

# EXPLORATION OF FUTURE COMPUTING PLATFORMS AT CMS

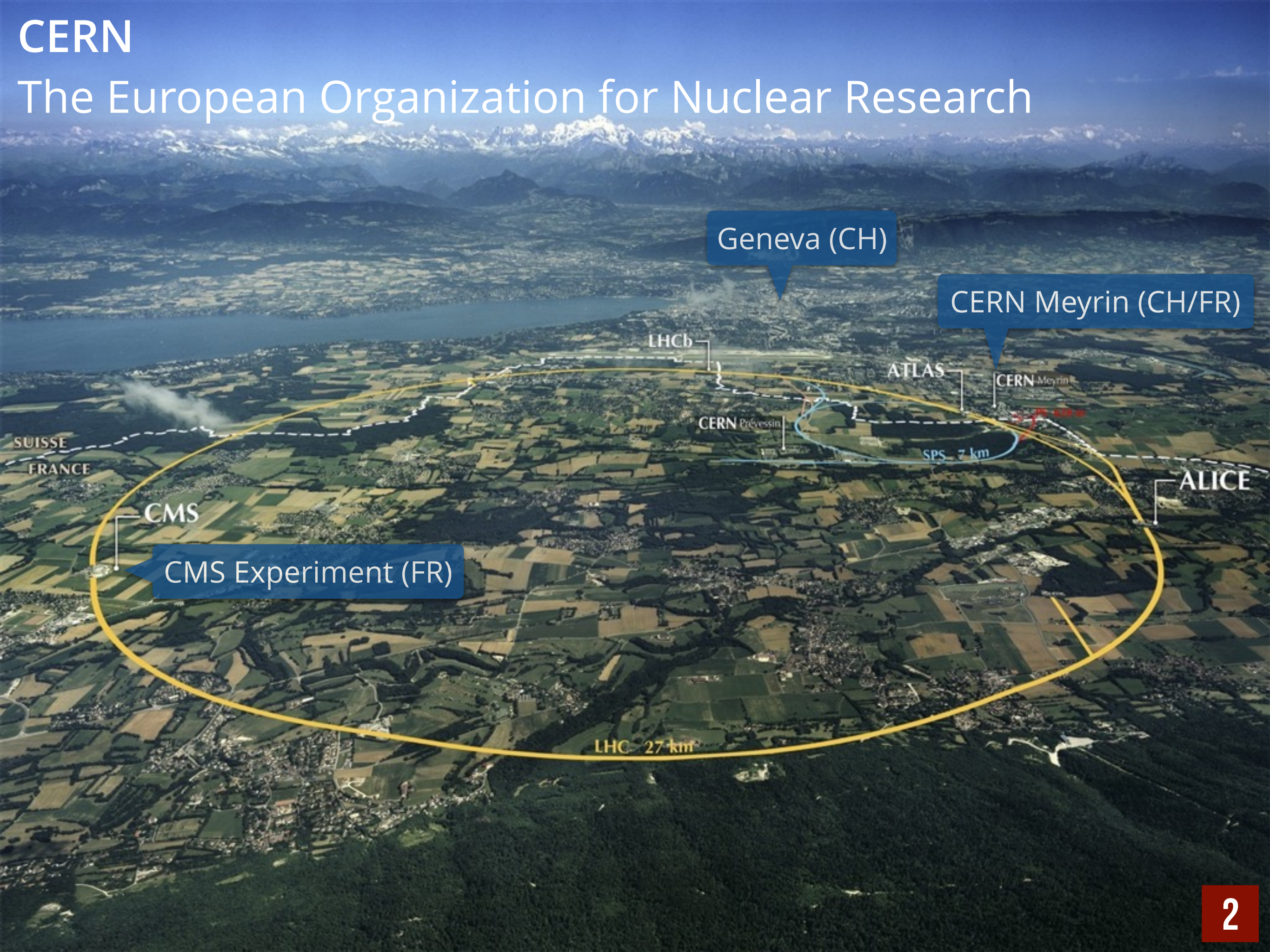
COMPUTING ON LOW-POWER ARCHITECTURES (COLA), 25.02.2016  
DAVID ABDURACHMANOV (FERMILAB)





