

# HEPiX CPU Benchmarking WG<sup>\*</sup>

## Status update

D. Giordano (CERN IT)

on behalf of the HEPiX CPU Benchmarking WG  
[hepix-cpu-benchmark@hepix.org](mailto:hepix-cpu-benchmark@hepix.org)

HEPiX Autumn 2019 Workshop

<sup>[\*]</sup> In charge of defining and maintaining a consistent and reproducible CPU benchmark to describe WLCG experiment requirements

Is HS06 still representative of the WLCG WLSs?

Which benchmark shall WLCG adopt after HS06?

# Is HS06 still representative of the WLCG WLSs?

## Which benchmark shall WLCG adopt after HS06?

- By end of day on **January 9, 2018 US Eastern Time**, SPEC will retire SPEC CPU2006.
  - After this day, further submissions not already under review will not be published by SPEC and **technical support for SPEC CPU2006 will end.**
  - For publication on SPEC's website by January 9th, 2018, results need to be submitted to SPEC by the **December 26, 2017 3AM US Eastern Time** submission deadline. Note that, per above, corresponding SPEC CPU2017 results are also needed.

<https://www.spec.org/cpu2006/>

# Follow the HEP sw evolution

## 1980's

MIPS (M Instr Per Sec)  
VUPS (VAX units)  
CERN units

## 1990's – 2000's

SI2k (SPEC INT 2000)  
INTEGER benchmarks  
200 MB footprint

## 2009

HS06 (SPEC CPU 2006 all\_cpp)  
INTEGER + FP benchmarks  
1 GB footprint  
32-bit  
x86 servers  
single-threaded/process on multi-core

## 2019

2 GB footprint (or more)  
64-bit  
multi-threaded, multi-process  
multi-core, many-core  
vectorization (SSE, ... AVX512)  
x86 servers, HPCs  
ARM, Power9, GPUs...?

- ❑ HEP software (and computing) evolves... so do HEP CPU benchmarks!
- ❑ As time goes by, *WLCG computing is becoming more and more heterogeneous*
- ❑ One of the challenges is how to summarize performance using a **single number**
  - Unfortunately, this is needed at least for accounting purposes

# A reminder about HS06

- ❑ Subset of SPEC CPU® 2006 benchmark
  - SPEC's industry-standardized, CPU-intensive benchmark suite, stressing a system's processor, memory subsystem and compiler.
- ❑ HS06 is suite of 7 C++ benchmarks
  - In 2009, proven **high correlation** with experiment workloads  
*<<CPP showed a good match with average lxbatch e.g. for FP+SIMD, Loads and Stores and Mispredicted Branches>> [1]*
  - Execution time of the full HS06 suite: O(4h)

Bmk	Int vs Float	Description
444.namd	CF	92224 atom simulation of apolipoprotein A-I
447.dealll	CF	Numerical Solution of Partial Differential Equations using the Adaptive Finite Element Method
450.soplex	CF	Solves a linear program using the Simplex algorithm
453.povray	CF	A ray-tracer. Ray-tracing is a rendering technique that calculates an image of a scene by simulating the way rays of light travel in the real world
471.omnetpp	CINT	Discrete event simulation of a large Ethernet network.
473.astar	CINT	Derived from a portable 2D path-finding library that is used in game's AI
483.xalanbmk	CINT	XSLT processor for transforming XML documents into HTML, text, or other XML document types

Correlation	Generation	Simulation	Reconstruction	Total
Atlas	0.9969	0.9963	0.9960	0.9968
Alice pp MinBias	0.9994		0.9832	0.9988
Alice PbPb	0.9984		0.9880	0.9996
LhcB	0.9987			
CMS HiggsZZ	0.9982		0.9987	0.9983
CMS MinBias	0.9982		0.9974	0.9974
CMS QCD 80 120	0.9988		0.9987	0.9988
CMS Single Electron	0.9987		0.9942	0.9981
CMS Single MuMinus	0.9986		0.9926	0.9970
CMS Single PiMinus	0.9955		0.9693	0.9955
CMS TTbar	0.9985		0.9589	0.9987

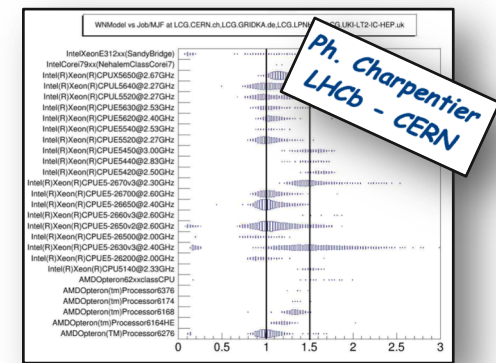
[1] Correlation of HEP-SPEC06 with several kinds of applications and different experiments

[1] "A comparison of HEP code with SPEC benchmarks on multi-core worker nodes" *J. Phys.: Conf. Ser.* 219 (2010) 052009 CHEP-09

# WG activities - short overview (1)

Fall 2015 (@ GDBs)

- Some experiments start reporting performance deviation respect to HS06
- Attention goes to **fast** benchmarks
  - LHCb DB12, Atlas KV, Root stress-test



<https://indico.cern.ch/event/319754>

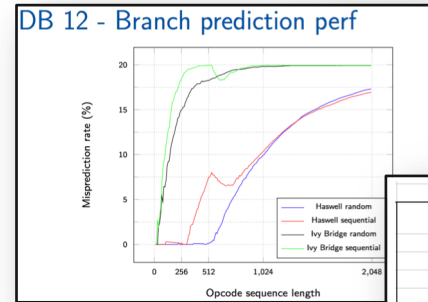
Spring 2016 (HEPiX)

- HEPiX CPU Benchmarking WG restarts coordinated studies



Up to Spring 2017 (WLCG workshop)

- Detailed analysis of fast benchmarks
  - In bare metal servers as well as VMs
  - ☞ Excluded DB12 for lack of robustness
- Understand systematics on HS06
  - ☞ E.g. 32-bits Vs 64-bits correction factor
- Increasing expectations in the future SPEC CPU benchmark



DB12 python		32 procs	
OS	version	ratio_to_pytho n2.6	ratio 32/16
slc6-base	Python 2.6.6	1	1.00
CC7	Python 2.7.5	1.09	1.02
cc7-base	Python 2.7.5	1.09	1.04
python:2.7	Python 2.7.13	1.18	1.01
python:3	Python 3.6.0	1.05	0.98

<https://indico.cern.ch/event/609911/contributions/2620190/>

# WG activities - short overview (2)

June 2017

- SPEC CPU 2017 finally available
  - Will it solve the HS06 “crisis”?

