

Experience from Monte Carlo Production in ATLAS

Wolfgang Ehrenfeld (Bonn)



Fast Monte Carlo Workshop in HEP
15 January 2014



➤ Motivation

➤ Monte Carlo production steps

- event generation, simulation, digitisation, reconstruction

Motivation

Produced MC Events

- mc11: 2.4×10^9 full and 2.1×10^9 fast simulation events
 - mc11a: 0.8×10^9 events
 - mc11b: 1.0×10^9 events (super seeds mc11a)
 - mc11c: 4.8×10^9 events (super seeds mc11b) → total: 4.8×10^9 events

- mc12: 3.8×10^9 full and 3.0×10^9 fast simulation events
 - mc12a: 5.9×10^9 events
 - mc12b: 0.5×10^9 events
 - mc12c: 0.2×10^9 events → total: 6.6×10^9 events

- total of 6.2×10^9 full and 5.1×10^9 fast simulation events

ATLAS Grid Resources

> grid resources

- Tier0: CERN
- Tier1: 10 (11) sites
- Tier2: ~70 sites
- Tier3: ~20 sites

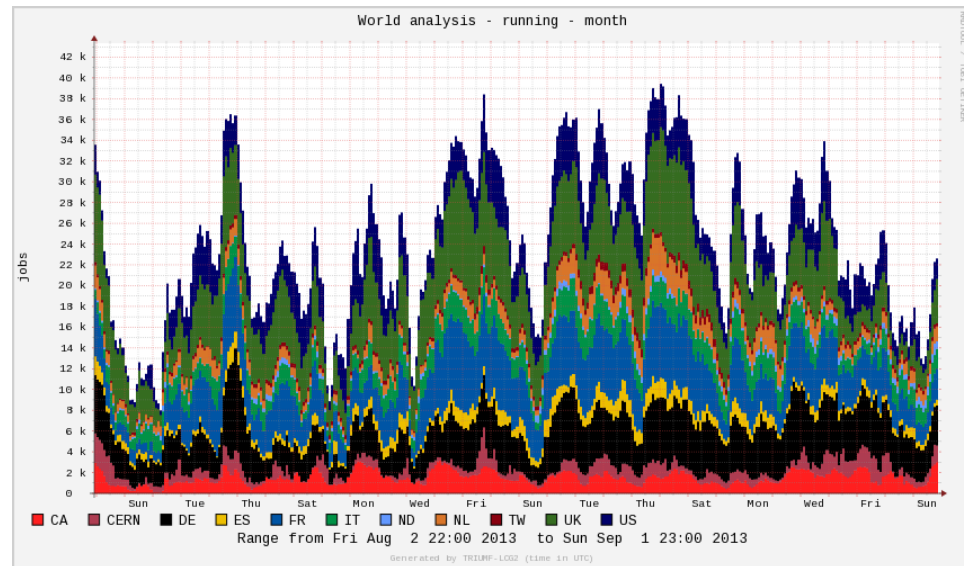
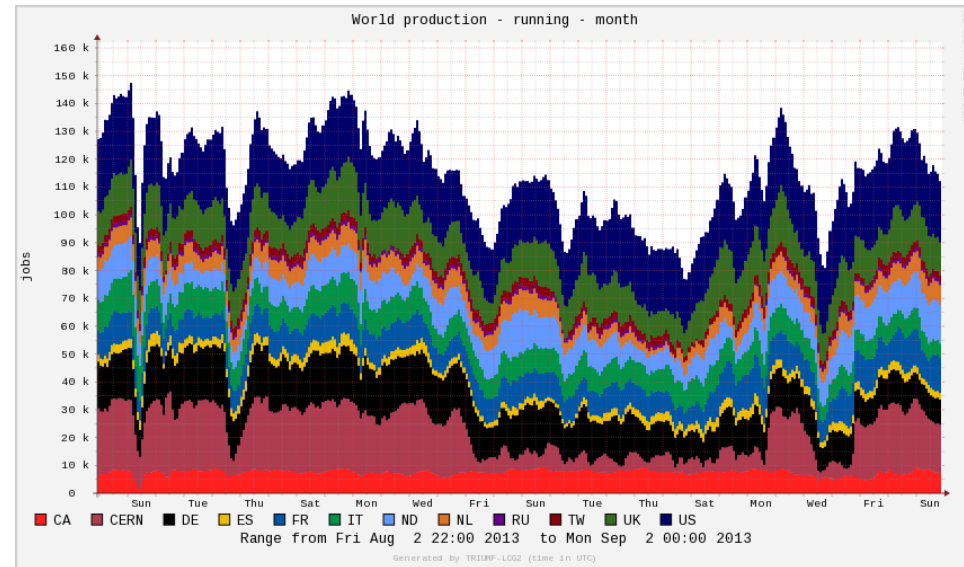
→ ~100 000 single core slots for MC production

