



HL-LHC Computing

Mikołaj Krzewicki, FIAS

on behalf of the

Trigger/Online/Offline Computing preparatory group



FIAS Frankfurt Institute
for Advanced Studies



TOOC members

ALICE

Latchezar Betev, Mikołaj Krzewicki, Pierre Vande Vuyvre.

ATLAS

Graeme Stewart, Benedetto Gorini, Nikos Konstantinidis, Imma Riu,
Stefano Veneziano.

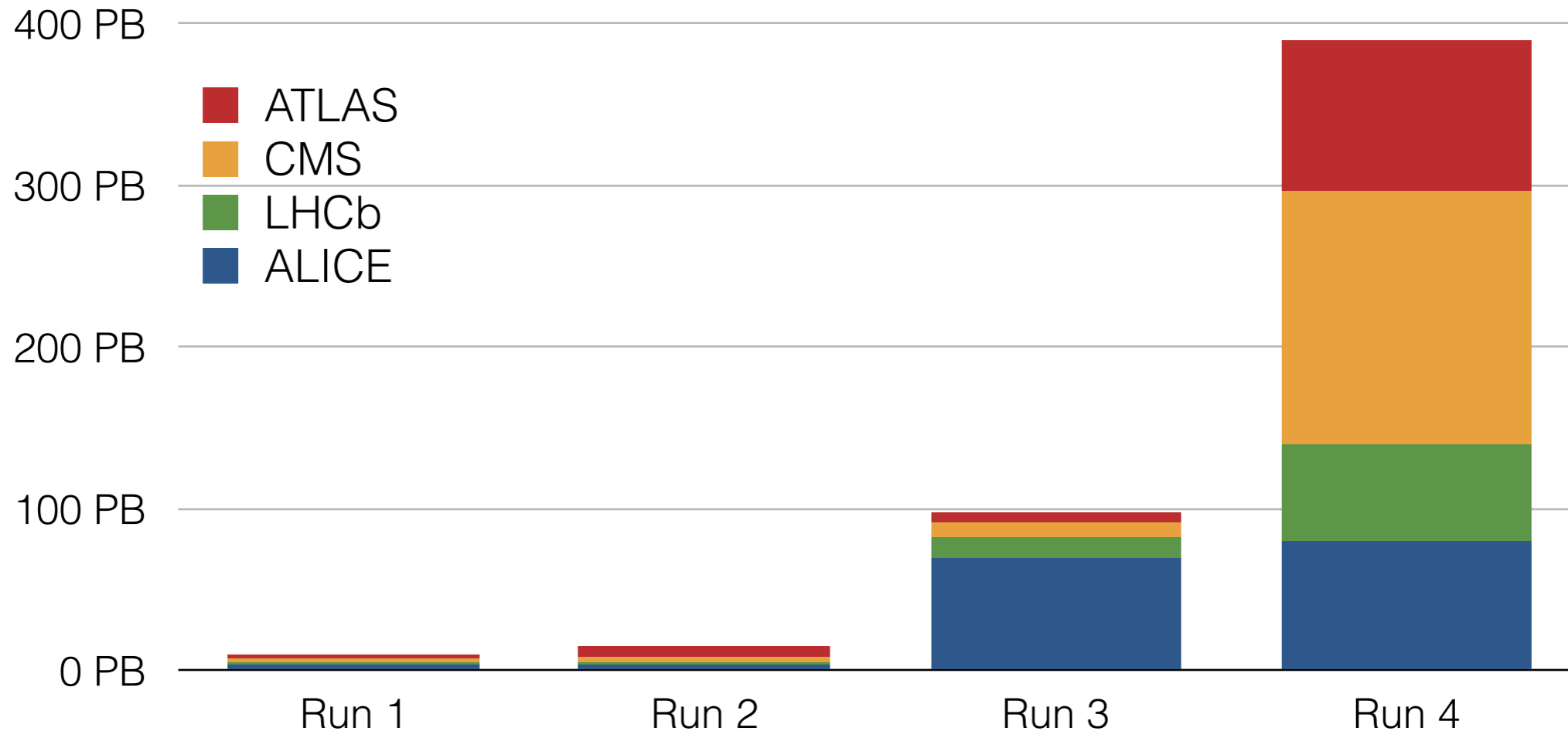
CMS

Wesley Smith, Maria Girone, David Lange, Frans Meijers.

LHCb

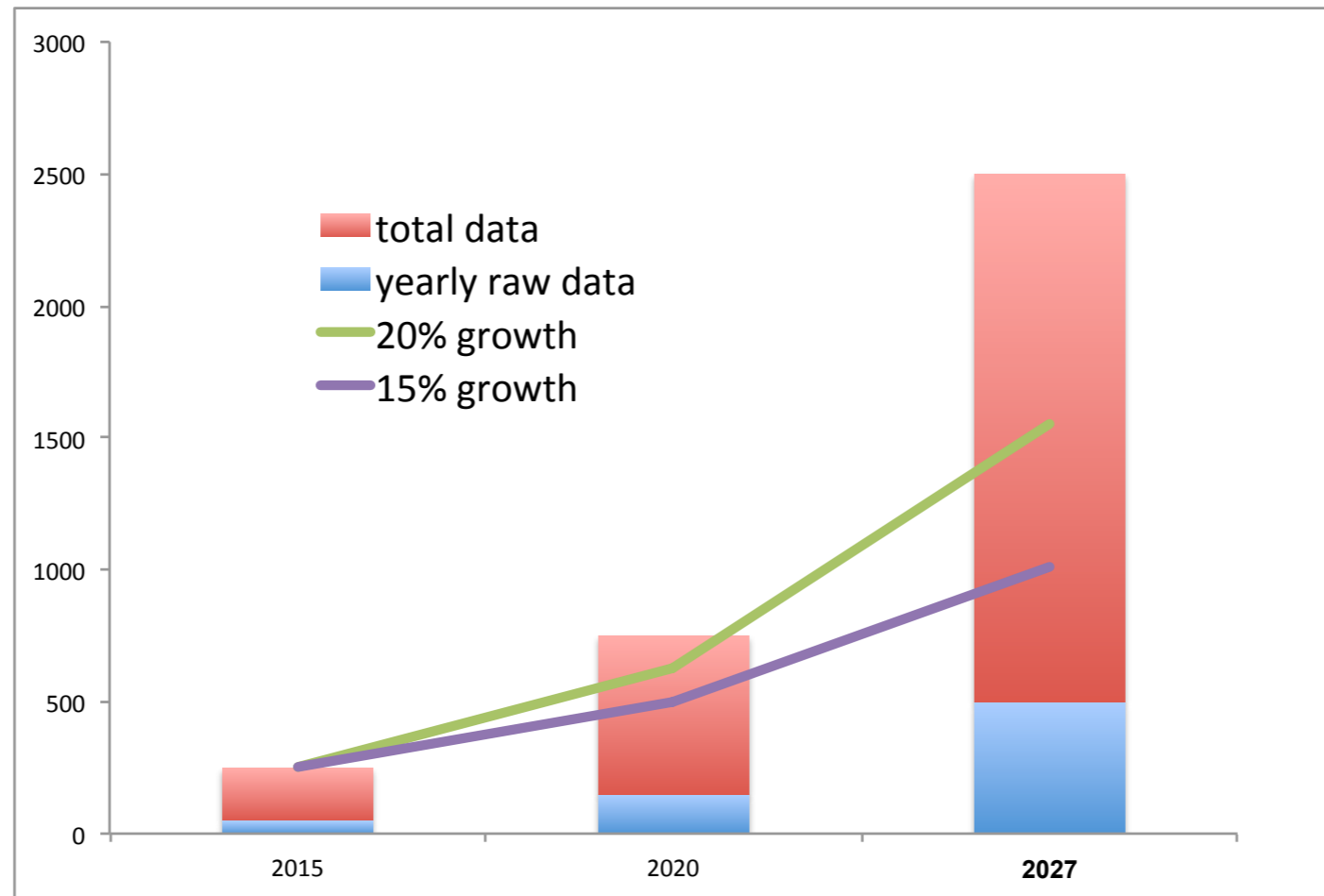
Peter Clarke, Vava Gligorov, Niko Neufeld.