Up to Summer 2018 (CHEP)

- Comprehensive comparisons of HEP Workloads and HS06

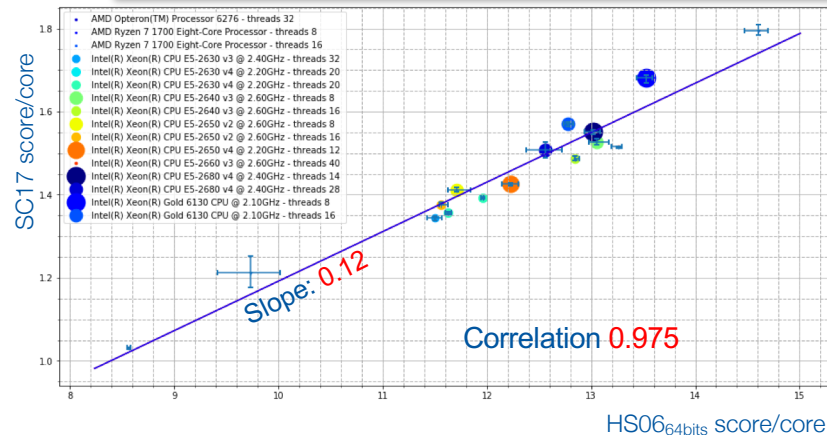
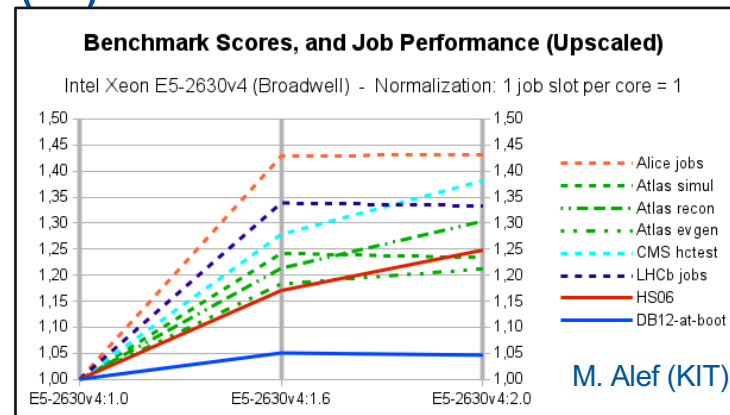
☞ *Confirmed discrepancies for Alice and LHCb*

- SPEC CPU 2017

- Detailed studies w.r.t. HS06

☞ *Extremely high correlation of the 2 benchmark suites*

☞ *No advantage is moving to SPEC CPU 2017*



HS06<sub>64bits</sub> score/core

# SPEC CPU2017

<https://www.spec.org/cpu2017/Docs/index.html#benchmarks>

## SPEC releases major new CPU benchmark suite

### The SPEC CPU2017 benchmark suite features updated and improved workloads, use of OpenMP to accommodate more cores and threads, and optional metric for measuring power consumption

Gainesville, Va., June 20, 2017 -- The Standard Performance Evaluation Corp. (SPEC) today released the SPEC CPU2017 benchmark suite, an all-new version of the non-profit group's software for evaluating compute-intensive performance across a wide range of hardware systems.

The SPEC CPU2017 benchmark suite is the first major update of the worldwide standard CPU performance evaluation software in more than 10 years. The new suite includes updated and improved workloads with increased size and complexity, the use of OpenMP to allow performance measurement for parallelized systems with multiple cores and threads, and an optional metric for measuring power consumption.

Current SPEC CPU subcommittee members include AMD, ARM, Dell, Fujitsu, HPE, IBM, Inspur, Intel, Nvidia and Oracle.

Larger suite, more complex code, shaped for multi-core and multi-threads

#### The Benchmarks

SPEC CPU2017 has 43 benchmarks, organized into 4 suites:

SPECrate 2017 Integer      SPECSpeed 2017 Integer  
SPECrate 2017 Floating Point      SPECSpeed 2017 Floating Point

Benchmark pairs shown as:

5nn.benchmark\_r / 6nn.benchmark\_s

are similar to each other. Differences include: compile flags; workload sizes; and run rules. See: [\[OpenMP\]](#) [\[memory\]](#) [\[rules\]](#)

SPECrate 2017 Integer	SPECSpeed 2017 Integer	Language <sup>[1]</sup>	KLOC <sup>[2]</sup>	Application Area
500.perlbench_r	600.perlbench_s	C	362	Perl interpreter
502.gcc_r	602.gcc_s	C	1,304	GNU C compiler
505.mcf_r	605.mcf_s	C	3	Route planning
520.omnetpp_r	620.omnetpp_s	C++	134	Discrete Event simulation - computer network
523.xalancbmk_r	623.xalancbmk_s	C++	520	XML to HTML conversion via XSLT
525.x264_r	625.x264_s	C	96	Video compression
531.deepsjeng_r	631.deepsjeng_s	C++	10	Artificial Intelligence: alpha-beta tree search (Chess)
541.leela_r	641.leela_s	C++	21	Artificial Intelligence: Monte Carlo tree search (Go)
548.exchange2_r	648.exchange2_s	Fortran	1	Artificial Intelligence: recursive solution generator (Sudoku)
557.xz_r	657.xz_s	C	33	General data compression

