# CHEP2010

## Rain and lecturing in Taipei
## Sam Skipsey

With thanks to Wahid Bhimji and Stuart Purdie for comments on parallel sessions I was unable to attend.

# Itinerary

- What is CHEP?

- Themes overview

  - A Summary of Summaries

  - Common themes

  - Presentation highlights

- Conclusion

# CHEP - 2010

- The big Computational High-Energy Physics Conference

- Hosted by Academi(c)a Sinica, Taipei, Taiwan

- Opened by the Vice-President of Taiwan, Vincent Siew (蕭萬長)

  - Good luck getting Nick Clegg in the UK!

# A brief slide about Taipei.

- Capital of Taiwan (Republic of China)

- Pronounced "Taibei" (臺北)

  

  - Has an international baseball team.
  - Has tons of Taoist, Buddhist and folk Temples.
  - Filled with "Night Markets"
  - Oh, and typhoon season is around October...

# A brief slide about Taipei.



- Capital of T
- Pronounced

We can't stop here; this is Typhoon Country!

# An overview of Summaries

- 7 Tracks of Parallel Sessions in total

  - Event Processing; Online Computing; Software Engineering & Data; Distributed Processing & Analysis; Computing Fabrics & Networking Tech.; Grid & Cloud Middleware; Collaborative Tools*

- Each was summarised in 30 minutes on Friday.

- We'll start with a summary of the summaries:

*we implicitly continue the agenda of Newton in pretending rainbows have 7 colours...

# Event Processing

- GEANT4 (simulation)

- FAIR plans (GEANT4, event analysis trains)

- Analysis and I/O (GPU, etc)

- Software standardisation

- Alignment and calibration (brilliant pictures)

# Event Processing tag cloud



- Fabio Cossutti (INFN) Ryosuke Itoh (KEK) Oliver Gutsche (Fermilab)

# Event Processing

## Conclusions

▶ LHC:

  ▶ In general, everything is working remarkably well

  ▶ Current data taking conditions are under control

  ▶ The future will show how the experiments will cope with the increasing PileUp conditions

▶ Beyond LHC

  ▶ FAIR experiments with non-traditional beam conditions mark a new frontier of challenges

  ▶ New experiments rewrite or develop new frameworks, software standardization helps them ramping up more quickly

▶ Performance

  ▶ **Multi-Core, GPUs and Vectorization**: buzz words of the processing world

  ▶ Applications show significant speed increase, but usage still very dependent on specific situations

  ▶ New experiments are designing their software for multi-core, multi-thread execution environments and also consider specialized hardware solutions like GPUs

Thursday, 4 November 2010

# Online Computing

- All LHC experiments ++good with real data taking. (DAQ eff > 90%)

- Emphasis on storage perf (at T0)

- Integration and automation!

- Data quality monitoring - and web access!

- Lessons learned: uniform stack (on and offline)

# Online Computing

## Summary

- The startup of the LHC experiments has been a tremendous success: DAQ efficiencies well over 90% and over all efficiencies for physics well over 80%
- Sophisticated tools for data quality monitoring allow remote and local experts to react and flag quickly
- Modern web-technologies make experiment info available everywhere
- But…
  - The LHC is improving and in some areas we are already beyond initial design parameters
  - Upgrades are coming: more data, new detectors, faster readout & storage

Thursday, 4 November 2010

# Software Engineering, Data*

- Heterogeneity of talks (also mentioned in other summaries)

- multithreading cool

- performance monitoring tools

- lots of "new" cool things: go, svn, html5, CVMFS etc

- software recycling for small exps.

- Data archiving and preservation

# Software Engineering, Data *



- Marco Cattaneo & Stefan Roiser

# Software Engineering, Data *

# Conclusions

- The software frameworks for LHC are in very good shape
  - Processes and tools are in place
  - A lot of efforts on the performance side are underway
    - We need more specialists in this area
- Other experiments should be able to profit from the work
  - We need for more collaboration across experiments
  - Implementation islands are getting bigger

# Distributed Processing and Analysis

- Successes

  - First year LHC stuff, organised processing, phenix, ATLAS DDM

- Future Architectures

  - FAIR, SuperB, Belle2, CDF, Fermi Space telescope

- Improvements

  - Global FS, ARC, VMs,

# Distributed Processing and Analysis



Thursday, 4 November 2010

# Computing Fabrics and Networking Tech

- Storage, Storage, vms, management, multicore

- No: lustre (as a sole topic), IPv6

- Fabric management (puppet)

- Data management architecture reworking

- NFs4.1, EOS, CPU scaling, Clouds (expensive)

# Computing Fabrics and Networking Tech.



- Tony Cass

# Grid and Cloud Middleware

- Operations and Monitoring

- Data Management - EMI plans

- Clouds -H1Data (ceph!), Boinc+CVM+CoPilot

- Virtualisation, messaging, integration

- Pilot jobs (improvement of)

# Grid and Cloud Middleware



- Markus Schulz

# Grid and Cloud Middleware

# Collaborative Tools

- Outreach Plenary

- Web2.0 internal/external communications (ATLAS), CMS collab. infrastructure

- Inspire, ATLAS Live, Glance information system

- EVO

- CERN Lecture archiving system,

- AbiCollab (like google Docs)

# Collaborative Tools

## Overview

- Increasing areas in the CT field
  - HD videoconferencing systems, Outreach and Inreach activities, Rich Media Content, Information systems, etc.

- Representative examples covering the activities in the HEP community
  - 1st session dedicated to Policies and New initiatives
  - 2nd session dedicated to SW systems and Collaborative Tools

- Plenary Talk (Lucas Taylor, FNAL/CMS)
  - Overview about the importance of the outreach activities for the HEP community
    - Contract with the Society
    - Need of a defined strategy: HQ messages, defined relation with the Media and use of latest multimedia technologies
    - Everyone needs to be involved

- Joao Fernandes – CERN, Philippe Galvez – Caltech, Milos Lokajicek – FZU Prague

# So, the summary of summaries:

- The LHC works!

- Data Management is hard

- Multi-core, GPGPU are exciting

- Virtualisation and Clouds are hot topics.

- Talking to people is good.
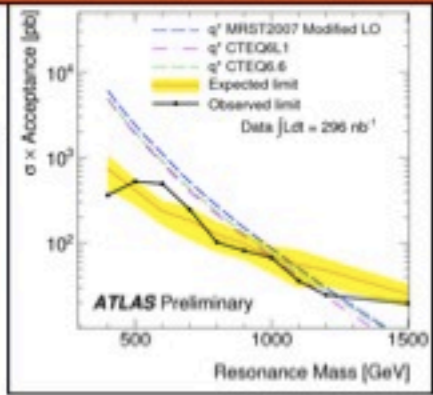
# Presentation Spotlights

