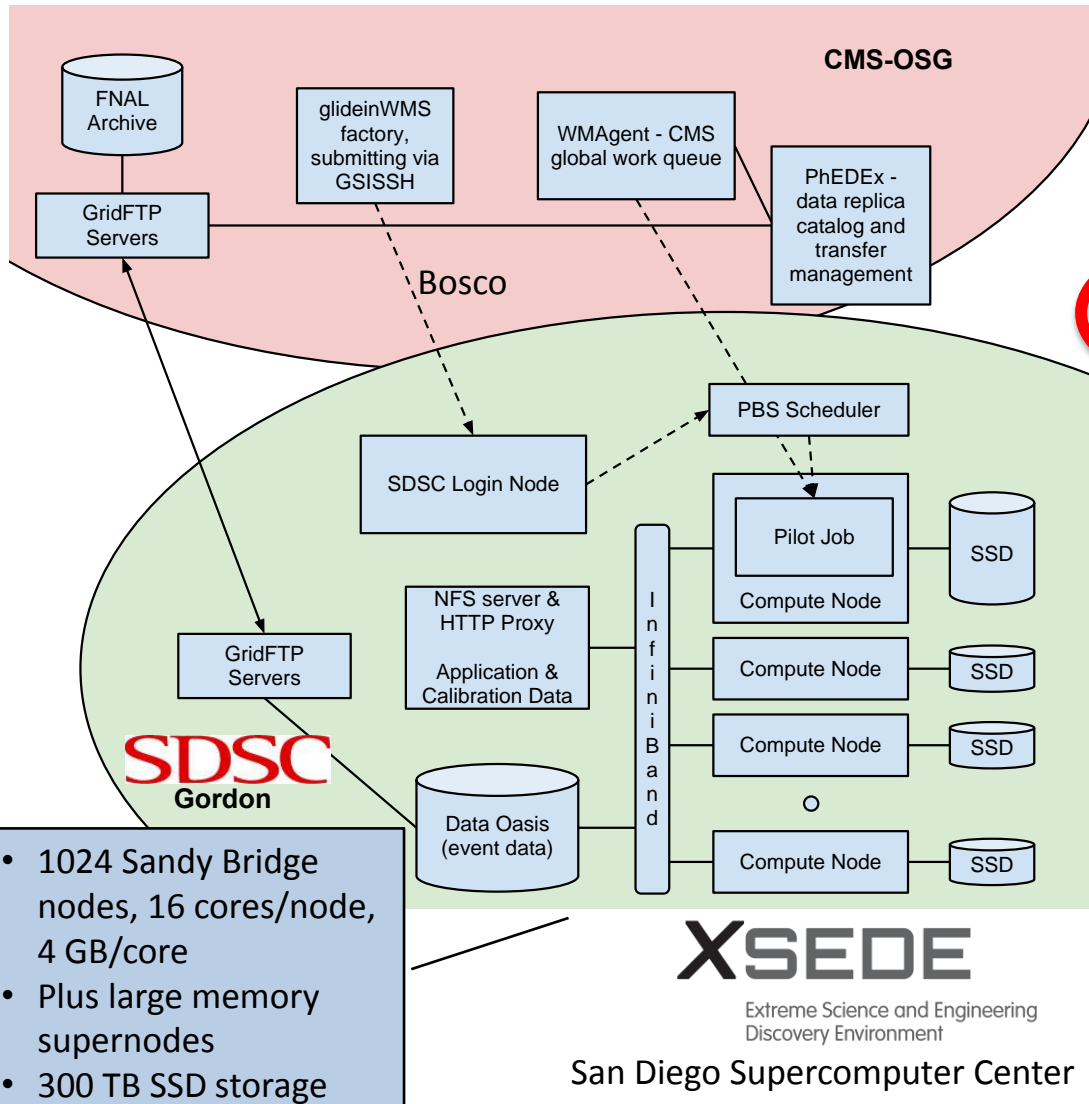
X86 (Intel+AMD) dominate –  
a familiar architecture