SPECrate 2017 Floating Point	SPECSpeed 2017 Floating Point	Language <sup>[1]</sup>	KLOC <sup>[2]</sup>	Application Area
503.bwaves_r	603.bwaves_s	Fortran	1	Explosion modeling
507.cactuBSSN_r	607.cactuBSSN_s	C++, C, Fortran	257	Physics: relativity
508.namd_r		C++	8	Molecular dynamics
510.parest_r		C++	427	Biomedical imaging: optical tomography with finite elements
511.povray_r		C++, C	170	Ray tracing
519.lbm_r	619.lbm_s	C	1	Fluid dynamics
521.wrf_r	621.wrf_s	Fortran, C	991	Weather forecasting
526.blender_r		C++, C	1,577	3D rendering and animation
527.cam4_r	627.cam4_s	Fortran, C	407	Atmosphere modeling
	628.pop2_s	Fortran, C	338	Wide-scale ocean modeling (climate level)
538.imagick_r	638.imagick_s	C	259	Image manipulation
544.nab_r	644.nab_s	C	24	Molecular dynamics
549.fotonik3d_r	649.fotonik3d_s	Fortran	14	Computational Electromagnetics
554.roms_r	654.roms_s	Fortran	210	Regional ocean modeling

[1] For multi-language benchmarks, the first one listed determines library and link options ([details](#))

[2] KLOC = line count (including comments/whitespace) for source files used in a build / 1000

<https://www.spec.org/cpu2017/press/release.html>

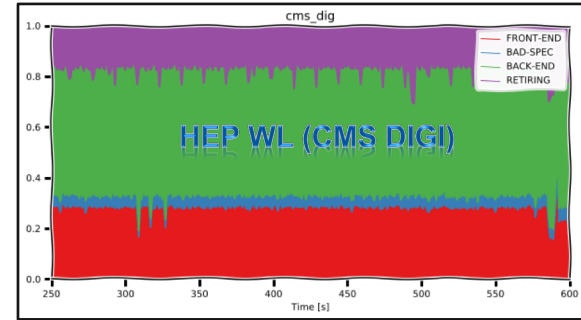
same application area as in HS06



# WG activities - short overview (3)

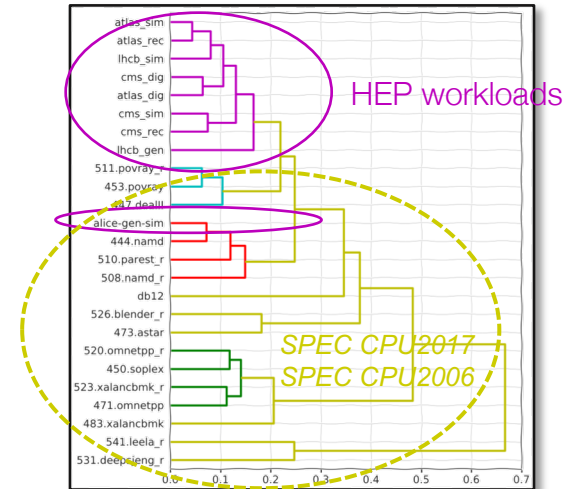
## Summer 2018

- Proposal of building a set of **HEP reference workloads** (WLCG MB)
  - Enable feature studies of the experiments' workloads
  - Build a HEP benchmark suite



## Fall 2018

- Collect instructions from LHC experiments to run reference WLS
- **Prototype** the build of HEP reference benchmarks in containers
- Studies on **hardware performance counters** (using Trident)
  - ☞ *HEP WLS have same characteristics and differ more respect to HS06 and SPEC CPU 2017 workloads*



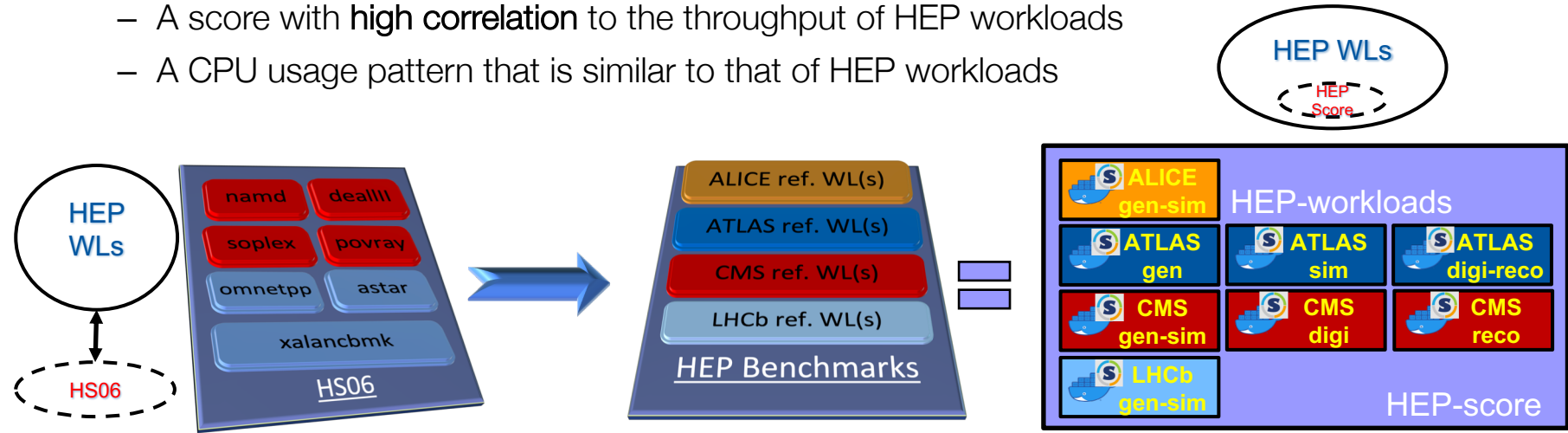
## 2019

- Start the **HEP Benchmarks project**

DENDROGRAM OF WL'S SIMILARITY

# Benchmarking CPUs using HEP workloads

- ❑ *By construction*, using HEP workloads directly is guaranteed to give
  - A score with **high correlation** to the throughput of HEP workloads
  - A CPU usage pattern that is similar to that of HEP workloads



- ❑ It is NOT a replacement of HS06 for today
  - This is the future approach to be adopted when the correlation with HS06 will be definitely broken

# HEP Benchmarks: project organization

## ❑ Current contributors

– M. Alef, J. M. Barbet, *O. Datskova*, *D. Giordano*, *C. Grigoras*, *C. Hollowell*, *M. Javurkova-Pagacova*, *V. Khristenko*, *D. Lange*, M. Michelotto, *L. Rinaldi*, *E. Santorinaiou*, *A. Sciabà*, *A. Valassi*

