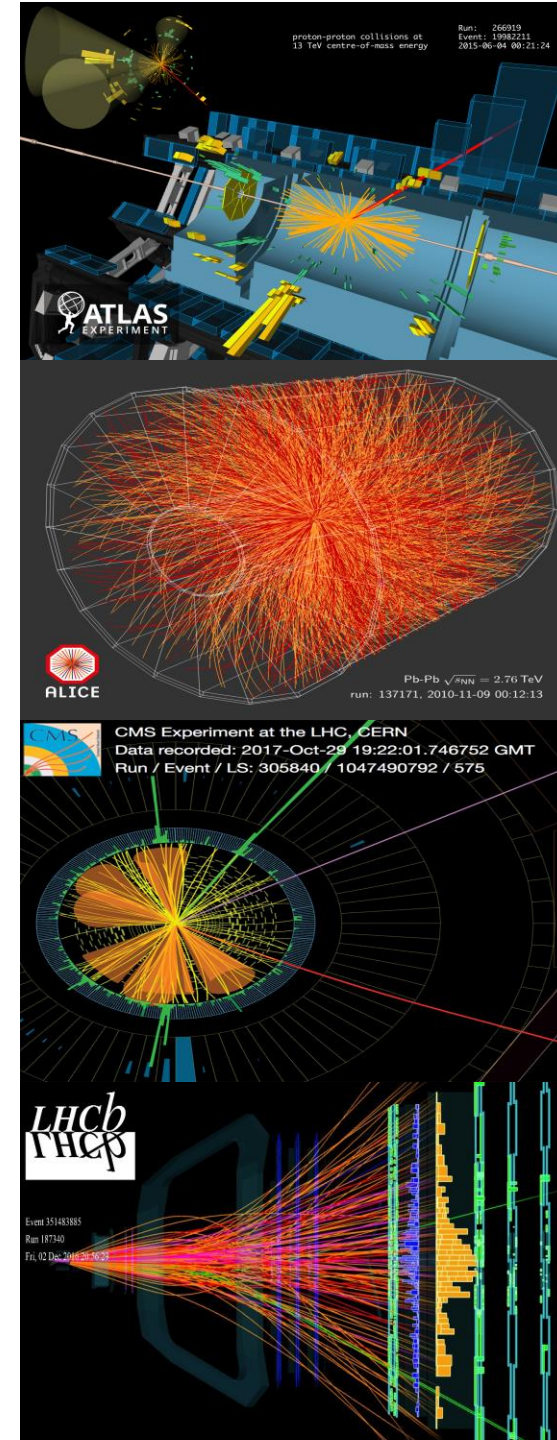# The Worldwide LHC Computing Grid
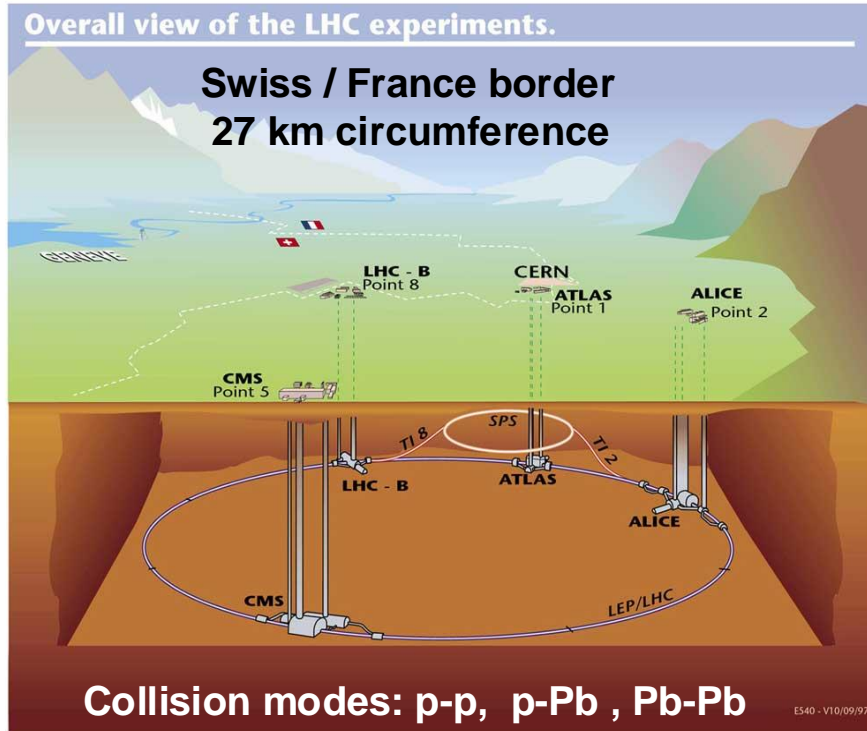
Reda Tafirout

