**Volker Guelzow**

**DESY**

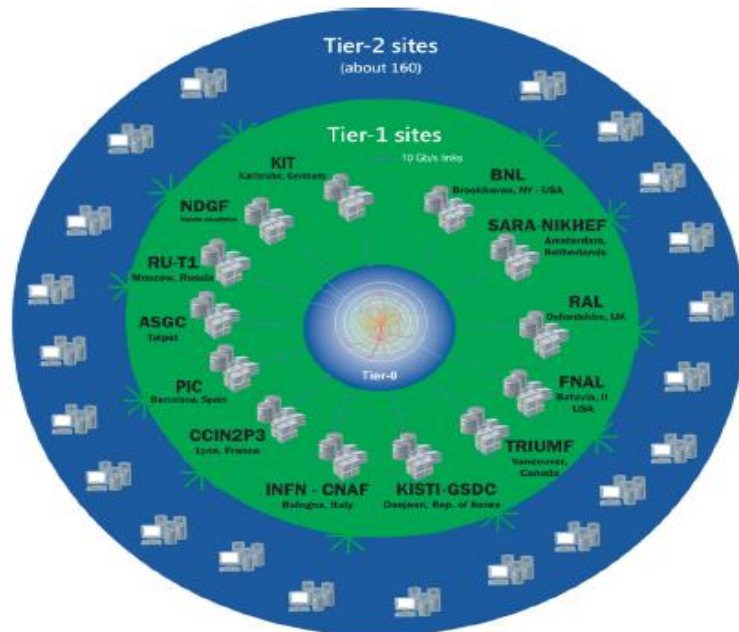**711. Heraeus Seminar, Bad Honnef**

# Worldwide LHC Computing Grid



**TIER-0 (CERN):** data recording, reconstruction and distribution

**TIER-1:** permanent storage, re-processing, analysis

**TIER-2:** Simulation, end-user analysis

nearly 170 sites, 40 countries

~350'000 cores

500 PB of storage

> 2 million jobs/day

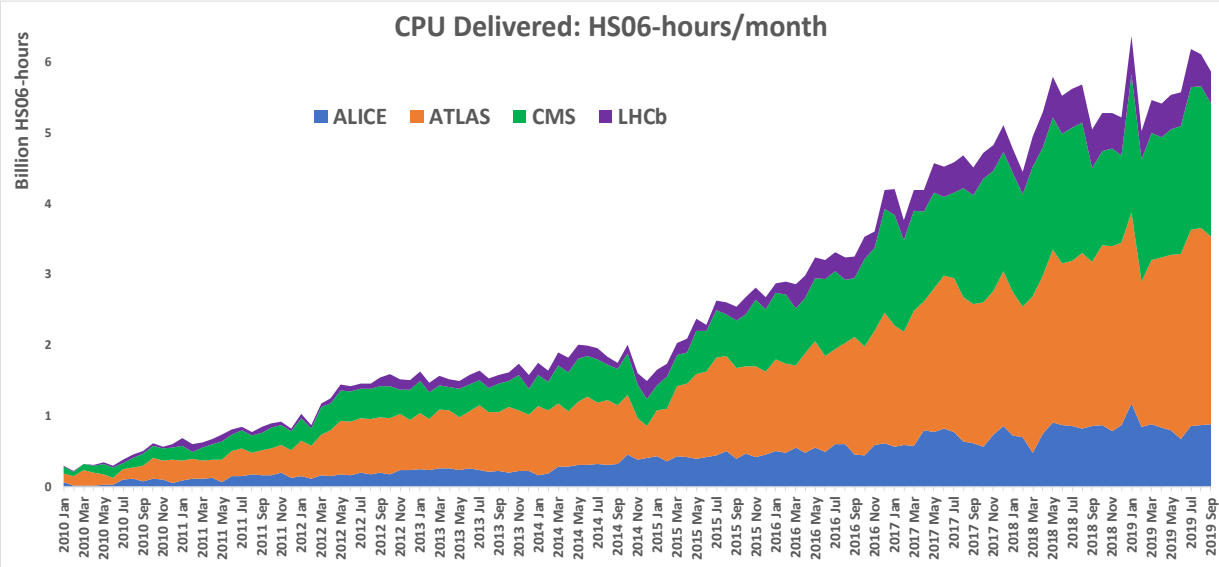10-100 Gb links

# WLCG Collaboration



October 2019:
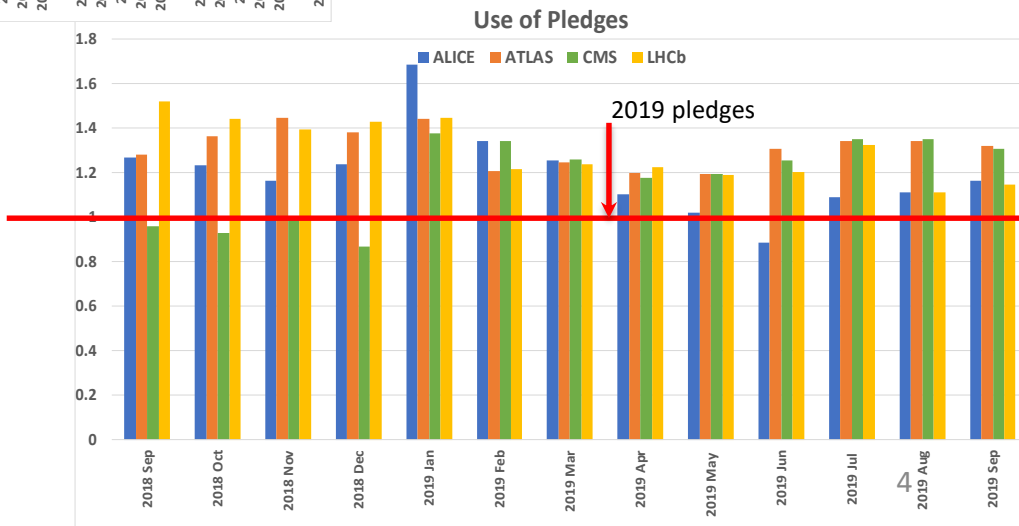- 65 MoU's
- 168 sites; 42 countries

Chinese University of Hong Kong Tier 2 ATLAS

CPU Delivered

New peak: ~270 M HS06-days/month

~ 860 k cores continuous

711. Heraeus Seminar

4

# Resource evolution



**Legend:**
- Actual installed
- 2021 request to C-RSG
- 2022 = 1.5*2018

**CPU Growth** — Pledge, 15% Growth from 2018

**Disk Growth** — Pledge, 15% Growth from 2015

**Tape Growth** — Pledge, 15% Growth from 2018

NB: Run 3 probably manageable overall, *but* constant budget growth until Run 4 is essential for HL-LHC

# Luminosity

$$L = \gamma \frac{n_b N^2 f_{rev}}{4\pi \, \beta^* \, \varepsilon_n} R; \qquad R = 1 / \sqrt{1 + \frac{\theta_c \, \sigma_z}{2\sigma}}$$

where $\gamma$ is the proton beam energy in unit of rest mass; $n_b$ is the number of bunches per beam: 2808 (nominal LHC value) for 25 ns bunch spacing; $N$ is the bunch population. $N_{nominal\ 25\ ns}$: $1.15\times10^{11}$ p ($\Rightarrow$0.58 A of beam current at 2808 bunches); $f_{rev}$ is the revolution frequency (11.2 kHz); $\beta^*$ is the beam beta function (focal length) at the collision point (nominal design 0.55 m); $\varepsilon_n$ is the transverse normalized emittance (no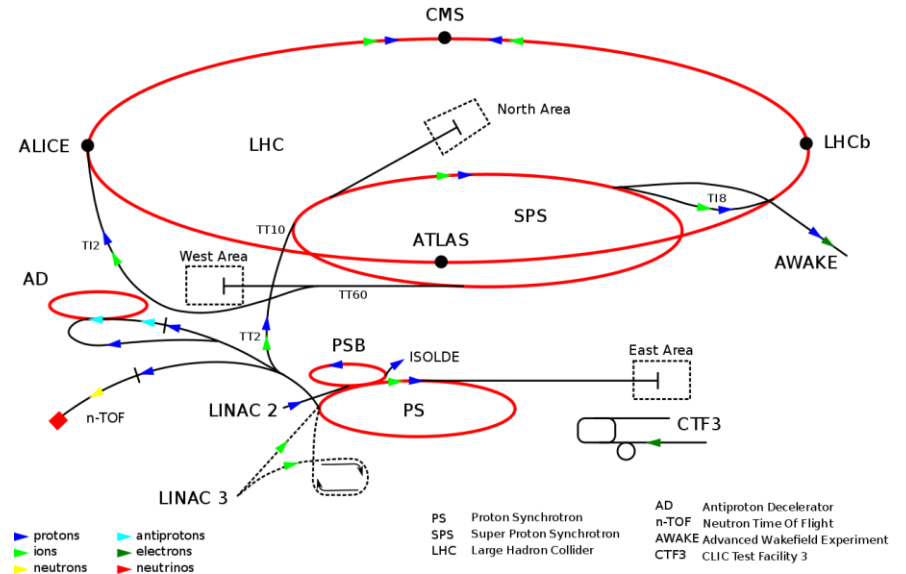minal design: 3.75 μm); $R$ is a luminosity geometrical reduction factor (0.85 at a $\beta^*$ of 0.55 m of, down to 0.5 at 0.25 m); $\theta_c$ is the full crossing angle between colliding beam (285 μrad as nominal design); and $\sigma$, $\sigma_z$ are the transverse and longitudinal r.m.s. sizes, respectively (nominally 16.7 μm and 7.55 cm, respectively)

# What is High Luminosity LHC (HL-LHC?)

- Let's start with some LHC numbers (for computing)
- LHC = Large Hadron Collider
  - Operating today @ 13 TeV, top $2 \cdot 10^{34}$ $cm^{-2}$ $s^{-1}$ instantaneous luminosity via pp collisions bunched @ 25 ns
  - Designed for a vast physics program; clearly the discovery / exclusion of the **Higgs boson was top in the list**
  - This means, given a total inelastic cross section of ~100 mb, **35 collisions per bunch crossing** averaged along O(10) hour fills
  - If we naively consider that the big detectors have ~100M acquisition channels (assume 1 byte/channel), the VIRGIN data rate of the big detectors (ATLAS, CMS) would be **4 PB/s**



| | |
|---|---|
| PS | Proton Synchrotron |
| SPS | Super Proton Synchrotron |
| LHC | Large Hadron Collider |
| AD | Antiproton Decelerator |
| n-TOF | Neutron Time Of Flight |
| AWAKE | Advanced Wakefield Experiment |
| CTF3 | CLIC Test Facility 3 |

protons   antiprotons
ions   electrons
neutrons   neutrinos

# Higgs boson production, expected mechanisms at LHC planning times

- Higgs production cross section (how probable to create one) increases very sharply with collider energy

- The actual number of produced events in a given process is proportional to its cross section, and the collider luminosity

- $N = \sigma \times L_{int}$

  How probable the process is "per collision" ($1 \, m^2 = 10^{28}$ barn)

  How many collisions we are trying $m^{-2}$

- **W**here $L_{int}$ is the integrated luminosity an experiment has been given

- Quite varying with the mass, but the typical Higgs production cross section is ~1-100 pb @ a 13 TeV collider

  - @ 1 TeV collider it would be ~ 100-1000 times lower, this is the reason why a direct positive discovery at TeVatron was not probable

# Access to rare Processes



proton - (anti)proton cross sections

$n_{Evt} = \sigma \cdot L$

$7 \times 10^{12}$ eV  Beam Energy
$10^{34}$ cm$^{-2}$ s$^{-1}$  Luminosity
2835  Bunches/Beam
$10^{11}$  Protons/Bunch

**7 TeV Proton Proton** colliding beams

Bunch Crossing  $4 \cdot 10^{7}$ Hz

Proton Collisions  $10^{9}$ Hz

Parton Collisions

New Particle Production  $10^{-5}$ Hz
(Higgs, SUSY, ....)

