# Understanding and modelling of the distributed infrastructure & computing models

With H-LHC in mind

# Input from Daniele and Eric

- Which I tried to integrate
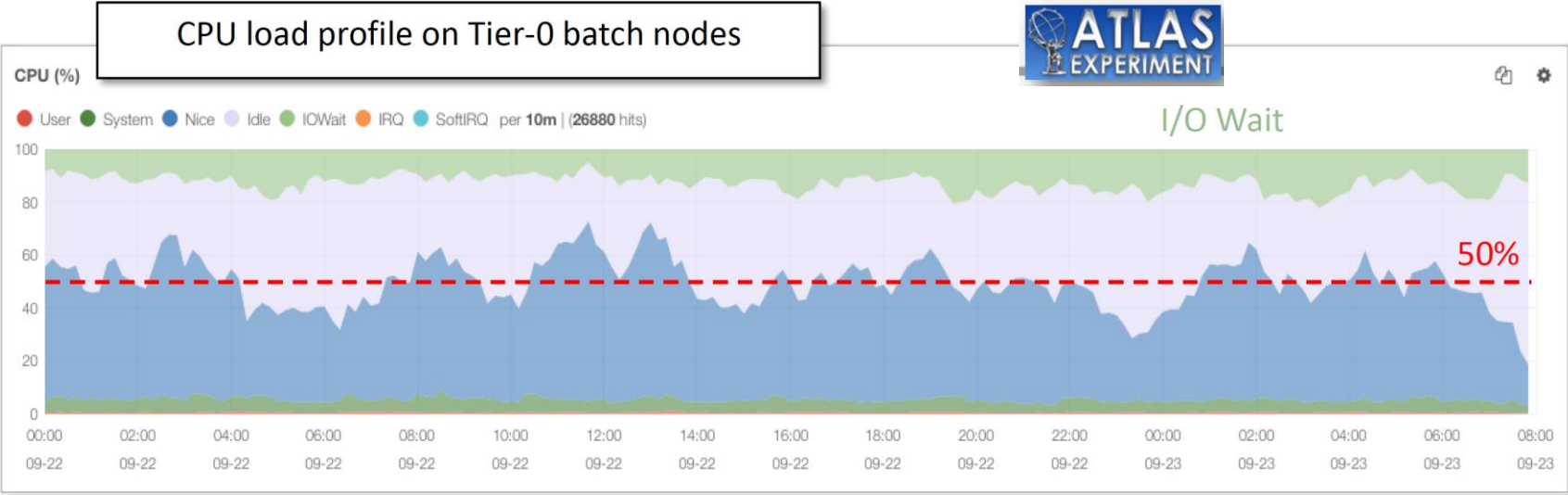- Attached to the Indico page

# From the Agenda

- *How well do we <span style="color:red">understand</span> our current workflows, their behavior and resource needs?*
    - *with respect to storage, remote access, networks, CPU, memory*
    - *How well do we <span style="color:red">understand</span> the behavior of our current infrastructure?*
    - *What can we do to improve this <span style="color:red">understanding</span> in an experiment independent way?*
    - *how independent can this be?*
- *What has been done already in experiments?*
- *What would be desirable? Ability to <span style="color:green">model</span> ideas of infrastructure to <span style="color:red">understand</span> performance, costs, etc.*
- *What is potentially common across experiments? What is specific?*
- ***Can we derive a cost <span style="color:green">model</span> for the infrastructure to <span style="color:green">explain</span> the full costs of computing and the relative costs of each component?***

# Information gathering (++)

- Monitoring data very fine grained for workflows and infrastructure
  - xrootd monitoring
  - PerfSonar
  - Data management monitoring
  - Fabric monitoring
  - Dashboards ..........
- Performance analysis tools
  - Detailed traces
    - Memory, cpu, storage etc.
- Analytics
  - Significant investment (people and hardware) using advanced tools
    - Machine learning etc.

