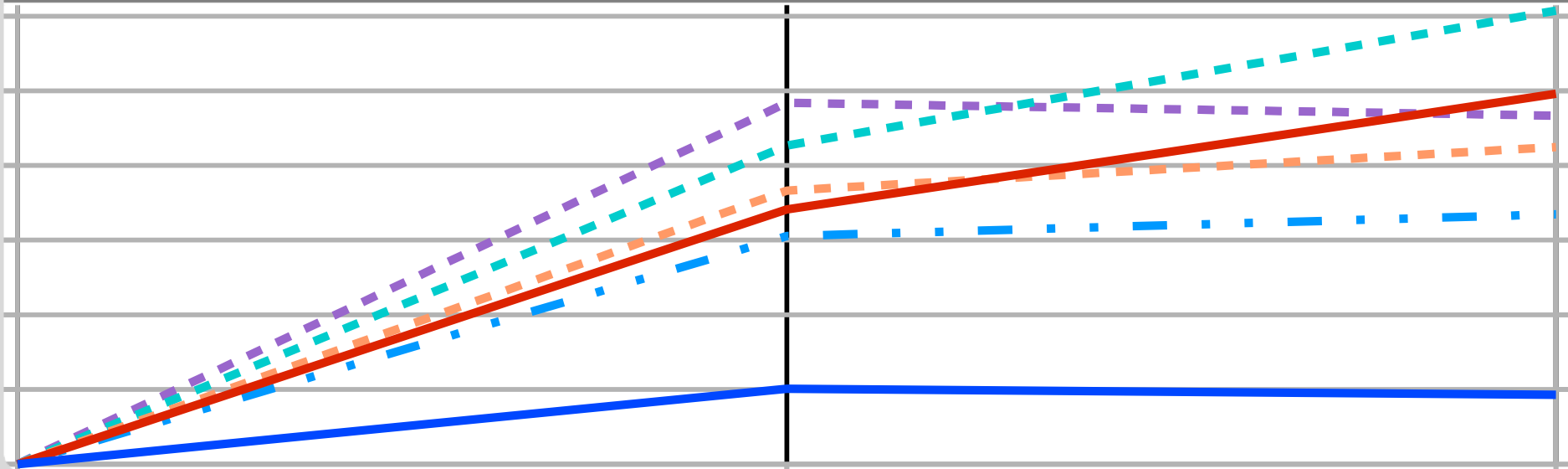# HEPiX Benchmarking Working Group Status Report Oct 2017

**Manfred Alef** (KIT), **Domenico Giordano (CERN), Michele Michelotto (INFN)**

STEINBUCH CENTRE FOR COMPUTING  (SCC)

www.kit.edu

# Mandate

## Benchmarking Working Group

- **Fast benchmark**
  - → Estimate performance of provided job slot or VM instance

- **Next generation of long-running benchmark**
  - → For installed capacities, accounting, procurements aso. (successor of HS06)

# Organization

- 60 subscribers of mailing list (hepix-cpu-benchmark@hepix.org)

- Biweekly Vidyo meetings

  → Kick-off at HEPiX Zeuthen (Apr 2016)

  → ~10 attendees per meeting

    - Site admins

    - Experiment representatives (Alice, Atlas, CMS, LHCb)

# Status of the Working Group

- Today:

  - ➔ Status update since last HEPiX meeting (Apr 2017), and April GDB

    - Talks by Domenico Giordano *

# Fast Benchmark

■ DIRAC Benchmark 2012 (DB12) is an attractive fast benchmark

➔ Python script running for around 1 min

➔ Very good correlation with Alice and LHCb jobs when running 1 benchmark copy ('DB12-in-job')

● However, DB12 doesn't show the stability and characteristics to probe all components of the CPU potentially used by HEP workloads; e.g. the limited instruction mix doesn't stress the memory subsystem

Manfred Alef et.al.:   HEPiX Benchmarking Working Group: Status Report Oct 2017   Steinbuch Centre of Computing

