

Software & Computing status

L. Poggioli, LAL

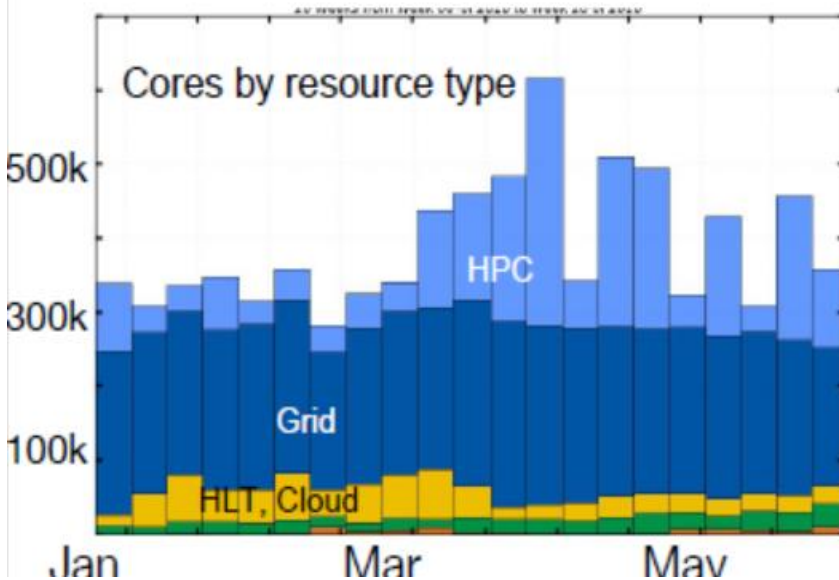
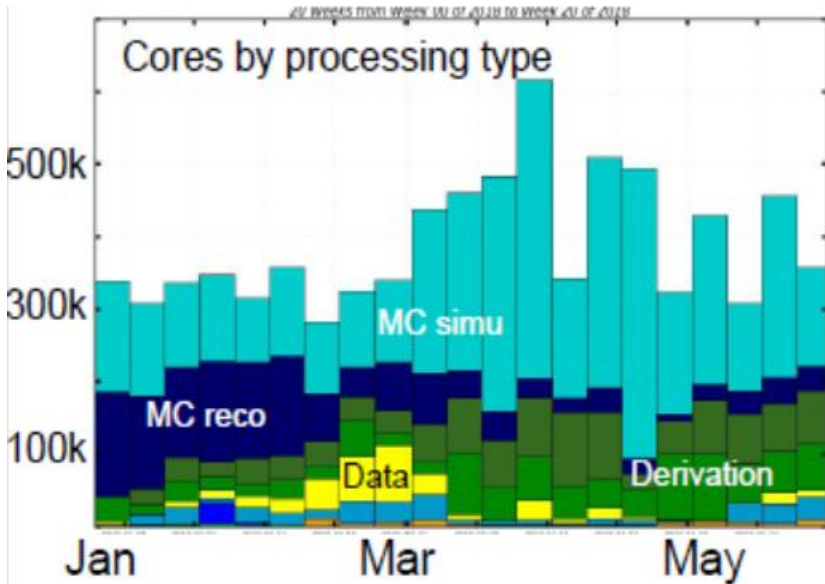
- Finalizing Run-2
- Preparing Run-3
- Towards HL-LHC

Based on

- Various S&C weeks
- HSF/WLCG in Naples March
- Tokyo ATLAS week June
- preGDB on networking 2017

Finalizing Run-2

2018 Operations

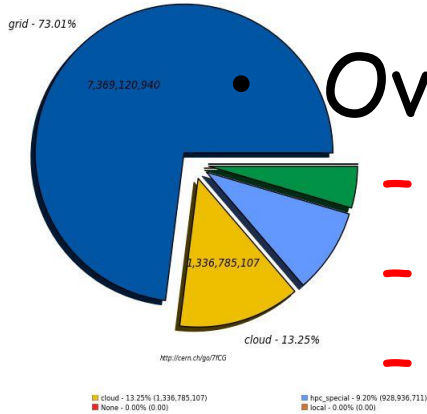


- T0
 - Since April +20% resource
 - 23k slots, 4GB/slot
 - Used by Grid if no T0 jobs
- Resource
 - Constant > 300-350k slots
 - Use HLT farm when no run
- Data
 - 2.5B events in AOD
 - 5B events simulated in 2018
 - Moving >1PB, >20GB/s, 2M files/day. Delete 10pB/wk
 - Automatization (replication, preplacement, balancing)