Scale of challenge: data



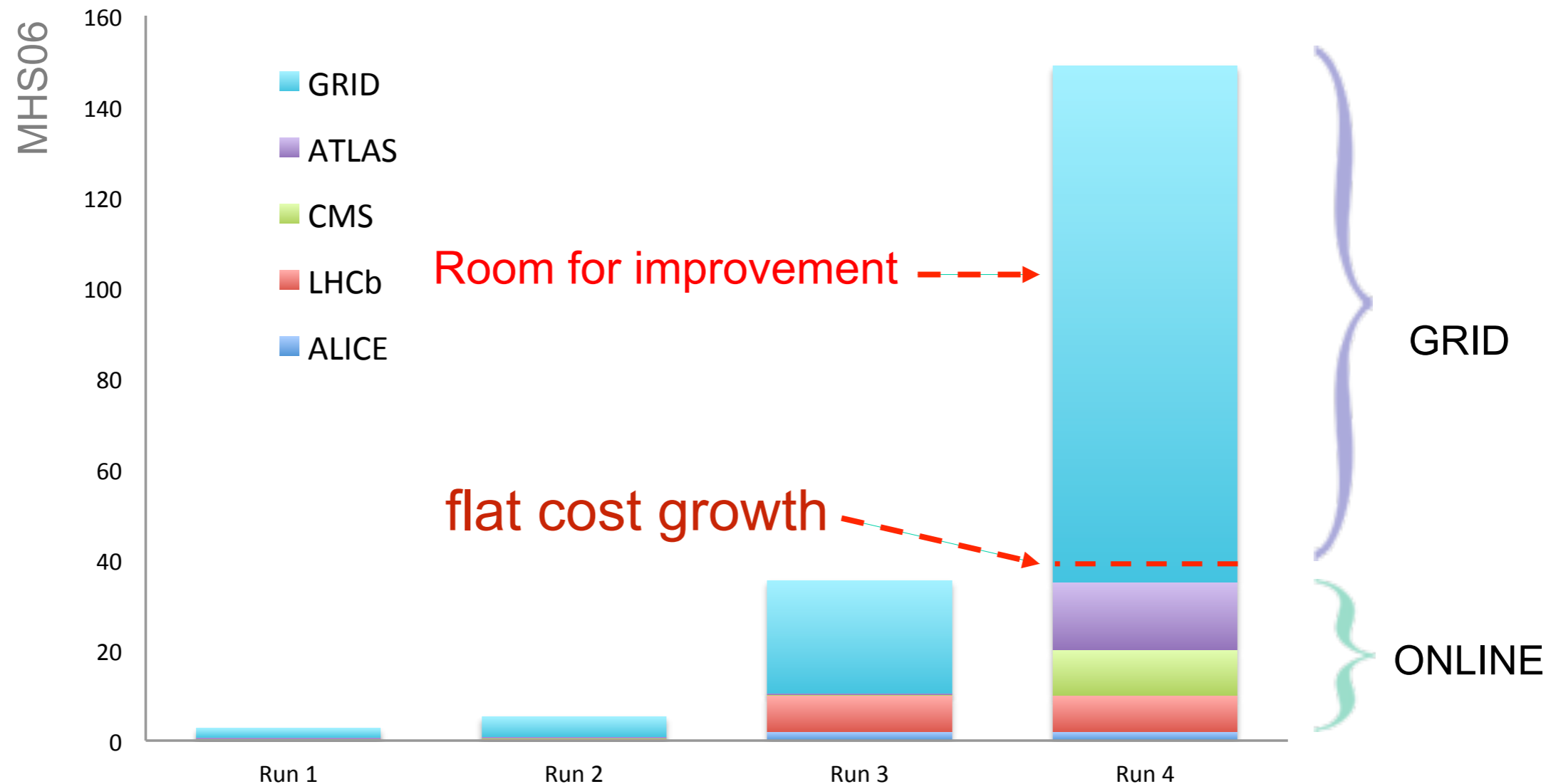
- Crude estimates based on the expected data rates (per annum).
 - ALICE: large part is a disk buffer in the online system, natural GRID evolution should provide the rest.
 - Data rates and event sizes vary within a run as much as factor 2.
- **EXCLUDES** derived data - typically factors more than RAW shown here.
 - ➔ Data volumes expected to grow dramatically.