- Core team for common infrastructure development and overall testing
- Experiment contacts for experiment-specific software, workloads, metrics

- *Developers*
- **Exp. experts**

– Contact with industry, access to new HW

- M. Girone (CERN IT-Openlab), L. Atzori (CERN IT-Procurement)

## ❑ Track work progress via [Jira Project](#) and [Twiki](#)

– Weekly meetings and Jira Sprint reviews

# Project repositories

Three repositories @ CERN Gitlab

– *hep-workloads*

- Common build infrastructure
- Individual HEP workloads

– *hep-score*

- “Single-number” benchmark aggregator from several WLS
- Steer the WLS’ run

– *hep-benchmark-suite*

- Automate execution of multiple benchmarks
- Publish results

<https://gitlab.cern.ch/hep-benchmarks>

The screenshot shows the GitLab group page for "HEP-Benchmarks" (Group ID: 19914). The page lists three subgroups and projects: "hep-benchmark-suite", "hep-score", and "hep-workloads". To the right of the list, a dependency diagram is overlaid, showing three boxes: "HEP-workloads" at the bottom, "HEP-score" in the middle, and "HEP-benchmark-suite" at the top. Arrows point from "HEP-workloads" to "HEP-score", and from "HEP-score" to "HEP-benchmark-suite", indicating a dependency flow.

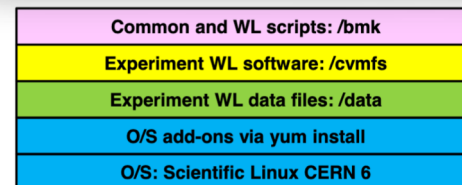
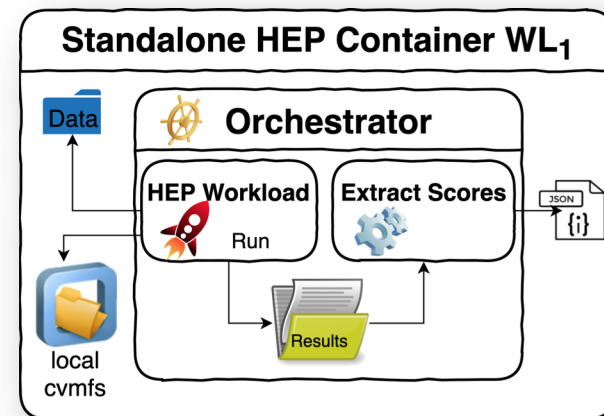
Project Name	Description
hep-benchmark-suite	Automates the programmatic execution of a number of benchmark
hep-score	Steer the HEP workloads' runs, collect results, compute the HEPsc
hep-workloads	Build standalone reference HEP workloads for benchmarking purp

# HEP Workloads

❑ **Standalone containers** encapsulating all and only the dependencies needed to run each workload as a benchmark

❑ Components of each HEP WL

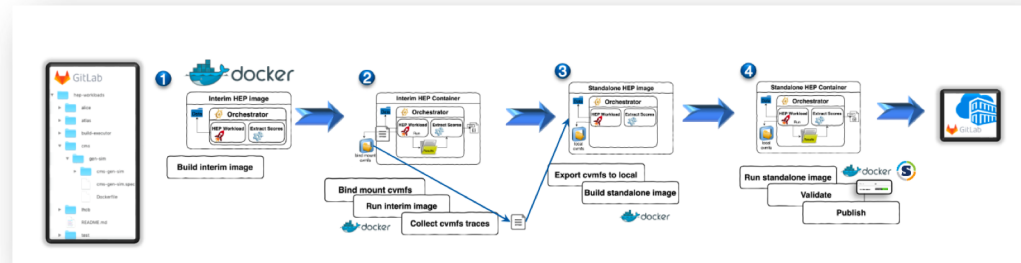
- SW repository (OS and CVMFS)
- Input data (event and conditions data)
- An orchestrator script (benchmark driver)
  - Sets the environment
  - Runs (many copies of) the application
    - Each copy may be multi-process or multi-threaded
  - Parses the output to generate scores (json)



Container images are made up of layers

# Solid Build Infrastructure

- Individual HEP workload container images are **built, tested** and **distributed** via gitlab
- Can be executed both via **Docker** and **Singularity**
- See Chris H. talk for more details



# Readiness of HEP Workloads

Summary (Status: pre-GDB 2019-10-08)



	ALICE				ATLAS				CMS			LHCb
	gen-sim	gen	sim	reco-digi	gen-sim	digi	reco	gen-sim	digi	reco	gen-sim	
Robustness	🚫 <sup>1</sup>	✅	✅	✅	✅	✅	✅	✅	✅	✅	✅	✅
Spread	🔧 <sup>2</sup>	🔧 <sup>3</sup>	✅	✅	✅	✅	✅	✅	✅	✅	✅	✅
Runtime	✅	✅	✅	✅	✅	✅	✅	✅	✅	✅	✅	✅
Memory	✅	✅	✅	🚫 <sup>4</sup>	✅	✅	✅	✅	✅	✅	✅	✅
Readiness	🚫	🔧	✅	🚫	✅	✅	✅	✅	✅	✅	✅	✅

✅ okay 
 🔧 fine tuning 
 🚫 to do / still in progress

Remarks:

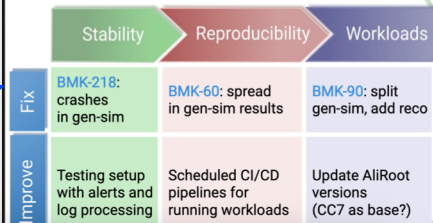
- Frequent crashes
- Use median of 3 or more runs to reduce spread
- Fine tuning in progress (e.g. more events, `-skipEvents`)
- Work in progress, improved container version will appear soon

14 pre-GDB 2019-10-08

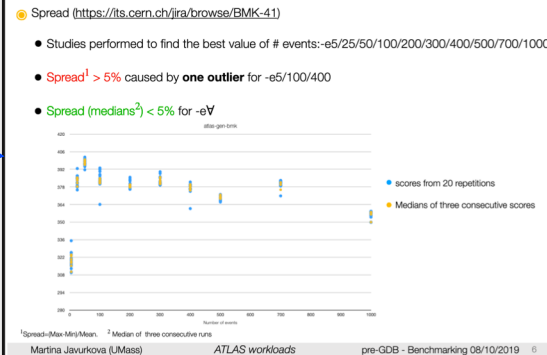
Manfred Alef: HEP Workloads as a Benchmark: Benchmark Validation (v2019-10)