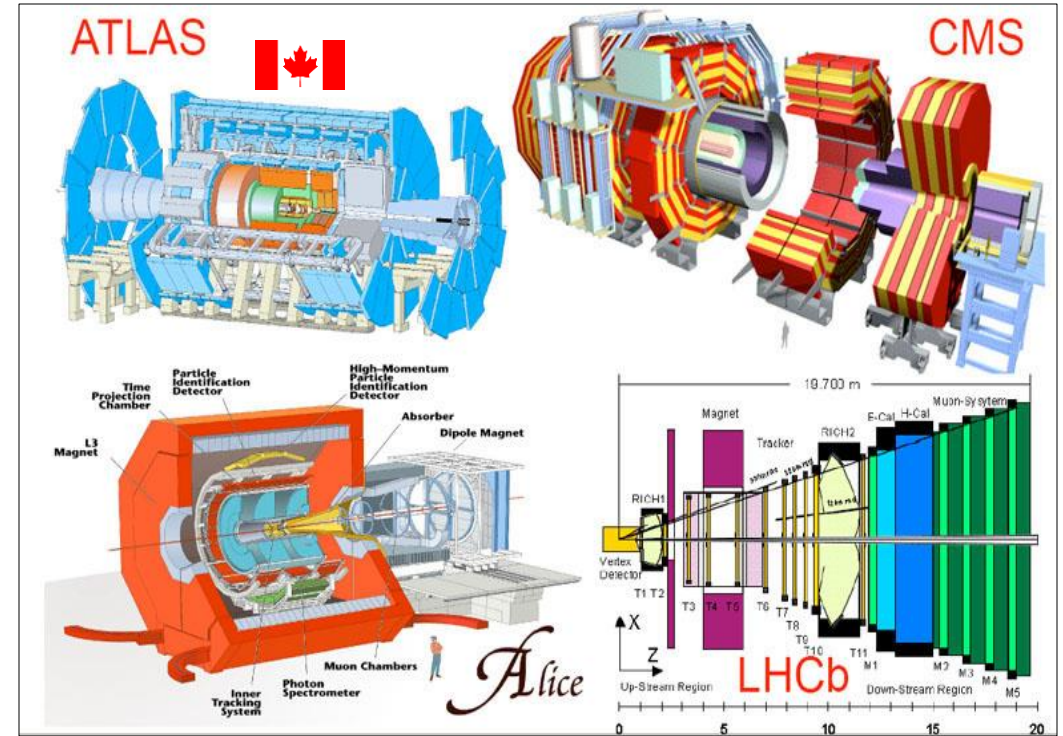TRIUMF Science Week 2024

2024-07-23

# LHC Scientific Program Tools

## Powerful Particle Accelerator
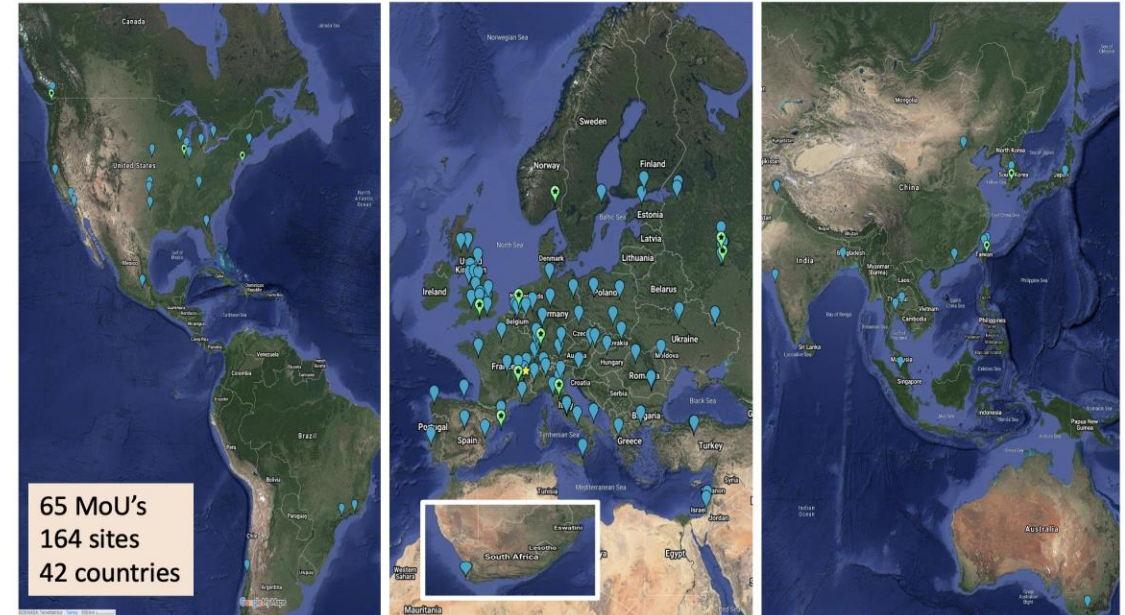


## Sophisticated Detectors



## Large-scale Distributed Computing

# WLCG: global collaboration & effort

- Coordinated resource sharing across multiple organizations

- 164 sites in 42 countries (tiered structure)

- MoU's between parties participating in the project

  - CERN (Tier-0), Tier-1 centres, Tier-2 centres

  - Baseline services & roles

  - Service levels & up-times

  - Management structure & Grid operations

- Pledged resources updated yearly based on experiments requirements.