Extra resource (1)



CPU HEPSPEC06 (Sum: 10,092,733,554)



Over last 3 months

- Grid 73%
- Cloud (standard, HLT, BOINC) 13%
- HPC_Special (Big US HPC) 9%
- HPC (Local, ie Grid-like) 5%

- General effort to handle heterogeneous resource
 - **Harvester** under development (common interface for ALL type of resource)
 - **Event Service** allows to work at event level (simulation). Also for sites (local cluster/T3)

Extra Resource (2)

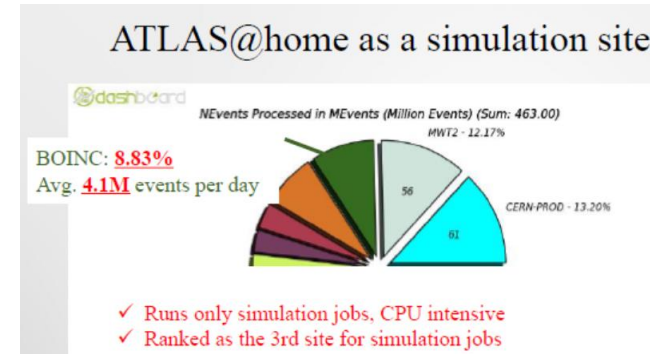
- Cloud
 - Stable but no real increase
 - Cost gain wrt grid? Manpower?

- ATLAS@Home: Increasing

- Free!!
- Used to backfill sites and optimize CPU usability

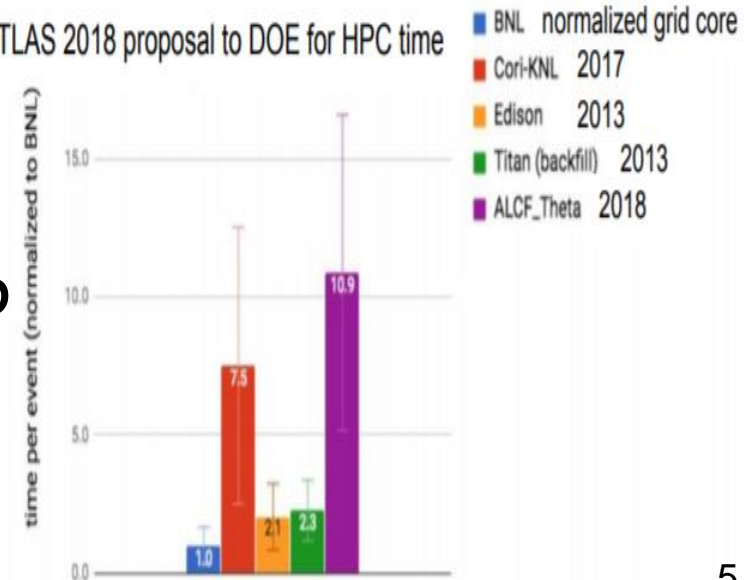
- HPC

- Not 'free' at prod scale
- HPC HS06 ~ 1/10 Grid HS06
- Test by CC@IDRIS.
OK but 2.5k/10k slots max