> clouds

- Amazon E2 cloud
- Google Computing Engine cloud
- Open Clouds

> opportunistic sites

- online trigger farm (16 000 slots)
- High Performance Computing

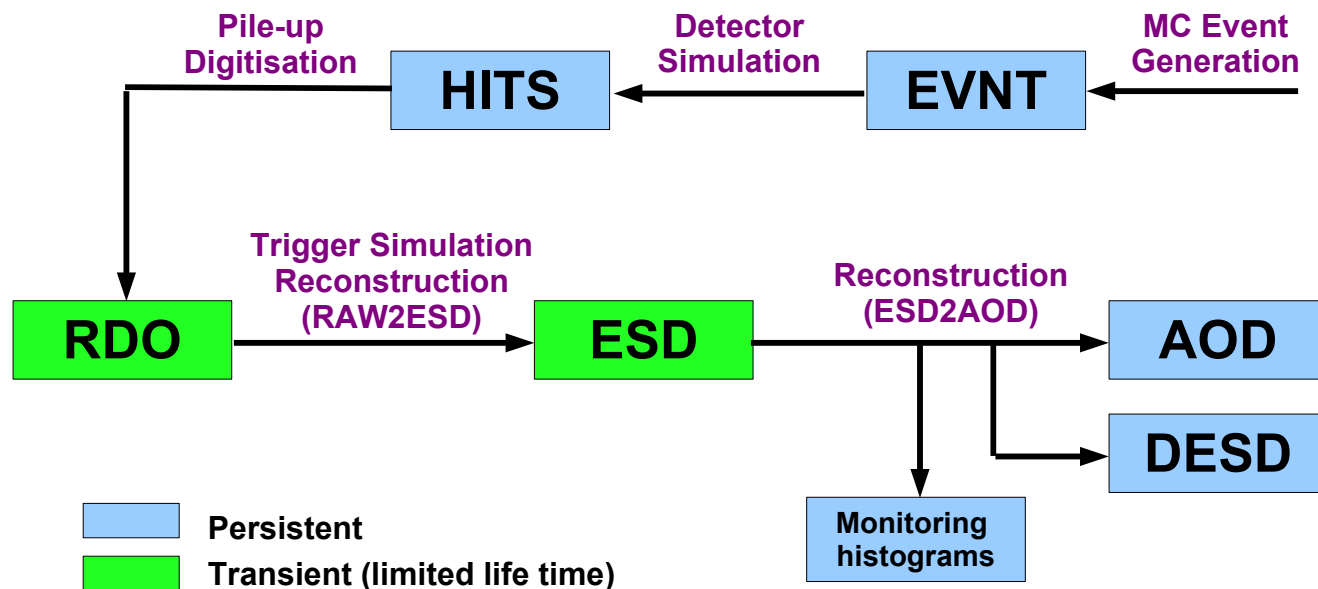


Monte Carlo Production

MC Production Steps

- > event generation
- > simulation
- > digitisation
- > reconstruction

ATLAS Monte Carlo Simulation Flow



Monte Carlo Production Chain Details

- the implementation of the Monte Carlo production chain is a compromise between different requirements/boundary conditions and evolves with time
- which data should be stored
 - event generation: store full event on disk for reuse (use same events in fast and full simulation)
 - simulation: keep HITS from time intensive simulation (full simulation) on tape for possible reuse (re-processing/updated pileup distribution)
 - reconstruction: keep AOD as the primary source for data analysis
- software development and release cycle
 - decoupled cycle for simulation and reconstruction (where to put digitisation?) reconstruction driven by data taking, simulation by the need to have sufficient number of events early in data taking
- grid
 - one configuration fits all
 - running time (optimal 8h, up to 2-4d available, factor 2 variation in CPU performance)
 - data size (optimise number of files and file size for data transfers and tape storage)

Event Generation

Event Generation

> ~30 generators used in ATLAS

- framework integrated generators
- stand-alone generators

> event generation work flows

- single step generation: Pythia6/8, Herwig(++), Sherpa
- two-step generation: parton level generator coupled via LHEF files to framework generator for hadronisation (Pythia(6/8), Herwig(++))
 - > default configuration: external, pre-made 4-vectors uploaded to the grid
 - > on-the-fly configuration: run external generator before hadronisation in the same job
- optional: run filter algorithms to populate dedicated phase space (for example number of stable leptons, pT sliced, invariant mass, B-decays, ...)