Active data - disk



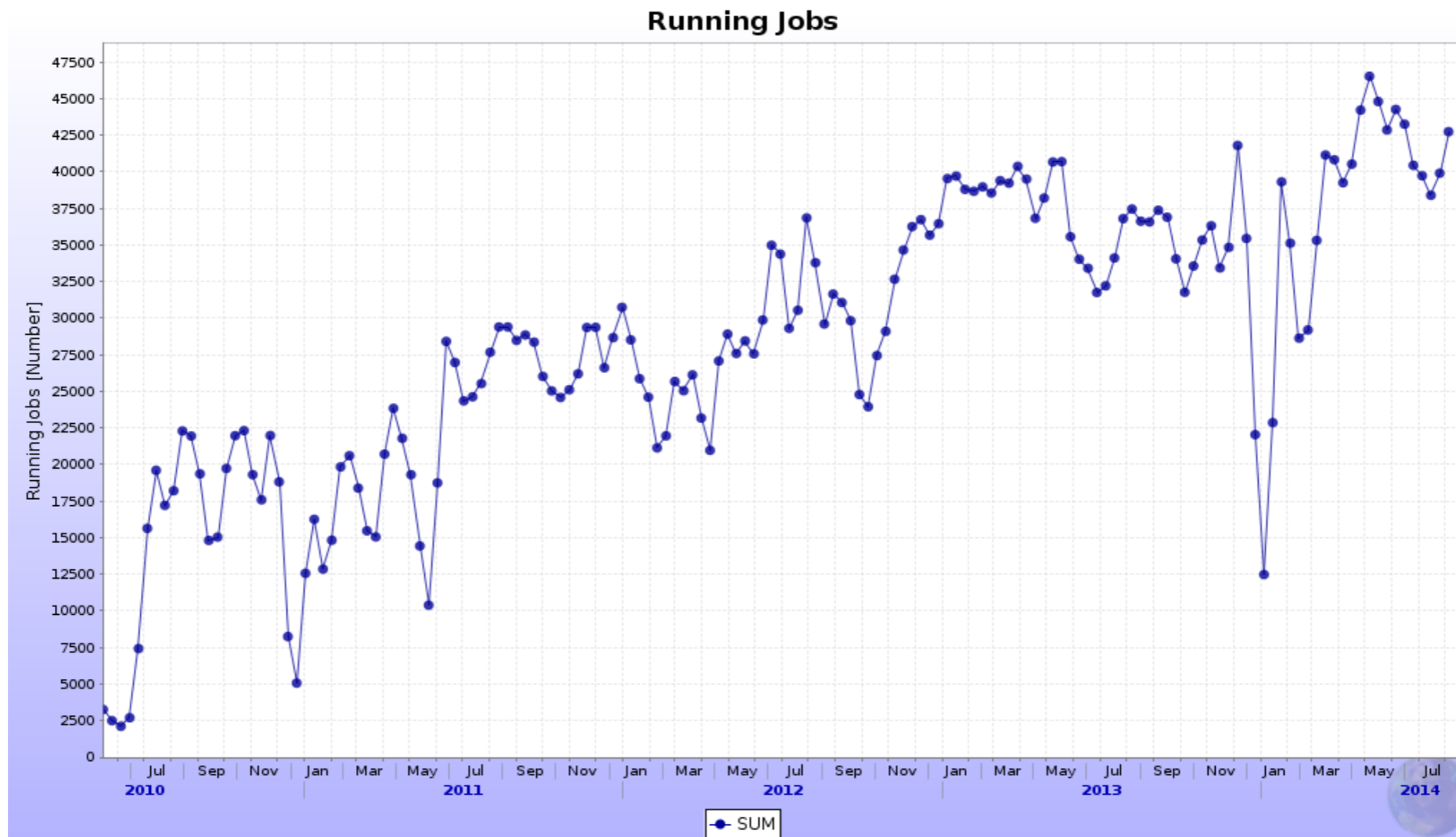
- Assumes ratio of disk to yearly raw data is as currently requested for 2015.
- Assumes flat budget annual growth remains at 15-20%.
- In 2025 cost is at least factor 2-3 above flat budget.

Scale of challenge: CPU



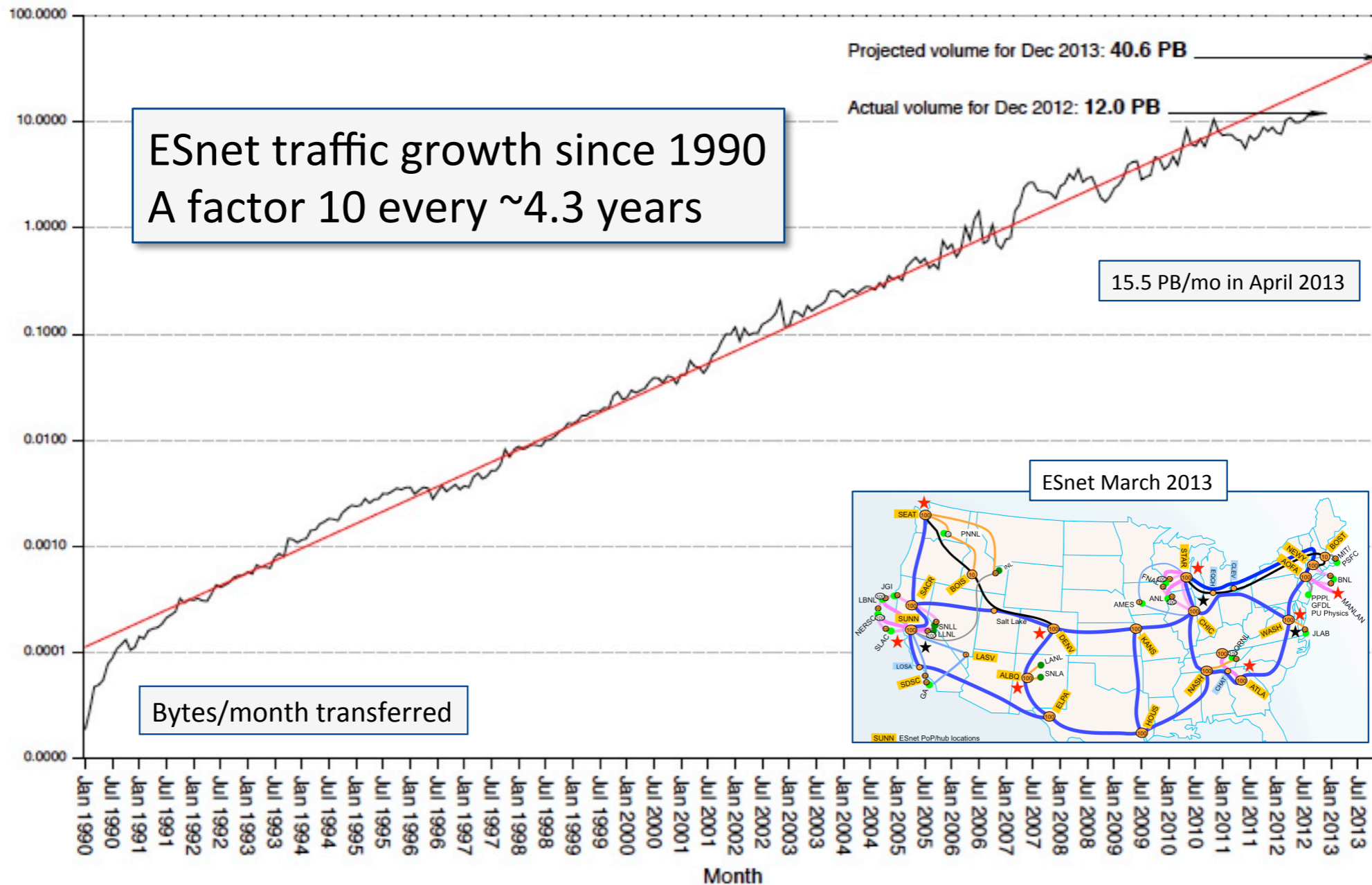
- Rough estimates of the CPU resources needed, based on extrapolations.
- It is clear CPU usage must be improved.