Steinbuch Centre for Computing

## Ongoing Work...



## atlas-gen-bmk WL: validation (# events)



- Robustness
  - Error detection, error handling
- Spread of results (total scores) when repeating the benchmark
  - Spread:  $(\text{score\_max} - \text{score\_min}) / \text{score\_mean} \leq 5\%$
  - Possible option: median of 3 or more consecutive runs
- Runtime
- Memory consumption:
  - Benchmark must run on default WLCG WN
  - 2 GB RAM (physical memory) per job slot

Remark: HS06 – runtime (median of 3 iterations per benchmark): ~3h, spread: ≤2%, memory footprint (32bit): 1 GB per copy

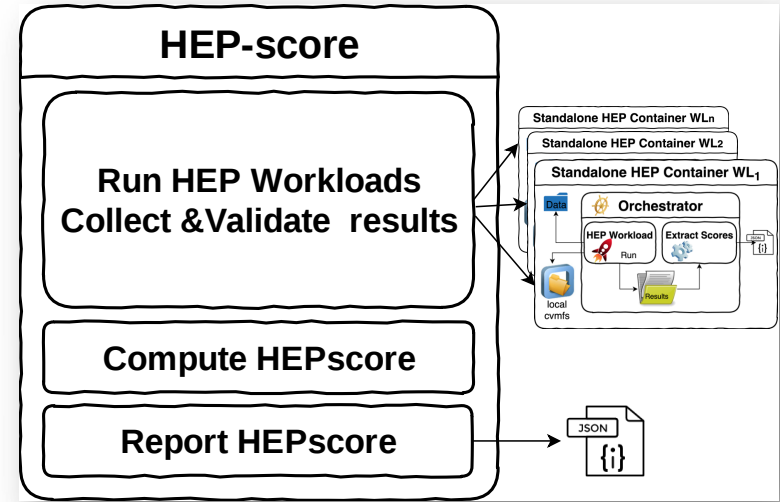
2 pre-GDB 2019-10-08

Manfred Alef: HEP Workloads as a Benchmark: Benchmark Validation (v2019-10)

Steinbuch Centre for Computing

# HEP Score

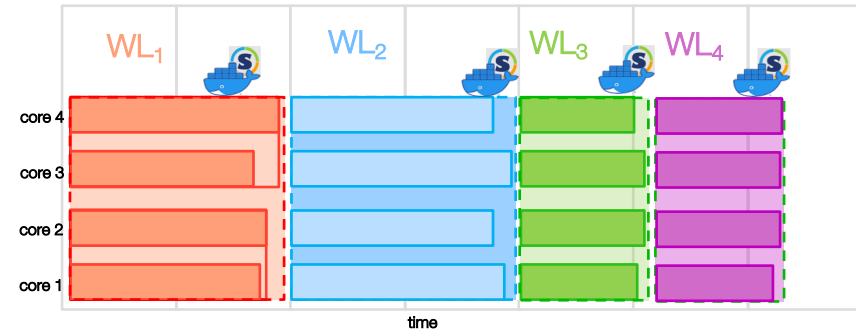
- ❑ Orchestrate the run of a series of HEP Workloads
- ❑ Compute & Report the HEPscore value
  - Default config. defines the HEPscore value
  - Other config. to perform specific studies
- ❑ HEP score does not include HEP Workloads' sw
  - HEP Workloads' sw is "*isolated*" in dedicated containers
  - Enable the utilization of **additional WLS**, as long as they comply with the expected API
  - Can be used by other domains





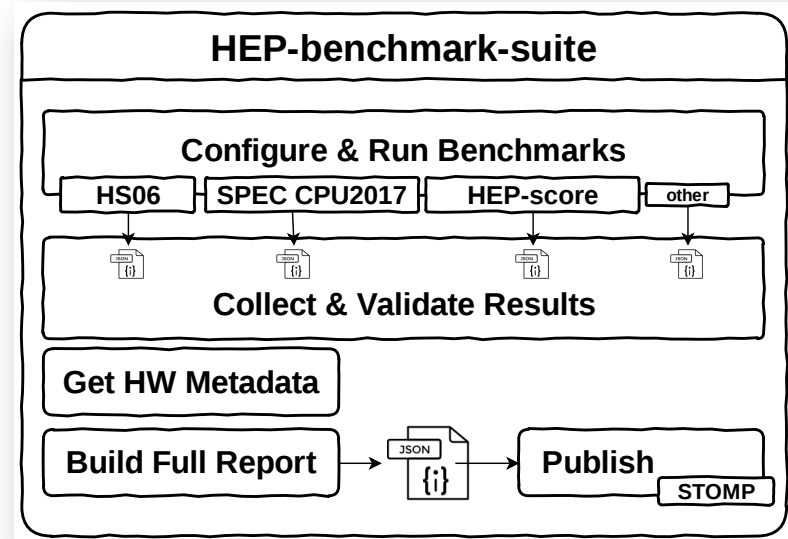
# HEP Score running mode

- ❑ HEP-score triggers HEP Workloads' runs in sequence
  - A **container** per WL
  - 3 times per WL, in sequence, and the **median** WL score is retained
- ❑ Each container runs the Experiment executable with a configurable number of threads (MT) or processes (MP)
- ❑ The available cores are saturated spawning a **computed** number of parallel copies
- ❑ The **score** of each WL is the cumulative event throughput of the running copies
  - When possible the initialization and finalization phases are excluded from the computation
  - Otherwise a long enough sequence of events is used
- ❑ A WL **speed factor** is computed as ratio of the WL score on a given machine w.r.t. the WL score obtained on a fixed reference machine
- ❑ HEPscore is the **geometric mean** of the WLs' **speed factor**