> ATLAS spends a significant amount of resources in event generation!

Event Generation - Performance

> requested samples very diverse

- 50 different generator combination in mc12 campaign
- ~34 thousand different samples produced in mc12 campaign

> job characteristics

- 5000 events per job → ~100 MB output file size
- low memory requirements: < 0.5-1 GB
- running time per job varies from
 - > a few minutes for simple final states/hadronisation of external 4-vectors
 - > hours or days for complex final states or low filter efficiencies
 - > number of events needs to be adjusted for optimal running time of 8 hours

> performance improvements:

- aiming for more automatisation and therefore improving turn around time
- on-the-fly generator setups: avoid storing 4-vector input files on the grid
- use pre-made integration files (Sherpa, Alpgen, MadGraph): reduce running time

Event Generation - Problems and Ideas

> long running jobs:

- avoid phase space integration by using integration grids → Sherpa, Alpgen, Madgraph
- integration grids need to be prepared manual → do this automatically (needs MPI for Sherpa)
- in case of large number of different samples (eg SUSY) use multi core for phase space integration and event generation (Madgraph)

> small number of output events

- jobs are running long and write out small files (dijet slicing, stable leptons, light/charm/b-jets, ...) → run on multi cores to get larger files or merge files

> Sherpa

- complex final states have a large initialisation time when using integration grids → can this be run parallel? Then also run event generation parallel

> Alpgen

- hard to tell Alpgen how many unweighted events to generate → run multicore and merge

> why multi core?

- multi core is another dimension of slicing and merging, but HPC resources might come for free

Simulation Flavours

- G4 full simulation:
 - every stable particle is tracked through the ATLAS geometry
 - one event takes ~5 minutes → major simulation time spent in calorimeters
- G4 full simulation with Frozen Showers (FS) in calorimeters: 25% speed up in mc12
 - showers are tracked down to very low energy by G4 → stop showering at a threshold and substitute each end particle by a pre-made list of energy deposits
 - residual difference between Frozen Shower and plain G4 is below 1%
 - frozen showers in the forward calorimeters as default in mc11/mc12 including upgrade production
- AtlFast-II (AF-II): factor 10 speed up in mc12
 - parametrise all particles except muons in the calorimeters
 - do not simulate particles except muons in the calorimeter
 - parametrise non-simulated particles before the digitisation step
 - in production since late mc10

Simulation Flavours - Technical Details

> common problems:

- where to store additional data and how to find it
- templates/parametrisation depends on detector geometry and needs to be prepared

> Frozen Showers

- where to store pre-made libraries?
 - > data area attached to full release build (at site, no version, slow update cycle)
 - > dedicated grid dataset (needs to be distributed to sites to avoid large number of transfers)
 - > store in release area
- choice depends on size and format (signal input ~100MB)
 - > ascii files in compressed tar ball (~360MB)
 - > ROOT files (~22MB)

> AF-II

- using conditions DB infrastructure for storing parametrisations
 - > use global conditions tag to introduce version and connection to payload data
 - > payload data from grid datasets or from CVMFS

Simulation - Performance

> job characteristics

- low memory requirement: ~1 GB
- run time per (averaged over grid cpus)
 - > G4 full simulation: 335 s/evt
 - > G4 full simulation with frozen showers: 250 s/evt
 - > AtI Fast-II: 20 s/evt

> simple speed ups

- tcmalloc: fast, multi-threaded malloc() from google perftools
- run in 64 bit → better performance while slightly increasing the memory
- Intel math library
- modern random number generator: SIMD-oriented Fast Mersenne Twister
- run on SL6
- new compiler
- ...