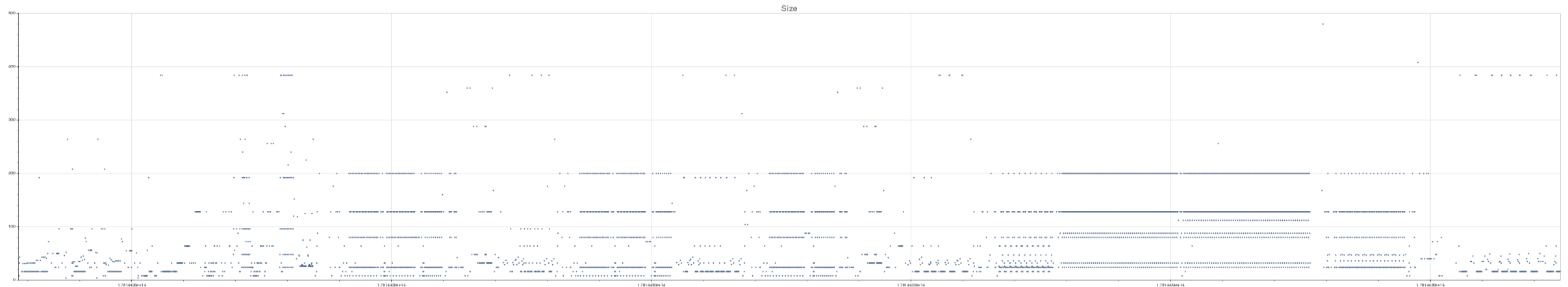# Examples



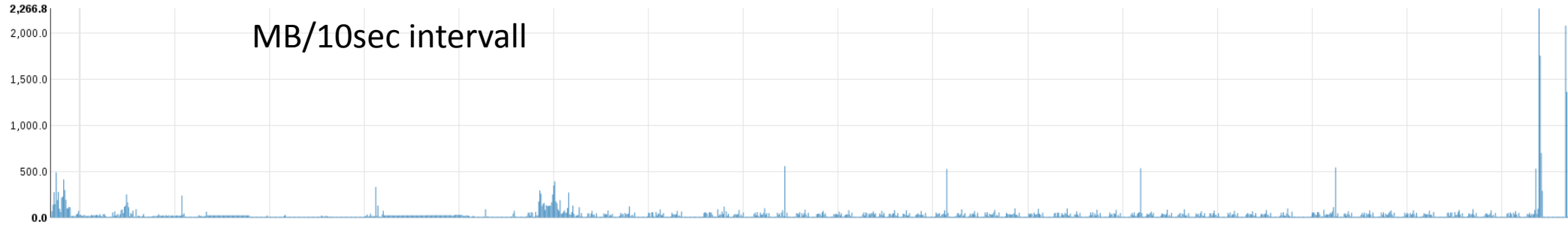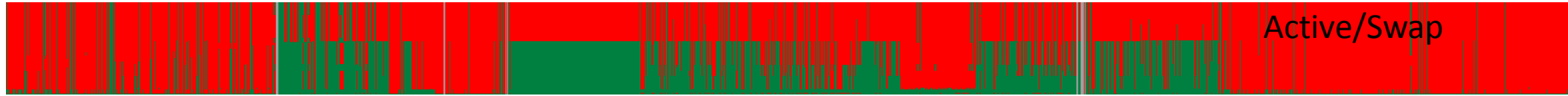CPU load profile on Tier-0 batch nodes

| Step | Setup | RAWtoESD* | RAWtoESD validation | ESDtoAOD setup | ESDtoAOD* | ESDtoAOD validation | DQHistogram Merge setup | DQHistogram Merge | DQHistogram Merge valid. |
|---|---|---|---|---|---|---|---|---|---|
| Wall time (*MP) | 6m 26s | 2h 47m 26s | 6m 29s | 7m 56s | 1h 0m 54s | 4m 13s | 19s | 2m 37s | 1s |
| CPU time, efficiency | N/A | 10h 20m 56s 92.7% | N/A | N/A | 2h 34m 16s 63.3% | N/A | N/A | 29s 18.5% | N/A |

| ESDtoDPD setup | ESDtoDPD* | POOLMerge Athena setup | POOLMerge Athena | ESDtoDPD validation | POOLMerge file validation | Finalisation |
|---|---|---|---|---|---|---|
| 3m 6s | 16m 22s | 1m 2s | 32m 58s | 5m 7s | 33m 15s | 8s |
| N/A | 20m 49s 31.8% | N/A | 5m 18s 16.1% | N/A | N/A | N/A |

Did we expect this?
Do we know why this workflow behaves this way (quantitatively) ?
What would change if we double/half the network/memory?

# Examples: Tracking Memory at the Nanoscale



Active/Swap

MB/10sec intervall

Size

# Examples: More



Allocation Statistics

Locality

**Incredible detailed information. Why do we see these patterns. Do we expect them? What is the impact on performance?**

# However….

- When things change we are often surprised:
  - CERN/Wigner performance differences
  - Virtualization performance differences
  - Move to multi threaded processing
- We are very good at <span style="color:red">noticing</span> and <span style="color:red">measuring</span> effects
  - Good monitoring and logging of "all and every thing"


- We are bad (quantitative) at answering:  **What if ?**
  - Can't predict well the effects of changes (workflows, infrastructure…)
  - Can't easily identify the main reasons and interdependencies
    - If more than one thing changes at the same time ( as it always does)
  - <span style="color:red">CERN-Wigner extension, single core → multi core, Spinning disks/SSDs</span>
- Not surprisingly given the complexity of the environment….
  - Very divers, many factors driving efficiency and performance.
  - Large phase space

# Knowledge / Understanding

- Understanding can be seen as a model based form of data compression *

  - *understanding* something means being able to figure out a simple set of rules that explains it.
  - Think about how the model of a rotating earth allows to predict data as brightness, temperature, and atmospheric composition during a day