**Selection of 1 event in 10,000,000,000,000**

# The Physics Drivers

- Electroweak Physics (incl. Higgs)
- Flavour Physics and CP violation
- Strong Interactions
- Neutrino Physics & Astroparticles
- Dark matter and Dark Sector
- Beyond the Standard Model

- Technology
  - accelerators
  - detector
  - computing
- Experiment
- Theory

## Goals

Exploration of the unknown at very short distances

Search for an understanding of the fundamental physical laws
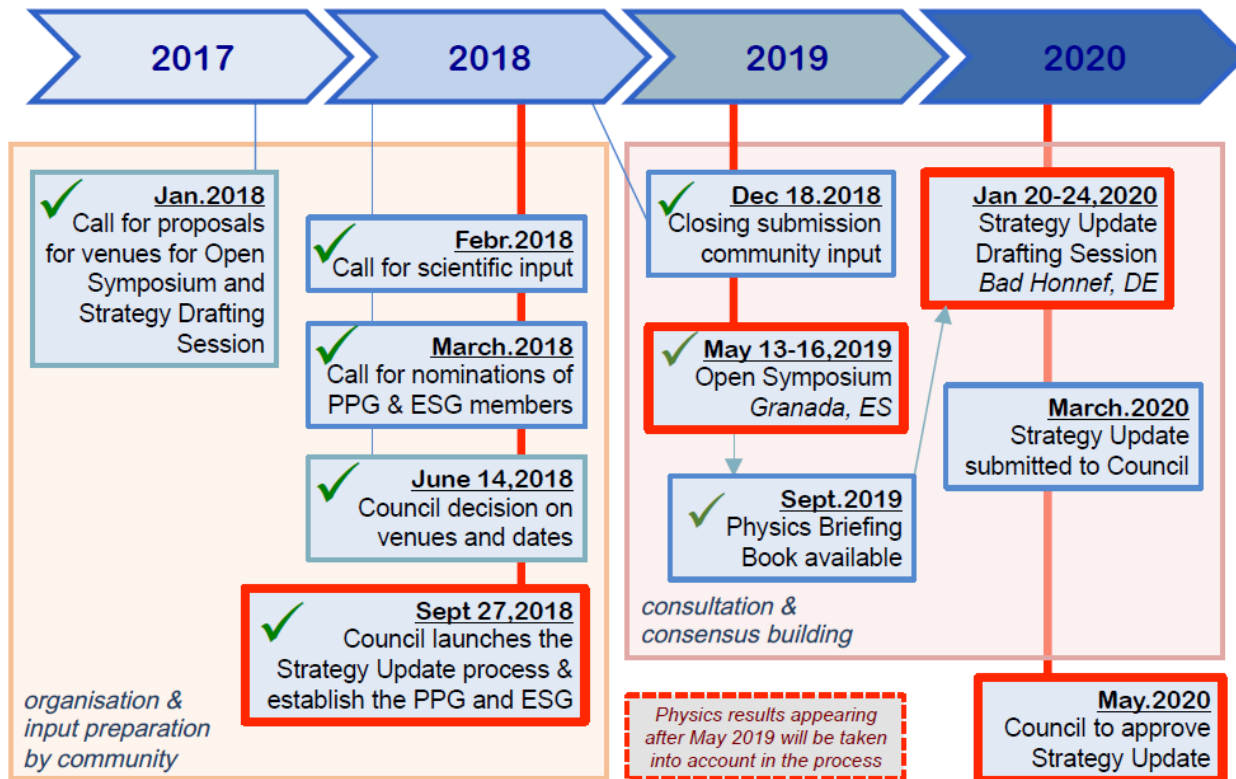
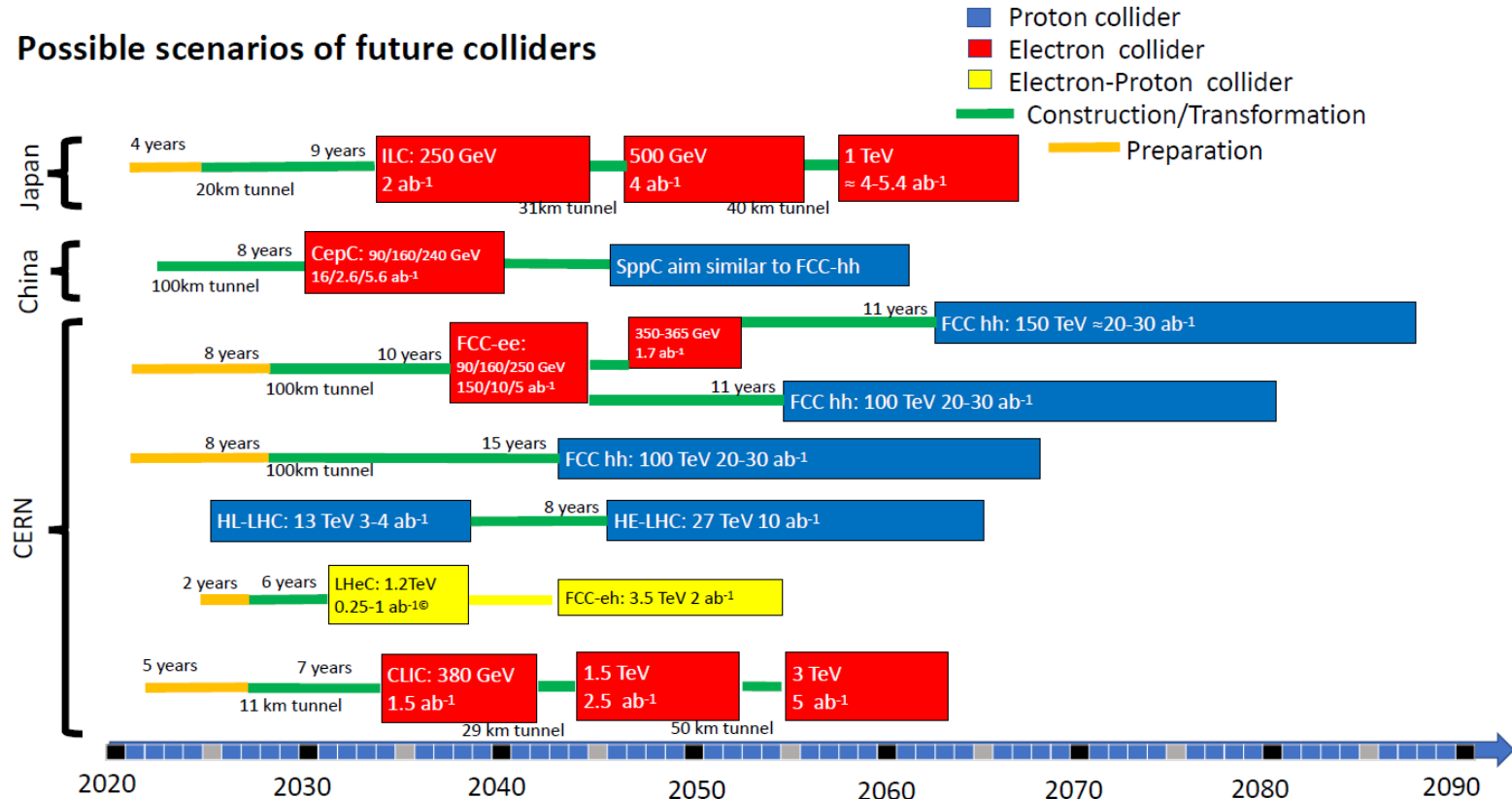## Instruments

diversity/variety/synergy

# European Strategy Update

## A bottom-up process
## to pave the near-term, mid-term and longer-term future



Timeline: 2017 | 2018 | 2019 | 2020

✓ **Jan.2018**
Call for proposals for venues for Open Symposium and Strategy Drafting Session

✓ **Febr.2018**
Call for scientific input

✓ **March.2018**
Call for nominations of PPG & ESG members

✓ **June 14,2018**
Council decision on venues and dates

✓ **Sept 27,2018**
Council launches the Strategy Update process & establish the PPG and ESG

✓ **Dec 18.2018**
Closing submission community input

✓ **May 13-16,2019**
Open Symposium
*Granada, ES*

✓ **Sept.2019**
Physics Briefing Book available

**Jan 20-24,2020**
Strategy Update Drafting Session
*Bad Honnef, DE*

**March.2020**
Strategy Update submitted to Council

**May.2020**
Council to approve Strategy Update

*organisation & input preparation by community*

*consultation & consensus building*

*Physics results appearing after May 2019 will be taken into account in the process*

# Future of HEP: Flagship Projects



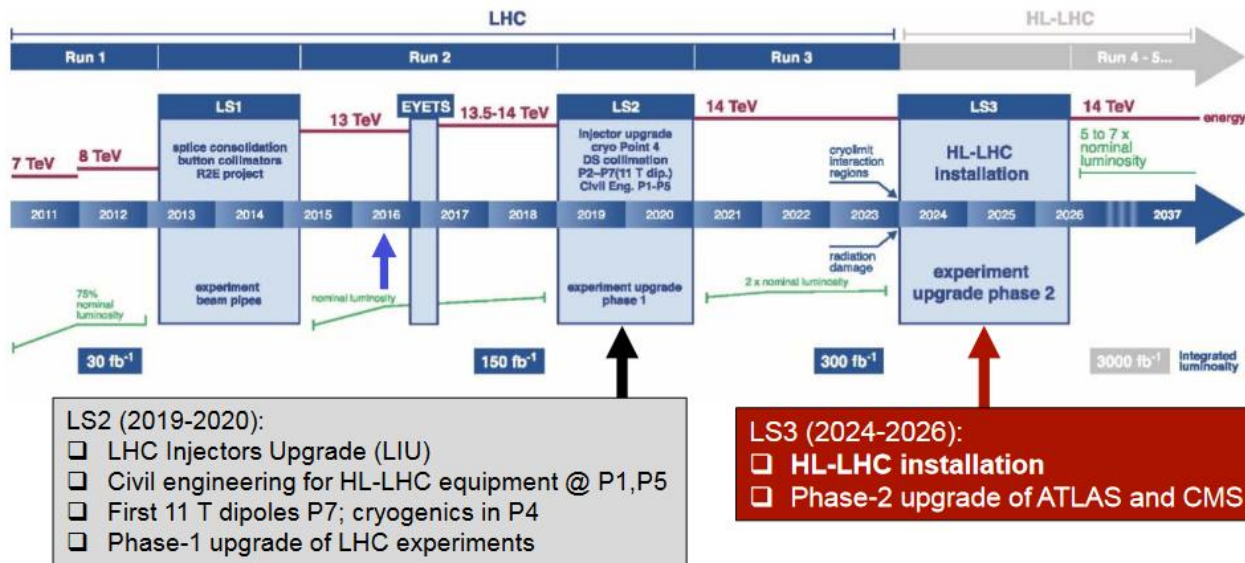Ursula Bassler, Granda, 13.5.2019

# LHC and HL-LHC

- LHC

  - 300 $fb^{-1}$ by 2023

    - 30 $fb^{-1}$ Run 1

    - ~40 $fb^{-1}$ (2015/16)

    - …

- HL-LHC

  - ~3000 $fb^{-1}$
    by ~2035

  - levelled luminosity



LS2 (2019-2020):
- ❏ LHC Injectors Upgrade (LIU)
- ❏ Civil engineering for HL-LHC equipment @ P1,P5
- ❏ First 11 T dipoles P7; cryogenics in P4
- ❏ Phase-1 upgrade of LHC experiments

LS3 (2024-2026):
- ❏ **HL-LHC installation**
- ❏ Phase-2 upgrade of ATLAS and CMS
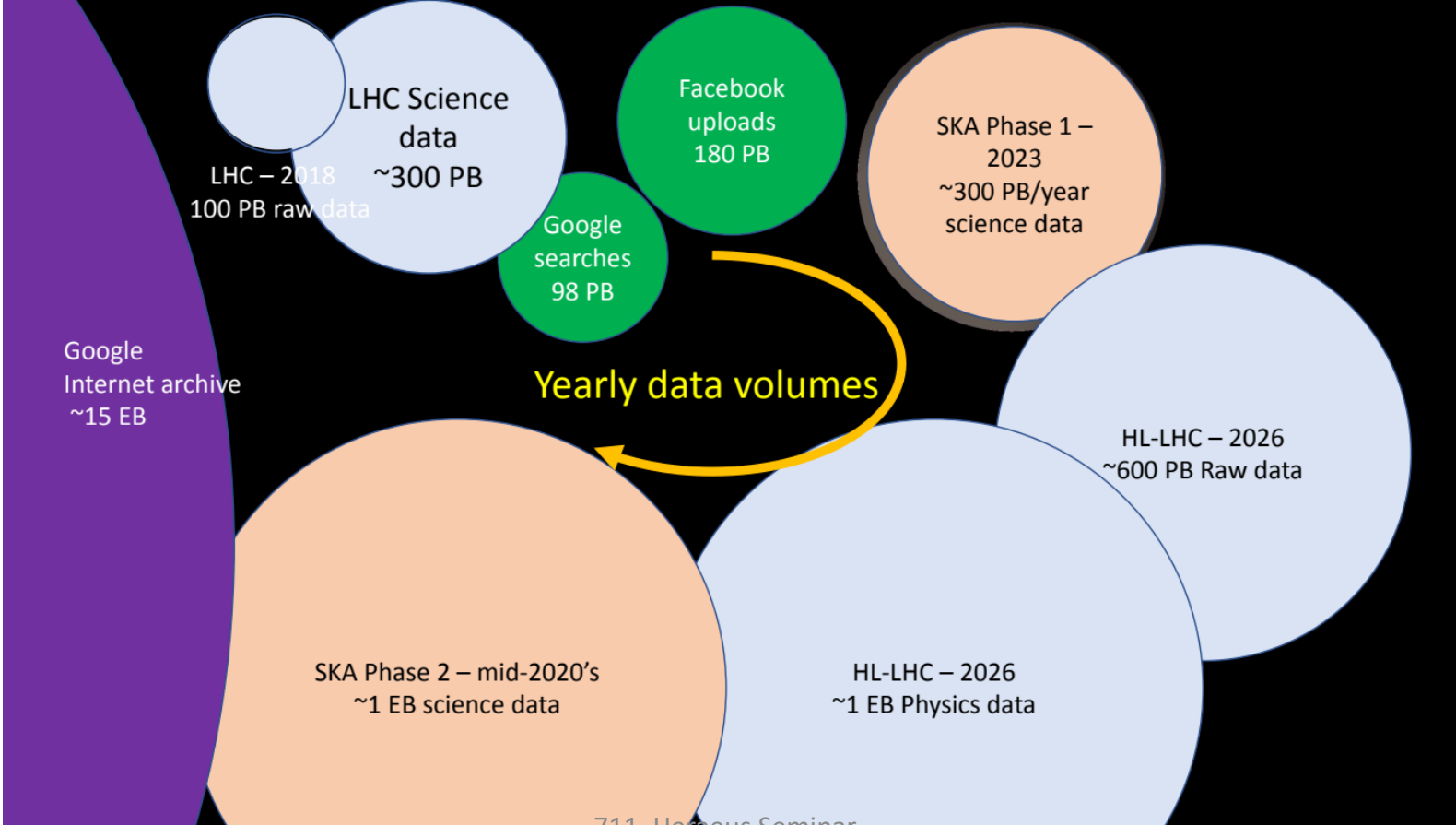
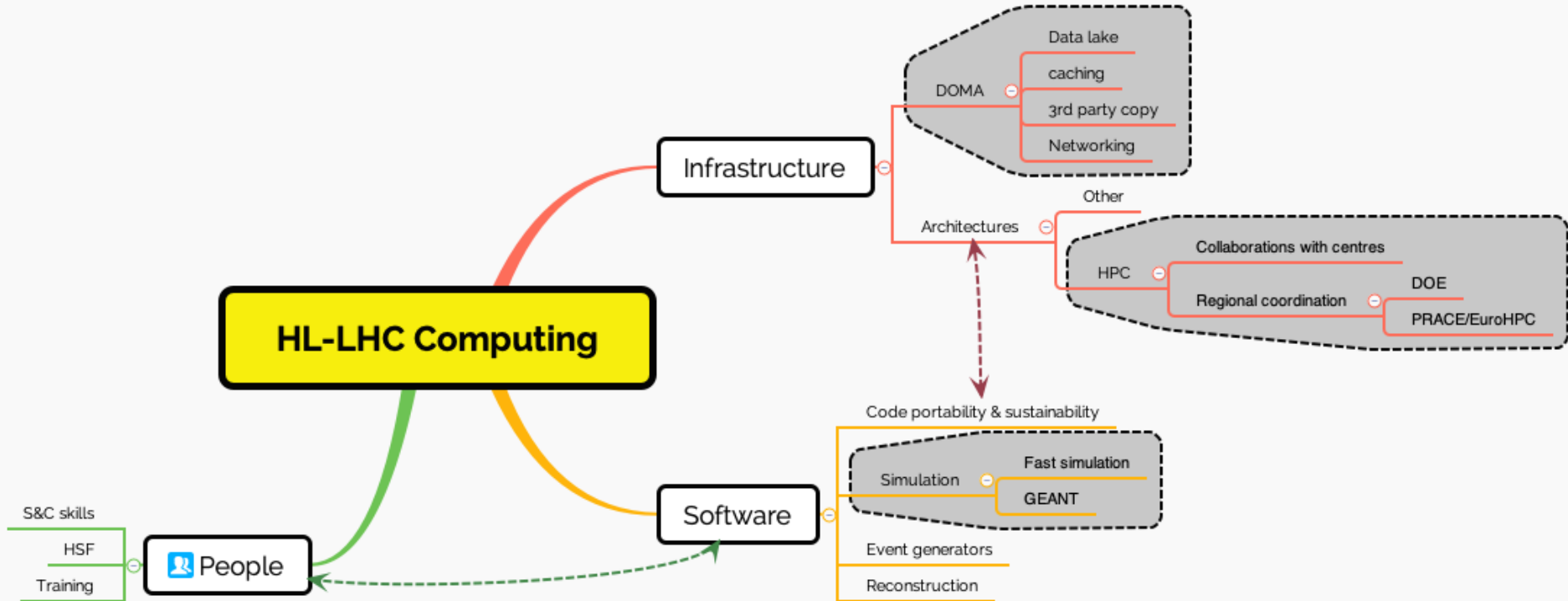# Proposed timeline for Run 4 computing