65 MoU's
164 sites
42 countries

**2024 pledge: ~2 EB tape, 1 EB disk, 550 kcores**

Tier-1: TRIUMF / SFU

Tier-2: Alliance sites (Victoria, SFU, Waterloo)

# WLCG: many key components

- Middleware services

- Information system (sites configuration)

- Security & authentication: IGTF (CA's)

- VO management service

- File transfer service & network provisioning

- Workload management

- Data management

- Accounting & monitoring

- Constant evolution while ensuring continuity

Standard Model Total Production Cross Section Measurements

- **A lot of data needed**
- **SM measurements**
- **BSM searches**
- **Large-scale simulations needed**

# ATLAS TDAQ Chain



- From 40 MHz beam crossing to ~3 kHz (storage at T0)
- HLT farm scale: ~100,000 cores (dedicated)

# ATLAS Computing Model



- Raw data (1st copy)
- Mass storage

- Calibration
- First pass processing

**TDAQ**

**Tier-0**

- Raw data (2nd copy)
- Mass storage (data + MC)

(10x)

**Tier-1** ... **Tier-1**

- Reprocessing
- Group prod.

- Simulation
- User Analysis

**Tier-2** **Tier-2** **Tier-2** **Tier-2**

**Tier-2** **Tier-2** **Tier-2**

(70x)

*chaotic scale*

(ATLAS wide resources / MoU)

- Code development
- Grid UI
- Small scale analysis & opportunistic production

**Tier-3** **Tier-3** ... **Tier-3**

(Private/local resources, Clouds, HPCs)

(~1000x)

7

# LHC Optical Private Network (LHCOPN)



8

# LHC Open Network Environment (LHCONE)



LHCONE L3VPN: A global infrastructure for High Energy Physics data analysis (LHC, Belle II, Pierre Auger Observatory, NOvA, XENON, JUNO, DUNE)

LHCONE Map Ver. 9.0, 2024-04-03 – WEJohnston, ESnet, wej@es.net

# WLCG: global data transfers

**Monthly data transfer throughput between WLCG sites (GB/s) – 3 years**



WLCG supports +40% more transfers since LS-2

Further scalability (x5) demonstrated in the Data Challenge

No strain to the services

# WLCG: cpu delivered globally

- Growing number of computing resources provided to the LHC experiments (+20% since LS2)

**CPU delivered: HS23 hours/month**
**(2010 – present)**



- Drop in winter 2022/23: energy crisis in Europe with high natural gas costs (supply issue from Russia)

- Back to "normal" now

- HEPScore 2023 (HS23) is a common cpu benchmark unit (used in WLCG pledges & accounting)

# ATLAS Distributed Computing Needs & HL-LHC Challenges

- ~10x data rate increase in Run-4
- Flat budget model "problematic"
- Significant development effort is required:
  - Improve software performance
  - Leverage modern architectures
- Data challenges planned

# ATLAS Software & Computing Roadmap for the HL-LHC

- Roadmap has several components dealing with various topics:

    - Network infrastructure ready for Run 4
    - Detector Description, Simulation and Digitization projects
    - HL-LHC datasets replicas and versions management
    - Core Software and Heterogeneous Computing / Accelerators
    - etc.

- ATLAS Heterogeneous Computing & Accelerators Forum established recently

- To tackle the combinatorics in a high luminosity environment, investigate tracking on GPU.  For this to succeed:

    - define a suitable Event Data Model,
    - develop a toolchain that supports e.g. CUDA kernels
    - provide GPU friendly implementations of the geometry and magnetic field.

ACTS - A Common Tracking Software



traccc
full scale demonstrator of an ATLAS-like track reconstruction chain for CPU/GPU

detray
geometry & propagation

algebra-plugins
encapsulation of algebra operation backend

covfie
generic vector field library

vecmem

# Dynamic network provisioning

- Demonstration of software defined networks and dynamic circuit provisioning based on demand (Supercomputing 2023 conference)

- Collaboration between TRIUMF, CERN, FNAL, KIT and network providers

- Traffic from both ATLAS & CMS

- Also work on network packets marking (scientific tags) in collaboration with HEPNet Canada





Canadian Tier-1 aggregate traffic

# ARM CPUs & GPUs

ARM CPUs process more events/Watt wrt x86

- Studies done with hardware Glasgow
- Differences between workflows
- Modern AMDs perform similarly to ARM



HEPScore/Watt



Older ——— time ——→ Newer

**GPUs (in HPCs): massive parallelization, Machine Learning toolkits**

Experiment workflows ported to ARM

- ATLAS, ALICE: validation successful
- CMS, LHCb: in progress



| | | madevent | | |
|---|---|---|---|---|
| CUDA grid size | | 8192 | | |
| $gg \rightarrow t\bar{t}ggg$ | MEs precision | $t_{TOT} = t_{Mad} + t_{MEs}$ [sec] | $N_{events}/t_{TOT}$ [events/sec] | $N_{events}/t_{MEs}$ [MEs/sec] |
| Fortran | double | 1228.2 = 5.0 + 1223.2 | 7.34E1 (=1.0) | 7.37E1 (=1.0) |
| CUDA | double | 19.6 = 7.4 + 12.1 | 4.61E3 (x63) | 7.44E3 (x100) |
| CUDA | float | 11.7 = 6.2 + 5.4 | 7.73E3 (x105) | 1.66E4 (x224) |
| CUDA | mixed | 16.5 = 7.0 + 9.6 | 5.45E3 (x74) | 9.43E3 (x128) |