- A selection of the presentations that were most interesting.
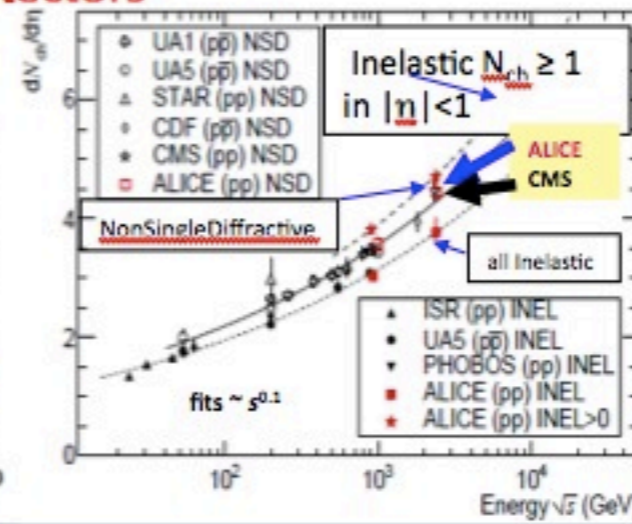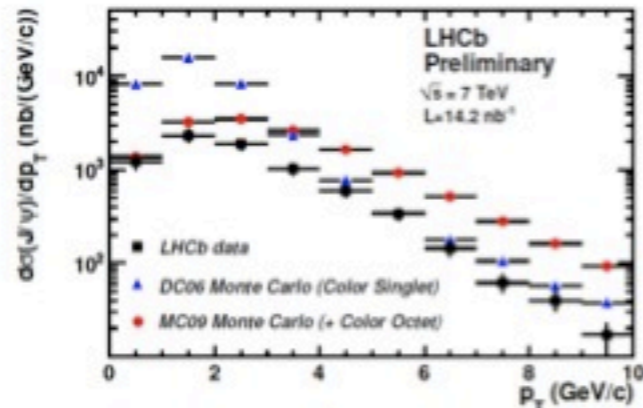
- Subjective! Caveat Auditor!
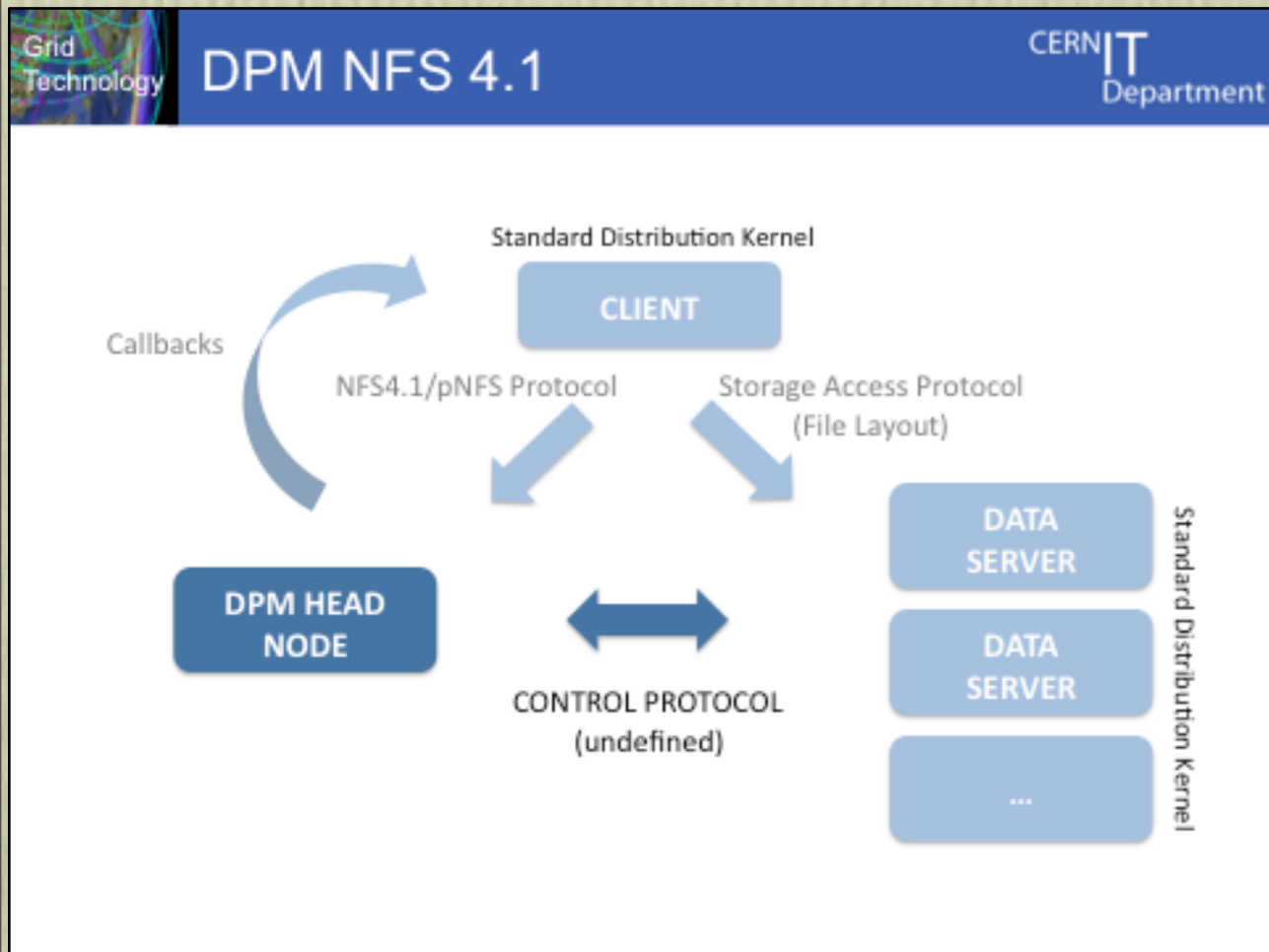
# The LHC works!

# The LHC works!

# Data Management

- WLCG Data Management meeting 16 June.

  - Many presentations sprung from this.

- New storage technologies evaluated.

  - SSDs, Ceph, CERNVM-FS

- Archiving and long-term storage.

# Amsterdam Themes

- Storage themes:

  - Dynamic data and Caches

  - Consistency - Messaging

  - Global filesystems (xrootd, mostly)
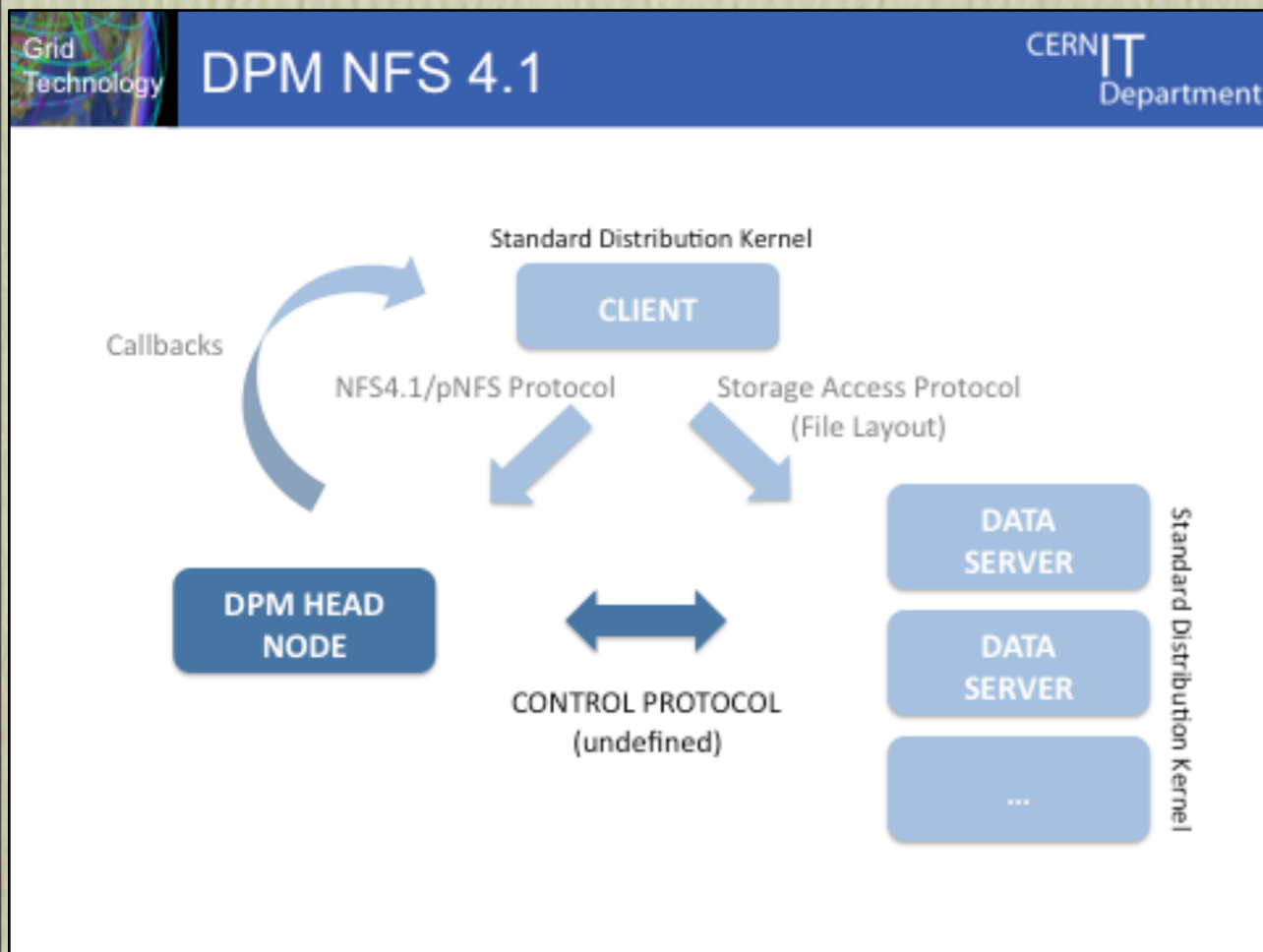
  - NFS4.1/pNFS

  - Archiving!