GRID growth



- Number of cores grows by 25% year on year (flat budget).
- Power/core ~constant.
- Storage growth at 20% per year.
- Projected at 2020 => ~3-4x the current power (storage and CPU, resp.).

Networking growth



- Dramatic growth, by example of ESnet.
- Factor 10 every 4.3 years.
- Could mean less data replication where appropriate (on demand data copy)?

Costs

Assuming similar computing models as today:

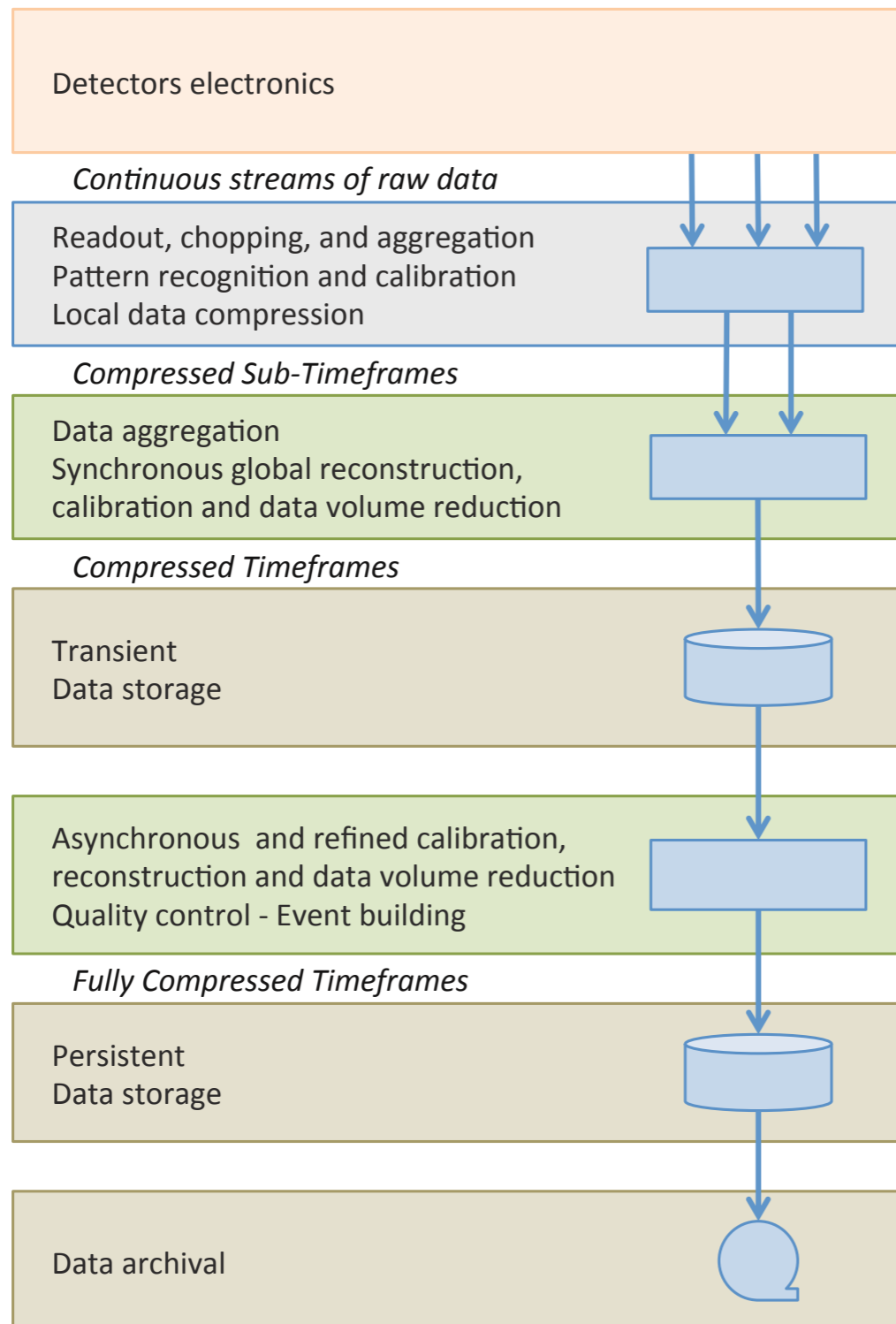
- Networks:
 - Technology growth will provide what we need;
 - Cost ? Affordable if today's trends continue.
 - Archive storage:
 - Tape (robotics, drives, media) – cost similar to today for full anticipated HL-LHC data growth.
 - Disk buffer cost will be much higher.
 - Active storage (data copies, caches, etc):
 - Costs factor 2-3 higher than flat budgets.
 - CPU:
 - Costs factor 3-5 higher than flat budgets.
- **Biggest impact on overall costs is disk storage.**

ALICE upgrade



- ALICE upgrade basic estimates:
 - Event rate 50KHz (Pb-Pb), 200KHz (p-p, p-Pb).
 - Event size 1.1TB/sec from detector; 20GB/sec average processed and compressed to storage.
 - Triggerless readout - basic data unit a “timeframe” instead of an “event”.
- RAW data rates and volume necessitate the creation of an online-offline facility (O2) for data compression, incorporating:
 - DAQ functionality – detector readout, data transport and event building.
 - HLT functionality – data compression, clustering algorithms, tracking algorithms.
 - Offline functionality – calibration, full event processing and reconstruction, up to analysis objects data.

The architecture of ALICE O²



synchronous with data taking - “online”

Internal O₂ data buffer @P₂

Calibration refinement @P₂

Custodial storage at T₀/T₁

- TDR: summer of 2015.

ALICE O² data reduction plans



- ‘Offline quality’ calibration critical for the data compression.
 - Compressed data allows reprocessing, i.e. finer-grain calibration is still possible (to a degree).
- Use of FPGAs, GPUs and CPUs in combination;
 - Software uses specific advantages of each.
 - A well-tested approach in production (current HLT).
- New framework to incorporate all tasks;
 - ALFA (ALICE-FAIR) being developed in collaboration with the FAIR collaboration at GSI Darmstadt.
 - Modular message based software framework
 - Very scalable, components communicate using a universal data/message transport (see Graeme’s talk).
- Run 2 is a test bed for many ideas, e.g. online calibration using the HLT.