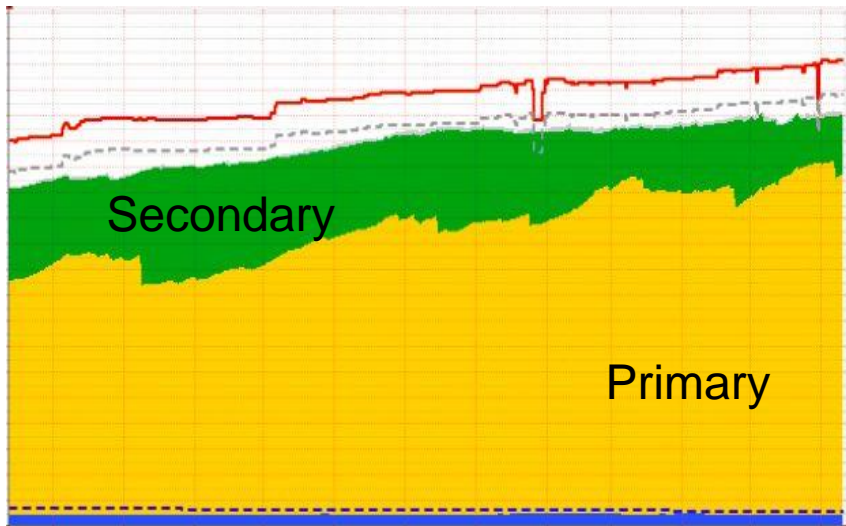
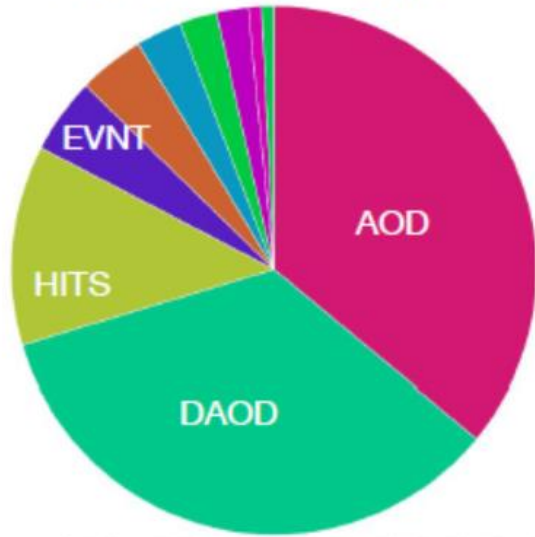
2018: BOINC 11kslots,
3rd site for Simulation

ATLAS 2018 proposal to DOE for HPC time



Disk space handling (2018 & LS2)

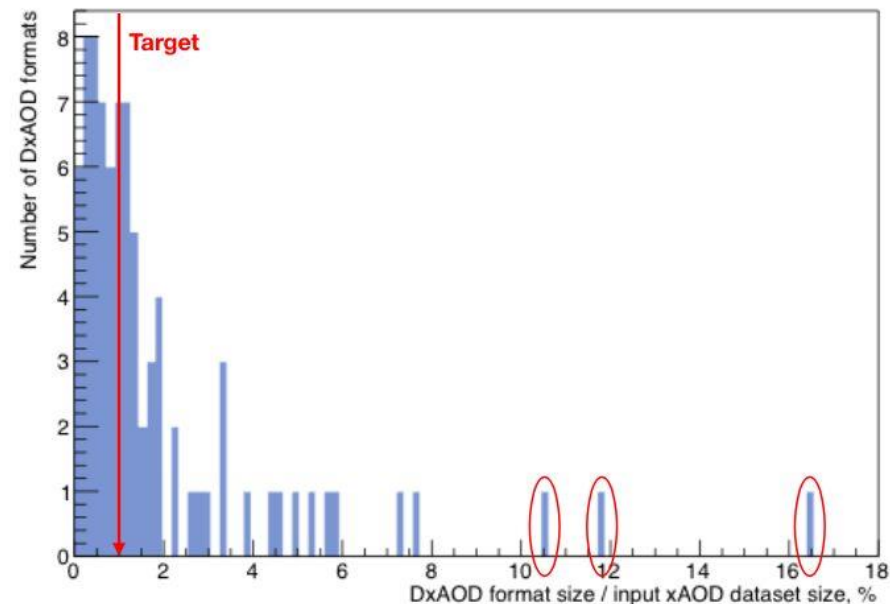
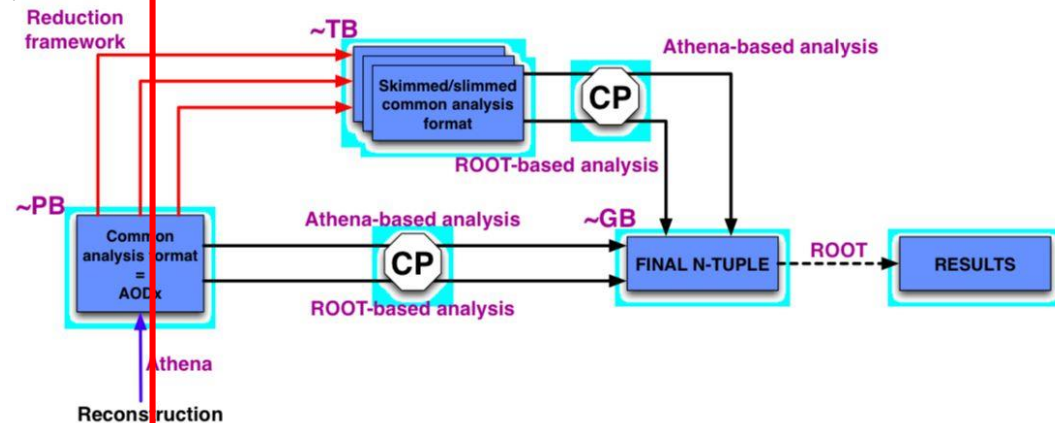
Disk usage by data type