# Amsterdam Themes



- Caches

- ssaging
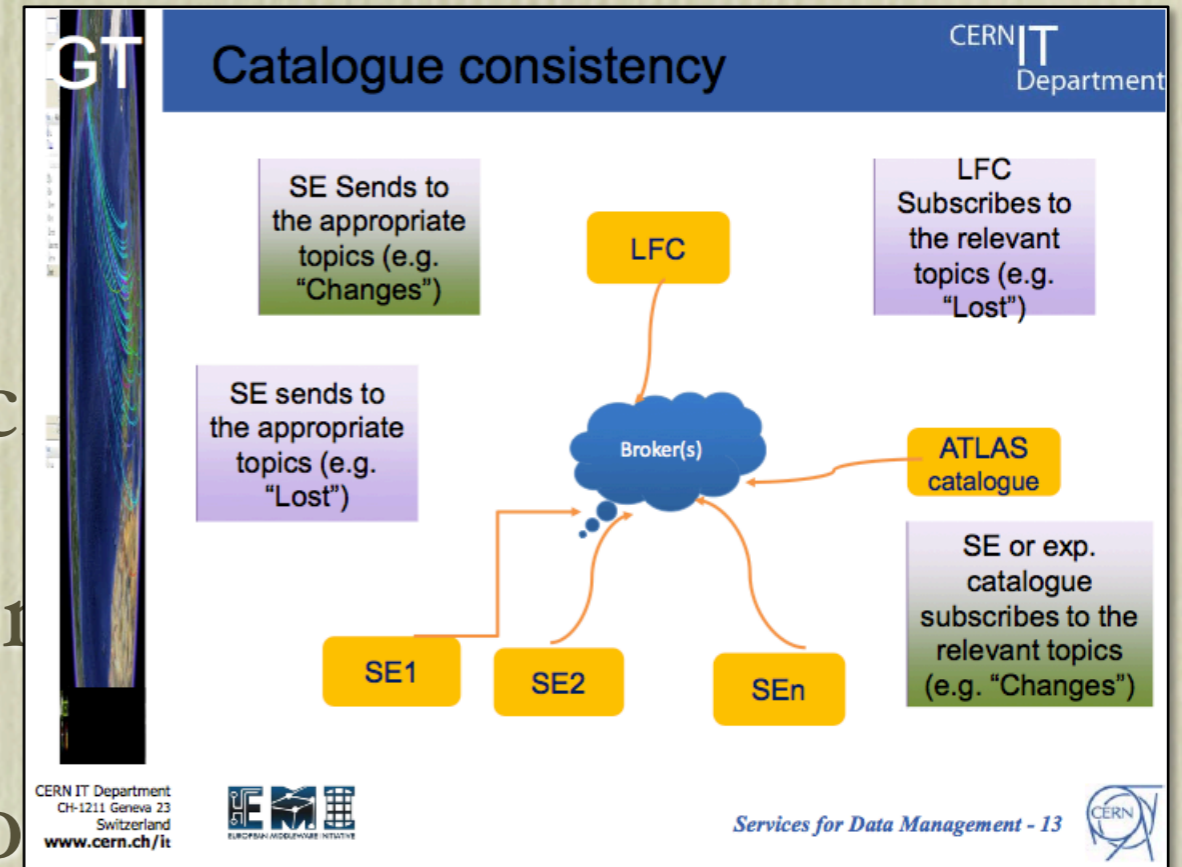
- (xrootd, mostly)

- NFS4.1/pNFS

- Archiving!

# Amsterdam Themes



- Cac...
- ...ssagin...
- (xro...
- NFS4.1/pNFS
- Archiving!

# Amsterdam Themes

# Amsterdam Themes



Grid Technology
## DPM NFS 4.1
CERN**IT** Department

Standard Distribution Kernel

GT
## Catalogue consistency
CERN**IT** Department

LFC

SE Sends to the appropriate topics (e.g. "Changes")

LFC Subscribes to the relevant topics (e.g. "Lost")