100+ PFlop plans:    NERSC: NERSC-8 2015-16    ORNL: OLCF-4 2016-17

The Register 6/2013



# Example of HPCs in Production – CMS on SDSC's Gordon Supercomputer



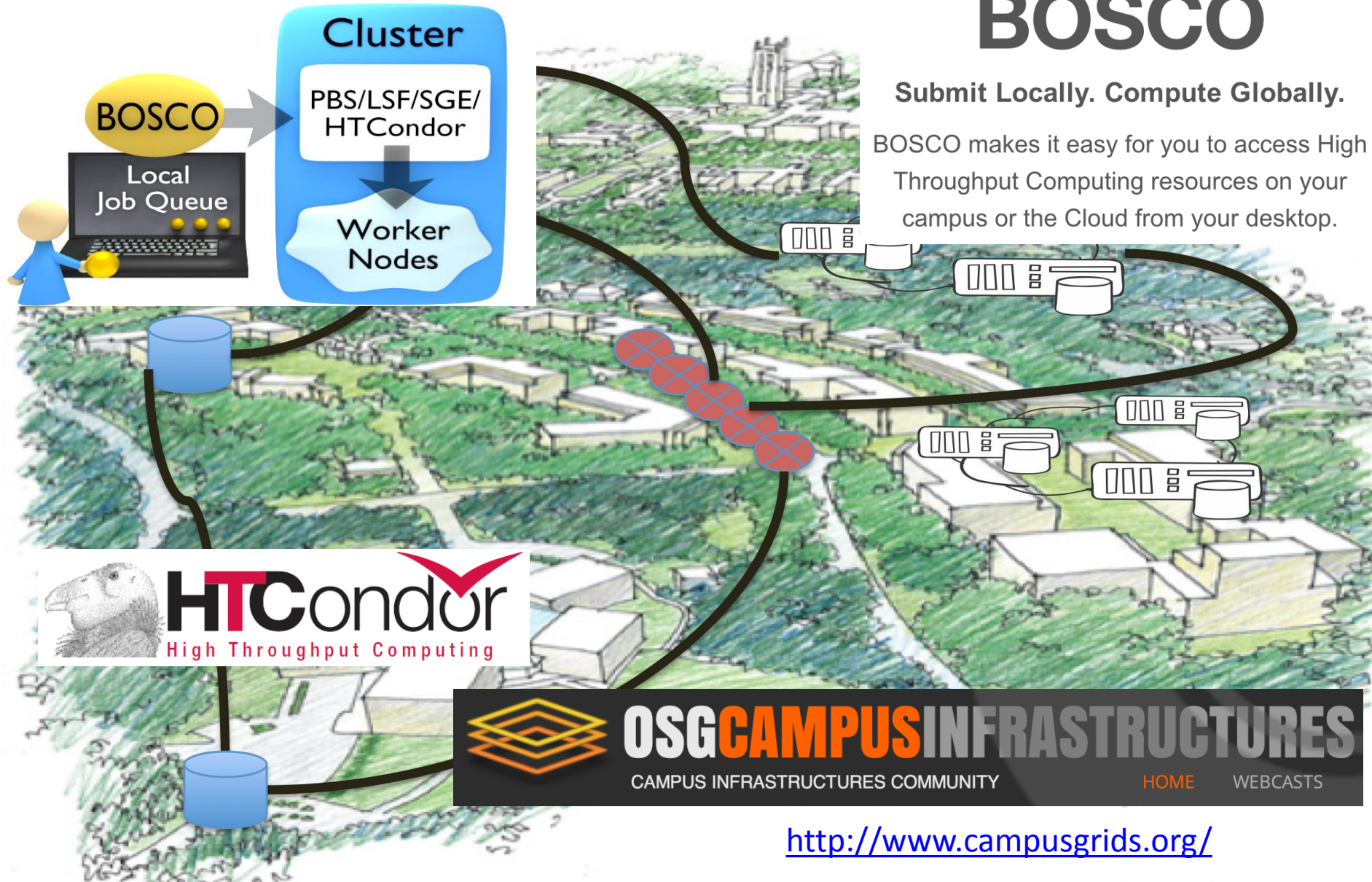
- CMS, OSG, SDSC teams
- 400M events processed in Feb-Mar 2013
- 125 TB in, ~150 TB out
- Enabled early processing of parked 2012 data to expedite physics analysis
- Good experience integrating an XSEDE (NSF) supercomputer via OSG

# BOSCO – Easy Access to Processing

# BOSCO

Submit Locally. Compute Globally.

BOSCO makes it easy for you to access High Throughput Computing resources on your campus or the Cloud from your desktop.