# Next-Generation Long-Running Benchmark

**Purpose of the 'long-running' benchmark is to measure installed and pledged compute capacities.**

**Hence it must scale (with a certain accuracy\*) with the average WLCG job mix, but it will probably not scale with any individual job type (simulation, event generation, reproduction, ...)**

**\*** Initial objective of HS06: spread ≤ 10%

# Next-Generation Long-Running Benchmark

- Current HS06 benchmark built on SPEC CPU2006

- New SPEC CPU2017 has been released Jun 20

  - ➔ Volunteering sites have already purchased the new benchmark suite, and they are now warming up

- Packaging Alice and Atlas reference workloads in Docker containers *

- HS06 scaling issues have been investigated in more detail

  - ➔ 64bit temporary workaround?

* https://indico.cern.ch/event/653573/contributions/2700565/attachments/1513184/2360433/HEPiX-workload-on-docker-container.pdf

# SPEC CPU2017 Benchmark Suite

■ Website: www.spec.org/cpu2017/

■ 43 single benchmarks

➔ Integer, and floating point

➔ Speed, and rate metric

➔ Many benchmark names are already known from CPU2006, but CPU2017 is coming with new releases, and running improved workloads

# SPEC CPU2017 Benchmark Suite

- SPECspeed and SPECrate metrics as before

  - ➔ Now different branches within the benchmark suite

    - SPECspeed:     20 benchmarks (10 integer + 10 fp)
    - SPECrate:       23 benchmarks (10 integer + 13 fp)

  - ➔ Memory requirements

    - SPECspeed benchmarks very memory-hungry (up to 16 GB), that's far too much for parallel copies as in HS06
    - SPECrate requires only 2 GB RAM per copy

- Current status:

  - ➔ Volunteering sites are warming up

  - ➔ First results at next HEPiX

# Scaling Issues of HS06 vs. HEP Applications

- 64 bit interim solution?

  - HS06 runs with mandatory -m32 compiler flag

  - Improved scaling with -m64?

    - Nearly linear increase by around 10...20% of 64bit benchmark scores

    - Double-checked SL6 + CentOS7

    - AMD Epyc: + ~33% (when running 1 benchmark copy per core)

  - Conclusion: migration to 64 bit doesn't fix the scaling issues

Manfred Alef et.al.:    HEPiX Benchmarking Working Group: Status Report Oct 2017    Steinbuch Centre of Computing

# Scaling Issues of HS06 vs. HEP Applications

■ Expanding to second dimension

➜ HS06 had been developed by the HEPiX Benchmarking Working Group from 2007 to 2008

➜ Typical WN hardware at that time without Hyperthreading feature:

● Intel: quad-core CPUs Xeon E53xx or E54xx

● AMD: 8...16-core CPUs Opteron 23xx or 61xx

➜ First servers with Hyperthreading feature (Intel E55xx) appeared on the market at the end of the project

➜ Variety of WN configurations at sites

● HT disabled

● HT enabled, more than 1 job slot per physical core

◆ E.g. ~1.5, or 2 job slots per core

Manfred Alef et.al.: HEPiX Benchmarking Working Group: Status Report Oct 2017 Steinbuch Centre of Computing

# Scaling Issues of HS06 vs. HEP Applications

- Expanding to second dimension

  - Experiment reports, for instance at several GDB meetings, have compared different hardware models

  - Only few reports taking into account the individual WN configuration, especially the number of job slots

  - Indications that this is important too

2017-10-18      Manfred Alef et.al.:     HEPiX Benchmarking Working Group: Status Report Oct 2017                Steinbuch Centre of Computing

# Scaling Issues of HS06 vs. HEP Applications

■ Expanding to second dimension

➔ Discrepancies in static benchmark scores (HS06, DB12-at-boot)

**Benchmark Scores (per Host)**

Intel E5-2630v4 (Broadwell) -- Normalization: 1 job slot per core = 1



➔ What about HEP applications?

