Big Data Technologies http://www.bigdatalab.nrcki.ru/ National Research Centre "Kurchatov Institute"





BigData and Computing Challenges in High Energy and Nuclear Physics

Alexei Klimentov



International Conference "Mathematical Modeling and Computational Physics, 2017" (MMCP2017)

National Research Centre "Kurchatov Institute"



Outline

- High Energy Physics and Nuclear Physics (HENP) scale of needs
 - BigData at HENP
 - Computing model, challenges, evolution
- Computing R&D projects highlights
 - Workflow and data management
 - Federated data storage
- Supercomputers role for HENP scientific program
- Summary and conclusion
- Future Challenges

Disclaimer 1 : This talk will have a "slight" bias towards ATLAS experiment @ LHC

The Context





Computing in HENP

- Has changed and evolved dramatically over the past decade
- Especially for the biggest experiments at the LHC (and soon for NICA, FAIR, DUNE BELLEII, LSST...)
- This talk is not about computing infrastructure in HENP but about how physicists make use of the infrastructure

The situation ~10-15 years ago

- Data processing was performed in large computing centers using local batch systems with dedicated shares
- A few satellite centers did simulations, occasionally reprocessing
- Users were mostly located near large computing centers, usually at the laboratory where the experiment was located, and used a combination of desktops and batch systems for analysis
- Final "Data Summary Tapes" versions were physically shipped for remote analysis

• Goals of this talk

- Present a new model of computing developed for LHC experiments
- Usage beyond the LHC
- Future evolution





The Science Drivers for Particle Physics

- Five intertwined science drivers, compelling lines of inquiry that show great promise for discovery :
 - 1. Use the Higgs boson as a new tool for discovery.
 - 2. Pursue the physics associated with neutrino mass.
 - 3. Identify the new physics of dark matter.
 - 4. Understand cosmic acceleration : dark energy and inflation.
 - 5. Explore the unknown : new particles, interactions, and physical principles.







Big Data Technologies http://www.bigdatalab.nrcki.ru/



Introduction. The ATLAS detector at the Large Hadron Collider



Physics 2013

ATLAS at CERN LHC is a flagship experiment in the High Energy Physics with multiple science drivers:

- Higgs Boson discovery in 2012. Nobel prize 2013 in Physics.
- Use the Higgs Boson as a new tool for discovery
- Identify the new physics of dark matter
- Explore the unknown: new particles, interactions and physical principles
- Current pace of research and discovery is limited by ability of the ATLAS distributed computing facilities to generate Monte-Carlo events - "Grid luminosity limit" and to process ALL LHC data in quasi real-time mode
 - LHC experiments use Grid computing paradigm to organize distributed resources
 - Currently ~300K cores available to ATLAS Experiment worldwide
 - Still not enough CPU power !
 - Many physics simulation requests have to wait for months
 - Supercomputers are rich source of computing power
 - ATLAS initiated R&D project aimed at integration of LeadershipClassFacilities and HPC resources (in general) into ATLAS distributed computing







Introduction. How Likely something interesting happen.



- Total Production Cross Section (== probability) vs Energy in pp collisions
- Notice the logarithmic scale on the Y-axis: it spans 11 orders of magnitude
- E.g. you produce 10 Higgs bosons out of 10¹¹ billions of collisions
- The probability increases logarithmically with energy

Alexei Klimentov

 Theory (lines) agrees very well with measurements (markers)





Introduction. How Likely something interesting happen.



New physics rate ~ 0.00001 Hz

Event Selection : **1 in 10,000,000,000**,000

Like looking for a single drop of water from the Geneve Jet d'Eau over 2+ days





Big Data Technologies http://www.bigdatalab.nrcki.ru/ Event Reconstruction

National Research Centre "Kurchatov Institute"



• 800,000,000 proton-proton interactions

- High efficiency
- Good resolution
- Low fake rate



- Robust against detector problems
 - Noise
 - Dead regions of the detector
- Be able to run within the computing resource limitations We are looking for this "signature"
 - CPU time per event
 - Memory use





LHC and Detector Upgrades









Excellent LHC Performance and Immediate Computing Challenges

Excellent LHC Performance in 2016 (Run2)

- Unprecedented peak instantaneous luminosity > 40% beyond LHC design
- Data accumulation ~60%
 beyond 25 fb⁻¹ goal for 2016
- High performance of the machine operation and data acquisition



Immediate challenge for the LHC Computing





Big Data Technologies http://www.bigdatalab.nrcki.ru/

LHC, Amazon & Google Computing Centers. Relative size of things.