## Half-Synthetic ROOT tests: Results

60MB TreeCache, all branches

0byte TreeCache, two

- ◆ NFS (original)
- ■ DCAP (original)
- ▲ NFS (flushed)
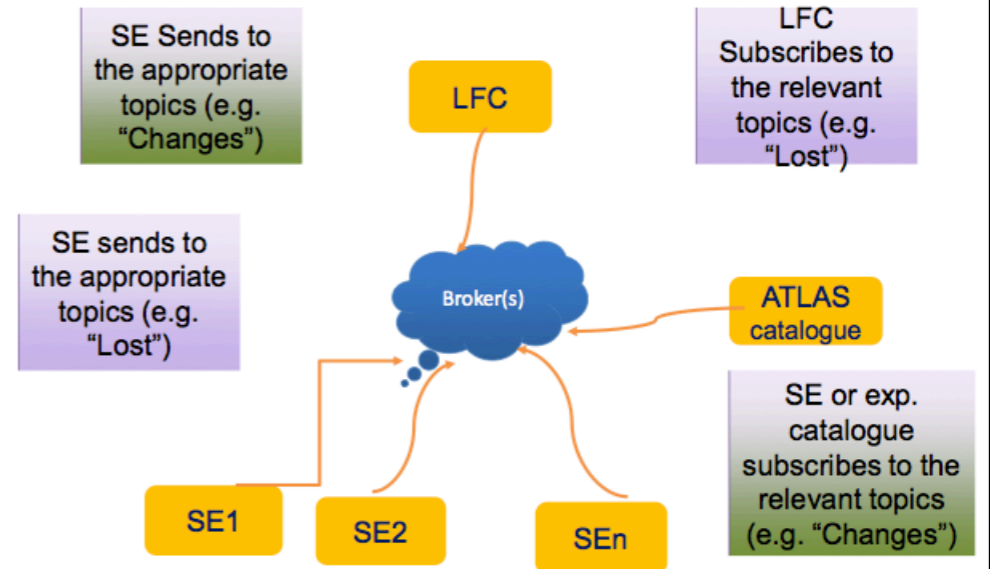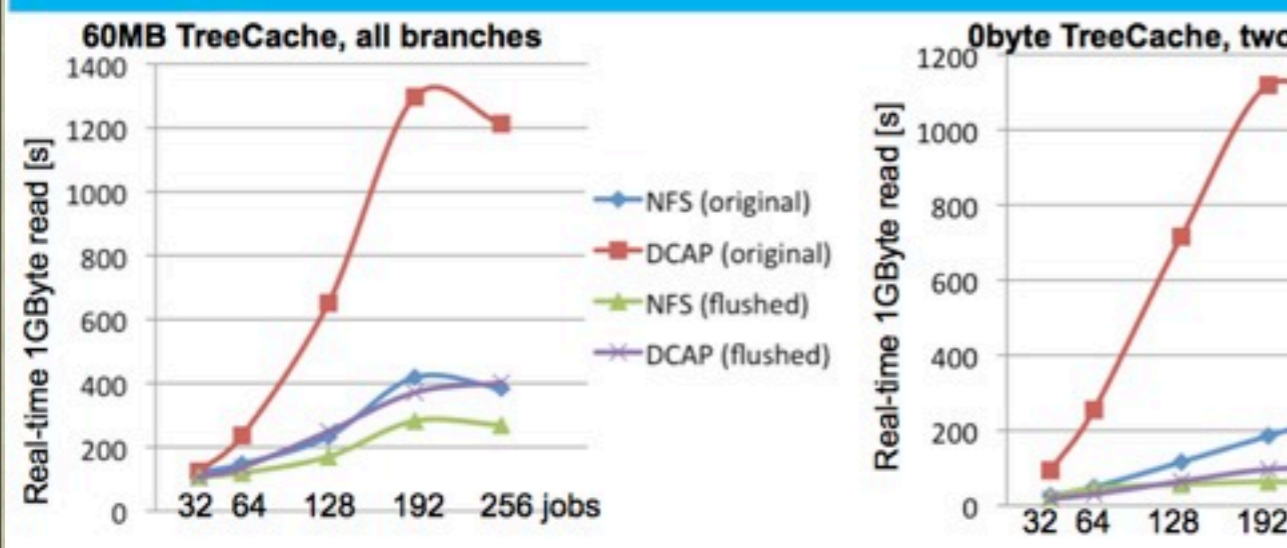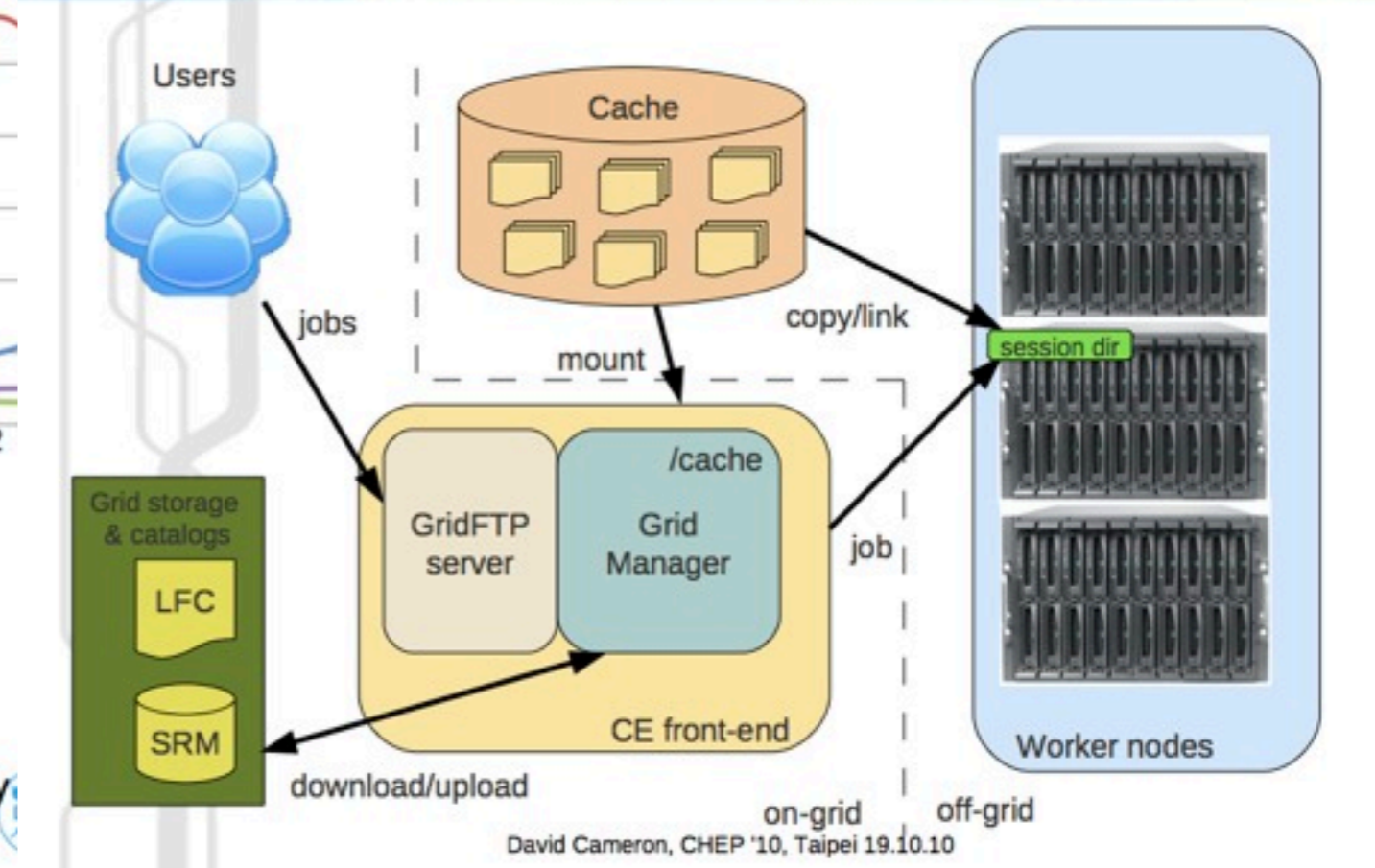- ✕ DCAP (flushed)

Real-time 1GByte read [s]

32 64 128 192 256 jobs

32 64 128 192

> NFS better for original and flushed files than dCap
  - Flushed: not much difference, original: Large difference
> TreeCache helps, NFS adds additional speed
> Peak at 192 clients not understood
> Remember: Just going through events and doing nothing … not really representative for analysis

Yves Kemp | LHC analysis uding NFSv4.1 (pNFS) | 10/20/2010 | Page 18

## NDGF
Nordic DataGrid Facility
## ARC Data Management Architecture

Users

jobs

Cache

mount

copy/link

session dir

Grid storage & catalogs

LFC

SRM

GridFTP server

/cache

Grid Manager

job

CE front-end

download/upload

on-grid    off-grid

Worker nodes

David Cameron, CHEP '10, Taipei 19.10.10

# "New" themes



Cloud Computing test configuration

- *Test-bed setup:*

Public network

Private switch

CC,CS,CLC

NC Nodes

NFS mount

CEPH Storage area

MON, MDS

MON, MDS, ODS

ODS

Bogdan Lobodzinski,     CHEP 2010, Taipei, Taiwan     9

# "New" themes



Cloud Comp[uting]

- Test-bed setup:

Public netwo[rk]

Private switch

CC,CS,CLC

CEPH Storage area

## Results (FileStager)

| Jobs:Cores | Storage | Efficiency | Throughput |
|---|---|---|---|
| | **Standard Node** | | |
| 8:8 | 1×Kingston Value SSD | 60% | 4.5 |
| 8:8 | 1×SATA HDD | 75% | 5.5 |
| 8:8 | 1×Intel X25 SSD | 80% | 6 |
| 8:8 | 2×SATA HDD (RAID 1) | 83% | 6.6 |
| 8:8 | 2×SATA HDD (RAID 0) | 90% | 7 |
| | **Magny-Cours Node** | | |
| 24:24 | 1×Intel X25 SSD | 50% | 12 |
| 24:24 | 2×SATA HDD (RAID 0) | 86% | 21 |
| | **Single Occupancy Efficiency (Measured)** | | |
| 1 | SATA HDD | 90% | 0.9 |

# "New" themes

# "New" themes

# "New" themes

# "New" themes

# Multi-Core, GPGPU and all that Jazz

## CHEP 2010

### How to harness the performance potential of current Multi-Core CPUs and GPUs

Sverre Jarp

CERN openlab

IT Dept.

CERN

Taipei, Monday 18 October 2010

CHEP 2010, Taipei

CERN openlab

## Today:
## Seven dimensions of multiplicative performance

- **First three dimensions:**
  - Pipelined execution units
  - Large superscalar design
  - Wide vector width (SIMD)

- **Next dimension is a "pseudo" dimension:**
  - Hardware multithreading

- **Last three dimensions:**
  - Multiple cores
  - Multiple sockets
  - Multiple compute nodes

Pipelining

Superscalar

Vector width

Multithreading

Nodes

Sockets

Multicore

SIMD = Single Instruction Multiple Data

Sverre Jarp - CERN

# CPU scaling



**Parallelizing Atlas Reconstruction and Simulation:**
Issues and Optimization Solutions for
Scaling on Multi- and Many-CPU Platforms

**Charles Leggett[1]**

Sebastien Binet[2], Paolo Calafiura[1], Keith Jackson[1], David
Levinthal[3], Mous Tatarkhanov[1], Yushu Yao[1]

[1]Lawrence Berkeley National Lab, [2]LAL, [3]Intel

# CPU scaling



## Memory Usage

- Single reco process uses ~1.5 Gb
- We can save significant memory through OS level sharing

**AthenaMP ~0.5 Gb physical memory saved per process**

- 8 core HT machine

6/23    CHEP 2010    10/19/10

Thursday, 4 November 2010

# CPU scaling

## Hyper-Threading

- Nature of instruction pipeline allows instructions to be interleaved
  - OS sees "virtual cores" as just another processor
  - Only effective until one of the threads saturates a shared resource and stalls



- Turn on HT when you have more processes than cores

# CPU scaling

## Memory Usage

## Hyper-Threading

## Conclusions: Longer Term Solutions

- We (HEP) are not the only ones facing these problems
  - Oracle, IBM, Google all have similar issues
    - They're only just now beginning to realize it

- We need new tools and analysis techniques
  - Current tools fail to show where the problem are, or are not suited to large scale deployment

- We need to drive these optimization techniques into the compilers and linkers themselves
- Changes at the hardware level would also improve the situation
  - It's already happening: new counters are being included in Intel's Westmere and Sandybridge chips which make profiling more useful
  - If we can show Intel exactly what's wrong, and what it will take to fix it in hardware, **they will listen**.