Multi Core Utilisation in the ATLAS Software Framework

- number of cores increases faster than memory size
- needed memory for digitisation and reconstruction increases with pileup
- athenaMP: multi core implementation of the ATLAS software
 - in a nut shell: fork main process into n clients which will process full events
→ optimising memory
 - performance improvement can come from better ratio between initialisation and processing time (needs additional (fast) output merging)
- threading of individual algorithms should bring some performance improvements

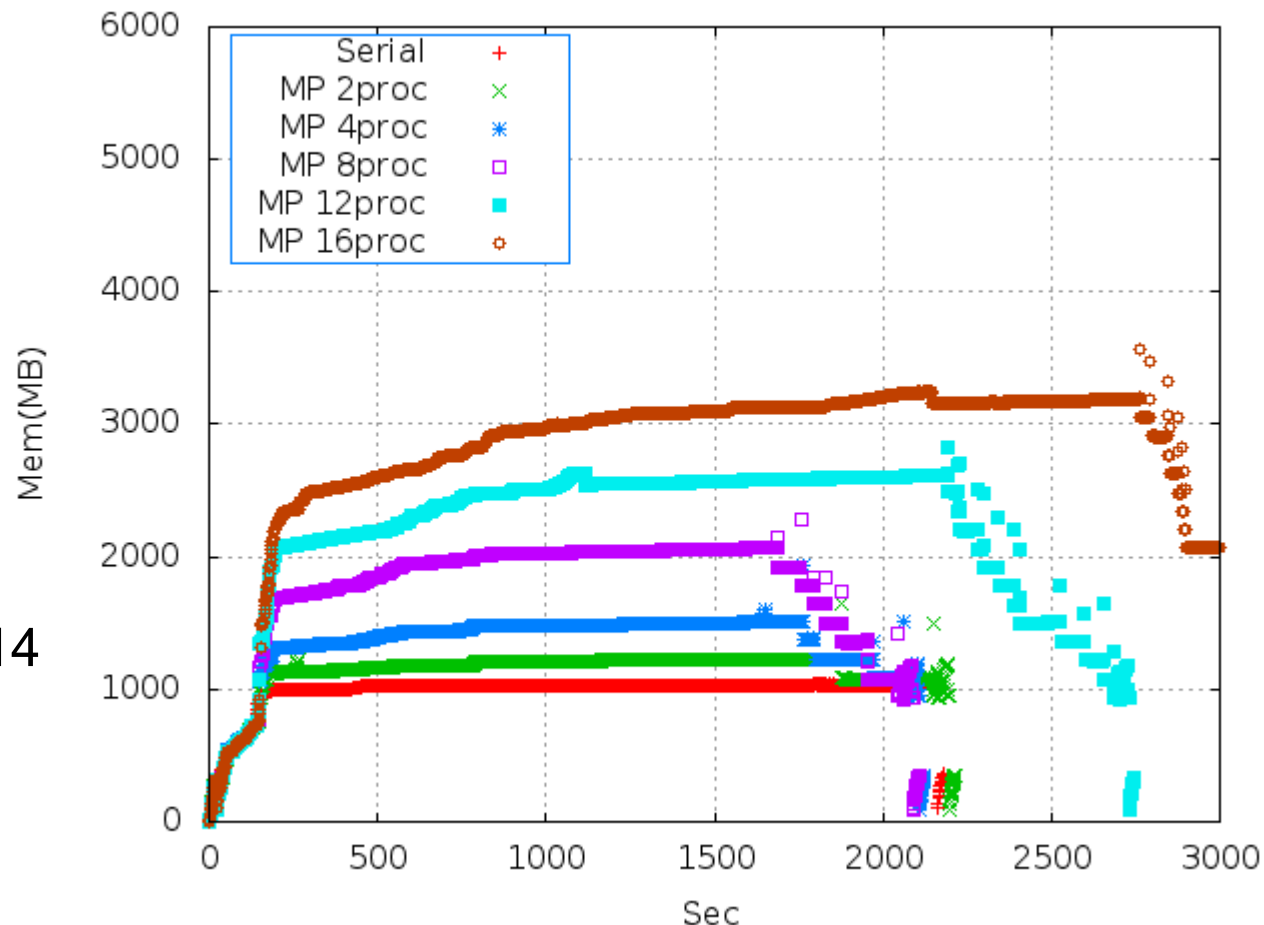
Simulation - Multi Core Utilisation with athenaMP

- single core: 1.0GB
- double core: 1.2GB
- ...
- 8 cores: 2.6GB
- 0.8GB + 0.16GB/core