NVidia V100, Cuda 11.7, gcc 11.2

- **WLCG discussing about pledging ARM CPUs. Resources already available at various sites**

# Conclusion & outlook

- WLCG: global infrastructure developed and operated over the last

  two decades

- Notable achievement: needs of LHC experiments successfully met

- Recent WLCG strategy document developed to tackle key areas:

  - Technical evolution

  - Financial sustainability

  - Heterogeneous grid infrastructure

  - Interaction with other communities with similar challenges

- The HL-LHC era will be a challenging computing environment

  - Need to ensure sustained innovation and development while

    ensuring continuing operations

# TRIUMF

# Thank you
## Merci

**www.triumf.ca**

Follow us **@TRIUMFLab**

**Discovery, accelerated**

# Additional Material

# Canadian ATLAS Tier-1 Centre

- Dedicated facility operated 24/7 (per WLCG MoU)
- Key player and contributor within ATLAS Distributed Computing Operations:
  - Availability, Reliability, Scalability & Performance
  - Critical user support for the entire ATLAS collaboration
- Data storage, data processing, simulations and user analysis in a highly secure environment
- Initially located at TRIUMF since 2007
- Transitioned to SFU in 2018, co-located with Cedar/Alliance; continues to be under the perview of TRIUMF
- Operated as federation (new + old site) since 2017
- Current capacity: 9,300 cores ; 17 PB disk ; 46 PB tape



TRIUMF-LCG2 Availability from 2018/09/01 to 2023/05/01



Tape usage (blue) (Canadian T1)

Disk usage (yellow & green) (Canadian T1)

# Energy efficiency

The electricity costs have been an unexpected development in the last couple of years. Environmental impact needs proper addressing!

What to do:

➢ Improve software performance

➢ Leverage modern architectures

➢ Invest in the facilities

There is no magic wand, however.

The peak of energy need happens in 2036 (start of Run-5): 400% higher than 2022 in the pessimistic scenario and 50% higher in the optimistic scenario.

# RNTuple

RNTuple is the successor of TTree, the ROOT columnar storage technology

Examples of recent commissioning progress:

ATLAS now capable to read/write all data formats in RNTuple, 20% saving in size for DAOD_PHYS. Substantial progress also for the other experiments

- Expected RNTuple speed-up improvements measured in a real environment at CERN using a community standard analysis benchmark

Take home message: **RNTuple progress well on schedule** thanks to a very good collaboration between experiments, CERN EP-SFT and IT

# Event Generators

A very good candidate for GPU acceleration with benefits for many experiments

CPU + GPU



only



Vectorisation

$gg \rightarrow t\bar{t}gg$
(float)

No Vectorisation

Sherpa gg->tt+ng

Matrix Element event throughput:
up to x10 gain when using GPUs

**Available for production**

Madgraph gg->tt+ng (n=2)

GPU-enabled Leading Order: being released to production.
By-product: enabling of CPU vectorisation: up to x8 gain in
ME event throughput (x6 global). Note: all CPUs in WLCG
provide vectorisation