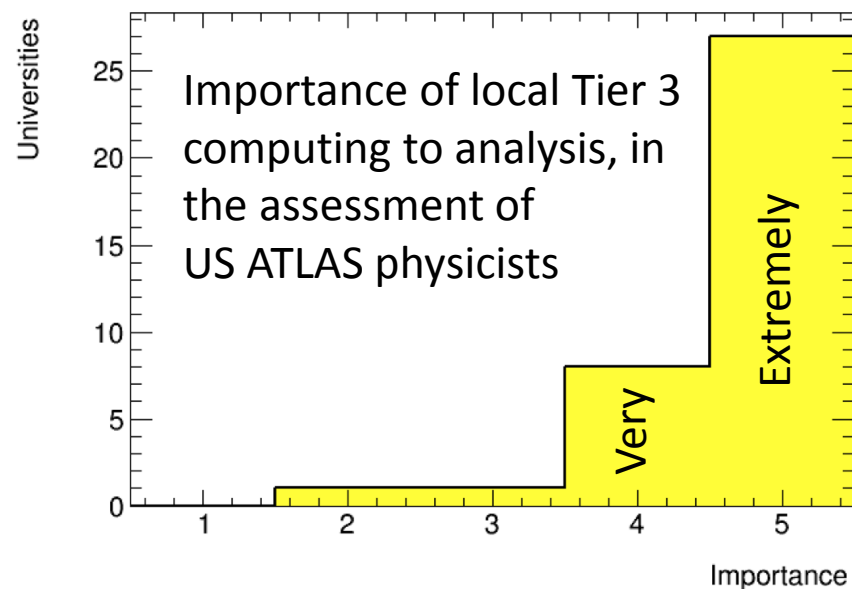


<http://www.campusgrids.org/>

# Ease of Use – Computing as a Service

- Operations and Tools technical exchange group (TEG) observed a strong desire for “Computing as a Service” especially at smaller sites
  - Easy to install & use, standard, packaged services requiring little manpower
  - Graceful reactions to fluctuations in load and user behavior
  - Straightforward and rapid problem diagnosis and remedy
  - “Cloud is the ultimate CaaS”, e.g. ATLAS is prototyping cloud-based Tier 3s
  - Many or most physicist users are at small sites, and they value their local resources as a key analysis tool

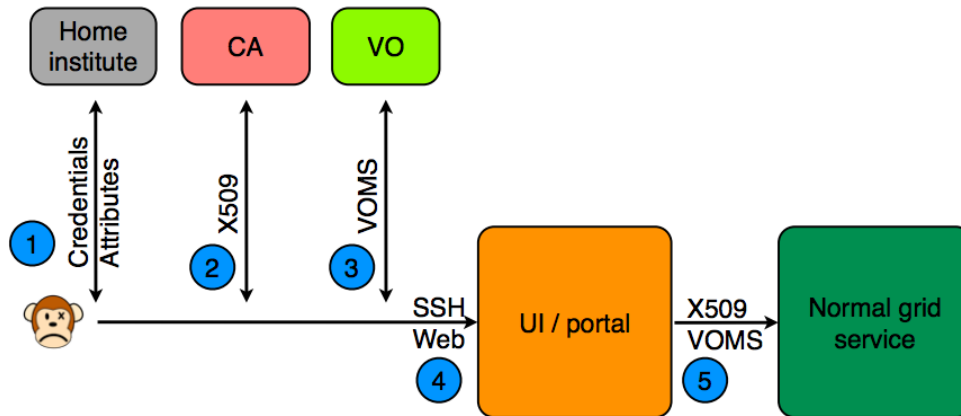
An interesting recent result from a survey of US ATLAS physicists...



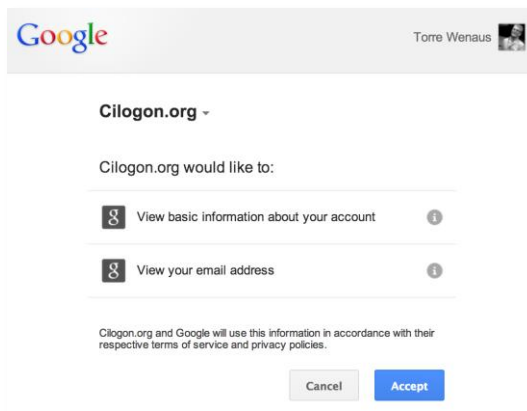
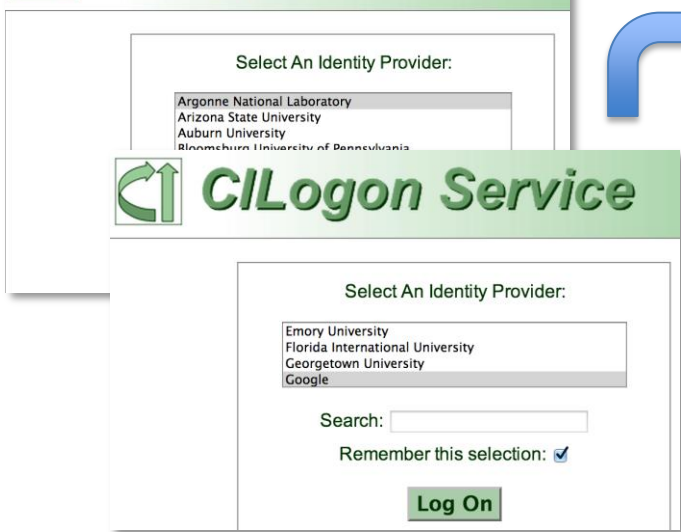
# Ease of Use – Improving on Grid Certificates

Universal authentication is at the root of the grid's success, and yet it's imperfect...