➤ athenaMP used in 2014 simulation

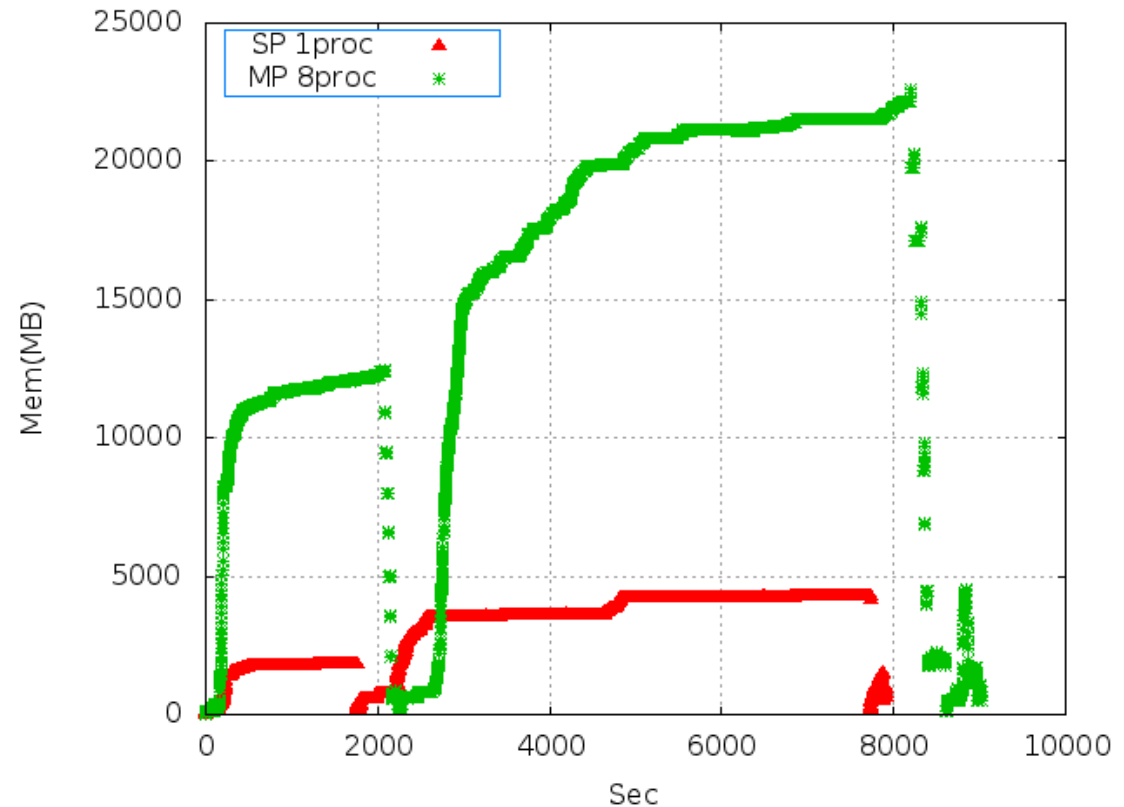
➤ production scenarios

- reducing number of job
- back filling multi-core slots
- high performance computing resources



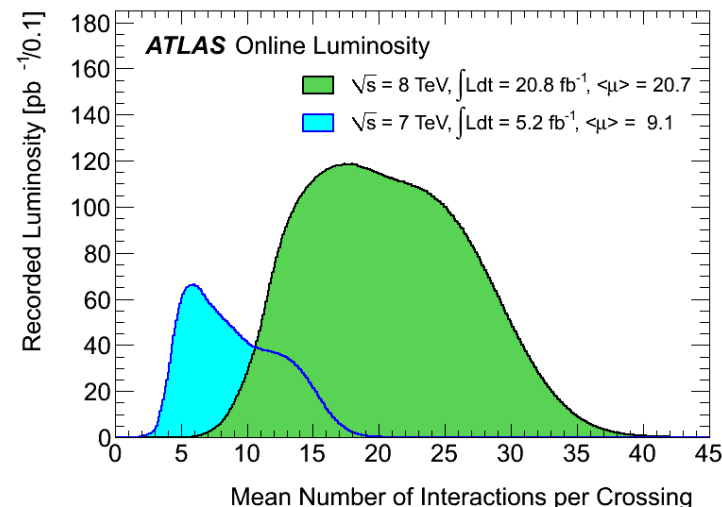
athenaMP - Memory Sharing in Digitisation+Reconstruction