- Currently commissioning new data placement and production system.
 - Typical lifetime > 5 years or so
expect new systems not before ~Run 4;
 - Run 2 & Run 3 similar in requirements for both only HL-LHC changes picture dramatically.
 - Need to learn from new system as well as need to know new requirements,
 - e.g. how to deal with accelerators; whole nodes scheduling should help.
- Future HW/SW technologies changes might offer completely new solutions.
- Work on optimising/modularising the software ongoing, e.g.:
 - dedicated EventServer for I/O running on same/different machine (enforcing all IO goes through the framework ...).
 - More speculative:
offload CPU intensive tasks to accelerators including parts of reconstruction (mostly tracking), file (de-) compression (on smallish GPUs/FPGAs!), Geant4 simulation, ideally these accelerators run on same or some other machine (incl. additional CPU cores for e.g. 'big.LITTLE' architectures).

ATLAS: disk usage

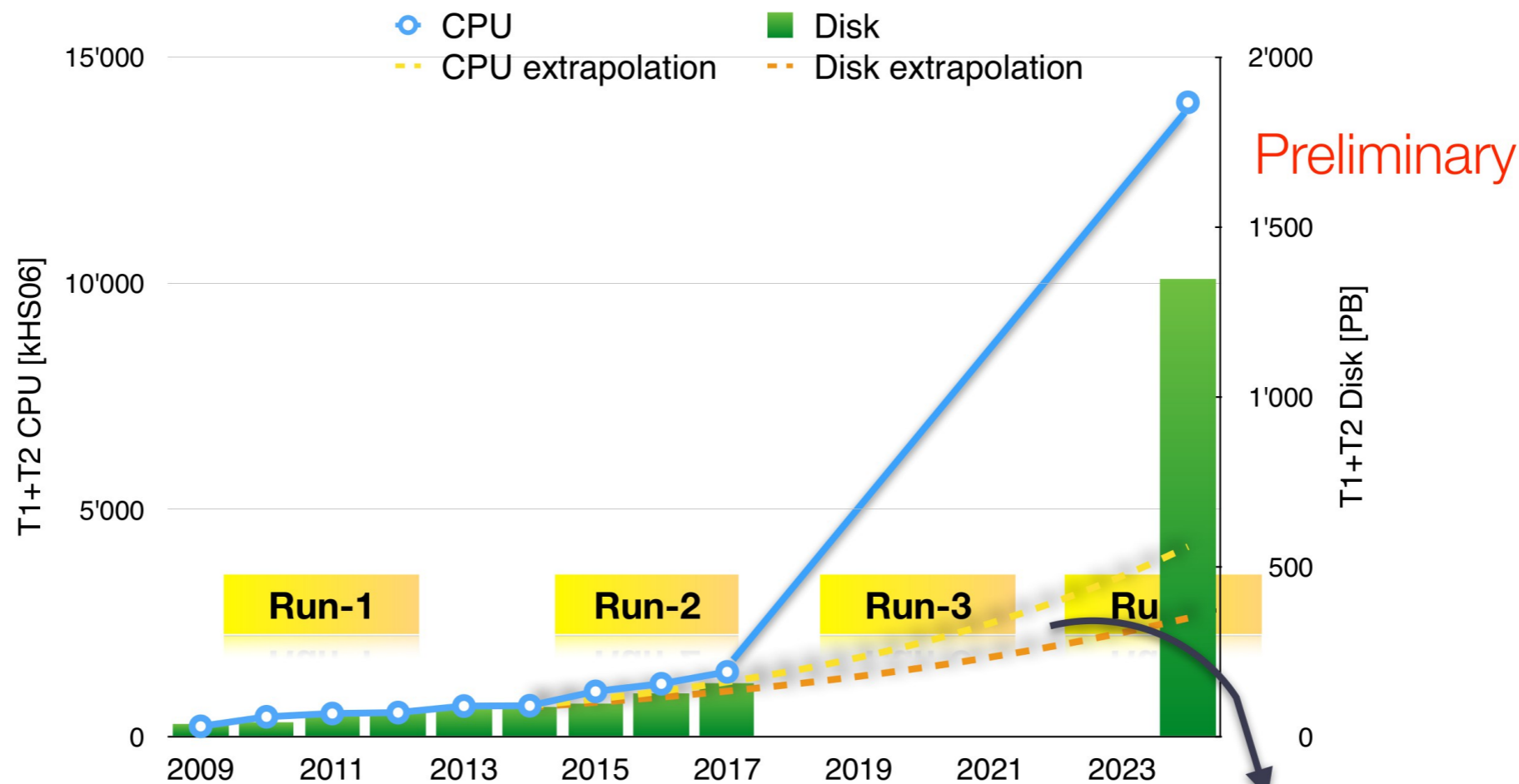


- Problem in Run 1: tuples often 1:1 copy of AODs (root readable); removing duplicated copies frees disk space for important new data -> one of main reasons for xAOD and new analysis framework.
- Resources will be even tighter with higher lumi/EF output rate -> need much more MC (2/5 billion planned for Run 2 for full/fast sim – how much is needed for Run 3 /4 ?).
 - Need to rethink what to store in xAOD files, and take a hit on what can be done with it ... ('redundant' information in AODs in Run 1 was used to apply some important fixes).
- Another application for fast reco/fast sim:
 - events directly to user ntuple to avoid storing large intermediate files never being looked at again.

ATLAS: projections

Run-4 (with 2014 performances)

ATLAS resource needs at T1s & T2s



Eric Lançon on behalf of the ATLAS collaboration

Extrapolation from 2014 (Moore law)

- Need to worry about disk and CPU usage for HL-LHC as well as access to disk (IO and capacity!).

CMS: resource needs



- CMS is planning for 5-7.5kHz of data in Run 4. In this scenario CMS would collect 25B-37B raw events per year.
- Estimating from the current software and using the upgrade simulation: events is more complicated to reconstruct and larger than the events we will collect in 2015.

Detector	Pile-up (Ave./crossing)	Reconstruction time (Ratio to Run 2)	AOD size (Ratio to Run 2)	HLT output rate (kHz)	Total
Phase 1	50	4	1.4	1	3
Phase-II	140	20	3.7	5	65
Phase-II	200	45	5.4	7.5	200

Scale of computing resource needs relative to Run 2 including the increase in projected HLT output rate

- Factoring in the trigger rate and taking a weighed average of the data and simulation tasks: computing challenge is 65-200 times worse than Run 2.
- Anticipating a factor of 8 in CPU improvements and a factor of 2 in code improvement: deficit of a factor of 3-15.
- Anticipating a factor 6 in storage improvements and having by Phase II events 4-5 times larger: deficit of 4-5 in storage.

CMS: targets



- Roughly 40% of the CMS processing capacity is devoted to task identified as reconstruction.
 - Prompt reconstruction, re-reconstruction, data and simulation reco.
 - Improving the number of events that can be reconstructed per computing unit per Swiss Franc is the single biggest savings.
- ~20% of the offline computing capacity is in areas identified as selection and reduction.
 - Analysis selection, skimming, production of reduced user formats.
- The remaining 40% is a mix.
 - Lot of different activities with no single area to concentrate optimisation effort.
 - Simulation already has a strong ongoing optimisation effort.
 - User analysis activities developed by many people.
 - Smaller scale calibration and monitoring activities.