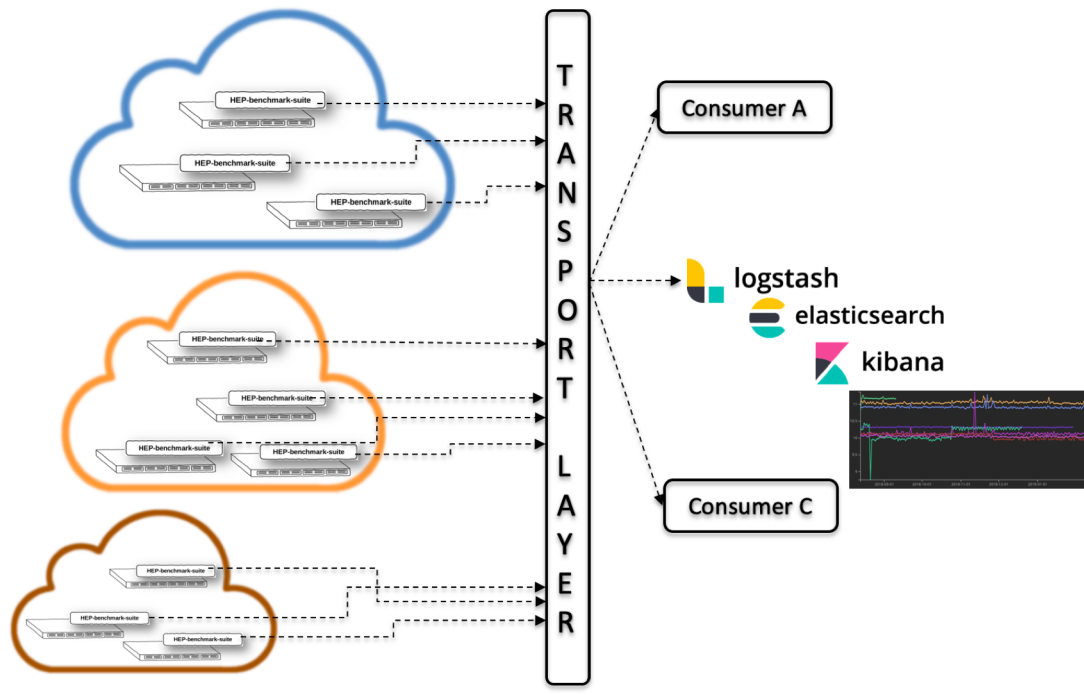
# HEP Benchmark Suite

- ❑ Control the execution of several benchmarks
  - HEP-score, SPEC CPU2017, HS06, KV, DB12, ...
  - NB: does NOT distribute sw under proprietary license (such as HS06), just requires that code to be pre-installed
- ❑ Simplify the sharing, tracking and comparison of results
  - Metadata track data-centre name, CPU model, host, kernel version, ...
- ❑ A sort of “Push the button” & Run & Get Results



# Centralise the benchmark data storage

- ❑ Global results “publishable” to a messaging system in json format
  - Ideal for monitoring and offline analysis
  
- ❑ Adoption
  - @ CERN to continuously benchmark available CPU models
  - Used in CERN commercial cloud procurements (HNSciCloud)
  - Tested by other site managers (GridKa, RAL, INFN-Padova, ...)



# Benchmarks for heterogeneous resources

- ❑ In the future WLCG resources will likely include HPCs with GPUs
  - “How to value pledged HPC resources”?
    - WLCG MB requested to investigate approaches
- ❑ First demonstrator of **standalone container** for GPU benchmarking available
  - Based on CMS reco with GPUs (Patatrack)
    - Pixel track reconstruction, Calorimeter reconstruction
  - Essential for us, to understand how to apply the approach used for **CPU HEP Benchmarks** to the **CPU+GPU** system

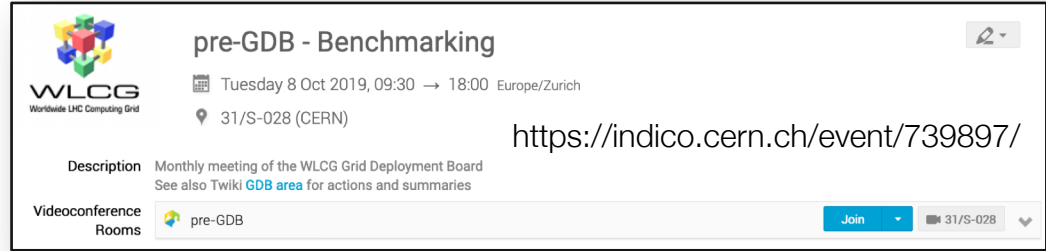
## MC and data event reconstruction

- **Event reconstruction** is where all experiments seem most advanced on GPUs
  - ALICE has a production-quality (IUUC) reconstruction for its online data processing in HLT
  - LHCb has an advanced prototype of HLT1 reconstruction, to be further reviewed internally
  - ATLAS is progressing on a heterogeneous framework to offload some algorithms
  - CMS has a functional heterogeneous framework, with several algorithms on GPUs
- Goals and GPU offload strategies differ in subtle and less subtle ways
  - LHCb has a standalone GPU application, all others have a heterogeneous framework
  - ALICE/LHCb focus primarily/exclusively on online reconstruction, unlike ATLAS/CMS
- *CMS event reconstruction may be IMO the first workload to reach GPUs on the Grid*
  - *We identified this workload as our first GPU candidate for the benchmarking suite*
  - *See the next two talks by F. Pantaleo and A. Sciaba*



# pre-GDB on Benchmarking

- ☐ Attendees
  - ~ 15 @ CERN, ~10 remotely
- ☐ Four dense sessions
  - Project Overview & License aspects
  - HEP reference Workloads (WL)
  - HEP Score and HEP Benchmark Suite
  - Reference Workloads running on GPUs
- ☐ Check the contributions for details!



The screenshot shows an Indico event page for "pre-GDB - Benchmarking". The event is scheduled for Tuesday, 8 Oct 2019, from 09:30 to 18:00 in Europe/Zurich. The location is 31/S-028 (CERN). The URL is <https://indico.cern.ch/event/739897/>. The description states it is a "Monthly meeting of the WLCG Grid Deployment Board" and refers to the "GDB area" on Twiki for actions and summaries. The "Videoconference Rooms" section shows a room named "pre-GDB" with a "Join" button and a room ID of "31/S-028". The WLCG logo (Worldwide LHC Computing Grid) is visible in the top left corner of the event card.