- running in 64 bit
- single core: 4.3GB
- 8 cores: 22.6GB
- 2.8GB/core
- better than 4GB/core but aim is 2GB/core
- athenaMP validated for digitisation+reconstruction
- production scenarios
 - reducing memory consumption
 - reducing number of jobs



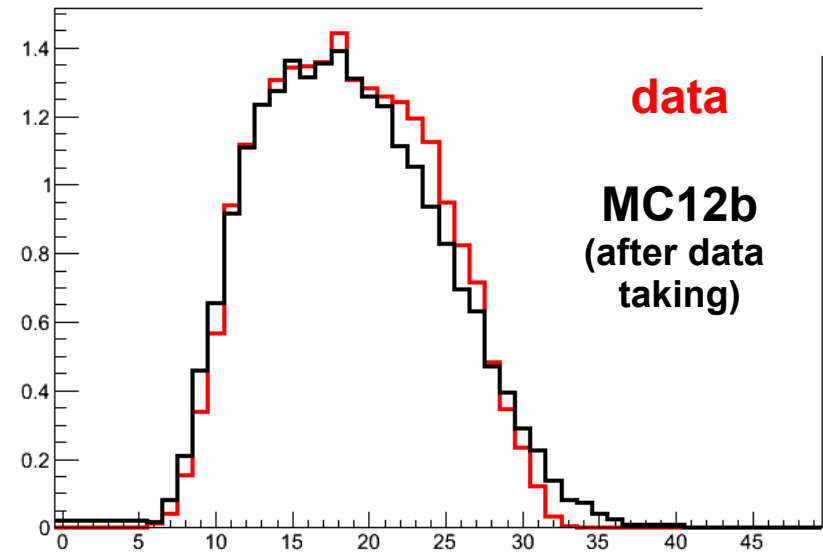
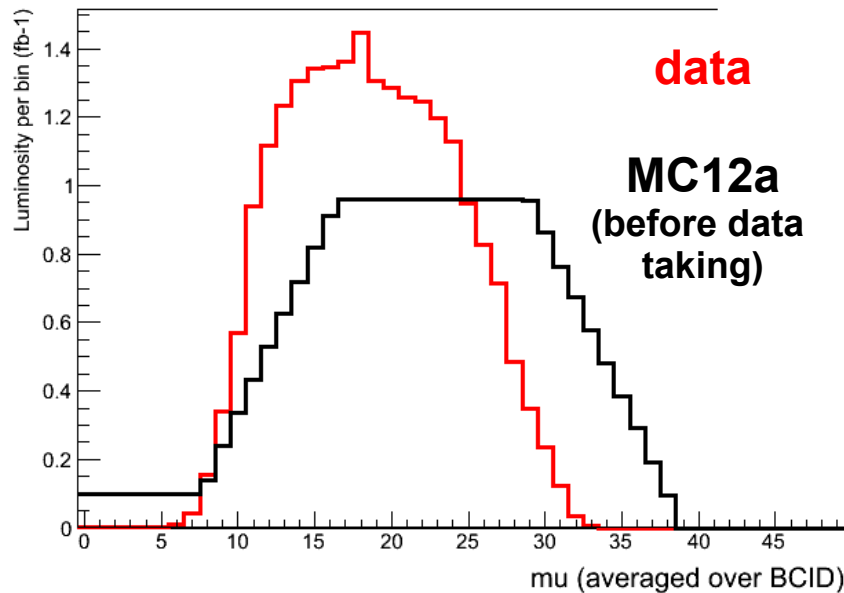
Digitisation

- simulate detector readout
- simulate pile-up contributions (multiple pp interactions on top of hard scatter event)
- overlay a number of pre-simulated minimum bias events on each signal event
 - $\langle\mu\rangle$ average number of additional pp collisions
 - fixed $\langle\mu\rangle$ (for performance studies)
 - pre-defined $\langle\mu\rangle$ profile (default for physics samples)
 - sample given $\langle\mu\rangle$ profile over 5000 events
→ small samples should be multiple of 5000 events
- optimise pile-up event storage and access
 - cache pile-up events in memory → memory intensive
 - flush memory early and re-load from disk on demand → I/O and CPU intensive
- minimum bias pile-up samples
 - separate into low-Q and high-Q ($Q=35\text{GeV}$) samples to allow for frequent re-use of low-Q events per job and limit re-use of within one sample