National Research Centre "Kurchatov Institute"





- One Google Data Center is estimated to cost ~\$600M
 - An order of magnitude more than the centre at CERN
- Amazon : 9 large sites/zones
 - up to ~2M CPU cores/site, ~4M total
 - 10 x more cores on 1/10 of the sites compared to our Grid
 - 500,000 users
- LHC Computing (WLCG)
 - 167 sites, 42 countries
 - 500+k CPU cores total
 - Disk 350PB, Tape 400+PB
 - ~5000 users



Big Data Technologies http://www.bigdatalab.nrcki.ru/

Relative Size of Things. Cont'd

- Storage :
 - Amazon supports millions of queries per second
 - Google has 10-15 exabytes under management
 - Facebook 300PB
 - eBay collected and accessed the same amount of data as LHC Run1
- Processing :
 - Amazon has more than 40 million processor cores in EC2
 - Google has ~1M servers so ~20M cores

HENP data and processing problems are about 1% the size of the largest industry problems, but we are still distribute more data and lead in the area of data and workflow management, and high-throughput computing in general.

HENP is good in distributed computing :

Datasets are large but custodially kept and protected

- We make dynamic use of tape systems
- We move to hundreds of sites
- We make effective use of global network links

We remain leaders in this challenging areas





ATLAS Production System Performance. Daily Completed Jobs.

National Research Centre

"Kurchatov Institute"





MC Simulation

MC Reconstruction

Data Processing

Big Data: often just a buzz word, but not when it comes to HENP...



Computing model evolution. Reducing Complexity



Wide area networks are very stable now

LHC Run1. Hierarchy

7/1cl62017

Network capabilities and data access technologies have significantly improved our ability to use resources independent of location Relaxing hierarchical model : Flat instead of Tiered Grid model

Now. Mesh

Workflow Management. PanDA. Production and Distributed Analysis System





Paradigm Shift in HENP Computing

New Ideas from PanDA

- Distributed resources are seamlessly integrated
- All users have access to resources worldwide through a single submission system
- Uniform fair share, priorities and policies allow efficient management of resources
- Automation, error handling, and other features in PanDA improve user experience
- All users have access to same resources

Old HEP paradigm

- Distributed resources are independent entities
- Groups of users utilize specific resources (whether locally or remotely)
- Fair shares, priorities and policies are managed locally, for each resource
- Uneven user experience at different sites, based on local support and experience
- Privileged users have access to special resources



Russian Science Foundation Award. «Machine Learning» algorithms to predict complex system behaviour

- Production System is a large, complicated, distributed system;
 - Hard to simulate;
 - Hard to detect anomalies
 - Hard to predict its behavior
- Very thorough logging;
- Machine learning (ML) algorithms are computationally intensive, using them on raw logs (database rows) is infeasible.
- However, it is possible to use machine learning algorithms if we limit their input to some aggregated metrics of Production System
- The most important metrics are :
 - Time To Complete for tasks/jobs
 - resource utilisation
 - percentage of failed tasks/jobs
 - Running/pending jobs ratio.









BigPanDA. Growing PanDA Ecosystem





BigPanDA in Genomics

- At NRC KI PALEOMIX pipeline was adapted to run on local supercomputer resources powered by PanDA.
- It was used to map mammoth DNA on the African elephant reference genome
- Using software tools developed initially for HEP and Grid reduced genomics payload execution time for Mammoths DNA sample from weeks to days.



Data Management Evolution

Storage and Compute loosely coupled but connected through fast network

- Heterogeneous computing facilities in and outside the cloud
- Different centers with different capabilities, for different use cases

We want to keep control of data

- Need to be able to deploy data to a diverse set of resources
 - Clouds, dedicated sites, HPC centers, etc
- Will need to a combination of real time delivery and advanced data caching

In order to replicate samples of hundreds TB in hours we will need the systems optimized end-to-end and a very high capacity network in between.











