

Performance?

- Summary of the panel discussions

HSF and its role in performance?

- V: This is the question that we should keep in mind throughout the panel discussion

What is Computational Efficiency for you?

- Holistic Performance Assessment?
 - Performance for a complete campaign!
 - Can we learn how to best spending money/effort?
 - Software (Frameworks)
 - CPU/DISK/NET
- How to measure?
 - Metrics
 - Phase space

- Graeme: The T0 and its use for data taking can serve as an example
 - Metric is event throughput rate in HZ
 - Different impedance for hardware and software
 - hardware requires 6 months to + 3 years to reach/leave production
 - software in a month(s) to production
 - Currently we are driven to adjust the software to the hardware (mostly)
 - In the long-term we have to tell the sites what hardware would be good
 - Currently we ask for terabytes, but not I/O operations/sec
- Olivier: The national infrastructure based on X86 has reached a plateau of efficiency and is power limited
 - Centres are going to move to GPUs
 - This is not our choice, x86 doesn't continue and we have to be able to work (well) on everything.
- D.Rohr
 - GPUs are here, but they also change. We have to change our data models. At the core they are similar.
 - In between we will run some loads on GPUs, you have a server and can mix.
 - A cluster at GSI provides only GPUs, but is dedicated to QCD
 - Gradual approach can be taken to gain experience. We(Alice) tried it in Run1 and have added more GPUs later
- Liz: We don't have a choice in the US. HTC has to join HPC.
 - HPC has a road map for 2020. We need to adjust our code and embrace heterogeneity.
 - Build and packaging systems need to improve to cover this
- Audience: GPUs change quickly, high performance, high prices, difficult to time, (buy tomorrow same performance for less money)

Can we have efficient Code/Algorithms in a heterogeneous world?

- Von Neumann
- Neuro Chips
- FPGAs
- GPUS...
- ????
- Possible, affordable?

- Liz: we have to invest into the build systems
- Audience: In the 90th all machine where von Neumann type, now vastly different concepts
- Liz: Building fat binaries and the manufacturers need to help using the new technologies
- Amir: New frameworks are needed to address this. HSF can coordinate this activity. Maybe we can define a common framework
- Jeff: Wider instead of faster. GPU optimised code runs also better on cpus. Not clear how different a framework needs to be.
- Sami: Problem is not the hardware
 - Need to define what we accept: Physics, effort, maintenance, etc.
 - Differences in physics and performance have to be expected.
 - What is acceptable??
- Liz: CMS and ATLAS improved without physics performance changes
- Pere: No clear that GPUs are the solution, 1st demonstrate how much you can gain.
 - If we haven't found any gain no effort is justified
 - We can't translate 1000s lines to work on GPUs well
- D.R.:It is easier to work with a concrete app like the HLT.
 - Greater diversity for frameworks like GEANT and ROOT
 - GPU vendors establish individual ecosystems to lock users in

Joining the HPC community and/or the next wave of commodity computing?

- Supercomputer and/or Mobile Phones
- How can we use HPC efficiently?
 - Beyond filling the opportunistic use
 - Real costs?
- Market shifts from desktop to mobile devices
- Commodity computing has been in the past the most cost effective approach
 - Will this continue?