# Scaling Issues of HS06 vs. HEP Applications

- Expanding to second dimension

  - Deeper analysis at KIT and at PIC

    - GridKa compute farm has been reconfigured
      - Default configuration: 1.5 (or 1.6) job slots per core
      - Latest hardware model (Intel Xeon E5-2630v4, Broadwell) with 3 different configurations:
        - 1.0 job slots per core                                                              (20 slots)
        - 1.6 job slots per core                                                              (32 slots)
        - 2.0 job slots per core (1 per logical processor)   (40 slots)
      - Correlations between job performance (events/s) and benchmark scores?
    - Dedicated benchmarking hosts at PIC

Manfred Alef et.al.:   HEPiX Benchmarking Working Group: Status Report Oct 2017     Steinbuch Centre of Computing

# Scaling Issues of HS06 vs. HEP Applications

- Expanding to second dimension

  - → Deeper analysis at KIT and at PIC

    - Performance results:

      - ◆ Benchmark scores (# copies == # job slots)
        - ■ HS06
        - ■ DB12-at-boot (MJF package)
      - ◆ Further benchmarks compared at PIC:
        - ■ Atlas KV
        - ■ CMS ttbar sim.

2017-10-18        Manfred Alef et.al.:        HEPiX Benchmarking Working Group: Status Report Oct 2017        Steinbuch Centre of Computing

# Scaling Issues of HS06 vs. HEP Applications

- Expanding to second dimension

    - Deeper analysis at KIT and at PIC

        - Performance results:

            - Performance of jobs run at GridKa (everyday job mix)
                - Alice (thanks to Costin Grigoras)
                - Atlas (values downloaded from Bigpanda, Tasks: simul=10944000, recon=11323845, evgen=11330855)
                - LHCb (thanks to Philippe Charpentier)
                - CMS:  n.a.
            - Alice and LHCb have also reported corresponding DB12-in-job scores (running 1 benchmark copy)
                - LHCb: DB16-in-job which is the same Python script as DB12 but with a modified internal calibration factor
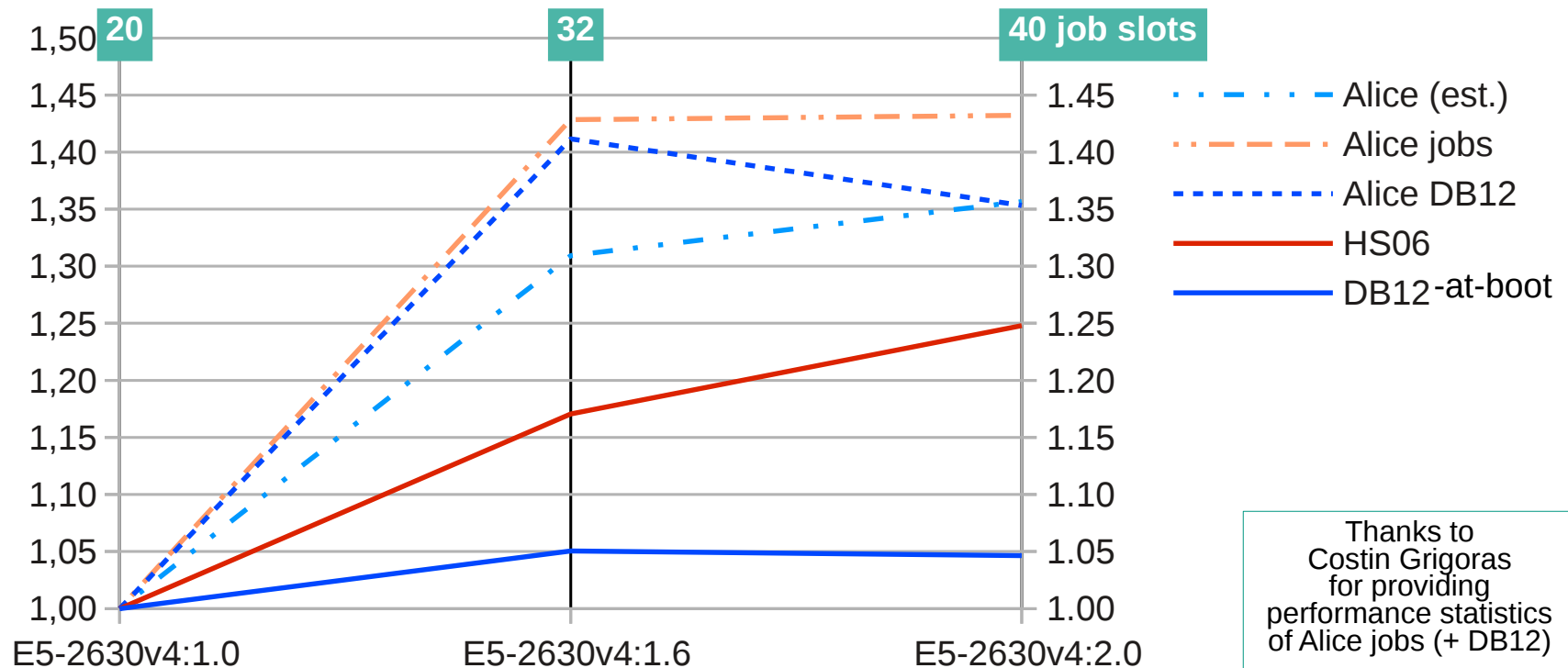
Manfred Alef et.al.:   HEPiX Benchmarking Working Group: Status Report Oct 2017   Steinbuch Centre of Computing