**GPU-related work brings immediate benefits also on CPUs**

# ATLAS Heavy Particle Searches* - 95% CL Upper Exclusion Limits

Status: March 2023

**ATLAS** Preliminary

$\int \mathcal{L} \, dt = (3.6 - 139)$ fb$^{-1}$        $\sqrt{s} = 13$ TeV

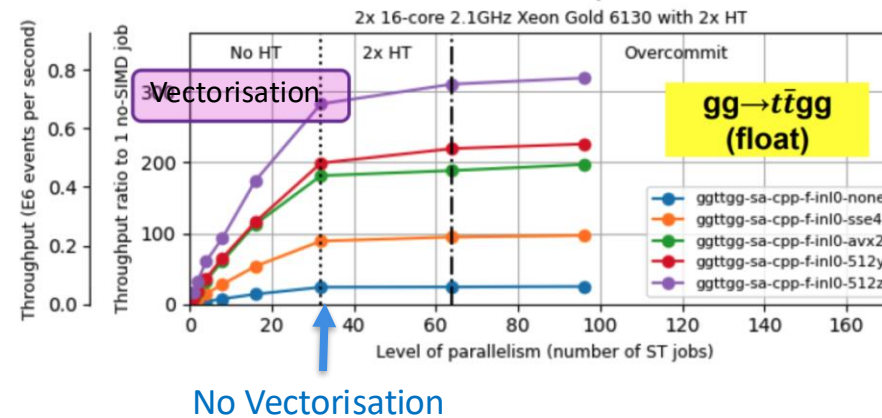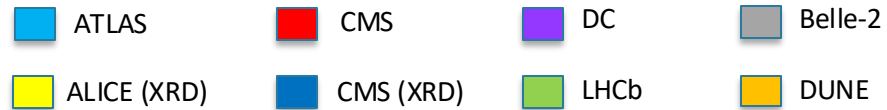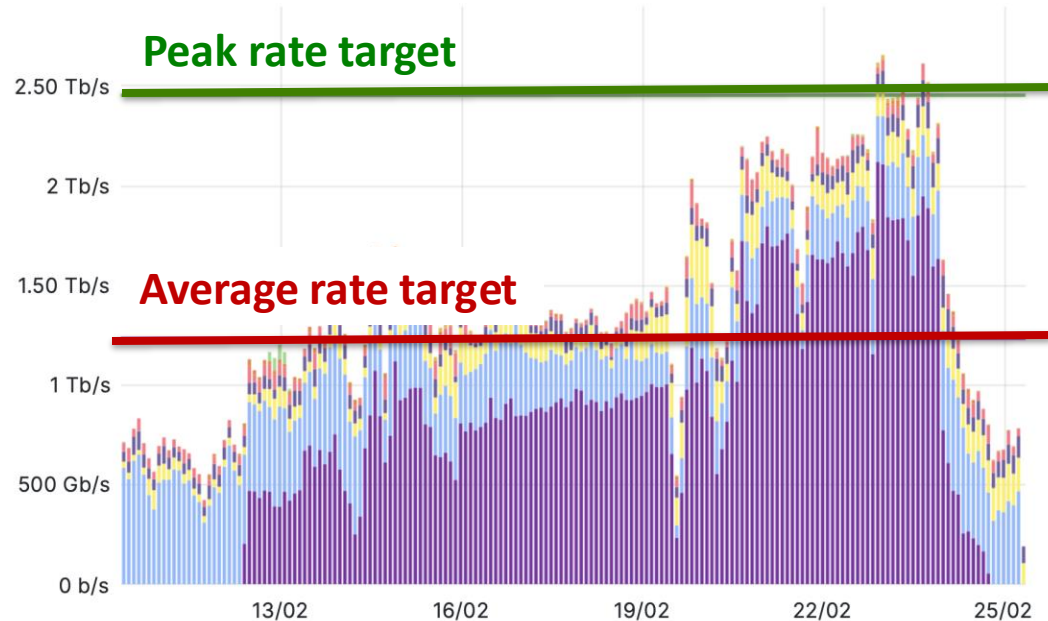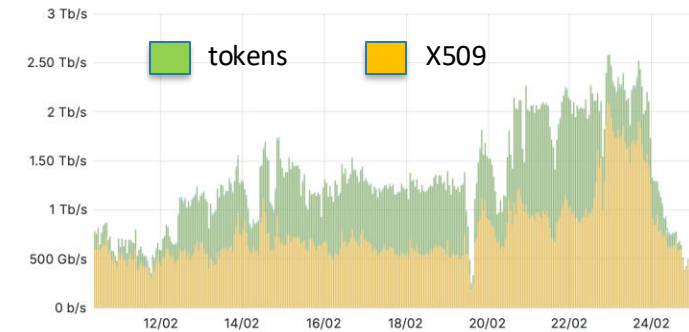| | Model | $\ell, \gamma$ | Jets† | $E_T^{miss}$ | $\int \mathcal{L} \, dt$[fb$^{-1}$] | Limit | | Reference |
|---|---|---|---|---|---|---|---|---|
| **Extra dimen.** | ADD $G_{KK} + g/q$ | $0\ e,\mu,\tau,\gamma$ | $1-4$ j | Yes | 139 | $M_D$ = 11.2 TeV | $n = 2$ | 2102.10874 |
| | ADD non-resonant $\gamma\gamma$ | $2\gamma$ | – | – | 36.7 | $M_S$ = 8.6 TeV | $n = 3$ HLZ NLO | 1707.04147 |
| | ADD QBH | – | 2 j | – | 139 | $M_{th}$ = 9.4 TeV | $n = 6$ | 1910.08447 |
| | ADD BH multijet | – | ≥3 j | – | 3.6 | $M_{th}$ = 9.55 TeV | $n = 6$, $M_D = 3$ TeV, rot BH | 1512.02586 |
| | RS1 $G_{KK} \to \gamma\gamma$ | $2\gamma$ | – | – | 139 | $G_{KK}$ mass = 4.5 TeV | $k/\overline{M}_{Pl} = 0.1$ | 2102.13405 |
| | Bulk RS $G_{KK} \to WW/ZZ$ | multi-channel | | | 36.1 | $G_{KK}$ mass = 2.3 TeV | $k/\overline{M}_{Pl} = 1.0$ | 1808.02380 |
| | Bulk RS $g_{KK} \to tt$ | $1\ e,\mu$ | ≥1 b, ≥1J/2j | Yes | 36.1 | $g_{KK}$ mass = 3.8 TeV | $\Gamma/m = 15\%$ | 1804.10823 |
| | 2UED / RPP | $1\ e,\mu$ | ≥2 b, ≥3 j | Yes | 36.1 | KK mass = 1.8 TeV | Tier (1,1), $\mathcal{B}(A^{(1,1)} \to tt) = 1$ | 1803.09678 |
| **Gauge bosons** | SSM $Z' \to \ell\ell$ | $2\ e,\mu$ | – | – | 139 | $Z'$ mass = 5.1 TeV | | 1903.06248 |
| | SSM $Z' \to \tau\tau$ | $2\tau$ | – | – | 36.1 | $Z'$ mass = 2.42 TeV | | 1709.07242 |
| | Leptophobic $Z' \to bb$ | – | 2 b | – | 36.1 | $Z'$ mass = 2.1 TeV | | 1805.09299 |
| | Leptophobic $Z' \to tt$ | $0\ e,\mu$ | ≥1 b, ≥2 J | Yes | 139 | $Z'$ mass = 4.1 TeV | $\Gamma/m = 1.2\%$ | 2005.05138 |
| | SSM $W' \to \ell\nu$ | $1\ e,\mu$ | – | Yes | 139 | $W'$ mass = 6.0 TeV | | 1906.05609 |
| | SSM $W' \to \tau\nu$ | $1\tau$ | – | Yes | 139 | $W'$ mass = 5.0 TeV | | ATLAS-CONF-2021-025 |
| | SSM $W' \to tb$ | – | ≥1 b, ≥1 J | – | 139 | $W'$ mass = 4.4 TeV | | ATLAS-CONF-2021-043 |
| | HVT $W' \to WZ$ model B | $0-2\ e,\mu$ | 2 j / 1 J | Yes | 139 | $W'$ mass = 4.3 TeV | $g_V = 3$ | 2004.14636 |