# Back of the envelope model …

- *«Computing @ HL-LHC would need a resource installation evaluated in 50-100x the current computing infrastructure»*
  - ( … if the processing model simply sales with inputs)
- *«Technology improvement helps in reducing the gap only partially»*
  - Moore's law @ 2x/18months is long gone
  - The same money buy you year-to-year 10-20% only more resources. In 8 years: **1.2\*\*8 = 4**
- *«Factors O(10-20)x are missing in order to be able to process HL-LHC data at the same cost»; otherwise*
  - Do less physics
  - Increase money on HL-LHC computing

# Data Volumes

# Towards HL-LHC
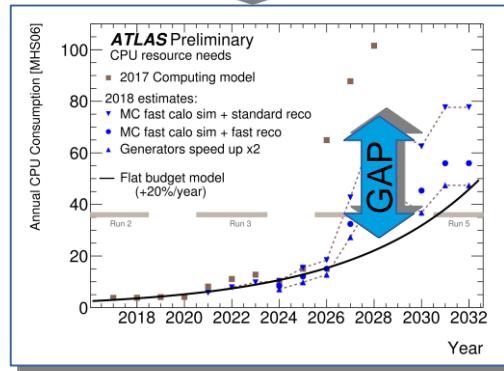
# Addressing the Resource Gap
## Examples

Access to resources
- Local farms
- Grid sites
- Cloud (private & public)
- HPC machines

Hardware trends
- GPU
- FPGA
- Multi-threading

Novel approaches Machine Learning



One effort likely not enough – a combination of many might do

# Evolution of WLCG

- **Community White Paper**
  - 1 year – bottom up review of LHC computing topics
  - 13 working groups on all aspects
  - Outlines how HEP computing could evolve to address computing challenges
  - https://arxiv.org/abs/1712.06982

- **WLCG Strategy Document**
  - Prioritisation of topics in the CWP from the point of view of the HL-LHC challenges
  - Set out a number of R&D projects for the next 5 years
    - Running global system should evolve towards HL-LHC
  - http://cern.ch/go/Tg79

# Strategy - Outline

The strategy develops around five main themes …

1) Software performance
2) Algorithmic improvements / changes (e.g. generators, fast MC, reco)
3) Reduction of data volumes
4) Managing operations cost
5) Optimizing hardware costs

It defines an R&D program with rough timelines, organized in sections:

- The HL-LHC challenge, hardware trends and a cost model
- Computing Models
- Experiments Software
- System Performance and Efficiency
- Data and Processing Infrastructures
- Sustainability
- Data Preservation and Reuse

This was discussed in depth in the WLCG/HSF workshop in Naples in March – many of the activities were started then

# HEP Software Foundation Community White Paper Working Group – Data Processing Frameworks

HEP Software Foundation: Paolo Calafiura[d] Marco Clemencic[a] Hadrien Grasland[b] Chris Green[c] Benedikt Hegner[a,e,1] Chris Jones[c] Michel Jouvin[b] Kyle Knoepfel[c] Thomas Kuhr[g] Jim Kowalkowski[c,1] Charles Leggett[d] Adam Lyon[c] David Malon[e] Marc Paterno[c] Simon Patton[d] Elizabeth Sexton-Kennedy[c,1] Graeme A Stewart[a] Vakho Tsulaia[d]

[a] *CERN, Geneva, Switzerland*
[b] *LAL, Université Paris-Sud and CNRS/IN2P3, Orsay, France*
[c] *Fermi National Accelerator Laboratory, Batavia, Illinois, USA*
[d] *Lawrence Berkeley National Laboratory, Berkeley, CA, USA*
[e] *Brookhaven National Laboratory, Upton, NY, USA*
[f] *Argonne National Laboratory, Lemont, IL, USA*
[g] *Ludwig-Maximilians-Universität München, Munich, Germany*
[1] *Paper Editor*

https://arxiv.org/pdf/1812.07861.pdf

- Throughput maximizing: here it is most important to efficiently move data through all the available resources (memory, storage, and CPU), maximizing the number of events that are processed. The workload management systems used by experiments on the grid work towards this goal.

- Latency minimizing (or reducing): online and interactive use cases where imposing constraints on how long it takes to calculate an answer for a particular datum is relevant and important. Dataflow and transaction processing systems work towards this goal.

# Chapter 11.2 ff  Computing

- It is also equally important to plan for an infrastructure that requires less hardware and less effort to maintain and operate as an experiment mature

# A European Data Science Institute for Fundamental Physics

Maurizio Pierini

CERN Experimental Physics Department

**ABSTRACT:** In order to facilitate the deployment of modern data science technologies (e.g., Deep Learning) into theoretical and experimental research in high energy physics, we suggest that the creation of a **European Data Science Institute for Fundamental Physics** is included among the recommendations of the European Strategy group. Such an institute would facilitate the development of cross-collaboration and across-border work on general-interest techniques, as well as yield knowledge transfer from and to other scientific communities (astrophysics, cosmology, computer science, etc.) and tech companies worldwide.

https://indico.cern.ch/event/765096/contributions/3295512/attachments/1785106/2906008/A_European_Data_Science_Institute_for_Particle_Physics.pdf

# Data & Computing Challenge
## HL-LHC Example

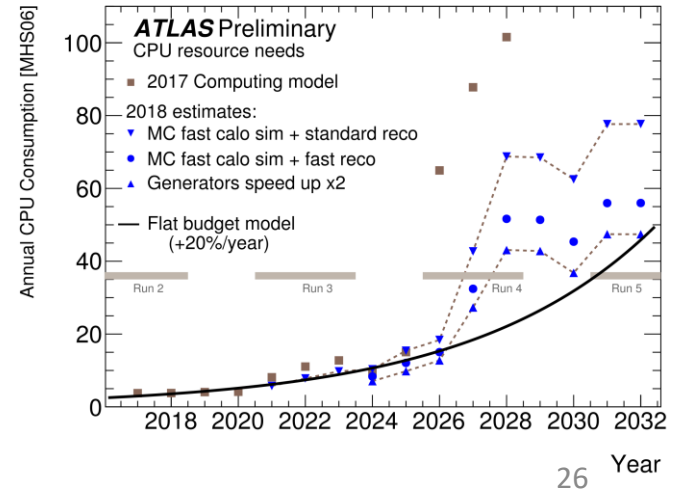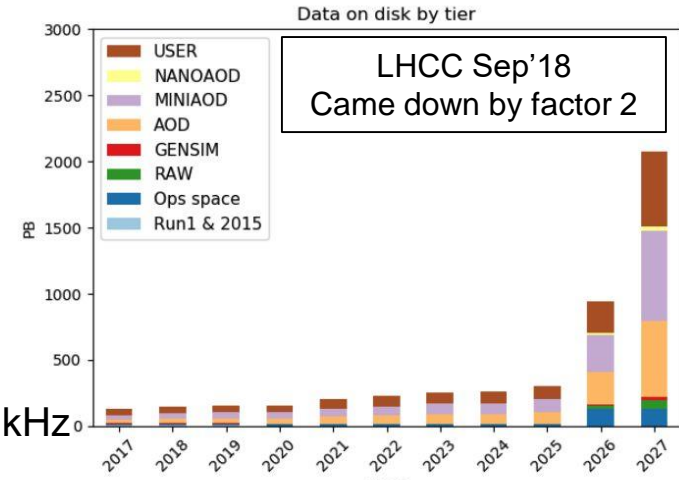**HL-LHC running conditions**

- Pileup goes to ~200 from ~35 in 2017

  – Event size increases by factor 10

  – Reconstruction CPU time demand increases by 10-15

- Logging rate goes to 7.5kHz (or even 10kHz) compared to ~1kHz

**Technology grows by 10-20%/year for the same investment**

- Recovers a factor 5-6 within 10 years

- Recent extrapolations favor lower values

**Adjustments to the computing models required**

- Some options:

  – Already active studied: Smaller data tiers, more use of "fastsim"

  – Optimize software and infrastructure

  – Unlikely: Increase Computing budget by factors (2 or 4…?)



LHCC Sep'18
Came down by factor 2

# So, by 2026….

- We can expect **Reconstruction** (on Data and Monte Carlo events) to be the dominant user of CPU cycles; **Geant4 simulation** following but somehow less important overall
- **Generation** will scale from today's fraction only if we start to need more precise simulations
    - LO → NLO → NNLO → … ?
    - V+ (1,2,3,4,5… N) Jets
    - The negative weights problem? A huge increase in resources if they are not solved
- We can expect the need to have sizeable Fast ("simplified"/"parametrized"/DL) Simulation; but this could clash with the need of more precise measurements
- ## How does analysis scale?:
    - Up: more precision needed, higher dimensional fits, …
    - Not as much: the number of users is ~ constant, «brain time» can be limiting

# TECHNOLOGY

# CPU Performance over last Decades
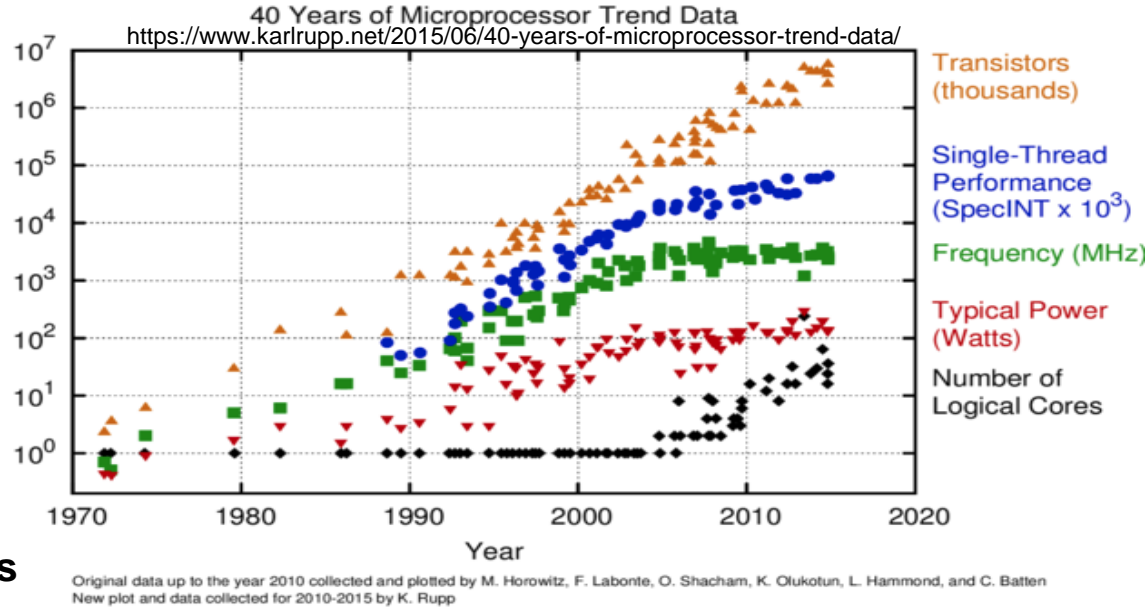
**Almost no performance increase for single thread**