# Scaling Issues of HS06 vs. HEP Applications

- Expanding to second dimension

  - Deeper analysis at KIT and at PIC

    - Performance results:

      - Job performance estimated by comparing runtime of top processes
        - Rough estimates, no high-precision accounting scores!
        - LHCb: n.a. (sophisticated autocalibrations)

Manfred Alef et.al.:    HEPiX Benchmarking Working Group: Status Report Oct 2017    Steinbuch Centre of Computing

# Scaling Issues of HS06 vs. HEP Applications

**Benchmark Scores  vs.  Alice Job Performance (Upscaled)**

Intel E5-2630v4 (Broadwell)  -  Normalization:  1 job slot per core = 1

Legend:
- Alice (est.)
- Alice jobs
- Alice DB12
- HS06
- DB12 -at-boot

Thanks to Costin Grigoras for providing performance statistics of Alice jobs (+ DB12)

Manfred Alef et.al.:     HEPiX Benchmarking Working Group: Status Report Oct 2017          Steinbuch Centre of Computing

# Scaling Issues of HS06 vs. HEP Applications



**Benchmark Scores vs. LHCb Job Performance (Upscaled)**

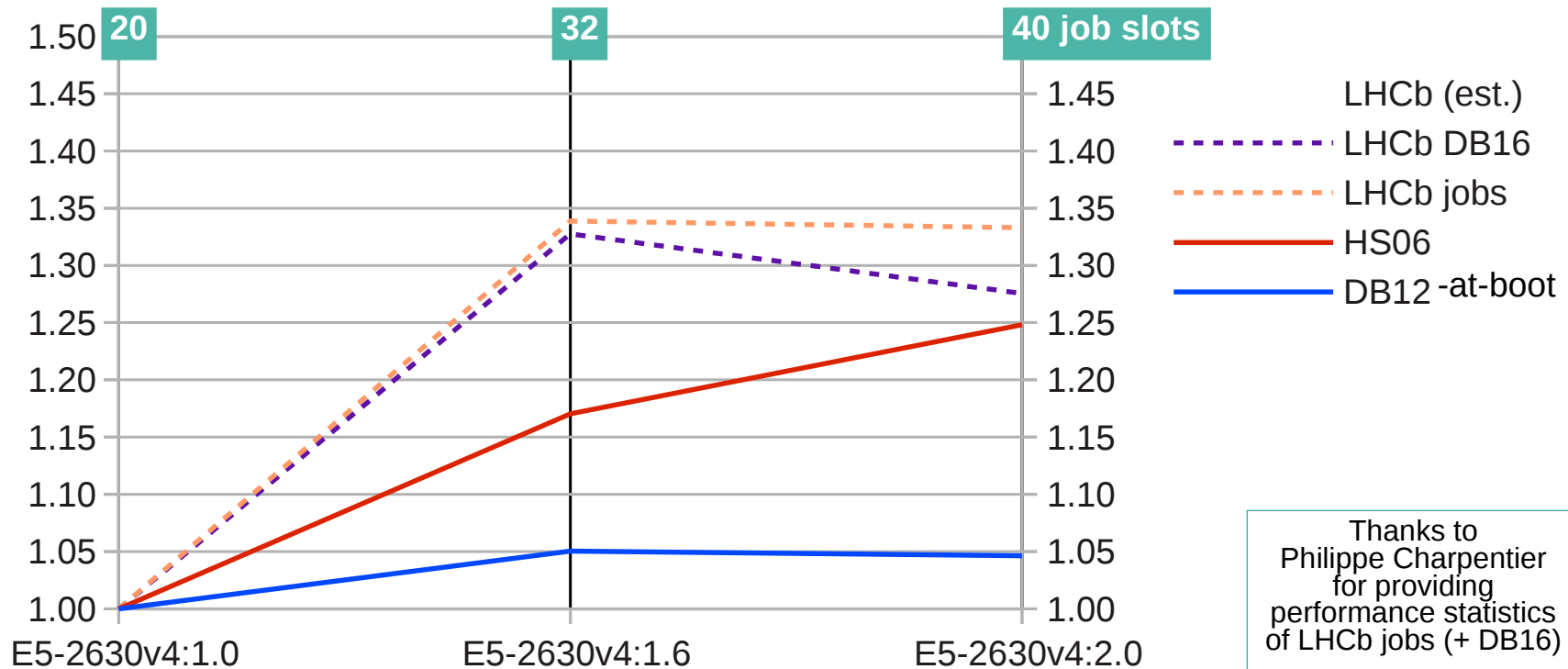Intel E5-2630v4 (Broadwell) -- Normalization: 1 job slot per core = 1

Legend:
- LHCb (est.)
- LHCb DB16
- LHCb jobs
- HS06
- DB12 -at-boot

Thanks to Philippe Charpentier for providing performance statistics of LHCb jobs (+ DB16)

**Benchmark Scores  vs.  Atlas Job Performance (Upscaled)**

Intel E5-2630v4 (Broadwell)  -  Normalization:  1 job slot per core = 1



Legend:
- Atlas (est.)
- Atlas simul
- Atlas recon