| | HVT $W' \to WZ \to \ell\nu\,\ell'\ell'$ model C | $3\ e,\mu$ | 2 j (VBF) | Yes | 139 | $W'$ mass = 340 GeV | $g_V c_H = 1$, $g_f = 0$ | 2207.03925 |
| | HVT $Z' \to WW$ model B | $1\ e,\mu$ | 2 j / 1 J | Yes | 139 | $Z'$ mass = 3.9 TeV | $g_V = 3$ | 2004.14636 |
| | LRSM $W_R \to \mu N_R$ | $2\mu$ | 1 J | – | 80 | $W_R$ mass = 5.0 TeV | $m(N_R) = 0.5$ TeV, $g_L = g_R$ | 1904.12679 |
| **CI** | CI $qqqq$ | – | 2 j | – | 37.0 | $\Lambda$ = 21.8 TeV | $\eta_{LL}^-$ | 1703.09127 |
| | CI $\ell\ell qq$ | $2\ e,\mu$ | – | – | 139 | $\Lambda$ = 35.8 TeV | $\eta_{LL}^-$ | 2006.12946 |
| | CI $eebs$ | 2 e | 1 b | – | 139 | $\Lambda$ = 1.8 TeV | $g_* = 1$ | 2105.13847 |
| | CI $\mu\mu bs$ | $2\mu$ | 1 b | – | 139 | $\Lambda$ = 2.0 TeV | $g_* = 1$ | 2105.13847 |
| | CI $tttt$ | ≥1 $e,\mu$ | ≥1 b, ≥1 j | Yes | 36.1 | $\Lambda$ = 2.57 TeV | $|C_{4t}| = 4\pi$ | 1811.02305 |
| **DM** | Axial-vector med. (Dirac DM) | – | 2 j | – | 139 | $m_{med}$ = 3.8 TeV | $g_q$=0.25, $g_\chi$=1, $m(\chi)$=10 TeV | ATL-PHYS-PUB-2022-036 |
| | Pseudo-scalar med. (Dirac DM) | $0\ e,\mu,\tau,\gamma$ | $1-4$ j | Yes | 139 | $m_{med}$ = 376 GeV | $g_q$=1, $g_\chi$=1, $m(\chi)$=1 GeV | 2102.10874 |
| | Vector med. $Z'$-2HDM (Dirac DM) | $0\ e,\mu$ | 2 b | Yes | 139 | $m_{Z'}$ = 3.0 TeV | $\tan\beta$=1, $g_Z$=0.8, $m(\chi)$=100 GeV | 2108.13391 |
| | Pseudo-scalar med. 2HDM+a | multi-channel | | | 139 | $m_a$ = 800 GeV | $\tan\beta$=1, $g_\chi$=1, $m(\chi)$=10 GeV | ATLAS-CONF-2021-036 |
| **LQ** | Scalar LQ 1st gen | 2 e | ≥2 j | Yes | 139 | LQ mass = 1.8 TeV | $\beta = 1$ | 2006.05872 |
| | Scalar LQ 2nd gen | $2\mu$ | ≥2 j | Yes | 139 | LQ mass = 1.7 TeV | $\beta = 1$ | 2006.05872 |
| | Scalar LQ 3rd gen | $1\tau$ | 2 b | Yes | 139 | $LQ_3^u$ mass = 1.49 TeV | $\mathcal{B}(LQ_3^u \to b\tau) = 1$ | 2303.01294 |
| | Scalar LQ 3rd gen | $0\ e,\mu$ | ≥2 j, ≥2 b | Yes | 139 | $LQ_3^d$ mass = 1.24 TeV | $\mathcal{B}(LQ_3^u \to t\nu) = 1$ | 2004.14060 |
| | Scalar LQ 3rd gen | ≥2 $e,\mu$,≥1 $\tau$ ≥1 j, ≥1 b | | – | 139 | $LQ_3^d$ mass = 1.43 TeV | $\mathcal{B}(LQ_3^d \to t\tau) = 1$ | 2101.11582 |
| | Scalar LQ 3rd gen | $0\ e,\mu$,≥1 $\tau$ $0-2$ j, 2 b | | Yes | 139 | $LQ_3^d$ mass = 1.26 TeV | $\mathcal{B}(LQ_3^d \to b\nu) = 1$ | 2101.12527 |
| | Vector LQ mix gen | multi-channel ≥1 j, ≥1 b | | Yes | 139 | $LQ_3^V$ mass = 2.0 TeV | $\mathcal{B}(\tilde{U}_1 \to t\mu) = 1$, Y-M coupl. | ATLAS-CONF-2022-052 |