# Conclusions

- ❑ The full **HEP Benchmarks** chain is in place
  - Already used by a number of beta testers



- ❑ Ongoing work (targeting Q1 2020)
  - Validation of the running WLs and produced “**scores**”
  - Inclusion of few remaining WLs
  - Consolidation of the WLs’ report
  - Consolidation of the HEP Score and HEP Benchmark Suite implementations
  - Compare Docker Vs Singularity HEP scores
  - Run at large scale on production nodes
    - Compare performance of HEP benchmarks and standard jobs
    - Compare with HS06 and SPEC CPU 2017



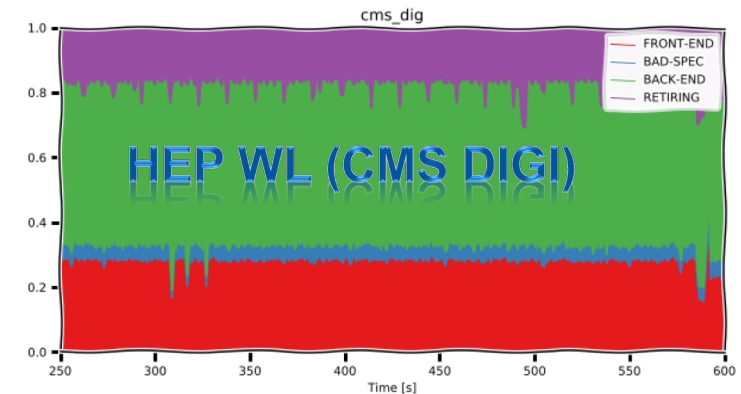
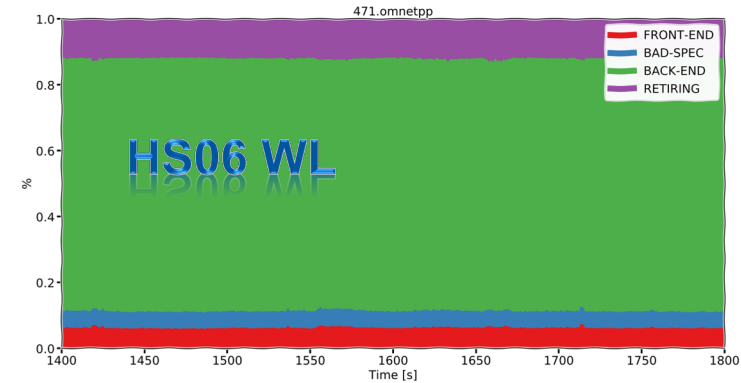
# Quantitative comparison with WLCG workloads

- Unveil the **dissimilarities** between HEP workloads and the SPEC CPU benchmarks
  - Using the **Trident** toolkit
    - analysis of the hardware **performance counters**

## Characterization of the resources utilised by a given workload

Percentage of time spent in

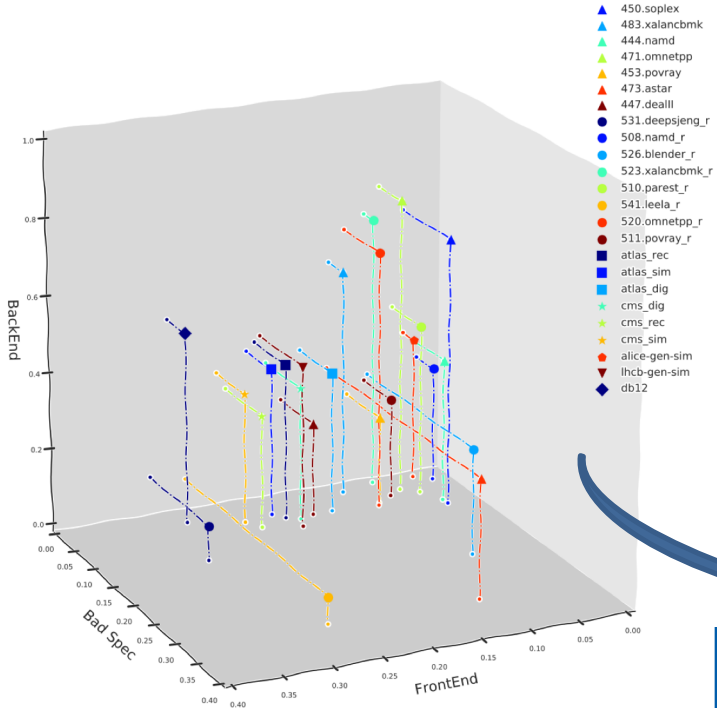
- **Front-End** – fetch and decode program code
- **Back-End** – monitor and execution of uOP
- **Retiring** – Completion of the uOP
- **Bad speculation** – uOPs that are cancelled before retirement due to branch misprediction



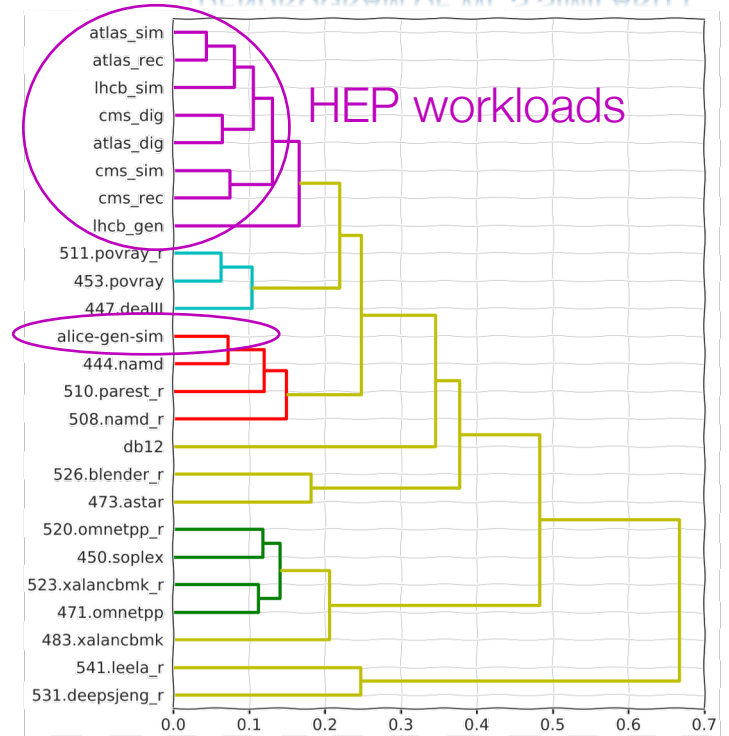


# Comparing all WLS together

- HEP WLS have same characteristics and differ more respect to HS06 and SPEC CPU2017 workloads