- Atlas evgen
- HS06
- DB12 -at-boot

Job performance statistics from Bigpanda, TaskIDs:
10944000 (simul)
11323845 (recon)
11330855 (evgen)

# Scaling Issues of HS06 vs. HEP Applications

■ CMS:

➔ ttbar sim. at PIC on Haswell host (J. Flix et. al. *):



CMS ttbar sim. - Intel Xeon E5-2640v3 @ 2.60GHz

Legend:
- DB12
- HS06
- KV (evts/sec)
- ttbar (evts/sec)

Normalization:
1 copy per core = 1

**DB12-at-boot**

Y-axis: Normalized power to phys. cores
X-axis: # parallel procs

➔ Estimates at GridKa similar to the Atlas ones

* https://indico.cern.ch/event/624830/contributions/2576000/attachments/1454803/2244865/20170505_CMS_Benchmarking_JFlix.pdf

Manfred Alef et.al.:   HEPiX Benchmarking Working Group: Status Report Oct 2017   Steinbuch Centre of Computing

# Summary

- **Fast benchmark:**

  → DB12 (in-job) scales with Alice and LHCb jobs

  - Runtime ~1 minute

- **Long-running benchmark (HS06 + successor):**

  → Not only the hardware model but also the configured number of job slots per physical core are important

  → Migration to HS06 64bit doesn't solve the issues

  → DB12-at-boot (multiple copies) is <u>not</u> a suitable candidate

  → Containerising reference workloads (Docker, CVMFS)

  → Investigating SPEC CPU2017

Manfred Alef et.al.:    HEPiX Benchmarking Working Group: Status Report Oct 2017    Steinbuch Centre of Computing

2017-10-18 Manfred Alef et.al.: HEPiX Benchmarking Working Group: Status Report Oct 2017 Steinbuch Centre of Computing