CMS: overview



- CMS is investigating ways to reduce the amount of computing spent on data reduction.
 - Event tags and catalogs can improve the selection speed and efficiency.
 - Big Data tools like Map Reduce can make scalable IO and reuse the selection criteria.
- CMS would like to investigate the scale of improvement in the cost per capacity of using specialised centres for dedicated workflows like reconstruction and event selection.
 - If this is the most efficient way of working, it could be a significant change in how computing services are supported and provisioned.
 - Not all services and capabilities will be at all sites.
 - It would introduce a more heterogeneous and complex system.
 - From an operations perspective and from a support and funding perspective.

Towards the LHCb Upgrade



- No revolution planned for the LHCb computing upgrade (Run 3).
- Rather an evolution to fit in the following boundary conditions:
 - Luminosity levelling at $2 \cdot 10^{33}$
 - Factor 5 c.f. Run 2
 - 100kHz HLT output rate for full physics programme
 - Factor 8-10 more than in Run 2
 - Flat funding for offline computing resources
- Computing milestones for the LHCb upgrade:
 - TDR: 2017Q1
 - Computing model: 2018Q3
- Therefore only brainstorming at this stage, to devise model that keeps within boundary conditions

LHCb: brainstorming for Run 3



- In Run 2, Online (HLT) reconstruction will be very similar to offline (same code, same calibration, fewer tracks).
 - If it can be made identical, why then write RAW data out of HLT, rather than Reconstruction output?
- In Run 2 LHCb will record 2.5 kHz of “TurboDST”.
 - RAW data plus result of HLT reconstruction and HLT selection.
 - Equivalent to a microDST (MDST) from the offline stripping.
- Proof of concept: can a complete physics analysis be done based on a MDST produced in the HLT?
 - No offline reconstruction.
 - No offline realignment, reduced opportunity for PID recalibration.
 - RAW data remains available as a safety net.
- If successful, can RAW data be dropped?
 - HLT then writes out ONLY the MDST.
- Currently just ideas, but would allow a 100kHz HLT output rate without an order of magnitude more computing resources.

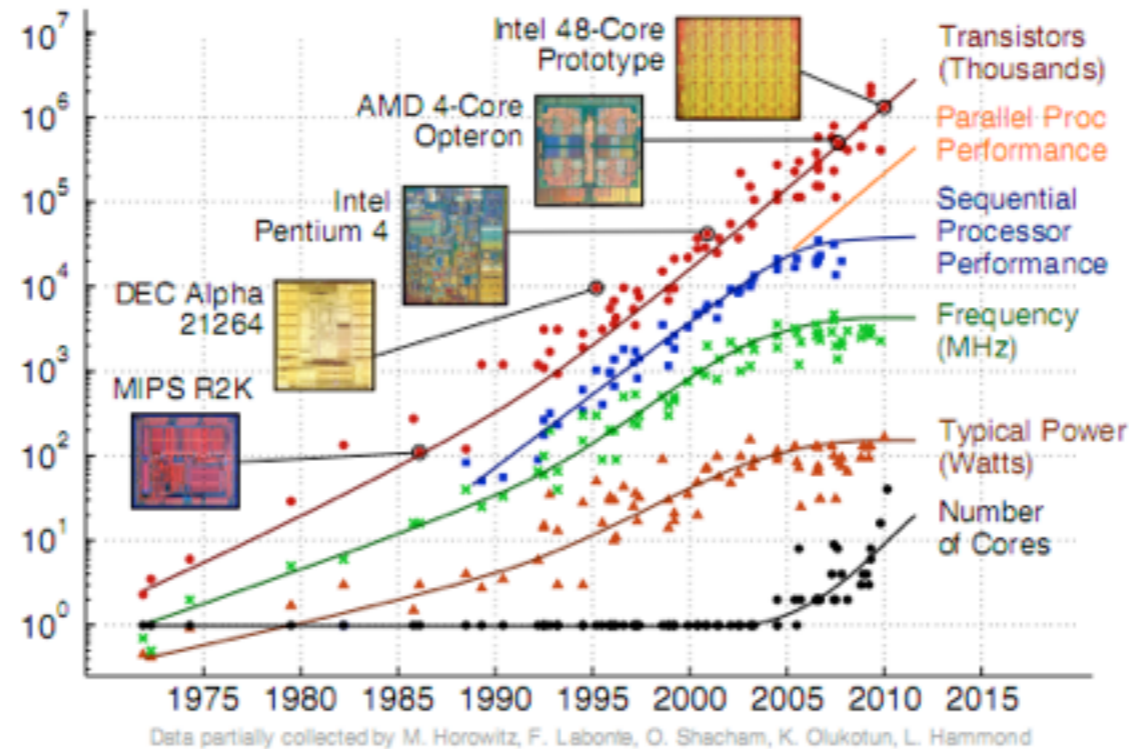
LHCb: simulation

- LHCb offline CPU usage is dominated by simulation (>60% of CPU already in 2016).
 - Many measurements start to be limited by simulation statistics.
- Simulation suited for execution on heterogeneous resources.
 - Pursue efforts to interface Dirac framework to multiple computing platforms.
 - Allow opportunistic and scheduled use of new facilities.
 - Extend use of HLT farm during LHC stops.
- Several approaches to reduce CPU time per event.
 - Code optimisation, vectorisation etc.
 - Contribute to and benefit from community wide activities, e.g. for faster transport.
 - Fast simulations.
 - Not appropriate for many detailed studies for LHCb precision measurements.
 - Nevertheless many generator level studies are possible.
 - Hybrid approach.
 - Full simulation for signal candidates only.
 - Fast techniques for the rest.
 - e.g. skip calorimeter simulation for out of time pileup.
- Avoid being limited by disk space.
 - Deploy MDST format also for simulated data.

What do we need to do?