- Disk space tight
 - Adding 2-3pB/week
 - Priority to **dAOD**
 - Old AOD&HITS deleted
- Actions
 - FastSim HITS deleted
 - Run Lifetime models
 - AOD compression (LZMA)
- Pledged disk 160pB
 - Primary cannot be deleted
 - Secondary can be deleted if needed

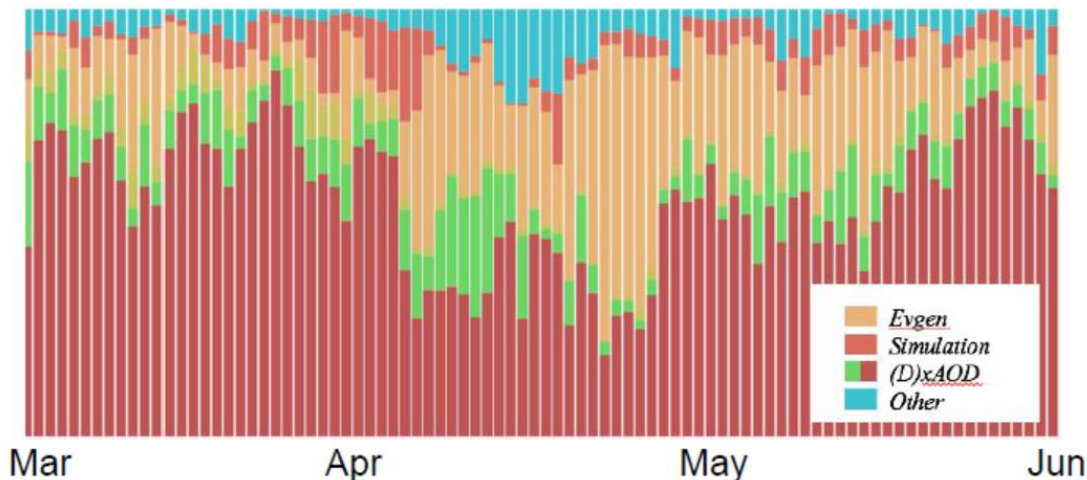
Distributed Analysis (1)

- Run-2 model successful
 - Still too long to reprocess all data
 - Analysis jobs tails
- AOD production for 2018 data
 - Running smoothly
 - Small CPU impact from AOD compression (LZMA)
- 2018
 - 27 trains, 84 derivations
 - dAOD still too big!