Russian Fund for Basic Research Award. Federated Storage

CERN, DESY, NRC KI, JINR, PNPI, SPbSU, MSEPhI,; ATLAS, ALICE (LHC), NICA (JINR) EOS technology : NRC-KI, JINR, T2 (ATLAS, PNPI, Gatchina), T2 (ALICE, SPbSU, Petergof), CERN dCache technology : NRC-KI, JINR, DESY

P.Fuhrmann, I.Kadochnikov, A.Kiryanov, A.Klimentov, D.Krasnopevtsev, A.Kryukov, M.Lamanna, A.Peters, A.Petrosyan, E.Ryabinkin S.Smirnov, A.Zarochentsev, D.Duelmann



R&D Project Motivation

Computing models for the Run3 and HL-LHC era anticipate a growth of storage needs.

The reliable operation of large scale data facilities need a clear economy of scale.

A distributed heterogeneous system of independent storage systems is difficult to be used efficiently by user communities and couples the application level software stacks with the provisioning technology at sites.

> Federating the data centers provides a logical homogeneous and consistent reliable resource for the end users

Small institutions have no enough people to support fully-fledged software stack.

 In our project we try to analyze how to set up distributed storage in one region and how it can be used from Grid sites, from HPC, academic and commercial clouds, etc.



Options for Future Computing & Collaboration

The ultimate question

- How will data be processed and analyzed in 7-10 years and beyond?
- Buy facilities
 - ✓ Pro : Own it! No impediment to running at full capacity when needed
 - ✓ Con : Must invest for peak utilization, even if not used
- Use services from other providers :
 - ✓ Pro : Others make capital investments
 - ✓ Con : Will usage be available/affordable when needed ?

We worked hard during last years to provide examples of infrastructure not owned by HENP and to integrate HPC with HTC

- Hybrid model
 - Own baseline resources that will be used at full capacity
 - □ Use service providers for peak cycles when needed



Big Data Technologies http://www.bigdatalab.nrcki.ru/

Impact on Data and Workflow Management



One of the biggest improvements in joining to a much larger pool of resources is breaking the idea we need to lay out our resources for average load

Workflows could be completed as they are defined and not over months In these processing models the workflow system needs to be able to scale to 5-10 times the average load

- We want to be able to burst to high values
- The least expense time to be delivered resources might be all at the same

If one is using commercially provided computing faults turn into real money

- Need to focus on potentially wasteful things
 - Infinite loops
 - Giant log output that trigger data export charges
 - CPU efficiency loss

All things we probably should have been worrying about with our dedicated systems, but somehow when you are directly paying for the resources you are a bit more careful

Running ATLAS jobs on HPC In opportunistic mode











10¹ 10⁻¹ 1940

1945 1950 1955 1960

1965 1970 1975 1980 1985 1990 1995 2000 2005

Top 500



- Large HPCs use a variety of architecture
- Half of computational power is concentrated in a small number of machines;
- Small HPCs use x86 architectures. Typically, these are ordinary server racks, with Infiniband interconnects. 94% of the bottom 400 of the Top 500 (including the last 130) are all x86

Supercomputers

2010 2015



High Performance Computing Scheduling

- supercomputer is full, means that the system have allocated all the cycles it is able to deliver
 - It is probably not all cycles it has
 - Just as there is room for sand in the jar of rocks, there's room for HEP jobs on even a "full" HPC







Big Data Technologies http://www.bigdatalab.nrcki.ru/

National Research Centre "Kurchatov Institute"





OLCF Titan Integration with ATLAS Computing



The largest supercomputer available for scientific applications



27 PFlops (Peak theoretical performance). Cray XK-7 18,688 compute nodes with GPUs 299,008 CPU cores AMD Opteron 6200 @2.2 GHz (16 cores per nodec 32 GB RAM per node NVidia TESLA K20x GPU per node 32 PB disk storage (center-wide Luster file system) >1TB/s aggregate FS throughput 29 PB HPSS tape archive

Distribution of queue wait time for ATLAS jobs





Supercomputers. Titan contribution to ATLAS MC simulations



ATLAS simulation jobs completed worldwide: Jan – Oct 2016

Titan backfill contributed
 7.8% of total simulation
 jobs CY16 to date.