* Gregory Chaitin

# What we could use Models for

- Understanding better the existing system
- Document and represent what we think that we have understood
  - Comparing measurements and model
- Spot gaps in our understanding
  - Guide analysis of infrastructure and workflows
- Exploring alternative approaches
  - Workflows and Infrastructures
  - More quickly and more cheaply
- Guide purchase decisions → meaningful cost model
  - When a cost model is included
  - X1 cores + Y1 disks + Z1 MB/core + R1 MB network  will me n events/hour of workflow D per Euro

# What kind of Model(s) might be useful?

- Model in the sense of "**simulation**" of the infrastructure elements and their interaction
  - Discrete Event Simulations ( used in real time systems and networks)
  - Hybrids
    - Like SimGrid (used for HPC, Grid, Cloud ..)
- Model in the sense of an **analytical model** describing the behavior of infrastructure and applications
  - Could be a set of rules to do back of the envelope calculations
  - In the most basic case an Excel table
  - Probably several already around
- In HPC it is standard practice to use modeling of workloads and machines during the design phase
  - Maybe we can profit from their expertise

# Cost model?



Cost of increase & renewal

- Legend: T2 - Disk renewal, T2 - Disk increase, T2 - CPU renewal, T2 - CPU increase, T1 - Tape increase, T1 - Disk renewal, T1 - Disk increase, T1 - CPU renewal, T1 - CPU increase

2017 pie chart: T1 - Tape 12%, T1 - CPU 16%, T1 - Disk 31%, T2 - CPU 28%, T2 - Disk 13%
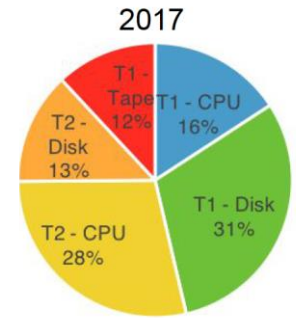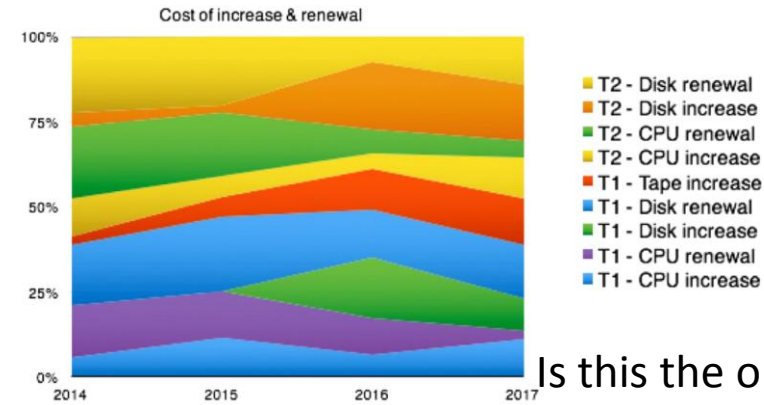
- **What is the metric that we want to look at**?
  - HepSpec/Watt?
  - HepSpec/CHF ?
  - **Workflow-A/B/C/D events per week /CHF ?**
    - Best in a mixture representing our needs
  - …….

Is this the optimum?

- Goal: Understand better on what we should spent our budgets on
  - Answer questions like: Is it better to move to from 2 to 3.5 GB memory per core and buy 10% less disk?
  - Needs good model for cost prediction

- **Complications:**
- How do we account for human effort?
  - Cloud, many/few sites, etc.

- Budgets are often not fluid

- Funding agencies are driven by additional/different motivations
  - But they are not brain dead and having quantitative arguments might help

# Common/Specific

- Common:
  - Models of the **infrastructure**
    - Global, local, generic storage….
  - ***Framework*** to model experiment workflows
  - Tools and "survey program" to analyze the workflows and their impact on the infrastructure
    - Like the FOM tool (HSF) and the allocation tracer
  - Metric to express the parameters of models
- Specific:
  - Models for specific storage systems/sites
  - The analysis and modeling of the different workflows
    - Using common tools as far as possible
  - Comparing model and monitoring data for a specific experiment/site
  - Exploration of new concepts

# What happened to MONARC ?

- LCG Modeling effort in the late 1990s early 2000
- Guided the design towards the hierarchical T0/T1/T2 mdel
- Then was used very little

- Was done *before* we had any infrastructure or established workflows
  - No rapid feedback loop
  - WLCG's focus quickly shifted to make it work all….

- Why not the obvious it as a starting point?
  - Tech evolved
  - Infrastructure evolved

# First Steps?
## Mostly speculative …

- Do not start a large (EC funded) project!
- Many different ongoing activities
  - Concurrency Forum
  - HSF
  - Experiments
    - Like RUCIO modeling based on SimGrid
- →**First step**: collect and document ongoing activities
- Modeling needs a home: WLCG / HSF ?
  - WLCG better connection to the infrastructure
  - HSF better linked with the workflows/applications
- Start with a very primitive model!
  - Maybe estimating on paper what the theoretical best efficiency of workflows could be
    - Academic exercise, but will collect input data that can be re-used later.
  - Be aware of the limitation of models! 20% precision == outstanding
  - Don't get lost in details ( packet level modeling of WLCG