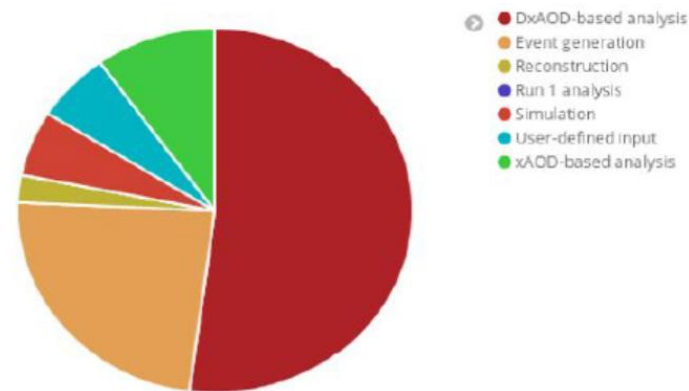


Distributed Analysis (2)

Fraction of distributed analysis jobs vs. time



Wallclock time usage

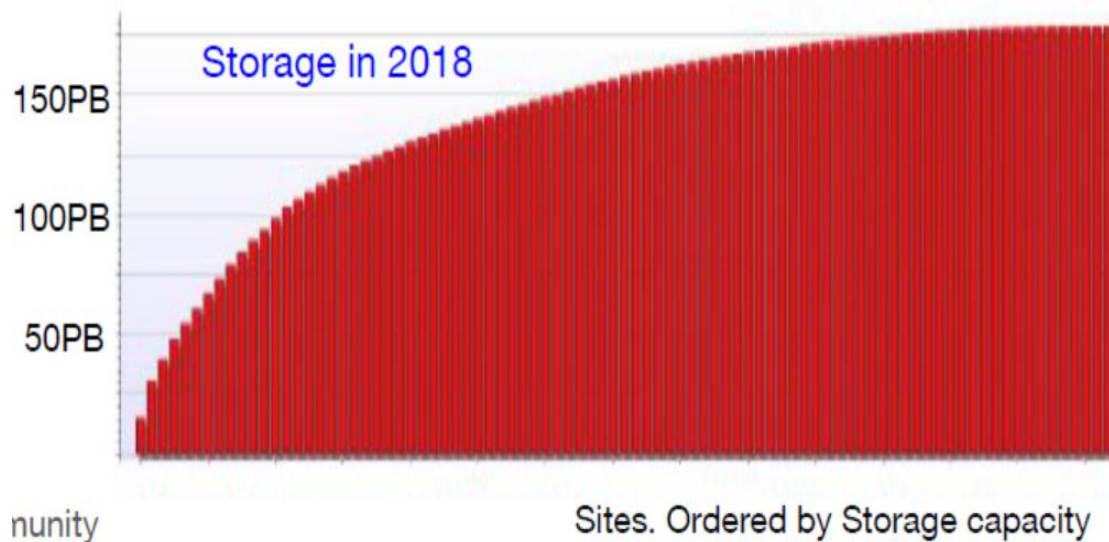


Typical analysis GRID usage over last 3 months

- Most Grid jobs (60% WT) use dAOD (&AOD) as inputs. In agreement with Train model
- 40% WT used by private Simu, Reco & Evgen (incl. Toy MC)