- No "trivial" performance gain (Which used to have until ~10y ago)
- Number of transistors still growing (Moore's law still partly holds)

**Number of CPU cores increasing**

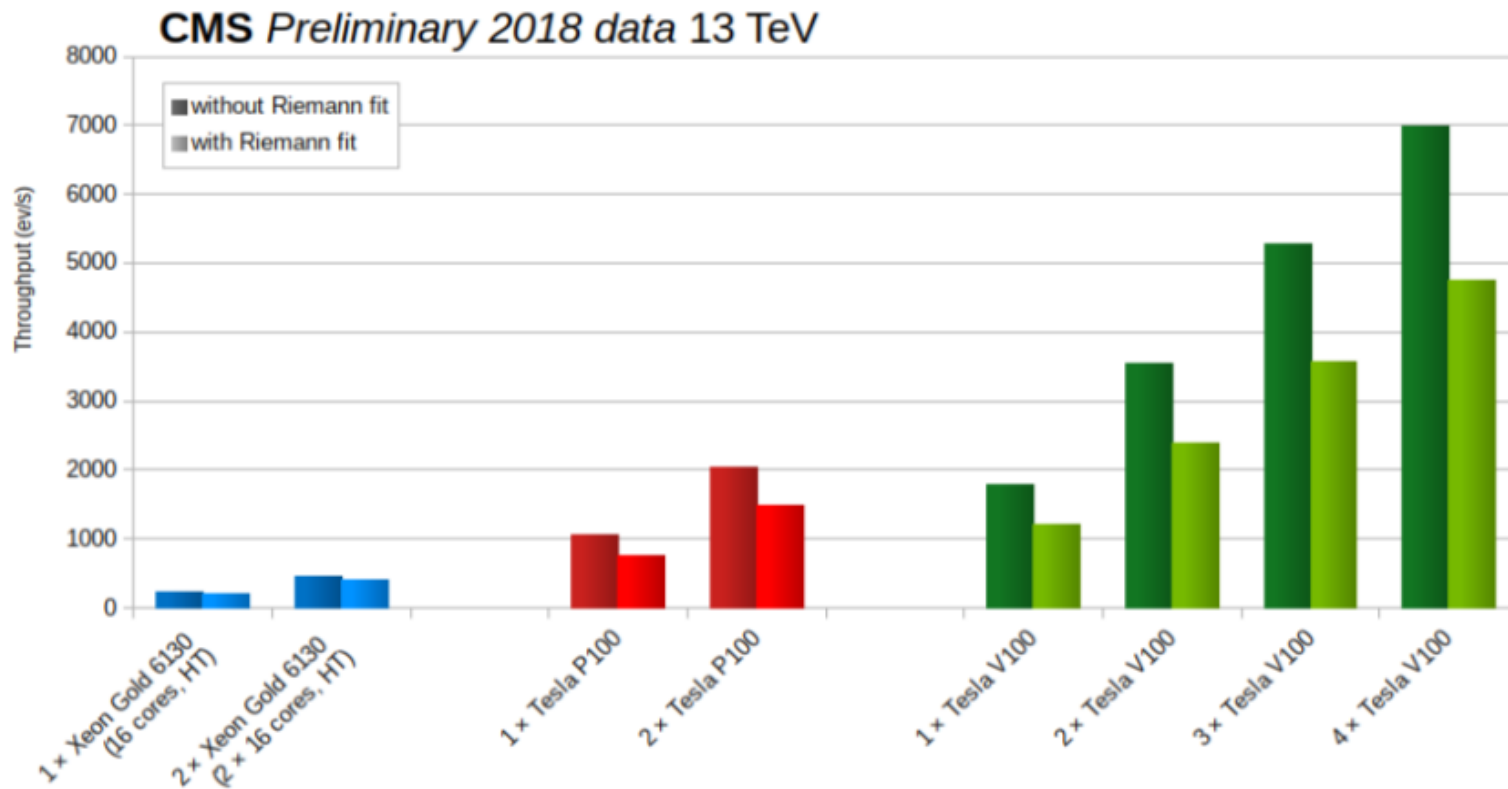- Requires <u>multi-thread</u> enabled and <u>thread-safe</u> applications

**GPU and other special (co-)processors**

- Very fast for specialized applications
- Require dedicated code development with special tools
- Code validation and workload management challenging
  - Present tools origin from Linux/x86 mono-architecture

### 40 Years of Microprocessor Trend Data
https://www.karlrupp.net/2015/06/40-years-of-microprocessor-trend-data/



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

# Fast Tracking on GPUs
## CMS Patatrack Example

# Utilizing multi-core Resources
## Multi-threaded Frameworks

### Advantages

- Significantly reduced memory usage

  - Most code components loaded only once

- Follow trend in hardware developments

  - Single core performance not increasing since ~10 years

  - Number of cores per machine constantly increasing
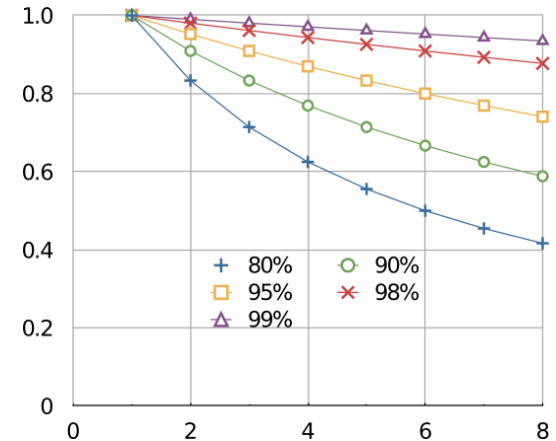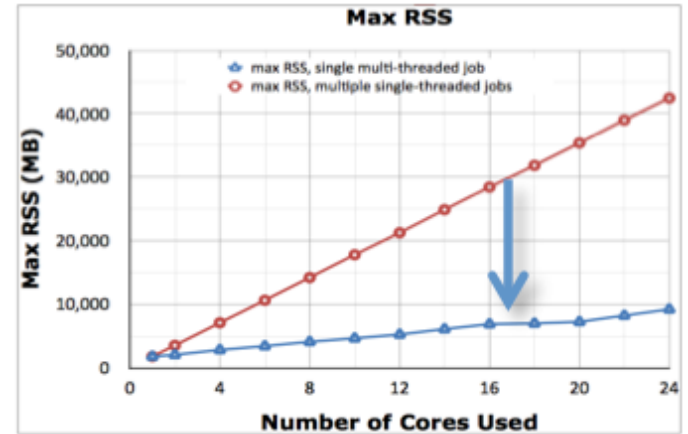
### Challenges

- CPU efficiency driven by fraction of thread-safe routines (Amdahl's law)

- Achieving expected (good) efficiency in the distributed infrastructure

  - Remote data access, big spectrum of workflows….

**CMS implementation based on Intel TBB in production**

- Can utilize GPUs, FPGAs transparently

**Efforts in other experiments ongoing (ATLAS slides in the backup)**

- Some approaches employ forking of single threaded processes

During the last years HEP used the assumption of a ~20% improvement rate for CPU and disk resources ($/HS06, $/GB) to extrapolate the future costing of computing equipment under the boundary condition of flat budgets.

CERN and several T1 sites now report deviations from the ~20% improvement rates.

It looks like the ~20% number is too optimistic and needs to be revised.

From: Bernd Panzer, Cern 2016

The WLCG Cost and Performance Modeling  working group (Markus) received numbers from several Tier 1 sites
and the they show a similar picture:   20% is too optimistic
The large variance of these numbers also point to strong site dependencies.


**Proposal for the future assumptions of cost improvements (just a starting point for discussions) :**

➢   **~10 % for CPUs**

➢   **~15 % for disk space**

➢   **~20 % for tape space (stays the same as before), BUT the future of tape per se is problematic !**
        **There are strong tendencies in the computing models to 'replace' disk space with tape**
         **needs very careful attention !**


**Need more input and discussions from the T1/T2 sites**
**Yearly adjustment of the figures !?**
**Weighted average !?**

From: Bernd Panzer, Cern (2018)

# Access to Cloud and HPC Resources

**Extending beyond 'classical' Resources**

**Classical resources likely not enough**

- Farm at host laboratory
- Grid sites for HEP

**Access "any" kind of resources**

- Clouds provided by institutes or commercially
- HPCs are special and each is different
  - Sometimes no outbound networking
  - Way to handle software/container

**Integration is often challenging**

- Interaction with data management and workflow management

**HEP is involved in a number of projects**

Examples!
(incomplete)



Working with commercial often has other than technical challenges

# So, how to gain back the E..?

1. Very easy solution: decrease some LHC/Experiment[...]
   parameters, like selection rate. **If 5x less data[...]**
   → **problem solved**
   – With a large price on Physics
   – It is like buying a Ferrari an[...]
     gasoline. Not too sma[...]
2. Try approaches[...]
   1. Be sm[...]
      [...]
      [...] **10M lines of code**
      [...]ns not present in today's
      [...]ity
      [...]s to be DeepLearning (training?)
3. Anyth[...]t the edge of technology?

Could QC be a solution to some specific parts of the problem????

# Is QC another "weapon" we should study?

- **Disclaimer: we are here mostly in the initial learning phase; our understanding of QC possibilities is not necessarily adequate**
  - **A very honest answer would be "we do not know yet"**

- Bird's eye evaluation:
  - **Quantum simulation** could in principle take the place of algorithmic generators, at least for some specific processes
  - **Quantum computing** could be used in principle for generic minimizations, or in order to speed up combinatorial algorithms
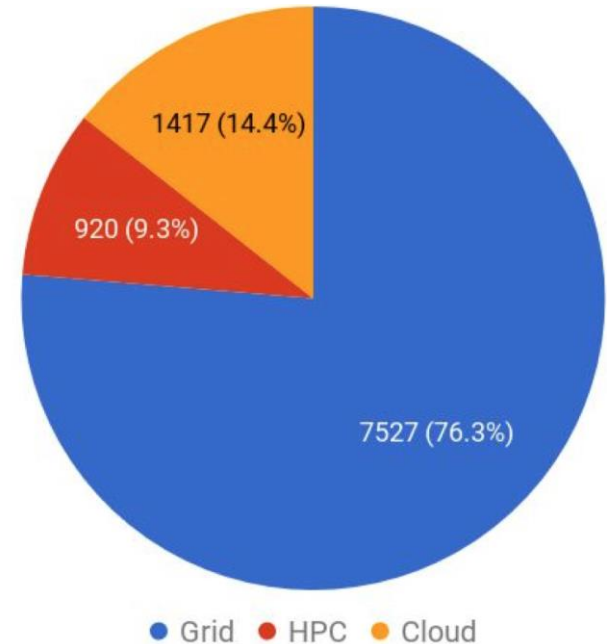    - Or in principle ANY algorithm via a Grover approach

# HPC Challenges

- Draft discussion document on challenges related to being able access and use large HPC
  - Policy & technical
- Working group on how to value HPC cycles for pledges and accounting
  - Very complex
  - Hand-in-hand with next round of benchmarking using suite of experiment codes
  - "HPC" here means GPU and non-x86
- Heading for a future where not all workloads will be efficient on some architectures → complexity and inefficiency
- In addition there is the software portability and sustainability challenge

# Opportunistic Resource Usage by ATLAS

- Opportunistic (=non-Grid) resources continue to play a significant role in ATLAS MC Production

- HPC is composed of specially prepared, dedicated jobs as well as running such resources as if they were additional grid sites

- Cloud is composed of jobs on volunteer computing and smaller clusters at institutes using `BOINC` as a lightweight submission mechanism as well as the **HLT (High Level Trigger) farm** at Point 1

- A new method of job submission is under development to allow tasks to **seamlessly** run with jobs on all resources: (grid, cloud, HPC,..)

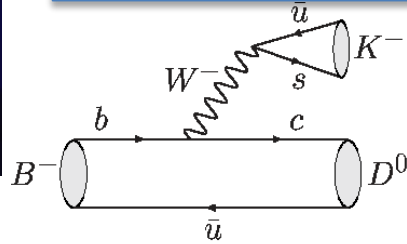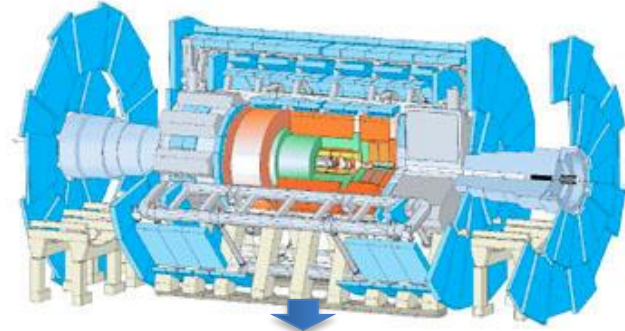M events per resource for 2017 (Full Simulation)



1417 (14.4%)

920 (9.3%)

7527 (76.3%)

● Grid  ● HPC  ● Cloud

# Software

# Reality



LHC collisions

Decay of unstable particles

ATLAS

Detector electronics

Trigger (selection)      SW

Reconstruction      SW

Analysis      SW

# Simulation - all SW

Theoretical model ("generators") → Simulation of decays of unstable particles → Simulation of interactions particle-detector → Simulation of detector electronics → Trigger Simulation → Reconstruction → Analysis



**GEANT4** A SIMULATION TOOLKIT

**LCIO**

**PYTHIA**

$$L_{QCD} = \sum_q \overline{\psi}_q \left( i \gamma_\mu D^\mu - m_q \right) \psi_q - \frac{1}{2} Tr \left[ \overline{G}_{\mu\nu} \overline{G}^{\mu\nu} \right]$$



- Data 2011+ 2012
- SM Higgs Boson $m_H$=124.3 GeV (fit)
- Background Z, ZZ*
- Background Z+jets, t$\overline{t}$
- Syst.Unc.