# CERN

## The European Organization for Nuclear Research



Geneva (CH)

CERN Meyrin (CH/FR)

LHCb

ATLAS

CERN Meyrin

CERN Prévessin

SPS 7 km

SUISSE  
FRANCE

CMS

CMS Experiment (FR)

ALICE

LHC 27 km



# Power In Data Centers

## **An Inconvenient Truth**

- ▶ Energy-related costs account for approximately 12 percent of overall data center expenditure and are the **fastest-rising cost in the data center**, according to Gartner, Inc. (September 29, 2010)
- ▶ CMS for 2012 data used ~100K x86\_64 cores from ~350K cores at Worldwide LHC Computing Grid (WLCG)
- ▶ Scaling up from the mix of machines at FNAL we estimate WLCG aggregate power consumption for machines at 10MW
- ▶ CMS expects 2 to 3 orders of magnitude increase in data produced in 15 years

## **Think Green**

- ▶ Local green or/and cheaper power source, e.g., Princeton energy plant (15MW) combines electricity, heat and cooling. When electricity cost increased gas, diesel or/and bio-diesel fuel is used to power local generators. Hot water and steam is provided from waste energy.
- ▶ Low-power and / or highly efficient hardware, e.g., Intel Atom, X-Gene (ARMv8 64-bit), GPUs, Xeon Phi, FPGA, etc.

# Computing & CMS

**CMS Detector (HLT)**



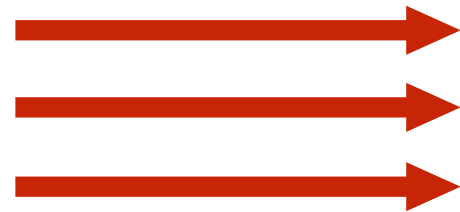
**Worldwide LHC Computing Grid (WLCG)**



**Full ownership** ✓

**Single "customer"** ✓

**High-bandwidth interconnect** ✓



**Partially owned**

**Multiple "customers"**

**Bandwidth varies**

**A virtual super computer (WLCG) is used to store, distribute and process LHC data**

Based on 170 computing centres in 42 countries

Distribute and analyse ~30PB of data annually generated by LHC

Experiments produce >15PB of new data annually

# Why new architectures?

Distributed computing in **HEP before ~2000** had multiple vendors involved, and incl. special workstations and heterogeneous computing