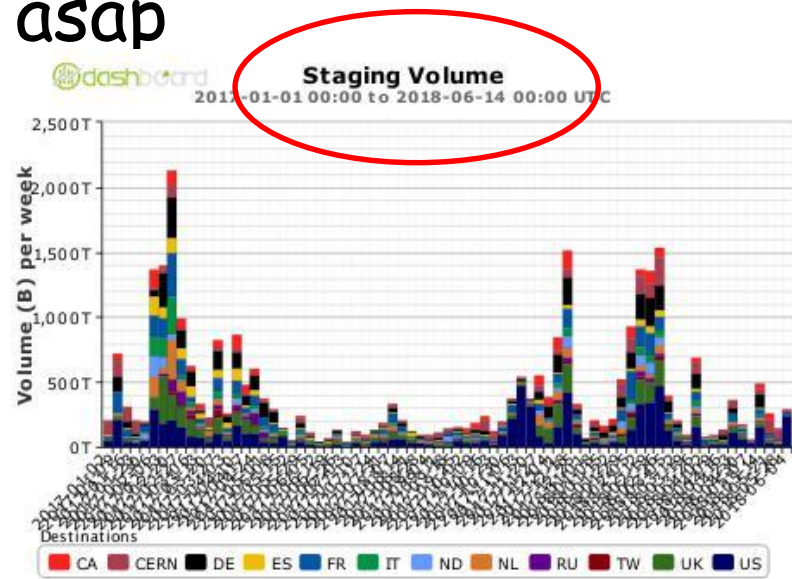
T2/T3 Consolidation

- 2010: All sites required to have storage
- 2016: smallest 50% sites ~ 10% storage
 - Operations heavy for sites and ADC
- Sites w/ small storage: **invest in CPU wrt disk**
 - Cutoff is **440 TB** (0.5% total T2s storage)
 - With +10-15% increase/year (follow flat budget)
- 2018: smallest 40% sites ~ 10% storage
 - Decommissioning in progress (**diskless sites**)
- Future: decouple storage and CPU
 - **"Data lake"** model



Distributed Computing status

- **Harvester** software to interface various platforms
 - eg HPCs: more homogeneity & automation
 - Possible use beyond ATLAS
- **Data carousel** mode of operation with **tapes**
 - R&D started to use tapes more efficiently
 - Inputs can be processed asap
 - Processed inputs eviction asap
 - Overall number of copies on disk < 1
- Data management
 - **Rucio** workshop: Interest among other experiments (eg CMS)



R&D projects

- **Google Ocean:** Use Google CPU & storage
 - User analysis
 - Store analysis copies output for user access & serves as cache
 - Data placement, replication & popularity
 - Store final derivation of MC and repro data
 - Google Network to make data available globally
 - Data streaming
 - Evaluate necessary compute for generation of sub-file products (branches/events from ROOT files)
- **CSCS** (3rd biggest HPC, Switzerland)
 - Evaluate T0 spillover & HLT reprocessing
 - Actions (with HL-LHC challenges in view)
 - Port ATLAS code to GPUs (eg HLT code), and machine learning (ML) apps and tracking algorithms
 - Pioneering EOS evaluation at CSCS w/ CERN & ATLAS

Lines of effort

- Software
 - Leverage additional resources (HPC, Boinc, ...)
 - Improve software and efficiency (SPOT group)
 - Run less full-simulation (and more **fast sim**)
 - Promote support for **software development**
- Workflow
 - T1s continue to exercise and improve perf. of DAOD production from **tape** inputs
 - **Harvester**, **Event service** (ES), **Overlay** (pileup handling),
 - New: Event Streaming service (**ESS**)
 - What ES is to computing, ESS is to input data transfer
- Computing Model
 - **Nucleus/satellites** model
 - **T2/T3** consolidation

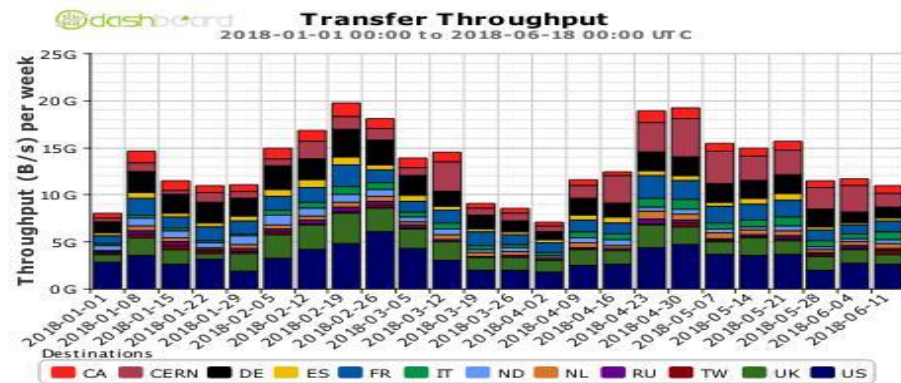
Network: end of Run-2

- 1PB transferred/day@10 GB/s (peak at 20)
- LHCOPN
 - Saturation CERN to T1 for some T1s (e.g. RAL) .OK for end Run-2

- LHCONE
 - OK w/ > 10Gb/s for end Run-2

- Remarks

- Use dynamic Data Distribution & job brokerage based on Networking (Perfsonar not integrated enough into DDM)
- Use Network matrix (closeness) to optimize job brokering



Preparing Run-3

Pledges 2019 (prel.)