MC12 Pileup Simulation

- pile-up profile in MC matched to observed distribution in data if possible

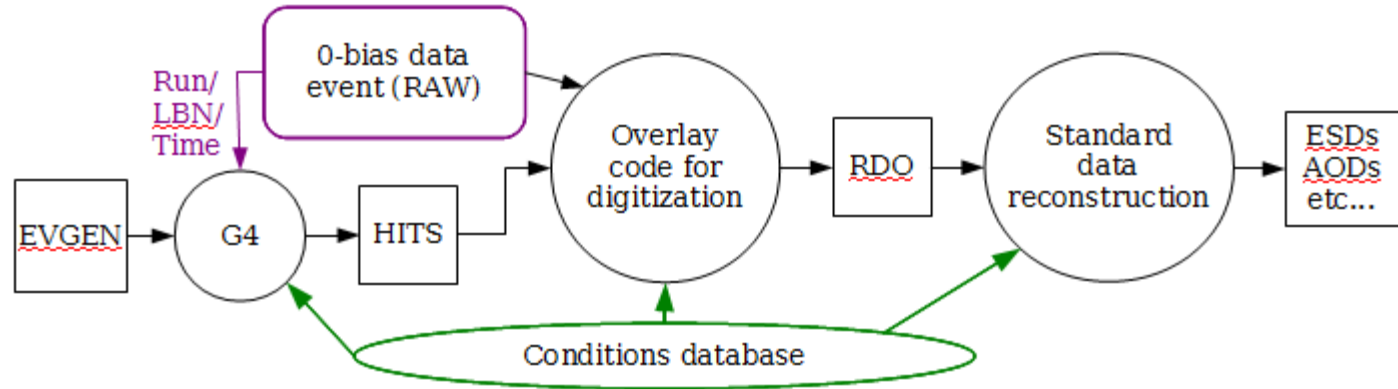


- mc12 pile-up sample configuration

- $\langle \mu \rangle$ profile samples from 0 to 40, with a mean of $\langle \mu \rangle = 20$
- 10M low/high-Q G4 simulated events (1.5/4.8 TB = 6.3 TB) → 5000/500 events per file
- 500 events per job: one signal file, 5 low/high-Q files → 4.8 GB of input files per job (100 events per job: one signal file, 1 low/high-Q file → 1.1 GB of input files per job)
- distribute minimum bias pile-up sample to T1 and larger T2 sites → 0.3-0.4 PB total

Zero-Bias Data Overlay (Embedding)

- improve pile-up simulation by using zero-bias data events



- conditions and beam spot need to be adjusted for each signal event to the corresponding zero bias event → run simulation and overlay in one job

- mc12 overlay configuration

- ensure a representative pile-up sampling in sets of 50 000 events
- 100 events per job: one signal file, 1 overlay file → 0.4 GB of input files per job
- 2012 pp 8 TeV zero-bias sample contains 50 million events (160 TB)
→ grid distribution needs to be improved

- overlay heavily used by heavy ion analysis for PbPb and pPb collisions as simulating heavy ion collisions is difficult and very CPU intensive

Zero-Bias Data Overlay (Embedding)

> performance improvements

- only merge zero bias data (already digitised)

> usage

- 8 TeV pp: under validation for general use, first specialised use cases
- HI running at 2 and 5 TeV: heavily used (~200M events for 5 TeV)

- different simulation flavours need different energy calibrations
 - mainly distinguish between full simulation and fast simulation
- what do to when mixing sim flavours (or data) in one event?
 - full sim min bias pileup events and fast sim signal event
 - data zero bias event and full sim signal event
 - different sim flavours within one events (ISF)