**ATLAS**
H→ZZ*→4l
√s = 7 TeV ∫Ldt = 4.6 fb$^{-1}$
√s = 8 TeV ∫Ldt = 20.7 fb$^{-1}$

Events/5 GeV

$m_{4l}$ [GeV]

# What are the typical algorithms doing?



- **Generation** is the simulation of a single particle collision, hence it has some modelling of a quantum system (be it via explicit matrix element calculation, or sequential steps, …)
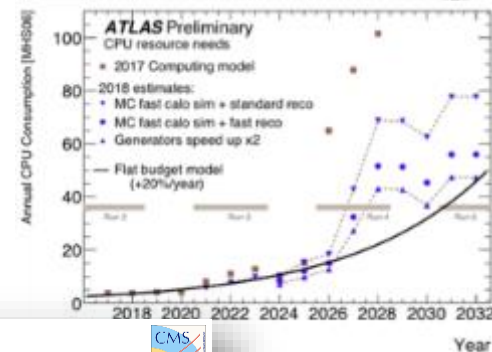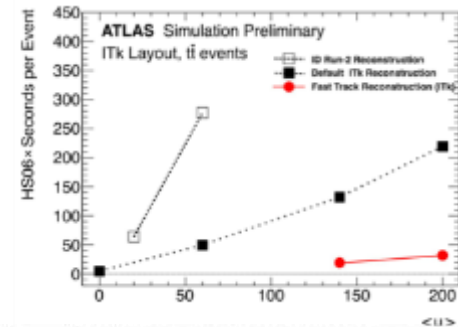  - Currently, done via approximations (perturbative orders, resummations, more and more loops and legs, …)
- **Simulation** in **Geant4** is mostly a transport problem, in which subsequent interactions particle/matter take place
  - Some of them only drive to energy loss, some others to decays / hard processes, …
  - The more the particles and the volumes (number, size), the more the time
- **Reconstruction** is an algorithmic problem, in general most of the time is spent in combinatorial algorithms (nested *for* loops)
  - Searching for doublets, triplets, quadruplets not atypical (N^2, N^3, N^4 …)
- **Analysis** is … anything!
  - In general, there is a selection step followed by a minimization (likelihood, …) step

# Software topics



- Several active HSF working groups
  - Event generators
    - Several workshops and meetings
  - Reconstruction and software triggers
    - Common topics: GPUs, real time analysis, links to other communities
  - Data Analysis working group
    - From DOMA to final analysis
    - Future analysis models, role of ML, etc.
  - Software frameworks
    - Just set up, conveners nominated
- Lots of work in experiments on software portability and performance
  - Use of HPC
  - Lots of work on tuning simulation; fast simulation (and where it is appropriate)
  - Performance and portability:
    - Adaptation of frameworks to accommodate heterogenous code (CPU+accelerators)
    - Portability libraries: Kokkos, Alpaka, SYCL, etc
      - Can there be one codebase for all architectures?



Software Portability

- Use same codebase for multiple backends (CPU, GPU, FPGA, …)
- Ongoing study of solutions (Kokkos, Alpaka, SYCL)
- Need to gain more experience to make sensible choice
- Collaboration with ATLAS and HSF

alpaka
kokkos
SYCL

14

43

# Processing needs by workflow - 2018

- Generators range between 1% and 10% of the total CPU needs;
  - Difference depends on the perturbative level (LO, NLO, NNLO), different choices on the market, ...
- Geant4 is currently the most demanding application
  - CMS: ½ of the CPU time for a simulated event
  - ATLAS: >50%
- Physics Object Reconstruction is 30-40% of the CPU budget
- Analysis depends critically on the experiment decisions
  - Some 10-30% of the overall budget

- But scaling with event complexity (so to 2026) is largely different

**Generators and Geant4 do not scale with LHC luminosity; the total time scales with the # of events to be processed**

**Reconstruction scales with the # of events processed, and scales more than linearly with the LHC luminosity**

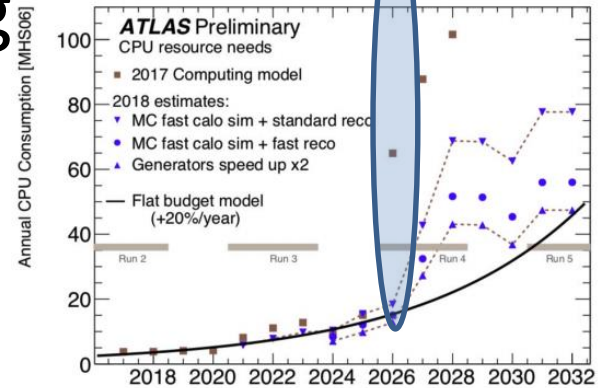**Analysis scales with the # of events, and mildly with their complexity**

# Simulation

- Is a major cost driver (~50% total computing cost)
- Long term supportability/portability/performance is essential
  - Must ensure code modernization & long term supportability, adaptability to changing computing landscape,
    - In a sustainable way - Not as one-off to e.g. GPU-version-x
  - Lot of effort in the world on portability to new architectures
  - Need a major effort on simulation for the future to tie in all of these R&D efforts
  - This is going to be a many-year effort
- This is where we really need to invest effort in the future
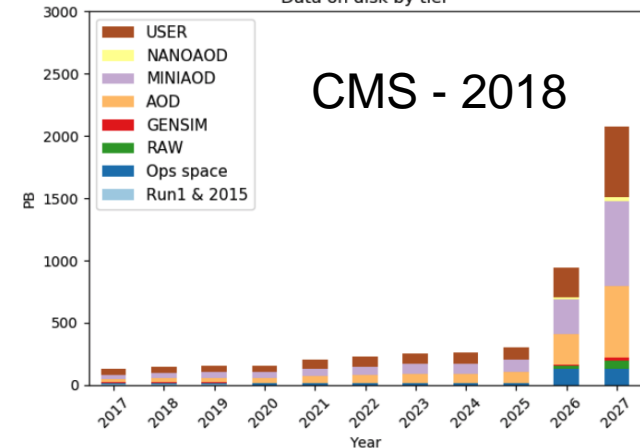  - And is a significant opportunity

# Clearly we can improve processing



CPU projections for HL-LHC

- Organize data in a better way
- Do less mistakes (avoid re-processing of the data, for example)
- use ML-driven algorithms to speed up reconstruction
- Slow down the offline system: publish (many years?) later …

- Some of these already implemented into experiments computing models: still, factors to go
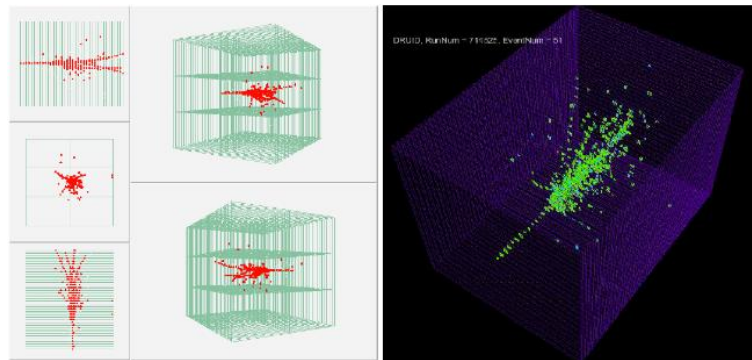  - Around 5x, probably



CMS - 2018

# AMALEA

**Helmholtz Innovation Pool Project**
**ML techniques for HEP, Photon Science and Accelerators**
**Sustainable Infrastructure Hardware and Software**
**Broad field of developments**

- Fast simulation and reconstruction for 3D images

- ultra-fast feedback algorithm for data reduction, compression and classification

- fast diagnosis and control systems



**Fast -but detailed- simulation of showers**

- Try generate shower images with Wasserstein GANs

- Order(s) of magnitude faster than classical particle propagation

- Detailed studies need to achieve competitive implementation

# DATA

# Rucio: A cross-community Tool for Scientific Data Managment?

**Rucio originally developed in ATLAS for LHC Run2**
**Operates on top of FTS3 (File Transfer Service)**
**Organization of files in Datasets or Containers**
**Policy engine**
**Manages ~200PB of ATLAS data**
**Selected as data management tool for other experiments:**

- CMS, Xenon1T

**Very likely used by**
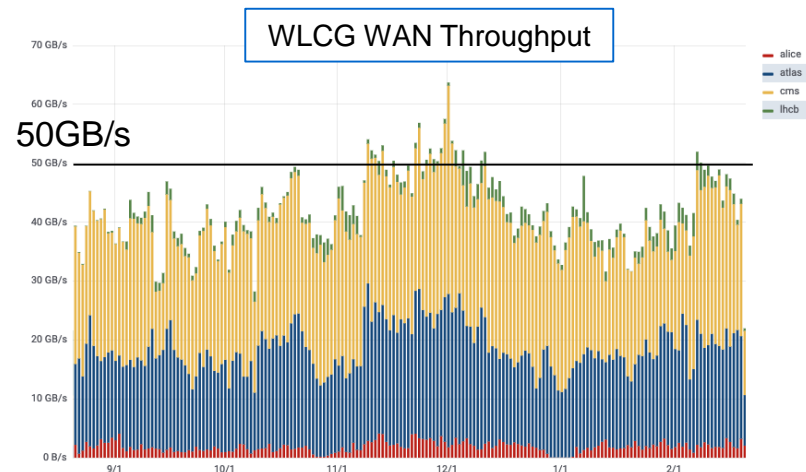
- Dune, Belle-II, IceCube, CTA

**Evaluated also by non-(astro)-particle groups**

- SKA, LSST, NSLS-II, LCLS-II

Astro-physics

Photon-science !!



WLCG WAN Throughput

50GB/s

# DOMA in a nutshell

**DOMA project**
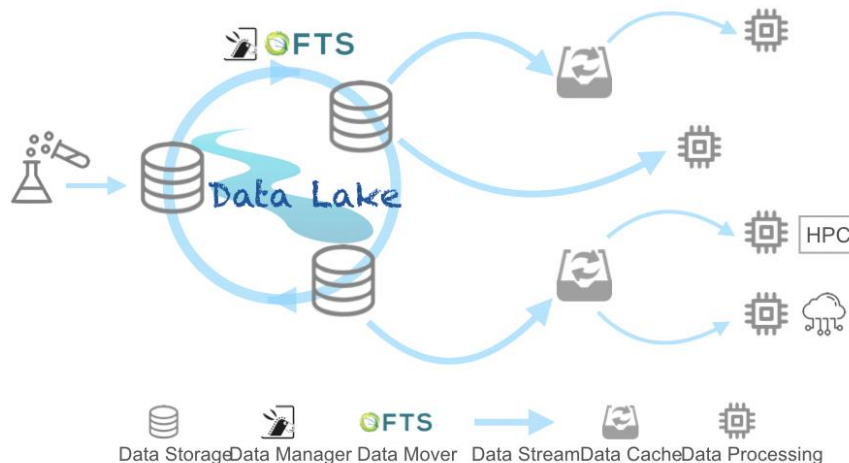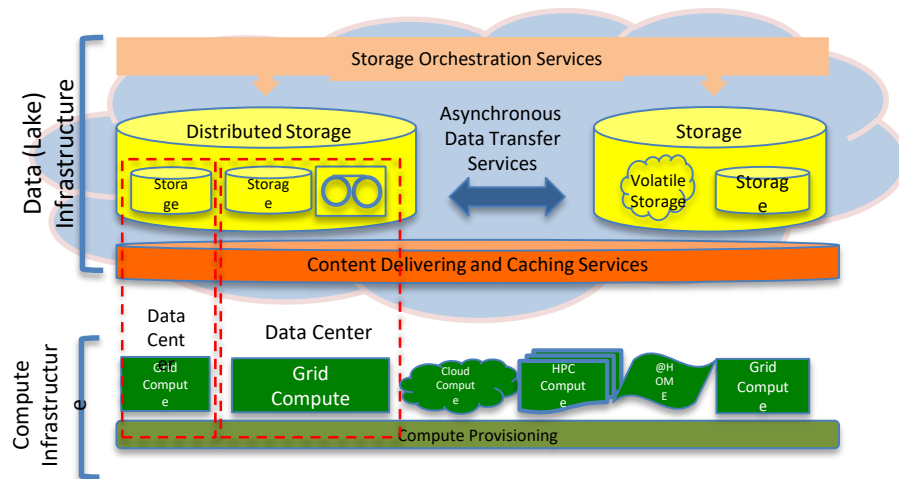(Data Organization, Management, Access)

A set of R&D activities evaluating components and techniques to build a common HEP data cloud

Three Working Groups

- ACCESS for Content Delivery and Caching
- TPC for Third Party Copy
- QoS for storage Quality of Service