- C-RSG April: No changes to 2019 requests Oct->now
- Preliminary questions from C-RSG
 - Coping w/ +20% lumi, mcore eff., tape usage, T0 spillover impact to grid, tape based workflows (prestaging 'data carousel')
- Preliminary pledges as presented at April C-RSG

	2018 (k)	2019 (k)	Delta (%)
T0 CPU	411	411	0
T0 Disk	26	27	4
T0 Tape	94	94	0
T1 CPU	949	1057	11
T1 Disk	72	88	22
T1 Tape	195	221	13
T2 CPU	1160	1292	11
T2 Disk	88	108	23

- Modest increase
- Disk @ T1&T2
 - Slightly high but essential wrt Model
- Tape
 - Occupancy below pledges
 - To catch up in 2018/2019
- Not approved yet!

Run-3 scenario (prel.)

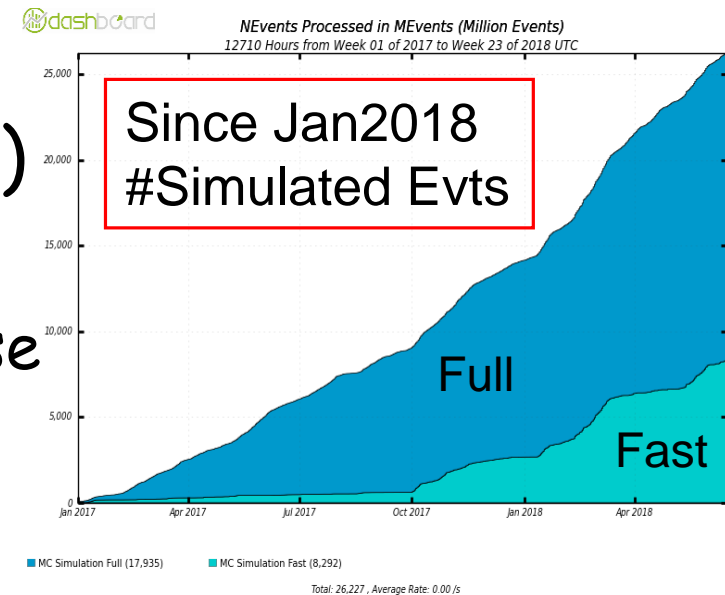
- Run-2
 - In 2018: 2×10^{34} , $\langle \mu \rangle \sim 39$ expect 50 fb^{-1}
 - Total $\sim 150 \text{ fb}^{-1}$
- Run-3
 - 3 years 2021 (shorter), 2022, 2023
 - $\langle \mu \rangle \sim 60-80$ 3.5×10^{34} expect $> 200-300 \text{ fb}^{-1}$
 - Trigger rate stays at **1kHz**
- A priori
 - Assumption resources will be **1.5x(resources in 2018)** Consistent with flat budget
 - Remember in 2017: **+50% data \leftrightarrow +20% resource**

Software for Run-3

- **AthenaMT**: Move towards a multithreaded framework to use modern architectures
- **FastCaloSim**: High priority for ATLAS
- Add new detectors to simulation and reconstruction (NSW)
- **ACTS** (A Common Tracking Software) for tracking. Streamlined ATLAS software, MT by construction. Recommendation to use some ACTS at end of June
- Severe lack of developers **15FTEs missing!!**

Progress on simulation(1)

- Validation ongoing but **manpower** is short!
- Full simulation
 - Dominates CPU usage (45% T)
 - Adds systematic to physics studies (eg VHbb)
- **Fast** Simulation
 - Uses FastCaloSim (version 2)
 - 73% G4 spent in calorimeters
 - Use param'zed shower response
 - As fast as V1 & more accurate
 - Physics validation underway
 - Gain $O(100)$ wrt G4 full simulation



Progress on simulation (2)