| | Vector LQ 3rd gen | $2\ e,\mu,\tau$ | ≥1 b | Yes | 139 | $LQ_3^V$ mass = 1.96 TeV | $\mathcal{B}(LQ_3^V \to b\tau) = 1$, Y-M coupl. | 2303.01294 |
| **Vector-like fermions** | VLQ $TT \to Zt + X$ | $2e/2\mu/\geq3e,\mu$ ≥1 b, ≥1 j | | – | 139 | T mass = 1.46 TeV | SU(2) doublet | 2210.15413 |
| | VLQ $BB \to Wt/Zb + X$ | multi-channel | | | 36.1 | B mass = 1.34 TeV | SU(2) doublet | 1808.02343 |
| | VLQ $T_{5/3}T_{5/3}|T_{5/3} \to Wt + X$ | 2(SS)/≥3 $e,\mu$ ≥1 b, ≥1 j | | Yes | 36.1 | $T_{5/3}$ mass = 1.64 TeV | $\mathcal{B}(T_{5/3} \to Wt) = 1$, $c(T_{5/3} Wt) = 1$ | 1807.11883 |
| | VLQ $T \to Ht/Zt$ | $1\ e,\mu$ | ≥1 b, ≥3 j | Yes | 139 | T mass = 1.8 TeV | SU(2) singlet, $\kappa_T = 0.5$ | ATLAS-CONF-2021-040 |
| | VLQ $Y \to Wb$ | $1\ e,\mu$ | ≥1 b, ≥1 j | Yes | 36.1 | Y mass = 1.85 TeV | $\mathcal{B}(Y \to Wb) = 1$, $c_R(Wb) = 1$ | 1812.07343 |
| | VLQ $B \to Hb$ | $0\ e,\mu$ | ≥2b,≥1j,≥1J | – | 139 | B mass = 2.0 TeV | SU(2) doublet, $\kappa_B = 0.3$ | ATLAS-CONF-2021-018 |
| | VLL $\tau' \to Z\tau/H\tau$ | multi-channel | ≥1 j | Yes | 139 | $\tau'$ mass = 898 GeV | SU(2) doublet | 2303.05441 |
| **Exctd ferm.** | Excited quark $q^* \to qg$ | – | 2 j | – | 139 | $q^*$ mass = 6.7 TeV | only $u^*$ and $d^*$, $\Lambda = m(q^*)$ | 1910.08447 |
| | Excited quark $q^* \to q\gamma$ | $1\gamma$ | 1 j | – | 36.7 | $q^*$ mass = 5.3 TeV | only $u^*$ and $d^*$, $\Lambda = m(q^*)$ | 1709.10440 |
| | Excited quark $b^* \to bg$ | – | 1 b, 1 j | – | 139 | $b^*$ mass = 3.2 TeV | | 1910.08447 |
| | Excited lepton $\tau^*$ | $2\tau$ | ≥2 j | – | 139 | $\tau^*$ mass = 4.6 TeV | $\Lambda = 4.6$ TeV | 2303.09444 |
| **Other** | Type III Seesaw | 2,3,4 $e,\mu$ | ≥2 j | Yes | 139 | $N^0$ mass = 910 GeV | | 2202.02039 |
| | LRSM Majorana $\nu$ | $2\mu$ | 2 j | – | 36.1 | $N_R$ mass = 3.2 TeV | $m(W_R) = 4.1$ TeV, $g_L = g_R$ | 1809.11105 |
| | Higgs triplet $H^{\pm\pm} \to W^\pm W^\pm$ | 2,3,4 $e,\mu$ (SS) | various | Yes | 139 | $H^{\pm\pm}$ mass = 350 GeV | DY production | 2101.11961 |
| | Higgs triplet $H^{\pm\pm} \to \ell\ell$ | 2,3,4 $e,\mu$ (SS) | – | – | 139 | $H^{\pm\pm}$ mass = 1.08 TeV | DY production | 2211.07505 |
| | Multi-charged particles | – | – | – | 139 | multi-charged particle mass = 1.59 TeV | DY production, $|q| = 5e$ | ATLAS-CONF-2022-034 |
| | Magnetic monopoles | – | – | – | 34.4 | monopole mass = 2.37 TeV | DY production, $|g| = 1g_D$, spin 1/2 | 1905.10130 |