Thursday, 4 November 2010

# CPU scaling

**ng: What Next-Generation Languages Can Teach Us About HENP Frameworks in the Manycore Era**

Author: Sebastien Binet
Institute: LAL/IN2P3
Date: 2010-10-19
Conf: CHEP-2010

Memory Usage

Hyper-Threading

Conclusions: Longer Term Solu

- We (HEP) are not the only ones facing these p
  - Oracle, IBM, Google all have similar issues
    - They're only just now beginning to realize it

- We need new tools and analysis techniques
  - Current tools fail to show where the problem are, or
    to large scale deployment

- We need to drive these optimization technique
  compilers and linkers themselves
- Changes at the hardware level would also improve the
  situation
  - It's already happening: new counters are being included in Intel's
    Westmere and Sandybridge chips which make profiling more useful
  - If we can show Intel exactly what's wrong, and what it will take to fix
    it in hardware, **they will listen**.

October 19, 2010    1 / 23

Thursday, 4 November 2010

# CPU scaling

## Next-gen (and not so next-gen) languages

- C1X + GCD/libdispatch (closures + work queues)
- C++0x (lambda functions, std::thread)
- Python/Cython + PyCSP + multiprocessing + mpi4py + ...
- Vala/Genie
  - http://live.gnome.org/Vala
  - http://live.gnome.org/Genie
- Haskell, Erlang
  - is HEP ready for functional programming ?
- go
  - http://golang.org

0                                                    October 19, 2010    3 / 23

---

Memory Usage

Hyper-Threading

**Conclusions: Longer Term Solu**

- We (HEP) are not the only ones facing these p
  - Oracle, IBM, Google all have similar issues
    - They're only just now beginning to realize it

- We need new tools and analysis techniques
  - Current tools fail to show where the problem are, or
    to large scale deployment

- We need to drive these optimization technique
  compilers and linkers themselves
- Changes at the hardware level would also imp
  situation
  - It's already happening: new counters are being inclu
    Westmere and Sandybridge chips which make profil
  - If we can show Intel exactly what's wrong, and what it will take to fix
    it in hardware, **they will listen**.

# CPU scaling

ng: What Next-Generation Languages Can Teach Us About HENP Frameworks in the Manycore Era

Next-gen (and not so next-gen) languages

Results

Memory Usage

Hyper-Threading

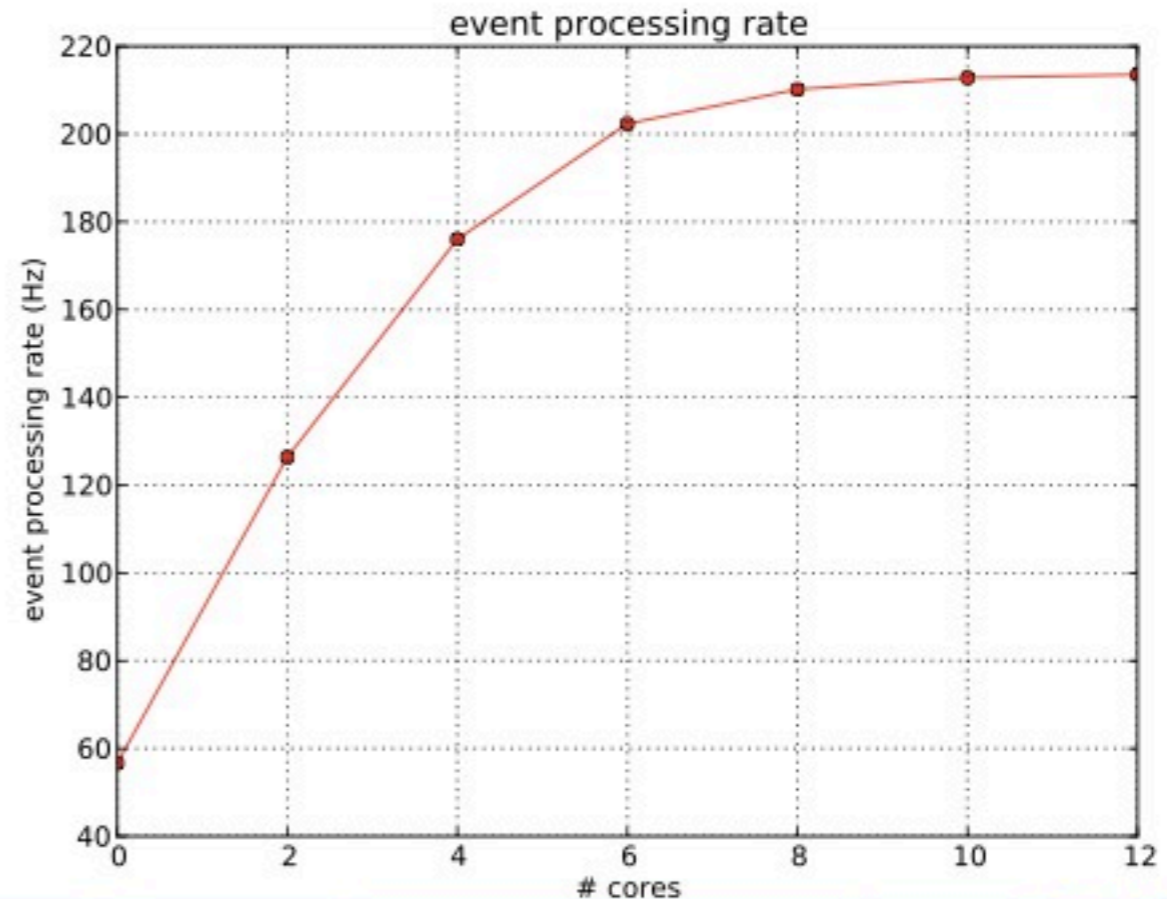Conclusions: Longer Term Solutions

- We (HEP) are not the only ones facing these p
  - Oracle, IBM, Google all have similar issues
    - They're only just now beginning to realize it

- We need new tools and analysis techniques
  - Current tools fail to show where the problem are, or a to large scale deployment

- We need to drive these optimization technique compilers and linkers themselves
- Changes at the hardware level would also impr situation
  - It's already happening: new counters are being inclu Westmere and Sandybridge chips which make profili
  - If we can show Intel exactly what's wrong, and what it in hardware, **they will listen**.

*event processing rate* (chart: event processing rate (Hz) vs # cores, ranging from ~55 Hz at 0 cores rising to ~215 Hz plateau near 10–12 cores)

Thursday, 4 November 2010

# GPGPU use

# GPGPU use

# GPGPU use



## Kalman Filter for CUDA
### D. Emeliyanov

- Standalone version successfully ported to C.

- Structs of arrays used to store track data.

- Vector data types (e.g. *float4*) for compact representation of data.

- One GPU thread per track.

- Modification of smoothing algorithm required for single precision arithmetic.

Muon tracks, $p_T$=10GeV, full MC simulation

- Over *5x* speed-up seen at 3000 tracks.

16/22

# Virtualisation and Clouds

- Virtual machines promising

    - but there are security and config implications

- Use of Clouds - works, but you still have to pay more (and aren't they just VMs in a datacenter?).