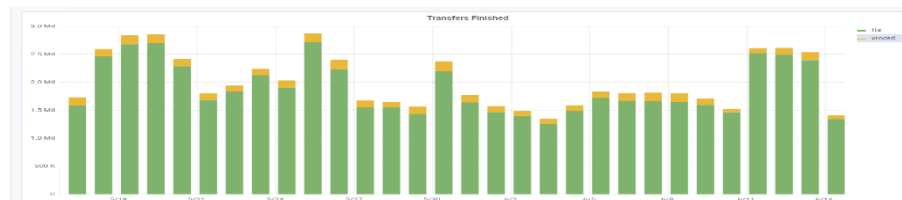
- **Pileup** handling
 - Standard (hits merging) , Data event overlay, MC event overlay
 - Overlay are less CPU&IO requiring but put higher constraint on Condition Databases
- **FastChain**
 - After calorimeter biggest CPU usage is Tracker
 - FATRAS: Simulation in Tracker (simpler Geometry & interactions with material model)
 - Fast digitization algorithm for Inner Tracker
 - Reconstruction uses truth-based pattern recog.

Analysis for Run-3

- Run 2 model very successful
 - Many derived AOD (DAOD) formats O(100)
 - AOD use 55 PB of disk / DAOD use 52 PB of disk
- Two task forces for AOD & DAOD
 - Reduced overall size-#versions used-Smaller evt sizes?
- Scrutiny group at last RRB
 - ATLAS uses more disk than CMS. Difference is growing
 - Encouraged to look into smaller data formats
- Run 3 Analysis Model Study Group
 - Run-3: More MC (FastSim), Bigger evts (μ), Same #data
 - Reduce disk use by at least 30% for same sample sizes
 - -> AOD, DAOD smaller, store less formats/evt, better tape usage

Network: Run-3

- Data growth faster than disk (under flat budget): -> smaller #replicas, -> more network usage
- Already exploiting federated storage & **remote access**. Sites need to dimension to handle remote data access. Expect 10-20% (today 5%)



- Remarks

- Up to now not needed to pledge netwk resource since sites able to properly dimension disk/cpu/network
- Today for most countries, 10-100Gb/s from NREN
- ATLAS would profit of big storage sites with **100Gb/s** (target) which will be able to run all the workflows

- Next

- **Event Streaming Service** (minimize load on network)
- More **tape** usage to store & reread more often -> more disk caches, and more network bandwidth (x2?)

Towards HL-LHC

Towards HL-LHC: Challenges

- Inputs
 - Trigger rate **10kHz**. Increase total evt numbers
 - $\langle \mu \rangle \sim 200$. Increase in CPU & storage needs (time in Reco, bigger evts)
- Today
 - **X 3 missing in CPU**. Seems affordable (many ideas)
 - R&D inside HSF, Accelerators (GPU, FPGA), Extra-resources (HPC, R&D with Google,...)
 - FastSim, Detecor layout, Machine Learning
 - **X 7 missing in storage** critical (less solutions available)
- **HSF** well established, **WLCG** strategy doc. available
- R&D areas
 - Data Organisation and Management Access (**DOMA**)
 - Software upgrade, HSF technical forum

Possible gains for storage

- Disk usage today: $\frac{1}{3}$ AOD, $\frac{1}{3}$ dAOD, and $\frac{1}{3}$ others
 - Others: samples mostly on tape, rotating onto disk cache when needed for processing (e.g. simu hits)
- -> Extend **tape carousel** to AOD & ultimately dAOD
 - Also Possible: Make **AODs 10x smaller** à la CMS, Streamline some physics analyses
- Difficult
 - Tape means delay, delicate workflows handling, but (d)AOD workflows time critical and very complex already
 - Tape is a geographically limited resource at T1s, while processing resources are much more widely distributed
- Limitation of **# replicas**
 - ≥ 1 replicas on disk today, -> dynamic, managed availability of actively used data via **Data Lake**, replica count $\ll 1$

Summary

- ATLAS S&C in very good shape!
 - Now able to focus on refinements, performance, and look to future with R&D
- ATLAS should be front and center in common R&D (inside HSF community)
- Run-3 a priori OK within flat budget. Key issue is software: AthenaMT & FastSim
- HL-LHC
 - Trend lines are good in CPU (constant progress)
 - Plans in storage to be quantified (today critical)
 - R&D, 'Data lake' model, ~~non-flat budget?~~