Reconstruction

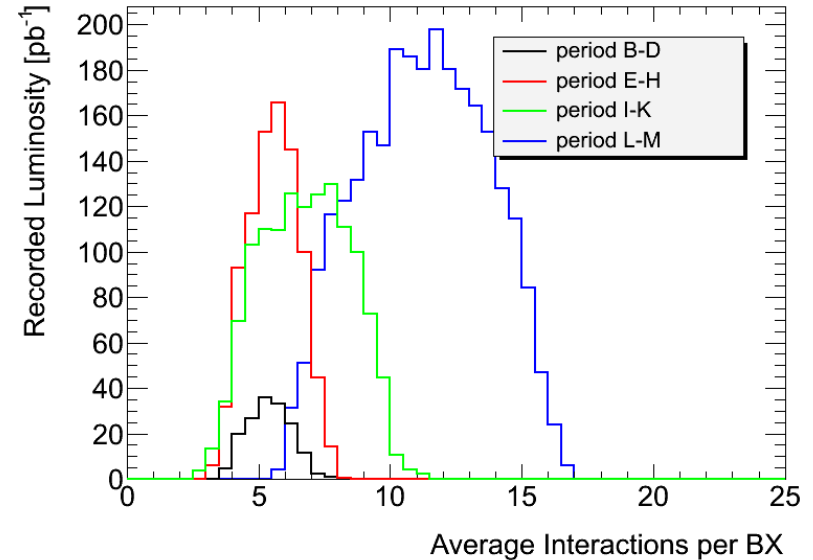
- reconstruct simulated events in the same way as data
- trigger simulation
- two step process:
 - RAWtoESD: main reconstruction → output is Event Summary Data (ESD)
 - ESDtoAOD: fast slimming process → output is Analysis Object Data (AOD)
- job characteristics
 - 500 events per job → ~220 MB output file size → merged up to 5000 events (~2.2 GB file size) for better grid transfer, processing and tape storage
 - high memory usage:
 - 3.6 - 3.8 GB in 32 bit
 - 64 bit would exceed the 4 GB (grid queue limits)

Upgrade Production

- preparations for Run 2, Phase 2 and beyond
 - planed detector upgrades
 - ATLAS+IBL (Insertable b-Layer (pixel detector extension) for Run 2)
 - ATLAS+ITK (silicon only inner tracker upgrade for Phase 2)
 - machine constraints: 50ns or 25ns bunch spacing and pile-up level
- ATLAS+IBL configuration: 25/50ns and $\langle\mu\rangle=20, 40, 60, 80$
 - simulation time increases due to higher centre-of-mass energy/more particles per event
 - higher pile-up level increases memory usage, especially in reconstruction and trigger
 - reduce trigger menu (<4 GB for 60@25ns and 80@50ns)
 - run on dedicated high memory queues (<6GB for 80@25ns)
 - for Run 2 simulate trigger between digitisation and reconstruction
 - running time: 100 s/evt for $\mu=20$; $\mu=40 \rightarrow \times 2.2$; $\mu=60 \rightarrow \times 1.8$
- ATLAS+ITK configurations: 25ns and $\langle\mu\rangle=80, 140, 200, (300, 400, 500)$
 - reconstruction stays (well) below 4GB
(trigger simulation not yet supported and no transition radiation tracker)
 - 25 events digitisation+reconstruction (top pairs): $\langle\mu\rangle=200: 14h \rightarrow \langle\mu\rangle=500: 48h$

Joining Steps in One Job

- joining two or more steps from the full MC production chain can be useful
 - digitisation+reconstruction (default)
 - avoid storing large digitisation output on the grid
 - easier for multi-period pile-up and trigger simulation (mc11)
 - avoids production system delays
 - fast simulation
 - avoid storing intermediate outputs
→ simplify data management
 - fast simulation
→ small loose in CPU in case of re-reconstruction
 - need to fit all steps into one job: secondary files, multi core processing, ..
 - as number of events per job is going down initialisation time might become significant



- real life experience with speed-ups and new trends in MC production
 - event generation, simulation, digitisation, reconstruction
 - usage of free resources (HPC) is faster production
 - improving production efficiency is faster production

Campaigns

- MC production campaigns correspond to data taking periods with same conditions
 - centre-of-mass energy, detector configuration, conditions, ...

- Major MC production campaigns
 - mc11: simulation configuration for 7 TeV in 2011
 - mc11a: digitisation+reconstruction configuration with Pythia 8 pile-up sample, estimated beam spot and pile-up profile based on three run periods
 - mc11b: same as mc11a with updated pile-up profile/conditions based on four run periods and two trigger menus
 - mc11c: same as mc11b with Pythia 6 pile-up sample
 - mc12: simulation configuration for 8 TeV in 2012
 - mc12a: digitisation+reconstruction configuration with Pythia 8 pile-up sample, estimated pile-up profile and beam spot based on 2011 data
 - mc12b: same as mc12a with beam spot and pile-up profile from data
 - mc12c: improved geometry description for precision measurements: simulation based on mc12 and digitisation+reconstruction based on mc12b