# Testing scale & usage growth

## Alternate Grids/Cloud usage scaling vs Time

Today, 1,000 jobs stable on Cloud or Hybrid ("virtualized Grids") possible (some challenges with stability / scalability at times)

**10 k to 100 k jobs needed for STAR – Should be "routine" by beginning of 2011**

Promising ... OSG ~ 13 M jobs at times … some way to go …

Jérôme LAURET – CHEP 2010, Academia Sinica Taipei/Taiwan – Oct 2010    10

# Models – 10-100 job scale series

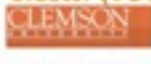| Model | Scale, Effic. & Scope | Observations | Key features & observations |
|---|---|---|---|
| Amazon/EC2 (amazon web services) Xen | Eff. > 99% Medium instance Raw simulation workflow (event generator) | • IO rather inadequate for large scale efforts (5 MB/sec per node) • Good for simulations and simple workflows: little "I", not that much "O" in IO • Normal instances unlikely suitable for HPC/HTC or large data mining (did not test the HPC instance. Has anyone?) | • Amazon has a concept of VM repository (needed) • Amazon AAA rudimentary (lacking?) • Amazon has a simple and competitive pricing model: $0.09 / hour - 300 jobs, week long cost ~ $5,600. A year long CPU @ 100 jobs saturation ~ 79k$ |
| Nimbus/EC2 NIMBUS Xen | Eff. ~ 85% (1st) and ~ 97% (2nd) Raw simulation workflow (event generator) | • Initial drop of efficiency due to batch system "inside" • More "natural" to handle submission (Cloud / Virtualization behind, Grid-like in front) • First REAL usage for Physics | • Contextualization needed at startup (GK not known a-priori, batch need to know topology) • OSG stack inside, GK+WN - virtual space looks like "another OSG site" – good attempt to unify |
| Virtual Org. Cluster (VOC) CLEMSON [ACAT 2010] KVM | Eff 100% (?) Did not lose a single job Raw simulation workflow (event generator) | • Looked like a "grid" site – submit job to it (pull or push), VM appeared on demand due to a "subscribe" mechanism for bacth • Transfer limited to University / National lab line • Required IPs (some scale problem; thought of overlaid network …) | • Contextualization remains a site specific (time) overhead • Interface is standard Grid – user is agnostic of technology • Performance improved by caching image locally OR directing changes to local disk – impossible on EC2. Final overhead < 1% (no NFS read) |

BROOKHAVEN NATIONAL LABORATORY

Testing scale & usage growth

# Models – larger scales …

| Model | Scale, Effic. & Scope | Observations | Key features & observations |
|---|---|---|---|
| Condor/VM<br><br>Condor | Scale 500+ jobs<br><br>Eff. unclear [80,85%]<br><br>*Raw simulation workflow (event generator)* | ▪ 10% of the VM never started, 15% stopped (crashed), 5% net loss for long simulation jobs (VM reboot every 24 hours). Need to be able to extend lease?<br>▪ Data transfer mechanism was through common SE and separate from workflow<br>▪ *Results used for an analysis PoP* | ▪ Interface remains grid-like but only starts VM; No real job get "inside" – "pull model" (via cron)<br>▪ As many VMs as one wants: nearly no contextualization (apart from SE) reduce overheads on local staff, condor steering<br>▪ IP space is local – no connection to outside – transfer of data out hard but SE & Cloud may be the path … |
| Clemson/Kestrel<br><br>CLEMSON | Scale 1000 jobs<br><br>Eff. Unclear, cruising average 90%<br><br>*Full simulation with track reconstruction, detector efficiency correction using the full STAR framework* | ▪ Job are "lost" if problem (communication, dead program, dying VM)<br>▪ FIRST time we mixed VMs from Clemson + CERN (true "Cloud" idea)<br>▪ Full database access as a service within the VM<br>▪ **Second massive usage of cloud targeting a conference, physics publication** | ▪ VM packaged everything<br>▪ Globus and MyProxy inside for transfer out – no real problems<br>▪ Job do not get "inside' – simple command trigger a script (external workload manager needed)<br>▪ Jobs start/stop were managed using a common Jabber-like client |

STAR framework btw is a single purpose framework for simu, reco, user analysis

BROOKHAVEN NATIONAL LABORATORY

Testing scale & usage growth
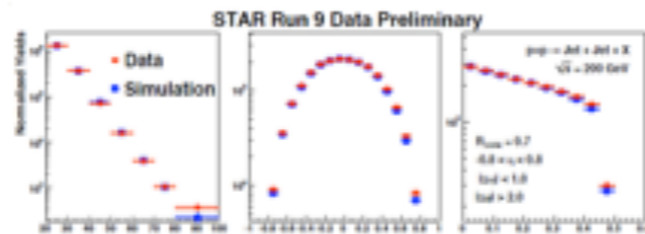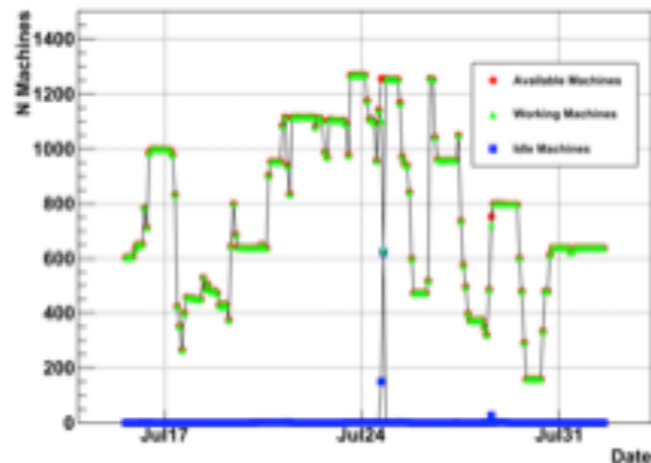
Clemson/Kestrel model

- Generation
  - 12 Billion PYTHIA events were generated – LARGEST sample produced we know off
  - Used over 400,000 CPU hours on 1,000 CPUs at Clemson (+CERN) over the course of a month

- Comparison to normal operation
  - Cloud allowed STAR to expand its computing resources by 25%. Student thesis work possible
    - Available #of CPU per users ~ 50
  - A year long science wait time.

- Achievement for this analysis
  - 4 orders of magnitude increase in number of events used in similar analysis in STAR
  - Near elimination of all uncertainties caused by statistics
  - Un-ambiguously demonstrated good agreement between our data sample and simulation
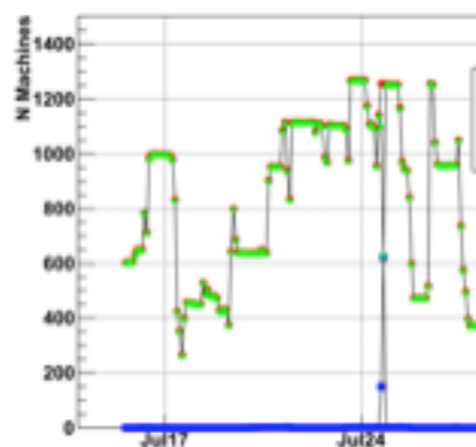  - Results presented at Spin 2010 conference (October)

## Testing scale & usage growth

Models – 10-100 job scale series

Models – larger scales

## Clemson/Kestrel model

- **Generation**
  - **12 Billion PYTHIA events were generated** – LARGEST sample produced we know off
  - **Used over 400,000 CPU hours on 1,000 CPUs** at Clemson (+CERN) over the course of a month

- **Comparison to normal operation**
  - Cloud allowed STAR to expand its computing resources by 25%. Student thesis work possible
    - Available #of CPU per users ~ 50
  - A year long science wait time.

- **Achievement for this analysis**
  - 4 orders of magnitude increase in number of events used in similar analysis in STAR
  - Near elimination of all uncertainties caused by statistics
  - Un-ambiguously demonstrated good agreement between our data sample and simulation
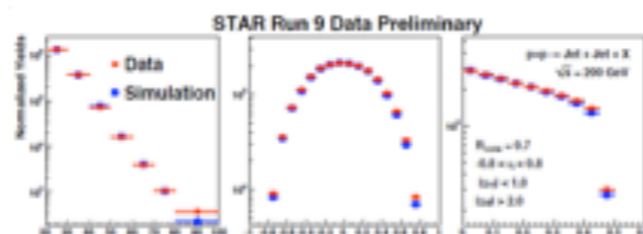  - Results presented at Spin 2010 conference (October)

STAR Run 9 Data Preliminary

- Data
- Simulation

STAR ☆   Jérôme LAURET – CHEP 2010, Academia Sinica Taipei/Taiwan – Oct 2010   14

### CernVM Users

**CernVM**
SoftWare Appliance

~3250 different IP addresses

CHEP 2010     Predrag.Buncic@cern.ch     Taipei, 19 October 2010 - 3

Testing scale & usage growth

Models 10-100 job scale series

Models larger scales

## Clemson/Kestrel model

- **Generation**
  - **12 Billion PYTHIA events were generated** – LARGEST sample produced we know off
  - **Used over 400,000 CPU hours on 1,000 CPUs** at Clemson (+CERN) over the course of a month

- **Comparison to normal operation**
  - Cloud allowed STAR to expand its computing resources by 25%. Student thesis work possible
    - Available #of CPU per users ~ 50
  - A year long science wait time.