Rest - 65.72% (15,084,631)
 MWT2_SL6 - 2.45% (561,738)
 TOKYO_MCORE - 1.68% (384,783)
 BOINC_MCORE - 1.54% (353,962)



Completed jobs (Sum: 22,954,360)

BOINC - 12.78% (2,932,678)
 BNL_PROD - 1.77% (406,023)
 MWT2_MCORE - 1.61% (370,151)
 TOKYO - 1.32% (302,391)

ORNL_Titan_MCORE - 7.82% (1,796,02
 CERN-P1_MCORE - 1.76% (403,267)
 FZK-LCG2 - 1.56% (358,715)





Workflow and Data Management Evolution

- Big centers for data reduction impacts workflow and data management
- Data selection workflow sits on top of "big data" tools
 - Focusing effort on reproducibility and shared selection criteria
- Data Management involves moving small samples to end sites
- Activity is triggered automatically
 - Needs throttling mechanisms
- The bulk of the data is placed at big sites
 - Reduced samples are moved and replicated
- Still a push to enable the processing on a variety of resources
 - Ability to burst to high capacity becomes even more important when access can trigger processing









Summary and Conclusions.

- 2000s The decade of the Grid
- At the first LHC Run(s) Distributed Computing has helped deliver physics rapidly
- Entering a phase of computing evolution
- Challenges for computing scale & complexity will continue to increase dramatically
- The distributed computing model allows us to incorporate clouds and supercomputing centers and to use them efficiently now and for LHC Run3+
 - Supercomputers offer important and necessary opportunities to the experiments
 - Great progress has been made to interface supercomputers to LHC Distributed Computing
 - It is demonstrated that we can run at scale
 - We show that we can use them opportunistically and in backfill mode
- Integrate more HPC into production environment
 - Many Supercomputing centers are very interested to collaborate with us.
 - Three technical thrusts
 - Integrate HPC into production environment
 - Port HENP code to each HPC system
 - Learn how to exploit accelerators where present
- Access to the supercomputers coupled with collaborative help in the transformation of HENP code would be a major scientific contribution to the physics discoveries of the next ten years



National Research Centre "Kurchatov Institute"



In meantime...

While we were developing the Grid, the rest of the world had other ideas.





The external world of computing is changing now as fast as it ever has and should open paths to knowledge in physics. HEP computing needs to be ready for new technical challenges posed both by our research demands and by external developments.







Summary

BUT

- the landscape has changed dramatically over the past decade
 - The Web, the Internet, powerful PCs, broadband to the home, ...
 - have stimulated the development of new applications that generate a massive demand for computing remote from the user
 - …. that is being met by giant, efficient facilities deployed around the world
 - and creates a market for new technologies capable of operating on a scale equivalent to that of HEP
- Whether or not commercial clouds become cost-effective for HENP data handling is only a financial and funding-agency issue

BUT

• Exploiting the associated technologies is an obligation



Challenges



- Technical challenges:
 - Optimization of the physics output vs cost
 - Software, algorithms, computing models, distributed infrastructure \rightarrow and implications (e.g. on networks needed)
 - Integration of *all* available resources: HPC, Cloud, opportunistic, traditional, etc.
 - Technology evolution will it be as much as we need?
 - Opportunity to re-think the computing models may be very different than today
- Sociological challenges:
 - Remove the "online-offline" boundary there is a computing challenge from detector to physics
 - Must ensure that Computing and Software careers are seen as Physics careers essential to build and maintain the skills we need
 - This requires change in the collaborations & in the Universities
 - Consolidation of resources (e.g. storage) must not be interpreted as removing the need for a global community and global contributions
 - Must find a path to reducing cost while maintaining the most broad and open contributing community
- Funding related challenges

HL-LHC will require more revolutionary thinking

Could there be a revolution here for physics computing





Thanks

- This talk drew on presentations, discussions, comments, input from many
- Thanks to all, including those I've missed

– I.Bird, P.Buncic, S.Campana, T.Childers, K.De,
I.Fisk, M.Grigorieva, M.Gubin, B.Kersevan,
A.Kirianov, T.LeCompte, T.Maeno, R.Mashinistov,
D.Oleynik, A.Patwa, A.Poyda, T.Wenaus,
A.Zarochentsev ...