The current bad old days:



WLCG and are pursuing an easy to use (and manage) CILogon.com based service. Objective: A certificate-less grid

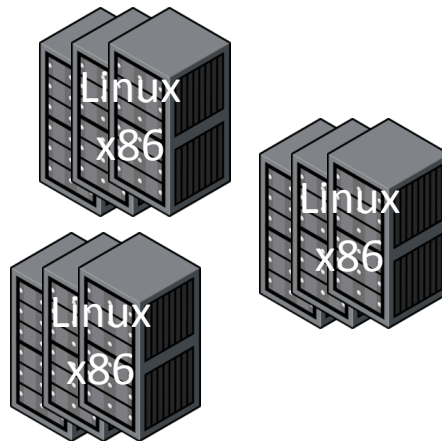


Welcome! Your new certificate subject is as follows.

/DC=org/DC=cilogon/C=US/O=Google/CN=Torre Wenaus A7321



1980s: Plethora of architectures & OSes

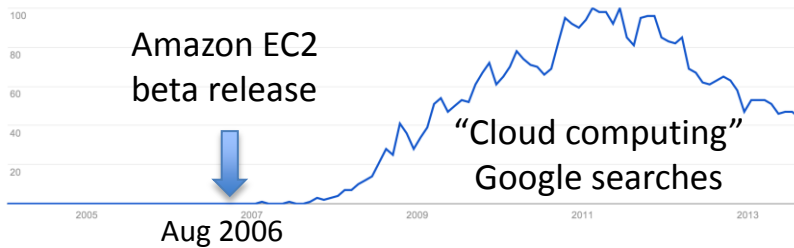


1990s: **Uniform OS/architecture**  
Linux/x86 standard for commodity cluster computing

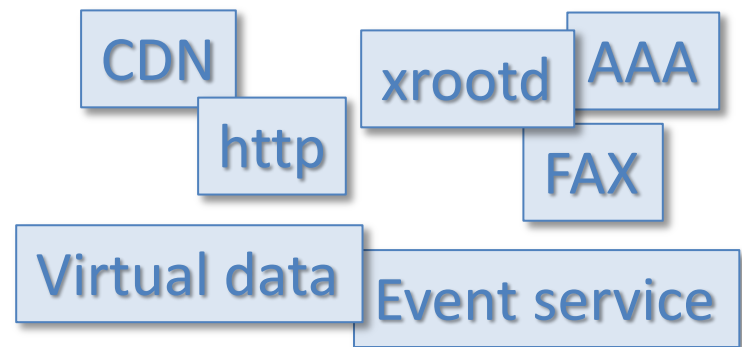


2000s: **Uniform fabric and access**  
Globally federated resources enabled by network and grid

**Computing evolution – growing uniformity counters growing scale and complexity**



2010s: **Uniform environment**  
VMs and clouds put the user in control of the environment – take it with you anywhere and everywhere



2010s: **Uniform data access**  
Working towards transparent distributed data access enabled by the network

# Conclusions

- LHC computing performed extremely well at all levels in Run 1 – we know how to deliver, adapting where necessary
- Excellent networks, flexible and adaptable computing models and software systems paid off in exploiting resources such as powerful, stable Tier 2s
- Charge for the future: live within limited budgets by **maximizing event throughput per (our) unit cost**
  - By utilizing resources we own as fully and efficiently as possible – despite the major development program required – which must be supported
  - By continuing and expanding our collaborative approaches to solving common problems
  - By using the resources others make available wherever possible – including non-traditional platforms like HPCs, clouds, volunteer computing
- Explosive growth in data and (highly granular) processors in the wider world gives us a powerful tool set and basis for success – our needs out to HL-LHC are orders of magnitude greater than today, but attainable with work
  - Industry is clearing paths for us... there's much to leverage
- The broad program of LHC computing upgrades underway is represented by many presentations at this conference

# Thank you!

- Much thanks to those who gave help and comments and whose materials I have drawn on
- Including but not limited to: L. Bauerdick, B. Bockelman, S. Campana, M. Ernst, I. Fisk, R. Gardner, A. Klimentov, W. Lampl, E. Lancon, F. Legger, H. Meinhard, R. Mount, S. Panitkin, D. Rousseau, M. Schulz, I. Vukotic, F. Wuerthwein
- Special thanks to Borut Kersevan, Predrag Buncic, Marco Cattaneo
- The talk benefited greatly from access to the comprehensive computing model update currently in draft