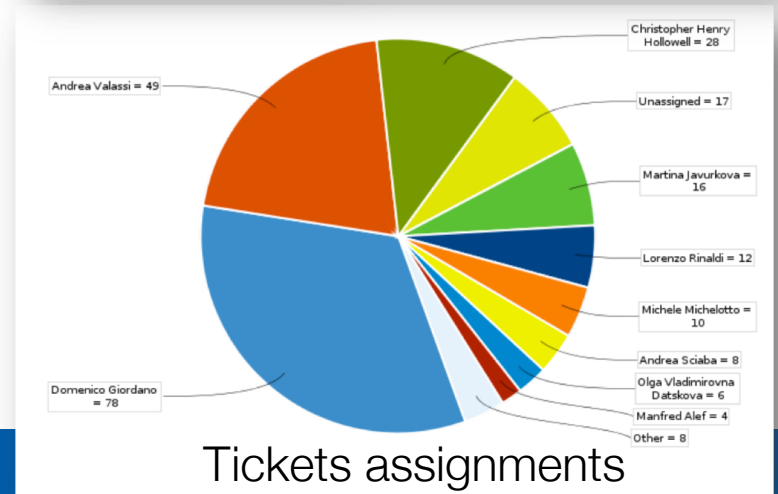
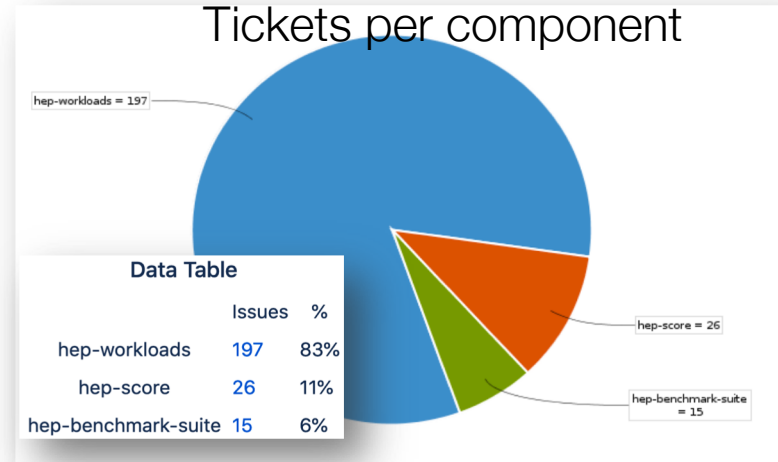
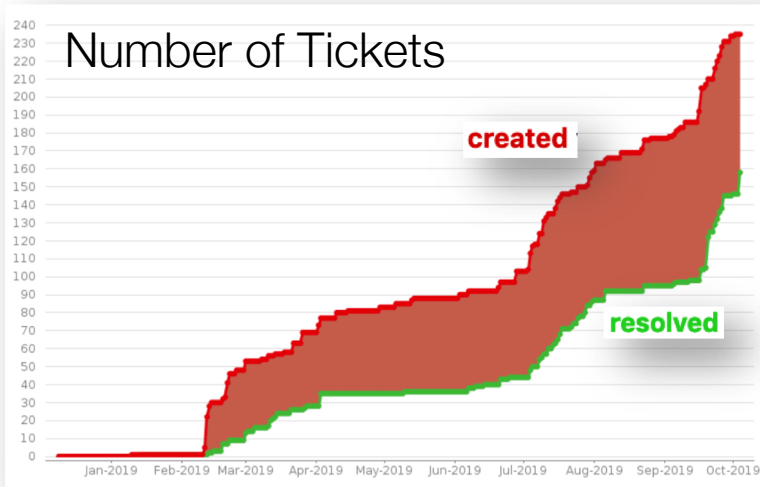
DENDROGRAM OF WL'S SIMILARITY



(\*) HS06 score is the geometric mean of the 7 HS06 apps

# Jira Statistics

- ❑ Work progression mainly in few months of activity
- ❑ Collective work and major attention gone to HEP-Workloads
- ❑ Main activities:
  - Develop the infrastructure
  - Develop the Experiment specific wrappers (drivers)
  - Validation



# CMS Patatrack in container

Felice Pantaleo: for the Patatrack Team  
for the CMS Experiment

*felice@cern.ch*



- A GPU-based full reconstruction of the Pixel detector from RAW data decoding to Pixel Tracks and Vertices determination, ECAL and HCAL local reconstruction has been implemented
- This reconstruction is fully integrated in the CMS Software
- Conversion to the legacy data formats and the standard validation can be run on demand
- Can achieve better physics performance, faster computational performance at a lower cost with respect to the baseline solution
- Working ongoing to develop the entire Phase-2 High Granularity Calorimeter Reconstruction directly with Heterogeneous Programming techniques
- Portable code is key for long-term maintainability, testability and support for new accelerator devices
  - Ongoing study and comparisons of solutions in Patatrack for CMS reconstruction
  - Starting from a CUDA code makes life **much** easier

14

## How to create the benchmark

1. Get all the executables and binaries
2. Get the input data and prepare it
3. Have a script that runs Patatrack and extracts a score
4. Package everything in a Docker image
5. Distribute it!

## Where are we?

- Simplified instructions to run the job by hand
  - <https://github.com/sciaba/patatrack-tests>
- CMS Open data set as input
  - In the working group EOS project space, takes 5 GB
- Scripts to build the Docker image
  - <https://github.com/vkhristenko/>
  - (the latest version is not yet committed)
- Image
  - Very large, about 50 GB – to become much smaller once the binaries are in CVMFS
  - All configuration parameters still hard-coded
  - Not yet publicly available
  - Still needs network connection for Frontier

Viktor Khristenko  
Andrea Sciabà