- Vincenzo: Is there still a place for commodity hardware? We are pressed to join HPC and there seem to be less opportunity with the now power saving centric new commodity hardware. Do we know that we can't achieve the same throughput/cost with ARM like systems than with HPC without the extra effort to parallelise our code?
- Graeme: Number of lines of code is the core of our cost. We have to rewrite as much as needed to cover most of the architecture phase space, but compilers have to do the heavy lifting
- J. Apostolakis: It is impossible to predict the dominant infrastructure in 5 years. We have to remain flexible and support multiple infrastructures at low cost.
- Olivier: We currently hide the detailed architecture aspects from the users in libs, a clear standard could simplify this problem, but currently there isn't one.
- Jeff: HPC will go away. Because the technology evolution slows down.
- Liz: No, they just make them denser. The manufacturer will address this.
- Amir: We can't depend on them, we have to change. The deep learning community is adapting. Maybe we should abstract on the level of machine learning. There is already now a split between hardware for HPC and gaming/machine learning. (double precision/ 16bit float). Offloading to GPUs is also intimately connect to I/O
- Liz: We can write code that can be handled by the compiler well, cooperation with compiler people necessary to vectorise. The exascale project addresses this.
- Pere: We have to prevent the users to deal with the low level aspects. If the users gives a high level description it is easier for the system to do optimisation internally.
- Geant: compilers might help in the future. We had to decide, but currently the compilers cannot do the work. Flexible to be adaptatoon for the future. Underlying layer has to perform. 20% rewrite for different architectures would be impossible.
- Pere: We need to assess the exact impact on the overall workflow of offloading parts to specialised HW and take into account the effort (human) to support them.
- D.R.: When ALICE rewrote the tracker it had the use of GPUs in mind, but the original version was not written for a GPU, but showed already massive performance gains.

Do we have to rethink I/O and data organisation?

- Distributed
- Local
- In memory
- Data structures
- Do we understand the impact on performance ?
- Will this impact the “Framework”?

- Audience: GEANT runs into bottlenecks with data reduction
- Pere: Several optimisations are under investigation, will progress over time
 - User should express high level things, optimise behind the scene, caches, map reduce
 - User should not deal with details.
- V: There is full distributed data indexing in the bio and astro world. Can be do the same?
- Audience: Question to the experiments, event indices has never been useful
- Pere: The user should not need to say in a loop that he wants a histogram to be filled, but express it at a higher level/
- Audience: Compression is the biggest bottleneck now, maybe the GPU can do this, the CPU should not be involved. DP-direct. This is what google does at the moment,
- V: This is similar to what our DAQ people do. Adding more intelligence into the frontends
- Pere: You have to see the impact if you optimise the I/O to see whether it makes sense. At least a back of the envelope estimate is needed before work starts.
- Liz: DIANA plans to work on root-i/o I hope it works out . It includes R&D.

Culture of performance? Lost?

- Generations of developers have been raised by:
 - “You shall code for clarity, they who did not trust in Moore’s law are forsaken to wander the desert among profilers and the multitude of compiler switches”
- How to (re) create a Culture of performance?
- Society of the Two Cultures?
 - One who comes up with ideas and prototypes (Science)
 - One who turns this into a refined code for production (Engineering)
 - Recipe for disaster?

- Not covered

Evolution or Revolution?

- Can we gain the (expected) needed factors in performance by either of them?
- What might be additional gains from rethinking computing for HL-LHC?
- Could we afford a revolution?
- Could we afford not having a revolution?

- Not covered

HSF and its role in performance

- Disseminate Information
 - Knowledge base, training
- Pere: Helping to link the different R&D projects and help them with visibility
- Audience: Knowledge of bottlenecks in code should be shared, best practices sharing, education
- Vincenzo: Collaboration with industry, like openlab?
- Amir: We need to define a list of what should be supported as R&D
 - Maybe whitepaper for the funding agencies
 - We should talk about the same topics together by creating workshops with clear focus
 - Need to be proactive
- Panel: The Concurrency Forum is the entity to cover these aspects
- M.Sokoloff: If the US NSF proposal moves forward we'll need a Community Whitepaper within the next 15 months. HSF the place to link activities, create working groups and task forces to define what we need to cover in R&D. A Community white paper from the HSF is very important for the NSF. We have to be proactive and the HSF is the natural organisation to do it.
- Jeff: The work has to be done by the people. We need a metric for HSF work so that we can make the best use of the effort, not require meetings or define roles.
- Liz: I would agree if we have to make a small community efficient, with NSF support we can get more people, but for that we need to write and agree on a roadmap. This is not a zero sum game. It has to cover all communities to attract funding.
- Dario: Lessons learned should be brought forward to others. Mainly by wiki based training we should increase the expertise of the existing people.