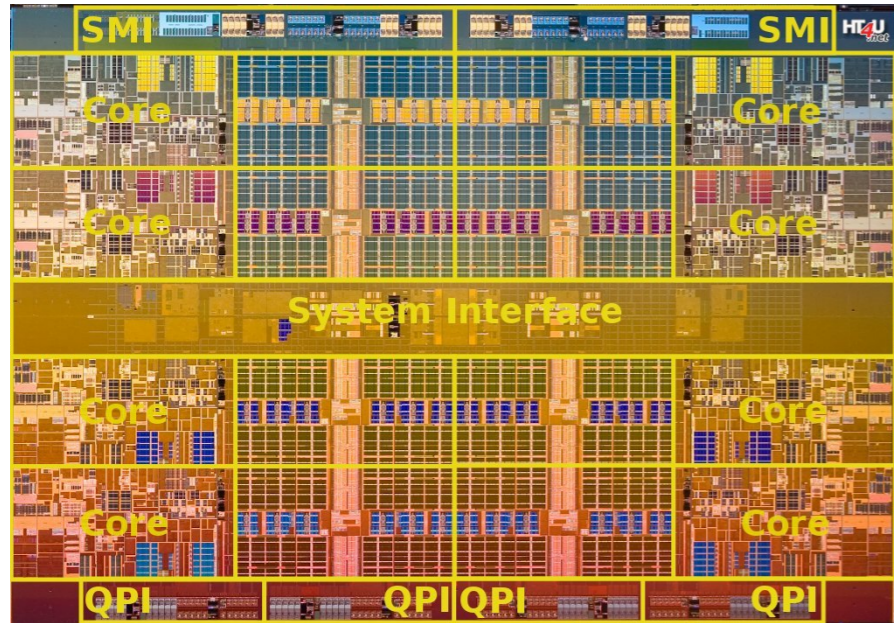
- >60% grid usage is MC.
 - Speedup existing frameworks.
 - Fast (parametrised) MC.
 - Optimize storage format.
- External HPC facilities (Titan, Mira), typically ran at ~90% efficiency.
 - for Titan it means ~300M core hours per year.
 - Frameworks to utilise this efficiently (e.g. PanDA).
- Use clouds for more optimised workflows.
- Mind IO performance on active storage for analysis.
- Rethink the data storage strategies?

Technology evolution



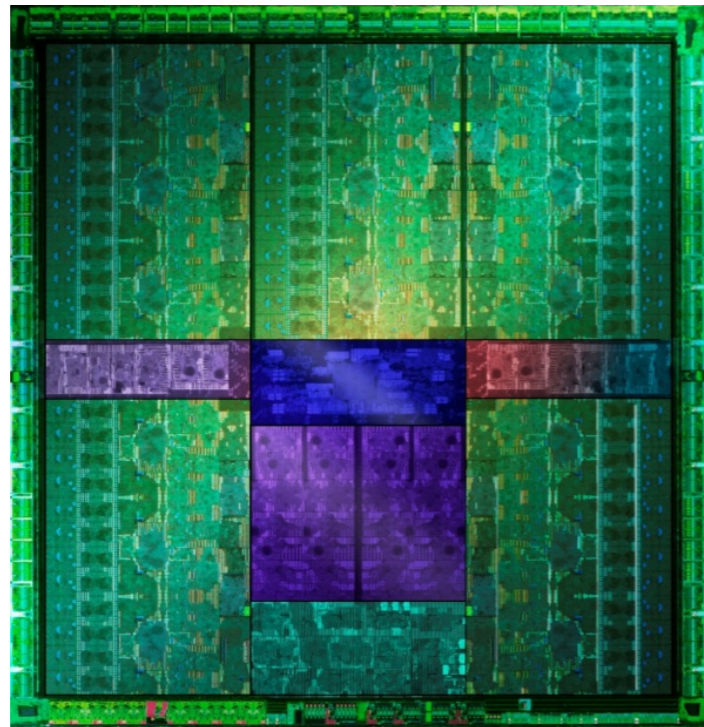
- Clock speed, power (per socket), performance per clock - flat since ~2006.
 - Issue: power dissipation/distribution.
- Number of transistors still growing exponentially (more cores added).
- Memory wall - see Graeme's talk.
- Disk capacity to performance ratio.

Silicon utilisation

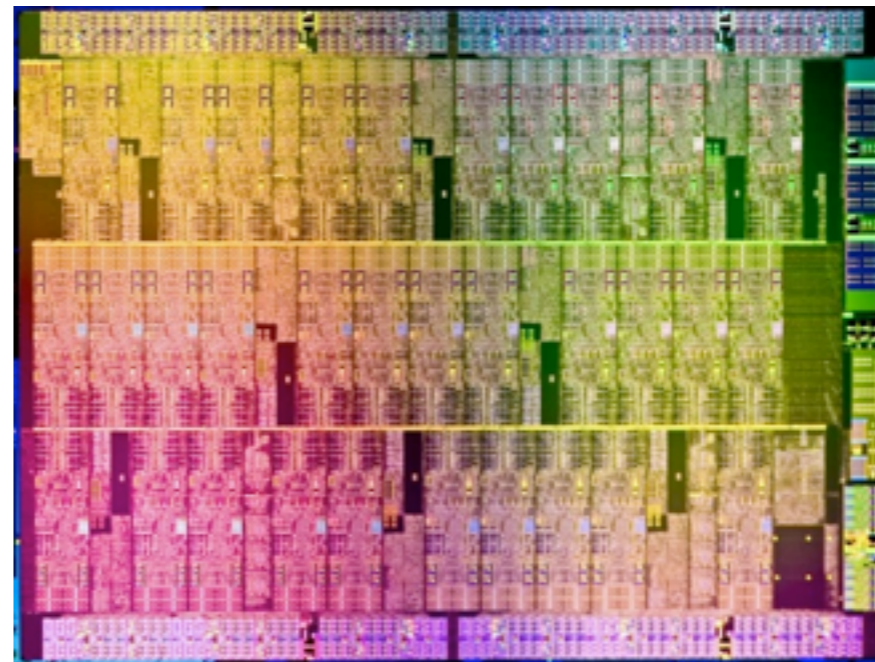


Intel Nehalem

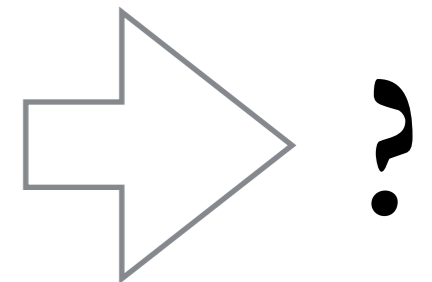
- CPU: only part of silicon used for ALUs.
- Trend to utilise more area, e.g. in accelerator boards (GPUs, etc...).
- Power dissipation (and distribution) problem also here.
 - Dark Silicon.