| $\sqrt{s} = 13$ TeV partial data | $\sqrt{s} = 13$ TeV full data |
|---|---|

$10^{-1}$        1        10        **Mass scale [TeV]**

*Only a selection of the available mass limits on new states or phenomena is shown.

†Small-radius (large-radius) jets are denoted by the letter j (J).

**ATLAS BSM searches**

# Data Challenge 2024 - Highlights

DC24 WLCG data transfers (Gbps) – 15 days: **all targets achieved**

**Peak rate target**

**Average rate target**



| ATLAS | CMS | DC | Belle-2 |
| ALICE (XRD) | CMS (XRD) | LHCb | DUNE |

New technologies (e.g. authentication tokens) introduced and validated



tokens    X509

WLCG services successfully supports DUNE and Belle-2 computing models
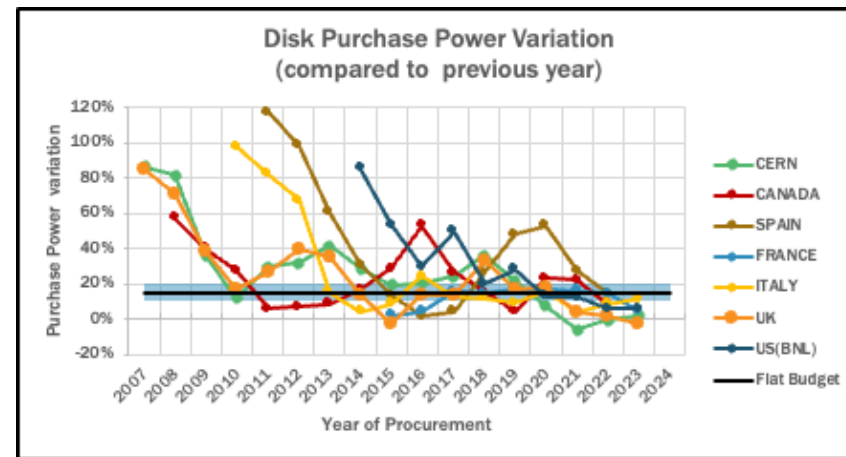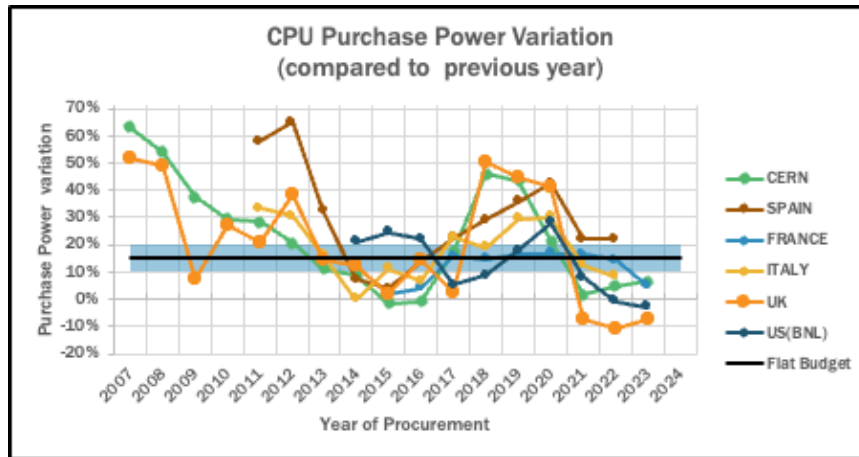
# Hardware Cost Evolution

The WLCG "flat budget model" assumption: +15% CPU, disk and tape every year **with the same level of funding**

We now monitor the HW trends in many countires. Last 5 years average is compatible with the 15% assumption but look at the first derivative …

CPU average variation (5 years): +14%

DISK average variation (5 years): +15%
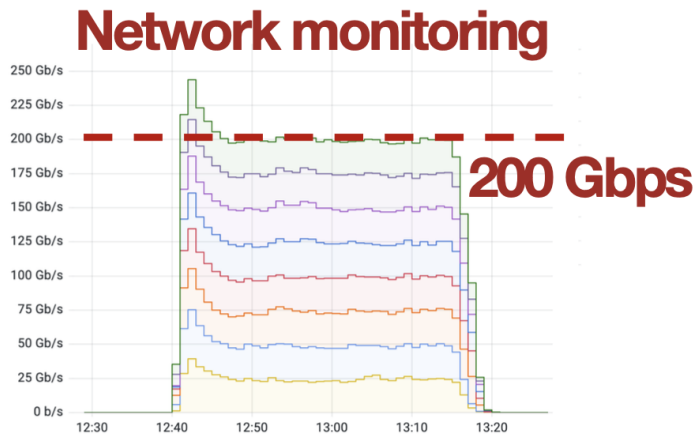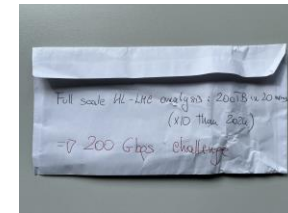
# IRIS-HEP: analysis 200Gbps (Grand) Challenge

Launched by IRIS-HEP to commission analysis capabilities for HL-LHC

- Commission services at increasing scale + introduce innovative aspects

$\Rightarrow$ Show readiness at 25% of HL-LHC scale (same as for data challenge)

Analysis models evolving => metrics of success hard to quantify (25% of?)



**Network monitoring**

**200 Gbps**

Based on the IRIS-HEP toolkit