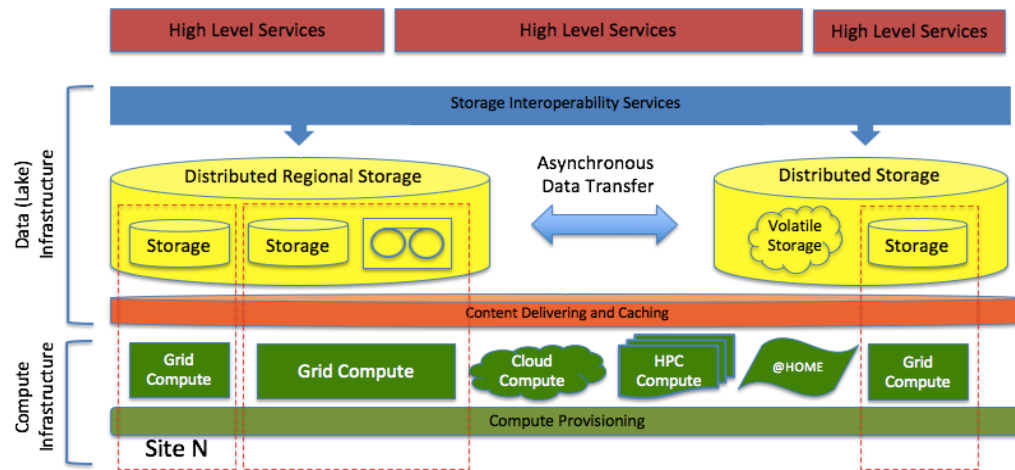And many activities, reporting regularly

https://twiki.cern.ch/twiki/bin/view/LCG/DomaActivities



From Simone Campana @ LHCC 10/09/19

# Data management ("data lake")

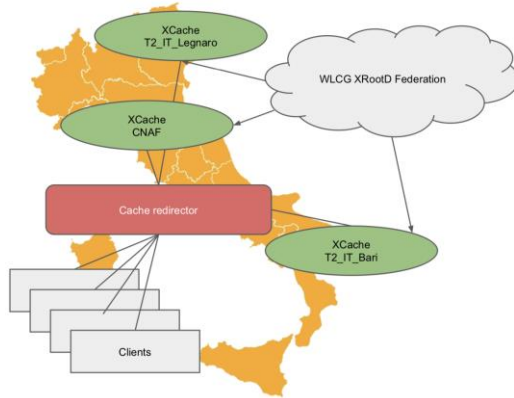Data Organisation, Management, Access (DOMA)

- Several activities and working groups
  - Storage consolidation
  - Caching and data access
  - Data transfer and access protocols:
    - 3$^{rd}$ party copy
    - Replacement of gridftp
  - Quality of Service
    - Performance/reliability vs capacity
    - Use of high-performance storage?
  - Use of networks and Investigation of low level protocols and optimization of data movement (with SKA, Geant, others)
    - Between parts of the data lake
    - Serving data
- A prototype "data lake" has been set up and can be used to explore technology and R&D questions
  - Several Tier 1s participating in the prototype



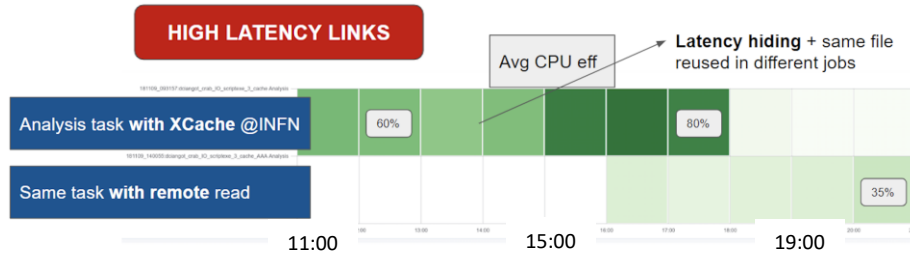- Idea is to localize bulk data in a cloud service (➔ data lake): minimize replication, assure availability
- Serve data to remote (or local) compute – grid, cloud, HPC, etc.
- Simple (unmanaged) caching is all that is needed at compute site
- Works at national, regional, global scales

# ACCESS: caching layer prototype

A distributed caching system in INFN



From Simone Campana @ LHCC 10/09/19

# Other active areas

- System performance and cost modelling
  - Very active group
  - Overall system optimization – detailed studies
  - Guidance on how to optimize costs at a site
  - This group can inform a lot of other work related to overall system design and optimization
- Technology and market tracking
  - Our cost estimates depend strongly on how technology evolves
  - Provide regular updates of cost evolution and likely technology directions
- Compute provisioning and access
- Open Access, open data
- Data preservation
- …

- AAI
  - Move to more modern token-based schemes (for end-users at least)
  - Lots of activity – WLCG, EC projects, OSG

# Towards a Computing TDR

**Goal:** WLCG Computing TDR recommended for approval to the LHCC by early 2024

*Initial meeting* in **May 2020** will focus on experiment specific issues (proposed: May 18-20)
- Charge to be delivered to WLCG management by Dec 6, 2019
- Establish a baseline computing model, data rates, computing and storage projections including the roles of the Tiers
- Establish anticipated cost drivers and infrastructure assumptions
- Outline technological risks and major areas of R&D

*Second meeting,* preliminary target **September 2021**, will focus on common tools and community software (Examples include Root, MC Simulation, Event Generators )
Mid 2022 begin formal TDR preparation

Referees comments:
Progress in studies shows that its the right moment now for the next step towards computing TDR

Referees recommendations:
Don't separate the discussion of experiment specific issues from common tools and community software

# This is the current masterplan

- Try preferentially to explore solution not impacting physics and not requiring more money

- Use the 8 years from now to 2027 to
  - Be prepared to use heterogeneous computing architectures, allowing to
    - Use the best performance/price ratio at any moment, following market
    - Enlarge the basis of potential resources (more HPC centers, more farms, more clusters, …)
  - Better understand analysis models, and reduce the needs for MC, processings, calibration steps, …

- Is this enough?
  - **Who knows for sure …**
  - In the communities, you can feel a mild optimism though …

# Summary

**Big computing challenges ahead for research in the program Matter**

- Needs for storage, compute capacity and network bandwidth growing by order of magnitude

- Approaches of today often not applicable

- Waiting for technology improvements will not be sufficient
**Recent trends from Industry could bridge a fraction of the gap**

- New specialized hardware architectures: GPUs & FPGAs, QC(?), Mem Driven, TPU's, …

- Opportunities to utilize (spare cycles of) HPCs or clouds

- Novel approaches in algorithms and methods

  – Machine learning has a huge potential
**A number of efforts already ongoing in the groups**

- Further cooperations with other communities, eg on SW developemnt, training etc.

# Data preservation, open access

- Currently many strands to these activities
  - Often started independently with different goals and interests behind them
- Main topics:
  - Data preservation – building on work of DPHEP
    - Several aspects from bit preservation, DC certification, metadata, knowledge retention, etc.
    - Bit preservation is what we already do at Tier 0, 1
  - Open data:
    - For example via open data portal
    - Experiments have different policies of what level, and how much data is made available; and for different intended purposes
    - This has a cost – today CERN provides ~5 PB disk for this – but growing
      - Currently not costed as part of pledges
  - Analysis preservation
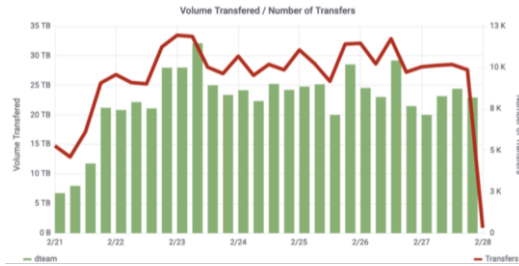    - E.g. via tools such as REANA, etc.

# Open access -  Concerns

- Active groups in IT, EP, experiments on all of these – but not necessarily coordinated
- No policy for how this should be funded
  - Today it is essentially CERN
    - Should it be a shared/distributed problem?
    - ESCAPE can provide a mechanism for this shared management (data lake)
- Scale and cost
  - If the scale increases it will need to be taken from the pledges
- Today cannot draft an overall coherent policy
- Propose to organize a workshop to address and coordinate these aspects
  - And to formulate a strategy for how this should be managed
  - This should become part of the overall strategy for the future and integrated with the other aspects
- Need feedback from funding agencies on what is mandated and affordable

# TPC

Goal: commission non-gridFTP protocols for asynchronous data transfer (Third Party Copy)

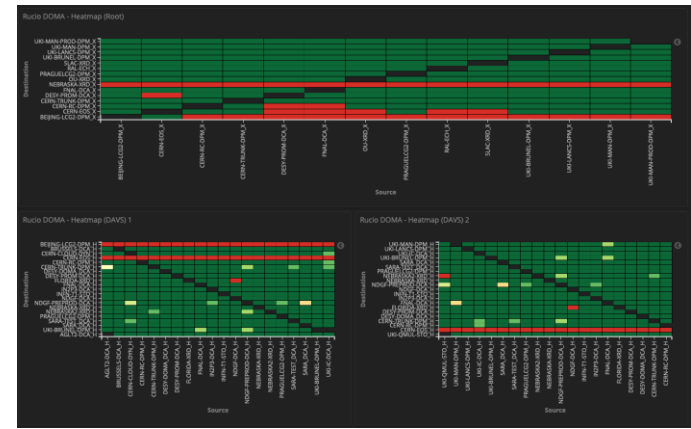- Phase-2 (deadline June 2019): all sites providing > 3PB of storage to WLCG should provide a non-gridFTP endpoint in production



Functional and Stress testing





Capable to fill available bandwidth

Point-point functional testing

- ## Phase-3 (Dec 2019): all sites to have a non-gridFTP endpoint

  NB: some features needed for TPC are available only in recent versions of storage

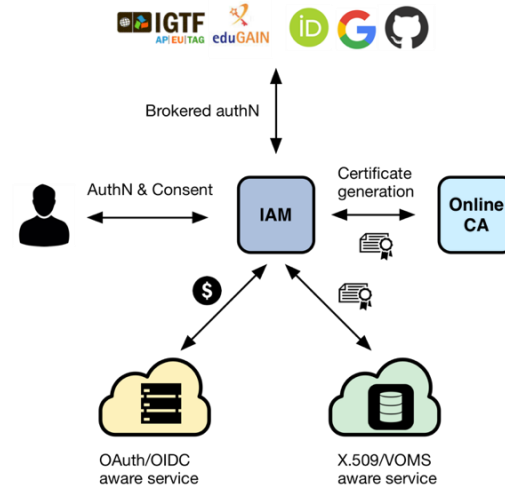From Simone Campana @ LHCC 10/09/19

# TPC and AAI

WLCG is planning to evolve AAI toward token based Auth/AuthZ and Federated Identities

The WLCG task force is finalizing the token profile as last item

While this is has a much broader scope than DOMA, TPC offers a well confined use case to start with

Rucio is integrating tokens. Storage is preparing to manage them.



From Simone Campana @ LHCC 10/09/19

# Process updated

- Strategy document delivered in May
- Discussed with LHCC
- Working groups active in many areas
- LHCC will organize a review of the strategy during 1H19 (tbd)
- Would also propose an update of the strategy document following the review
- Agreed that TDR for computing would be then on a timescale of 2022
  - Earlier does not make sense, and a review is a good checkpoint
  - There would be a general TDR with complementary experiment-specific documents
- Intend to provide ~yearly updates of estimated requirements vs anticipated budget
  - To show convergence

## Modernisation and Tuning of Software

▶ **Modernise and improve the performance of the CMS sw stack**

**Why?**

▶ Accommodate within computing resources an **ambitious Run3 Physics program**
▶ **Be ready for Run4**
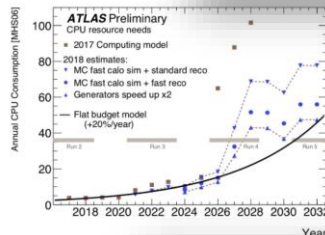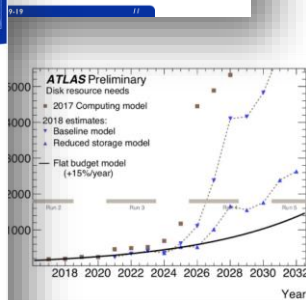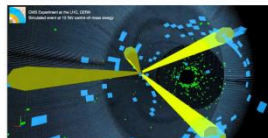  • Use Run3 also to test solutions targeting Run4

**How?**

▶ **Optimisations**: technical (e.g. compiler flags) and algorithmic in CPU code
▶ **Size reduction** of AOD(Sim) and RAW on storage media
  • E.g. compression settings/algorithm, precision, content of tiers, row Vs. columnar storage
▶ **Accommodate in CMSSW heterogeneous code**, i.e. CPU + accelerator (e.g. GPU)
  • Evolve respecting present CMSSW architecture
  • Identify the right tools for **performance portability: one codebase for all architectures**
  • Start from framework and high level trigger code

---

## Showcase: Improvements in Simulation

▶ Continuing efforts to improve performance of FullSim

  • **Preliminary result: 20-30% runtime reduction possible for Run3**

▶ Several elements to achieve this success:

  • Switch from Geant4 10.4 to 10.6
  • Tune energy-dependent propagation through EM fields (*smart tracking*)
  • Optimize usage of the VecGeom library

▶ Investigating technical solution to run simulation efficiently on HPCs with accelerators

---

## Upgrading the O&C Software Toolset

▪ Sustainability of software tools on the Run 4 timescale is a concern
▪ Strategy: turn to common solutions, put in production the products already for Run3

**CRIC** Computing Resource Information Catalogue (used by Atlas et al.)

▶ Access physical and CMS logical computing resources
▶ Replace Information System
▶ Already there

**DD4HEP** (used by ILC/CLIC, evaluated by LHCb)

▶ Detector description tool, EU financed (AIDA 2020)
▶ Review and optimize current detector description too!
▶ Steady progress, replacement planned next year

**Rucio** (originated in Atlas, rapidly growing adoption!)

▶ Data management solution replacing Phedex / Dynamo
▶ Steady progress, looking for power users this fall
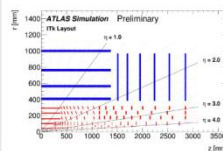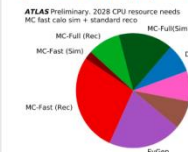▶ One big step forward: transfer ownership of NanoAOD to Rucio

Potential mitigation of costs and improved sustainability: common solutions with industry and other experiments
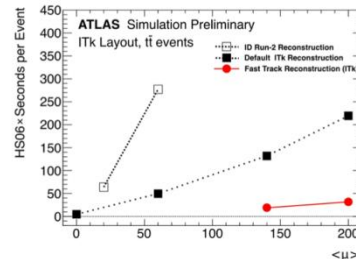
---

## Running on GPUs and other accelerators for HL-LHC

• Resources available to us on acceleration hardware (opportunistic)
  ○ Much debated topic in the past year or so
  ○ Complementary to non-x86 resouces that could be used by recompiling the SW stack

• Prototype of Fast Simulation in CUDA
  ○ Self-contained kernel, collaboration with computing scientists
• ACTS module for GPUs in initial design phase (IRIS-HEP)
• Cross experiment initiatives:
  ○ Prospects of running event generators on GPUs.
  ○ Geant4 GPU kernels (?)