NVIDIA Kepler GPU



Intel Xeon Phi

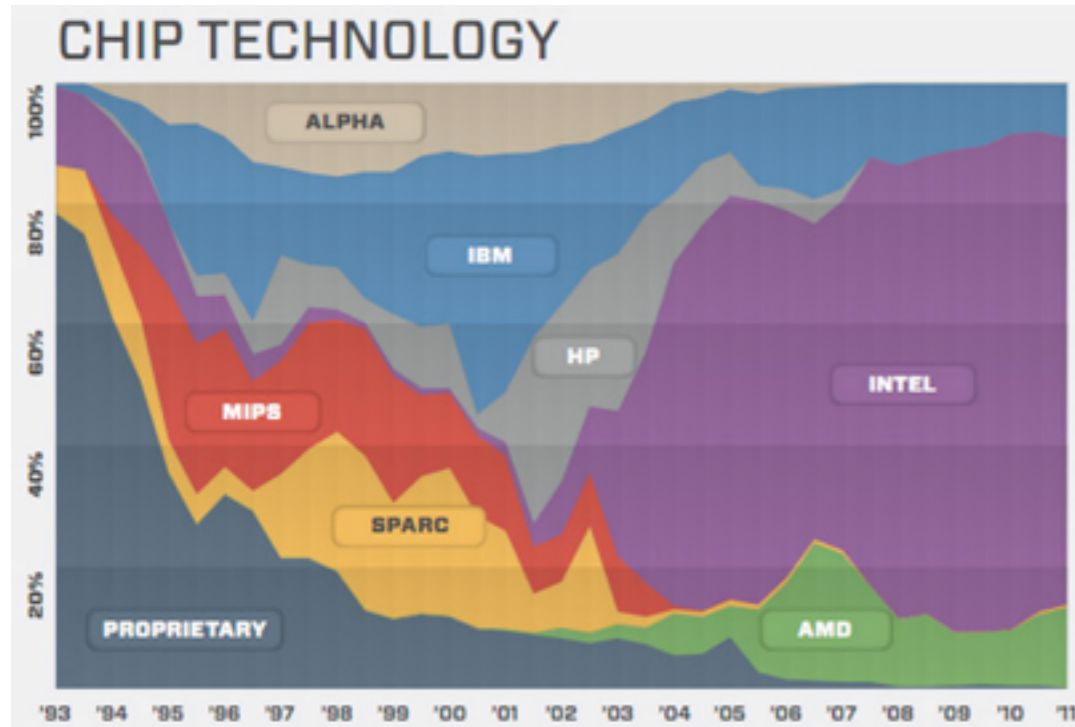


Life in a multi-core landscape

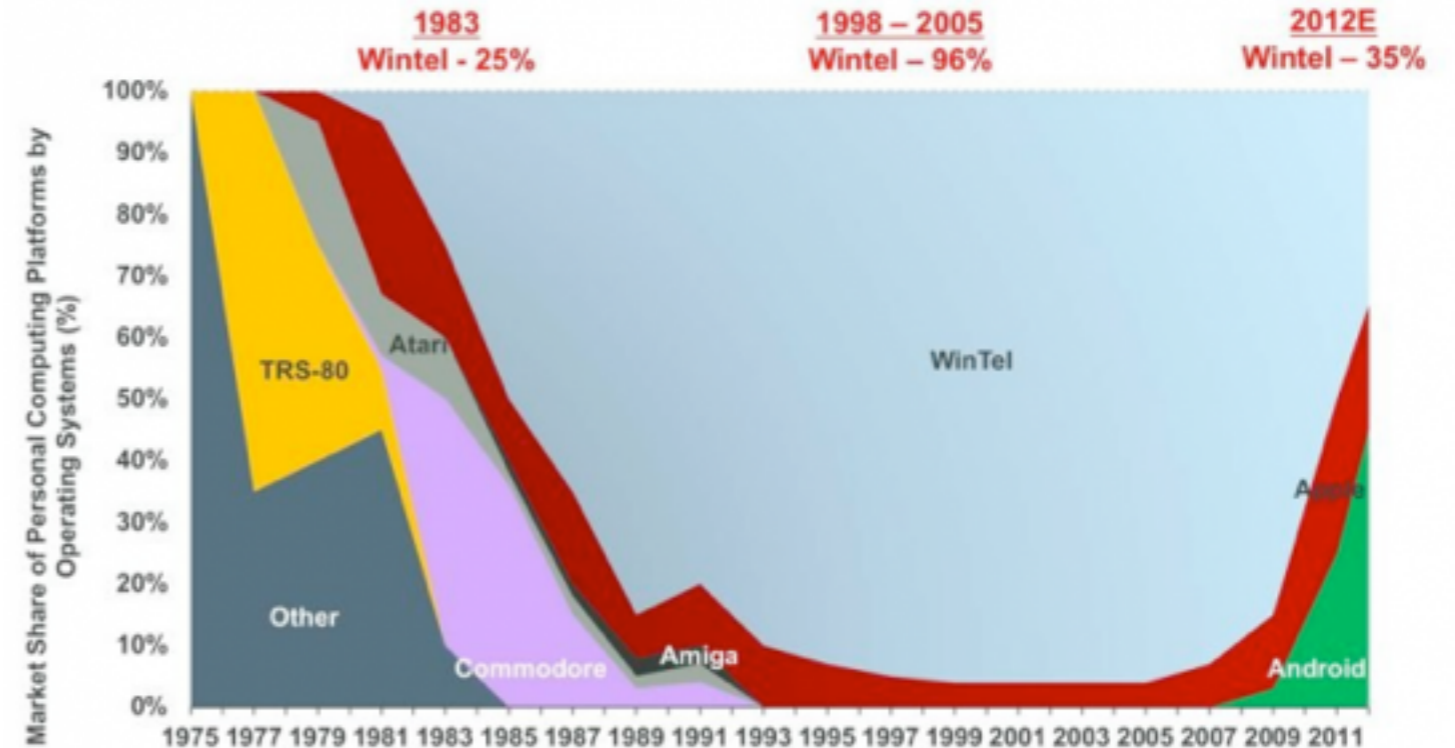
- Shift data processing paradigm to utilise the silicon more efficiently.
- Use heterogeneous systems: CPU+specialised coprocessors (FPGA+GPU).
- Adapt code where appropriate to use coprocessors.
- Multi-core utilisation.
 - Possible memory issues?
 - Multi-threading to relieve part of memory strain.
- Code optimisations:
 - e.g. vectorised code.

(see Graeme's talk)

Industry trends



Global Market Share of Personal Computing Platforms by Operating System Shipments, 1975 – 2012E



KPCB

Source: Asymco.com (as of 2011), Public Filings, Morgan Stanley Research, Gartner for 2012E data. 2012E data as of Q3:12.

24

- In the past: scientific computing dominated by specialised architectures.
- Industry “settled” on x86 at some point.
 - We followed suit - standardised on Intel/Linux - commodity hardware/software.
- Market trends nowadays: other architectures emerge.
- Big players (Google, Facebook, ...) do Big Data differently (few specialised HPC farms, etc.)
 - Synergy between architectures: mobile end (e.g. ARM) and big server backend (e.g. Xeon).
- Have to rethink again?

Summary

- Resource needs large.
 - Technology evolution alone will not close the gap.
- Efforts on the way or already undertaken.
 - Storage.
 - CPU.
- Large online farms for data compression or online triggering planned.
 - Usage for offline duty, simulation, etc.
- We need to investigate also other resources.
 - Spare cycles of BIG computing centres.
- Keep an eye out for new hardware developments...