- **Achievement for this analysis**
  - 4 orders of magnitude increase in number of events used in similar analysis in STAR
  - Near elimination of all uncertainties caused by statistics
  - Un-ambiguously demonstrated good agreement between our data sample and simulation
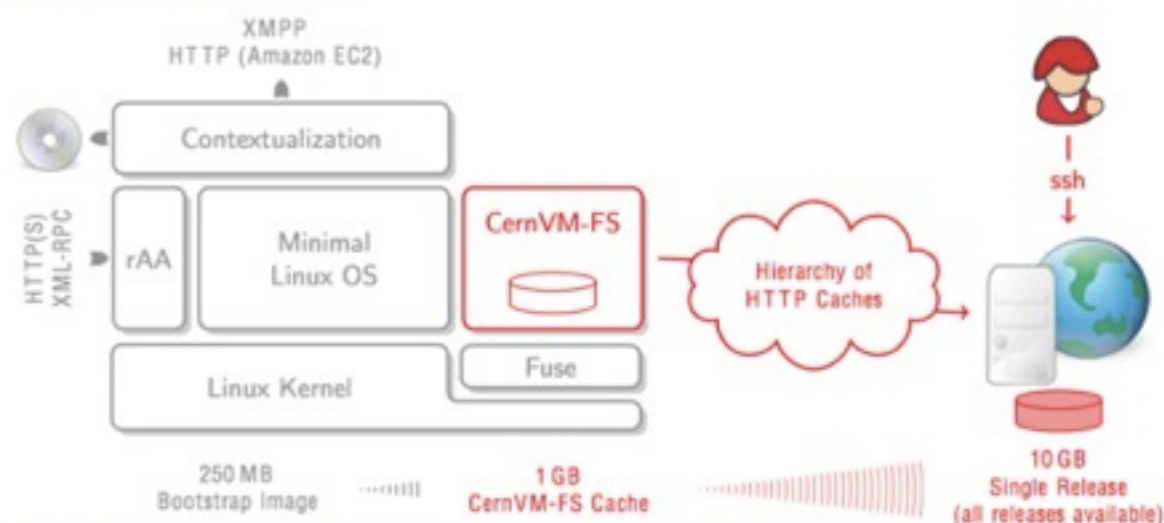  - Results presented at Spin 2010 conference (October)

STAR Run 9 Data Prelim
- Data
- Simulation

CHEP 2010

---

**CernVM Users**

## Part #1: Minimal OS image

XMPP
HTTP (Amazon EC2)

Contextualization

HTTP(S)
XML-RPC

rAA

Minimal Linux OS

CernVM-FS

Hierarchy of HTTP Caches

Linux Kernel

Fuse

ssh

250 MB
Bootstrap Image

1 GB
CernVM-FS Cache

10 GB
Single Release
(all releases available)

- Just enough OS to run LHC applications
- Built using commercial tool (rBuilder by rPath)
  - Top-down approach - starting from application and automatically discovering dependencies
- Small images (250MB), easy to move around

## Testing scale & usage growth

Models = 10-100 job scale series

Models = larger scales

## Clemson/Kestrel model

- **Generation**
  - □ **12 Billion PYTHIA events were generated** – LARGEST sample produced we know off
  - □ **Used over 400,000 CPU hours on 1,000 CPUs** at Clemson (+CERN) over the course of a month

- **Comparison to normal operation**
  - □ Cloud allowed STAR to expand its computing resources by 25%. Student thesis work possible
    - ▪ Available #of CPU per users ~ 50
  - □ A year long science wait time.

- **Achievement for this analysis**
  - □ 4 orders of magnitude increase in number of events used in similar analysis in STAR
  - □ Near elimination of all uncertainties caused by statistics
  - □ Un-ambiguously demonstrated good agreement between our data sample and simulation
  - □ Results presented at Spin 2010 conference (October)



STAR Run 9 Data Prelim
- Data
- Simulation

BROOKHAVEN NATIONAL LABORATORY

**STAR** ☆    Jérôme LAURET – CHEP 2010, Academia Sinica Taipei/Taiwan – Oct

---

**CernVM Users**

**Part #1: Minimal OS image**

**Part #2: CernVM-FS**

**As easy as 1,2,3**



1. Login to Web interface

2. Create user account

3. Select experiment, appliance flavor and preferences

CHEP 2010    Predrag.Buncic@cern.ch    Taipei, 19 October 2010 - 18

# Talking to people.

# Talking to people.

# Talking to people.
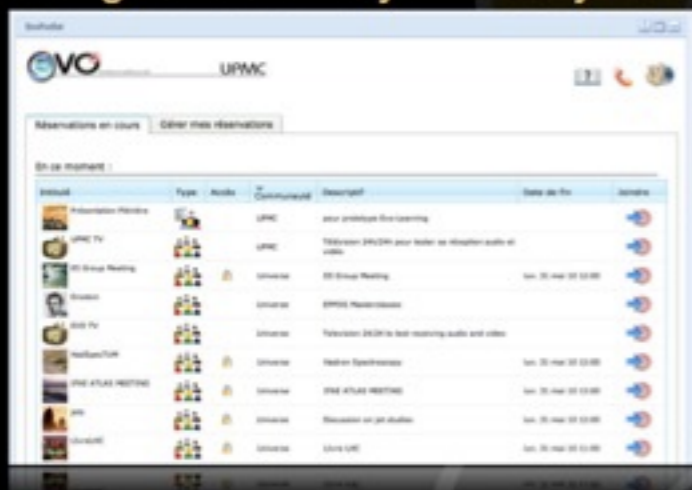
# Talking to people.

# Talking to people.

# Talking to people.
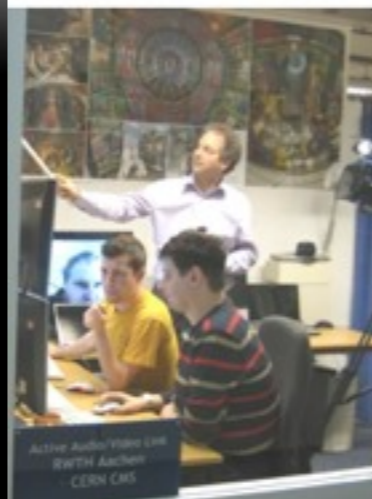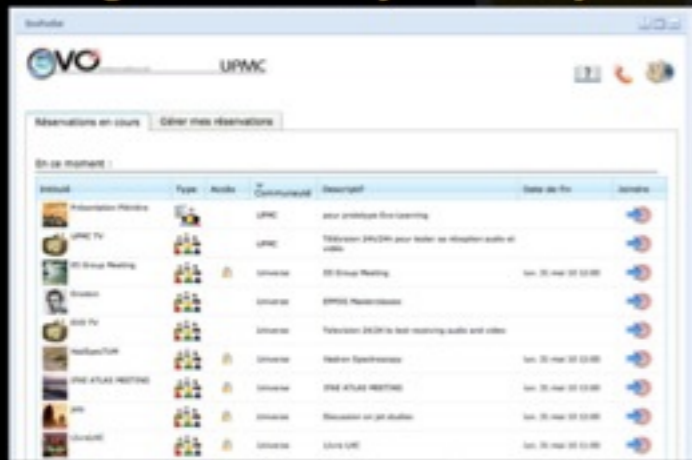
# Talking to people.

# Talking to people.



## EVO Usage

## EVO Web Portal

Web Interface to allow users to access information on EVO meetings and click to join directly from the portal



Integration with Shibboleth and others authentication mechanisms will be added as we deploy the portal

EVO.CALTECH.EDU

Origins

What are they used for ?