• Focus on frameworks for running on heterogenous resources
  ○ Two ATLAS senior developers charged with accelerators R&D
  ○ Current prototypes in CUDA running on NVIDIA GPGPUs.
    ▪ Issue with sustainability, code duplication and validation
    ▪ How practical is this outside Online or other contained environments?
    ▪ How do we keep both CPU and GPU busy?
  ○ We need to focus on **portability**, kokkos, SYCL
  ○ Not all HEP code suitable for GPUs
  ○ The technology is evolving (We will soon evaluate Intel's OneAPI beta)

---

## Tracking reconstruction improvements

Need to speed up reconstruction at high <mu>
Optimised tracker (ITk) with x10 more channels

• Optimised track selection
• Improved seeding algorithm (for ITk)
• Omission of ambiguity resolving (to be partly recovered by the new fitter)

# New Approaches: Deep Learning

- A full spectrum of "new" tools and libraries

- Change in "software culture"
    - HEP used to primarily use own grown tools: ROOT, GEANT, PAW…
        - Specially developed for HEP needs
    - Most modern Deep Learning tools come from data science industry
        - Developed by internet industry with billion dollar investments
        - Many software released as open source – *make money with data not with software*
        - Development priorities clearly decided outside the scientific community

- Already promising results using the new tools

# Machine Learning at Belle-II Example

**Belle II in Japan**



- **intensity frontier** flagship experiment at KEK
- precision measurement and (extremely) **rare B-decays**
- Almost any analysis employs machine learning!

## ECL cluster shape calibration



*more training*

*Fake (WGAN, Belle II)*

*Fake (WGAN, Belle II)*

*Fake (WGAN, Belle II)*

*Real (Belle II)*

Semi-supervised learning: Wasserstein GAN learns to create 'fake' images that look like real Belle II images.

Example: E1oE9 shower shape variable

use of a WGAN for generating *images of calorimeter showers needed to correct imperfect 'conventional' simulation*

# Summary

- LS2 is busy for the experiments & facilities
  - Ongoing processing, analysis, etc.
  - Preparations for Run 3 – simulations, software preparation, etc.
- Run 3 looks like an evolution of Run 2 for ATLAS and CMS
  - LHCb & ALICE major changes – but sw & computing preparations in hand
  - Resource outlook seems realistic
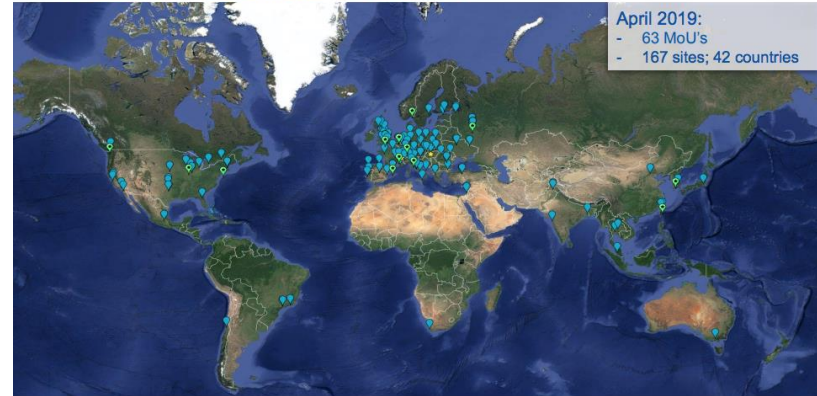- Data preservation , open access workshop to be held 26 Nov
  - Initial discussion to align and agree policies, strategies, goals, and resource needs
- Further outlook to HL-LHC
  - Many R&D topics progressing well
  - Significant work in experiments closing the gap between requirements and likely resources
    - Although the cost evolution of hardware is a major concern
  - Software challenges are potentially significant – but are opportunities for the longer term sustainability
- LHCC will hold a review of HL-LHC computing preparations in ~Spring 2020

# What is High Luminosity LHC (HL-LHC?)

**WLCG Collaboration**



April 2019:
- 63 MoU's
- 167 sites; 42 countries

- **4 PB/s/exp** is clearly unfeasible, hence the need for complex selection / suppression / compression algorithms
  - Various level triggers: hardware, software, …
  - LZMA compression, Zero Suppression, …
- Current data rates to offline ~ 1-3 GB /s
  - 1-3 kHz of O(1MB/ev) events
- Together with the LHC livetime (~7 Ms/y) this drives the computing requests
  - Collect ~ 10 PB/y of RAW data
  - You need at least as much MC simulation
  - You need to process both («CPU») to provide physicists with predigested samples

- **Well, it worked!**

1 «xeon-core» ~ 10 HS06
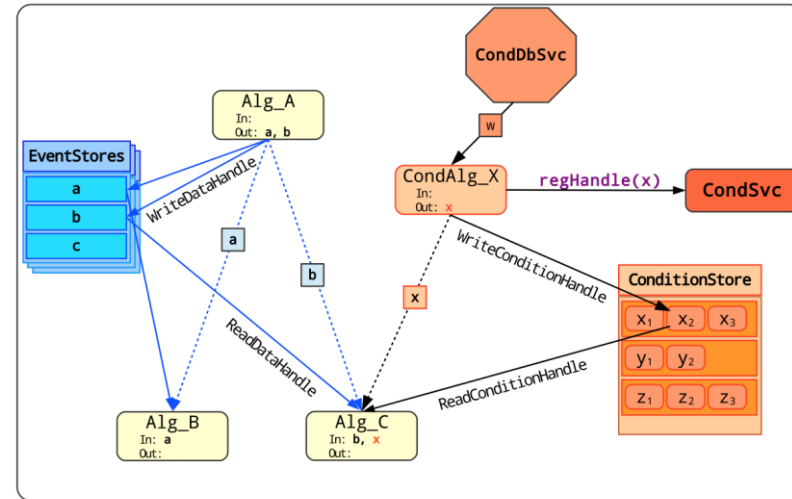
| Experiment | CPU (kHS06) | Disk (PB) | Tape (PB) |
|---|---|---|---|
| **ALICE** | 1000 | 100 | 85 |
| **ATLAS** | 2800 | 230 | 310 |
| **CMS** | 2000 | 160 | 280 |
| **LHCB** | 450 | 45 | 90 |
| **TOTAL** | **6250** | **535** | **765** |

# Multi-threaded Application
## ATLAS Experiment

> The Athena/Gaudi Atlas software framework was designed with serial processing in mind, one event at a time, on one thread, essentially using a single-core

> Emerging technologies require a concurrent, multi-threaded approach to be adopted, which is the aim of **AthenaMT**

> So called "*Shared Software Services*" such as conditions, that must be madethread-safe and be able to simultaneously process requests from different events, in an asynchronousdata stream

> One solution employed for example in accessing data conditions is to access data using smart references or **ConditionsHandles**, which store information pertinent to multiple data ranges in dedicated containers

# Software-related aspects

The software challenges are key to addressing the current mismatch between requirements and affordability (given expected technology)

- HSF – several ongoing activities called out in the CWP, addressing performance
  - New sub-groups being set up to work on:
    - Detector simulation;  Reconstruction and software triggers;  Data analysis
  - Workshop on physics generators and related computing challenges organized in November
- Many experiment-specific investigations on core software and key topics such as data models and data formats
- NSF funding awarded to IRIS-HEP project (the CWP was part of the proposal process)
  - $25 M over 5 years for a "software institute" to work on core software for HL-LHC
  - Institute for Research and Innovation in Software for HEP (IRIS-HEP)
- Important to understand that this requires community-wide investment in software
  - HSF was a good start
  - CWP outlined the problems
  - Many opportunities for funding …

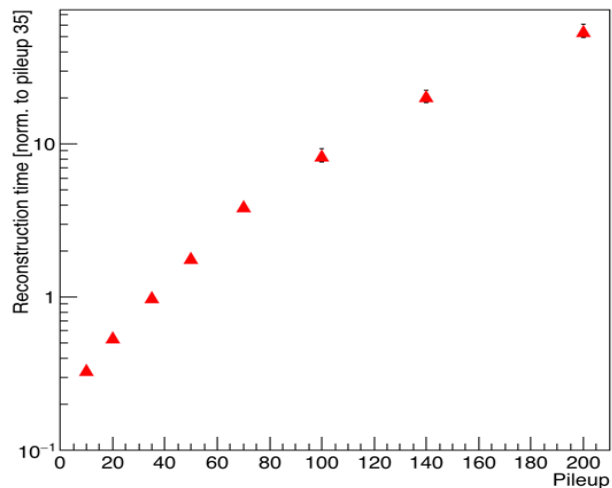# Resource Needs for HL-LHC

**Big change in computing needs for Run4**

- Logging from 1kHz → 7.5 (or 10)kHz: ~10x events to store and process
- Pileup ~35 → ~200: Reconstruction time increases by ~10-15
  - Number based on 2017 detector
- Event size increases ~10 times (for the same data tier)

Need about 100x more Computing in 2027

**There will be some technology improvement**

- How much?
  - Likely not to compensate for a factor 100
- Where to find the missing factors?



Reconstruction of t$\bar{t}$ at $\sqrt{s}$ = 13 TeV with CMS 2016 configuration

# Putting all together ...

- If your goal is to have **10.000.000 produced Higgs in 5 years** (per experiment)
- $L_{int}$ = 100 fb$^{-1}$ ($10^7$/(10000fb))  and then, scaling to the instantaneous lumi (assuming an efficiency factor ~5 for shutdown periods, vacations, repairs, etc)
- $L_{int\_max}$  = 100 fb$^{-1}$
- If you remember that 1 b = $10^{-24}$ cm$^2$ → $L_{int}$  = $10^{42}$ cm$^{-2}$

$$L_{INST} = 5 * 10^{42} \text{ cm}^{-2} / (5 \text{ y} * 3 * 10^7 \text{s/y}) =  O(10^{34}) \text{ cm}^{-2} \text{ s}^{-1}$$

- **SO: the extreme LHC parameters are the only way to "guarantee" LHC would have been able to discover / exclude the Higgs boson in the energy range where we were searching for him.**
- **Any machine with lower parameters could have not been able to close the issue on the Higgs (if you want, not well spent money)**
- **But: the very same parameters drive to the data flux O(PB/s) → we have a computing problem!**

# Multi-threaded Application

## ATLAS Experiment

> Another example in the multi-threaded approach is the asynchronous call-backs to the **IncidentService**, which registers calls such as "*BeginEvent*", "*OpenFile*" and so on



> Calls to the **IncidentService** outside of the event execution loop are now made schedulable, so call-backs are correctly executed

> Algorithms in **AthenaMT** are by design thread-safe in that they only process a single event, and whilst concurrent processing can be achieved by producing multiple instances this results in an calculable increase in memory requirements

> Therefore, new "re-entrant" algorithms are being developed, which so long as they are thread-safe and stateless (e.g employing **const** methods) may be executed concurrently in multiple events