Business model



▸ Franchise business model – each institute pays
▸ Standard design with commodity systems

15

# And another thing:
# EMI

# And another thing: EMI

# And another thing: EMI

# And another thing: EMI

# The best plenary

- Lucas Taylor's CERN Outreach talk.

# The best plenary

- Lucas Taylor's CERN Outreach talk.

# The best plenary

- Lucas Taylor's CERN Outreach talk



**EUROBAROMETER**

## A majority of European citizens agree

"Scientists do not put enough effort into informing the public about new developments in science and technology"

http://ec.europa.eu/public_opinion/archives/ebs/ebs_340_en.pdf

CMS Centre in the Austrian
Parliament, Vienna

# The best plenary

- Lucas Taylor's CERN Outreach talk.

## LHC Communications Strategy

1. Coherent and high-quality messages

2. Open engagement with traditional media

3. Exploitation of new media (Web 2.0)

CMS Centre in the Austrian
Parliament, Vienna

# The best plenary

# The best plenary



## How about with younger people ?

Poll says not enough effort is made to reach young people

EUROBAROMETER

cern.ch popularity by age (relative to general internet population)

| | |
|---|---|
| 18-24 | |
| 25-34 | |
| 35-44 | |
| 45-54 | |
| 55-64 | |
| 65+ | |

▸ 3,500 teachers went through CERN Teachers Programme and now teach O(100,000) school students at any time

CERN Teacher Programme Participants 1998 - 31 July 2010

# The best plenary

- How well are we doing?

## Language monitoring of online and print media

http://www.languagemonitor.com/news/top-words-of-2009/

| Top Phrases of 2009 | Top Words of 2009 | Top Names of 2009 |
| --- | --- | --- |
| 1. King of Pop | 1. Twitter | 1. Barack Obama |
| 2. Obama-mania | 2. Obama | 2. Michael Jackson |
| 3. Climate Change | 3. H1N1 | 3. Mobama |
| 4. Swine | 4. Stimulus | 4. Large Hadron Collider |
| 5. Too Large to Fail | 5. Vampire | 5. Neda Agha Sultan |
| 6. Cloud Computing | 6. 2.0 (next gen.) | 6. Nancy Pelosi |
| 7. Public | 7. Deficit | 7. M. Ahmadinejad |
| 8. Jai Ho! | 8. Hadron | 8. Hamid Karzai |
| 9. Mayan Calendar | 9. Healthcare | 9. Rahm Emmanuel |
| 10. God Particle | 10. Transparency | 10. Sonia Sotomayor |

with nothing at all LHC-related in 2008

# The best plenary

- How well are we doing?
  What can you do ?

  How ... with ...

  )09/

  **63% of respondents agree that**

  *"scientists working at a university or government laboratories are the best qualified to explain scientific and technological developments"*

  els

  | 10. God Particle | 10. Transparency | 10. Sonia Sotomayor |

  with nothing at all LHC–related in 2008

# The worst plenary (a moment of indulgence).

- ACER Marketing SSD talk.

- An object lesson in how marketing people will misuse graphs to attempt to mislead you.

- SLC (expensive) SSDs used for performance comparisons.

- MLC (cheap) SSDs used for price comparisons.

- Write performance carefully not explored much (15K HDDs and RAID0 both beat even some SLC SSDs at this).

# The worst plenary (a moment of indulgence).



...alk.

marketing people will
...t to mislead you.

...sed for performance

...l for price comparisons.

- Write performance carefully not explored much (15K HDDs and RAID0 both beat even some SLC SSDs at this).

# The worst plenary (a moment of indulgence).

...alk.

...eople will ... you.

...rmance ...mparisons.

**HDD vs. SSD Throughput Performance**

Price/performance?
What about writes?

I/O Per Second    Throughput

Better

IOPS
12,000
10,000
8,000
6,000
4,000
2,000
0

11,497
+1
3,994

24 Seagate® Savvio 15K SAS HDDs    3 Intel® X25-E Extreme SATA SSD

- Hardware: a shelf of 3 Intel® X25-E Extrem... cases, Principled Tech. used a Newisys N...
- Software: Jetstress; Performance results i...
- Average throughput in MB per second for t...

**HDD vs. SSD Read/Write Performance**

If you ignore 15K HDDs and RAID arrays

Sequential          Random

SSD is Better

Megabytes
2,300
160
150
100
89   59   30
89   59   30
100

2,300
70  60  51
60  60  53
100

Sequential Read MB/s    Sequential Write MB/s    Random Read IOPS    Random Write IOPS

■ WD 7200 HDD    ■ WD 5400 HDD    Toshiba 4200 HDD    ■ SanDisk SSD

- Hardware: Intel® Core™2 Duo processor E8400 (3 GHz, 6 MB L2 cache, 1333 MHz FSB), 2 GB DDR2 Non-ECC SDRAM, 800 MHz
- Software: IOMeter, measured as a secondary drive

- Wri... lored much (15K HDDs and RAID0 both beat even some SLC SSDs at this).

# The worst plenary (a moment of indulgence).

_...alk._

_...eople will_
_...you._

## HDD vs. SSD Throughput Performance

Price/performance?
What about writes?

**I/O Per Second**    **Throughput**

_Better_

- Hardware: a shelf of 3 Intel® X25-E Extrem cases, Principled Tech. used a Newisys N[
- Software: Jetstress; Performance results i
- Average throughput in MB per second for t

IOPS: 12,000 / 10,000 / 8,000 / 6,000 / 4,000 / 2,000 / 0

3,994 — 24 Seagate® Savvio 15K SAS HDDs
11,497 — 3 Intel® X25-E Extreme SATA SS
+1

## HDD vs. SSD Read/Write Performance

If you ignore 15K HDDs and RAID arrays

**Sequential**    **Random**

_SSD Better_

Megabytes: 2,300 / 150 / 100 / 50 / 0

89, 59, 30 — Sequential Read MB/s
89, 59, 30 — Sequential Write M
160

■ WD 7200 HDD    ■ WD 5400 HDD

- Hardware: Intel® Core™2 Duo processor E8400 (3 GH 2 GB DDR2 Non-ECC SDRAM, 800 MHz
- Software: IOMeter, measured as a secondary drive

## Trends in Selling Prices

$USD / per GB

● SLC Flash

—— SSD
—— HDD

Source: IDC, 2009

MLC Flash (original slide)

_The gap is getting closer!_

10 / 8 / 6 / 4 / 2 / 0

2008  2009  2010  2011  2012  2013

- Wri...
  much (15K HDDs a...
  some SLC SSDs at...

# The summarised Summary, concluded.

- Data Management is hard.

  - So we have to be cleverer about how we do it.

- But the LHC works.

- And lots of cool tech is waiting for us in the future, if we can use it.

# Conclusion

- 謝謝你們。你們要訊嗎？

- (Thank you. Do you have any questions?)

# List of talks referenced.

- EMI, the Future of the European Data Management Middleware (PS15-1-463)

- Standard Protocols in DPM (PS15-3-443)

- LHC Data Analysis Using NFSv4.1 (pNFS): A Detailed Evaluation. (PS35-4-289)

- ng: What Next-Gen Languages Can Teach Us About HENP Frameworks in the Manycore Era [114]

- Parallelizing Atlas Reconstruction and Simulation: Issues and Optimization Solutions for Scaling on Multi- and Many-CPU Platforms (PS18-3-126)

- Algorithm Acceleration from GPGPUs for the ATLAS Upgrade [273]

- Adaptive Data Management in the ARC Grid Middleware (PS21-1-156)

- When STAR Meets the Clouds – Virtualization & Grid Experience (PS44-1-322)

- Tests of Cloud Computing and Storage System Features for Use in the H1 Collaboration Data Preservation Model (H1 Collaboration) ( PS44-2-346)

- Establishing Applicability of SSDs to LHC Tier-2 Hardware Configuration (PS35-3-288)

- Distributing LHC Application Software and Conditions Databases Using CernVM File System (PS06-5-434)

- CernVM: Minimal Maintenance Approach to the Virtualization (PS29-4-432)

- ATLAS, CMS, New Challenges for the Public Communication (PL-08)

- One Small Step: A View of the LHC Experiments' Offline Systems After One Year of Data-Taking (PL -02)

- How to Harness the Performance Potential of Current Multi-Core CPUs and GPUs (PL-04)

- CMS Centres Worldwide - A New Collaborative Infrastructure (PS34-2-267)