# HL-LHC Parameters

| | Nominal LHC (design report) | HL-LHC 25 ns (standard) | HL-LHC 25 ns (BCMS) | HL-LHC 50 ns |
|---|---|---|---|---|
| n collision [TeV] | 7 | 7 | 7 | 7 |
| | $1.5 \times 10^{11}$ | $2.2 \times 10^{11}$ | $2.2 \times 10^{11}$ | $3.5 \times 10^{11}$ |
| | 2808 | 2748 | 2604 | 1404 |
| lisions in IP1 and IP5 | 2808 | 2736* | 2592 | 1404 |
| | $3.2 \times 10^{14}$ | $6 \times 10^{14}$ | $5.7 \times 10^{14}$ | $4.9 \times 10^{14}$ |
| [A] | 0.58 | 1.09 | 1.03 | 0.89 |
| e [μrad] | 285 | 590 | 590 | 590 |
| on [$\sigma$] | 9.4 | 12.5 | 12.5 | 11.4 |
| | 0.55 | 0.15 | 0.15 | 0.15 |
| | 3.75 | 2.50 | 2.50 | 3 |
| | 2.50 | 2.50 | 2.50 | 2.50 |
| pread | $1.13 \times 10^{-4}$ | $1.13 \times 10^{-4}$ | $1.13 \times 10^{-4}$ | $1.13 \times 10^{-4}$ |
| ngth | $7.55 \times 10^{-2}$ | $7.55 \times 10^{-2}$ | $7.55 \times 10^{-2}$ | $7.55 \times 10^{-2}$ |
| [h] | 80–106 | 18.5 | 18.5 | 17.2 |
| al [h] | 61–60 | 20.4 | 20.4 | 16.1 |
| neter | 0.65 | 3.14 | 3.14 | 2.87 |
| factor $R_0$ without crab cavity | 0.836 | 0.305 | 0.305 | 0.331 |
| factor $R_1$ with crab cavity | (0.981) | 0.829 | 0.829 | 0.838 |
| P without crab cavity | $3.1 \times 10^{-3}$ | $3.3 \times 10^{-3}$ | $3.3 \times 10^{-3}$ | $4.7 \times 10^{-3}$ |
| P with crab cavity | $3.8 \times 10^{-3}$ | $1.1 \times 10^{-2}$ | $1.1 \times 10^{-2}$ | $1.4 \times 10^{-2}$ |
| ty without crab cavity [cm$^{-2}$ s$^{-2}$] | $1.00 \times 10^{34}$ | $7.18 \times 10^{34}$ | $6.80 \times 10^{34}$ | $8.44 \times 10^{34}$ |
| sity with crab cavity, $L_{peak} \times R_1/R_0$ | $(1.18 \times 10^{34})$ | $19.54 \times 10^{34}$ | $18.52 \times 10^{34}$ | $21.38 \times 10^{34}$ |

| | | | | |
|---|---|---|---|---|
| Events/crossing without levelling and without crab cavity | 27 | 198 | 198 | 454 |
| Levelled luminosity [cm$^{-2}$ s$^{-2}$] | - | $5.00 \times 10^{34\dagger}$ | $5.00 \times 10^{34}$ | $2.50 \times 10^{34}$ |
| Events/crossing (with levelling and without crab cavities for HL-LHC) | 27 | 138 | 146 | 135 |
| Peak line density of pile-up event [event/mm] (maximum over stable beams) | 0.21 | 1.25 | 1.31 | 1.20 |
| Levelling time [h] (assuming no emittance growth) | - | 8.3 | 7.6 | 18.0 |
| Number of collisions in IP2/IP8 | 2808 | 2452/2524$^{\ddagger}$ | 2288/2396 | 0$^{**}$/1404 |
| $N_b$ at SPS extraction$^{\dagger\dagger}$ | $1.20 \times 10^{11}$ | $2.30 \times 10^{11}$ | $2.30 \times 10^{11}$ | $3.68 \times 10^{11}$ |
| $n_b$/injection | 288 | 288 | 288 | 144 |
| $N_{tot}$/injection | $3.46 \times 10^{13}$ | $6.62 \times 10^{13}$ | $6.62 \times 10^{13}$ | $5.30 \times 10^{13}$ |
| $\varepsilon_n$ at SPS extraction [μm]$^{\ddagger}$ | 3.40 | 2.00 | <2.00$^{***}$ | 2.30 |

# Key topics identified

- Software improvements
- Algorithmic improvements
- Event generators
- Reduce Data volumes
- Managing operations costs
- Optimizing  HW cost

**ESFRI Science Projects**

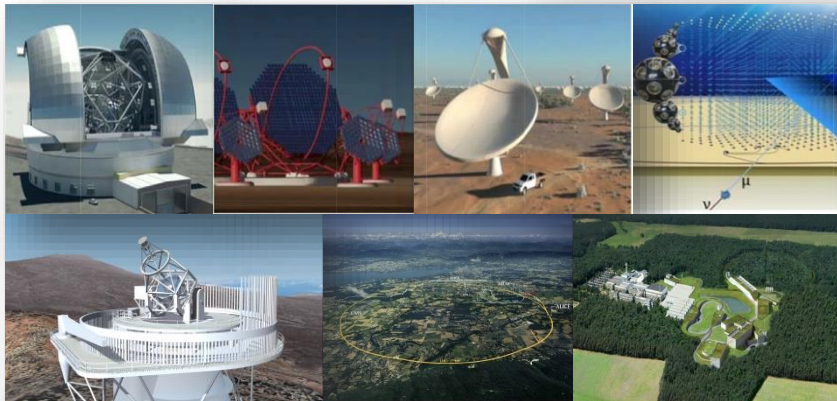| | |
|---|---|
| HL-LHC | SKA |
| FAIR | CTA |
| KM3Net | JIVE-ERIC |
| ELT | EST |
| EURO-VO | EGO-VIRGO |
| (LSST) | (CERN,ESO) |

# ESCAPE

**European Science Cluster of Astronomy & Particle physics
ESFRI research infrastructures**

**Goals:**
Prototype an infrastructure for the EOSC that is adapted to the Exabyte-scale needs of the large ESFRI science projects.

Ensure that the science communities drive the development of the EOSC.

Has to address *FAIR* data management, long term preservation, open access, open science, and contribute to the EOSC catalogue of services.

**Work Packages**
WP2 – Data Infrastructure for Open Science
WP3 – Open-source scientific Software and
　　　 Service Repository
WP4 – Connecting ESFRI projects to EOSC through
　　　 VO framework
WP5 –  ESFRI Science Analysis Platform

**Task 2.2 Content Delivering and Caching**

**Task 2.2 Storage Orchestration Service**

**Task 2.1 Storage Services**

**Task 2.1 Data transfer services**

**Task 2.3   Efficient Access to Compute**

HTC/Grid

HPC

Cloud/
commercial

citizen

**Task 2.4 Networking**

**Task 2.5 AAI**

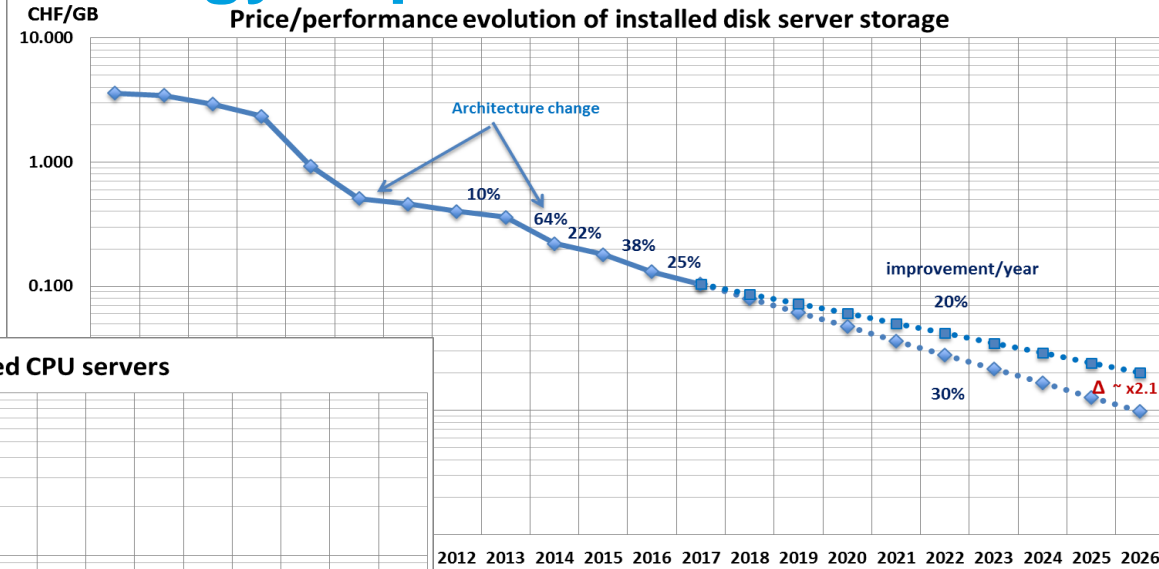**Data centres** (funded in WP2)
CERN, INFN, DESY, GSI, Nikhef, SURFSara, RUG, CCIN2P3, PIC, LAPP, INAF

# Extrapolation of Technology Improvements
## CERN Study

**Constant budget**

- ~15-20% increase per year

- Gain a factor ~5-6 in 10 years
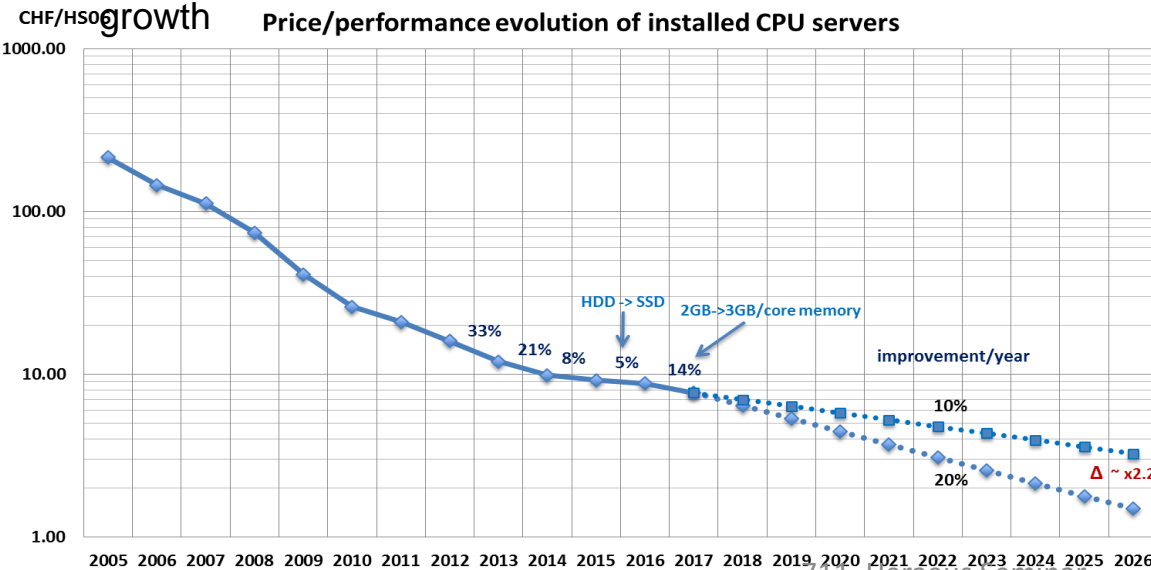
- Recent studies indicate lower growth



Price/performance evolution of installed disk server storage



Price/performance evolution of installed CPU servers

Material by Bernd Panzer-Steindel (CERN)
*Computing Evolution: Technology
and Markets, Jan-2017*
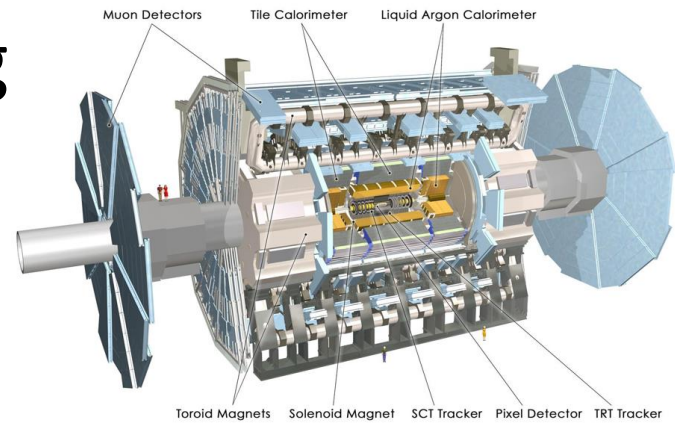https://indico.cern.ch/event/570249/contributions/2404412/

# Run 3 expectations

- Bunch intensities ramp up from 0 to 1.4e11ppb over the year
    - with limited availability of the injectors/LHC resulting in **only 20% machine efficiency**.
- For <u>contingency</u> planning, the machine efficiency assumed to reach normal value of 50%. This results in the following luminosity envelope:

|  | Baseline | Upper limit |
|---|---|---|
| ATLAS / CMS | 17 fb$^{-1}$ | 42 fb$^{-1}$ |
| LHCb | 3 fb$^{-1}$ | 7 fb$^{-1}$ |
| ALICE | 36 pb$^{-1}$ | 90 pb$^{-1}$ |

- ❑ NB. The upper limit is **contingency planning only** (i.e. raw data tape storage), **not physics**.
- ❑ Pb-Pb assumed to be a full production year: >2 nb$^{-1}$ for ATLAS, ALICE and CMS.
- ❑ 2022: We assume a full production year with 1.5 x 2018 resource levels

> ❑ To be updated once running conditions better specified (End Nov)
> - In particular different assumptions on pileup (55 vs 45) will make noticeable difference

# CMS and ATLAS computing scaling @ HL-LHC

- **# events collected/y** = Experiment live time * Experiment rate to offline
  - LHC RunII: 7 Ms/y * 1000 Hz = ~ 7 B events/y
  - LHC RunIV: 7 Ms/y * 7.5 kHz = ~ 50 B events/y
- **Bandwidth, total storage** = # events collected * $(1+ f_{MC})$ * typical_event_size
  - $f_{MC}$ ~ 1-2
  - Typical event size:
    - LHC RunII: 1 MB/ev
    - LHC RunIV: 5-10 MB/ev

~ 7.5 * 10 → O(50-100)x for storage

- **Computing power** = # events collected * $(1 + \alpha*f_{MC})$ * F(event_complexity)
  - F(event_complexity) usually superlinear in instantaneous luminosity
  - $\alpha$: how much more expensive is to process a simulated events with respect to a real data one. $O(2) < \alpha < O(20+)$
- Storage is also ~ integral with time
- $Storage_{YearN+1} = Storage_{YearN} + Delta_{NEW\ EVENTS}$

~ 7.5 * 10 → O(50-100)x minimum for CPU