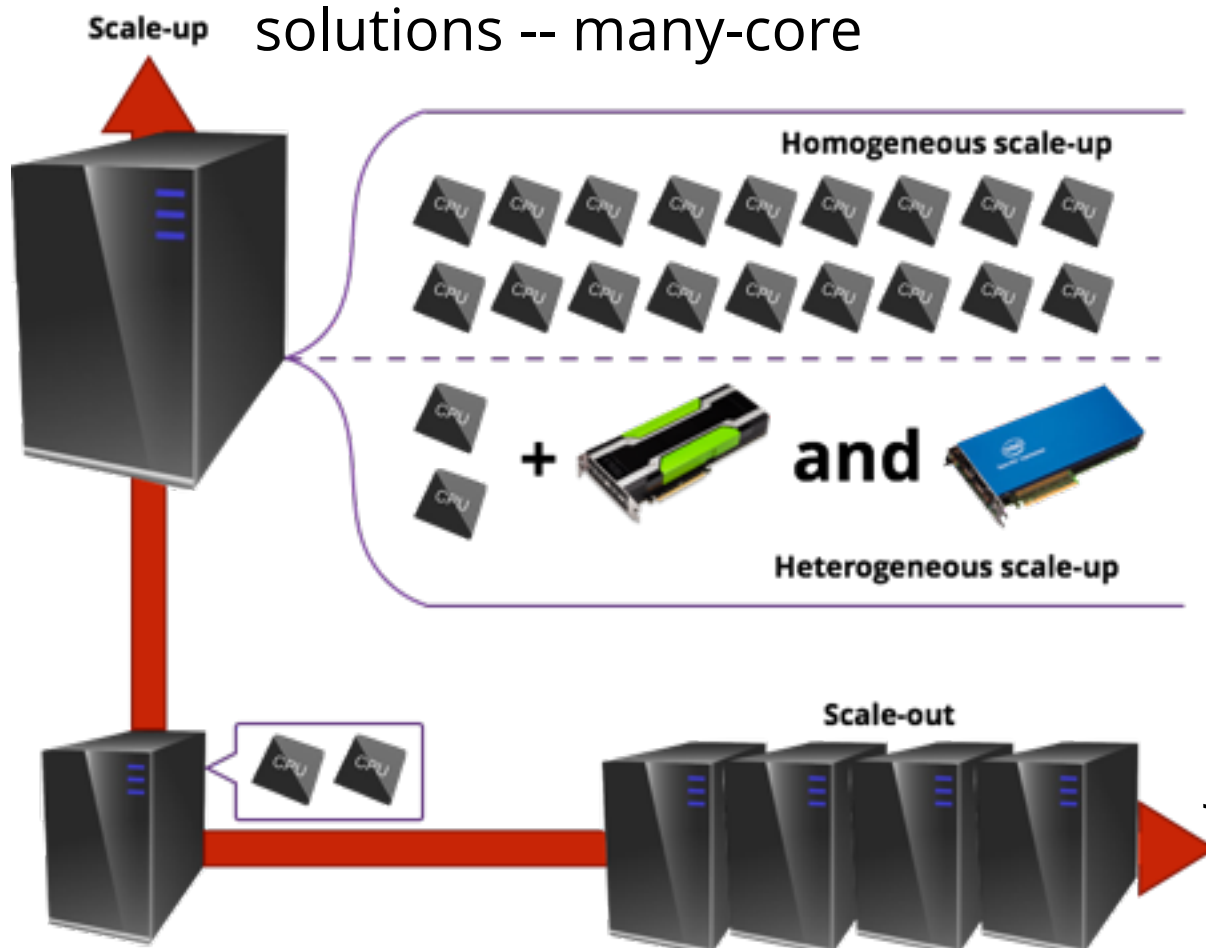
High Throughput Computing (HTC) converged on x86/Linux at ~2000

Commodity hardware enabled the current model of WLCG:

## Build Once, Run Everywhere

**Two vendors:** Intel (dominating) and AMD

The current commodity hardware itself is limited by power wall with stop-gap solutions -- many-core



Specialised processors and heterogeneous computing rise up

Lightweight general-purpose low-power high-density, vector units, GPUs, Xeon Phi (highly-parallel long-vector), etc

The focus is shifting to **performance/watt**, not just **performance/price**

# How we do it?

No single job batch submission system, incl. **LSF, HTCondor, Slurm, SGE, Torque/Pbs**

No single storage solution, incl. NFS, GlusterFS, **Hadoop** (popular in US)

Has 100+ different CPUs from the last **10 years**, most 4-5 years old

Common operating system: **RHEL/CentOS/Scientific Linux (SL)**

Dominated by **SL 6** co-developed by CERN and Fermilab

**CentOS 7 + CERN Special Interest Group** to follow **SL 6**

Software and essential precomputed data (e.g. LUT) distributed via **CernVM File System (CVMFS)**

**HEP SPEC '06** benchmark is used for accounting in WLCG and by experiments

Designed to represent worker node activity under full load

Based on CPU SPEC 2006 **all\_cpp** benchmark set

# CMS Software Bundle

The actual application software for "pattern recognition", "simulation", etc.

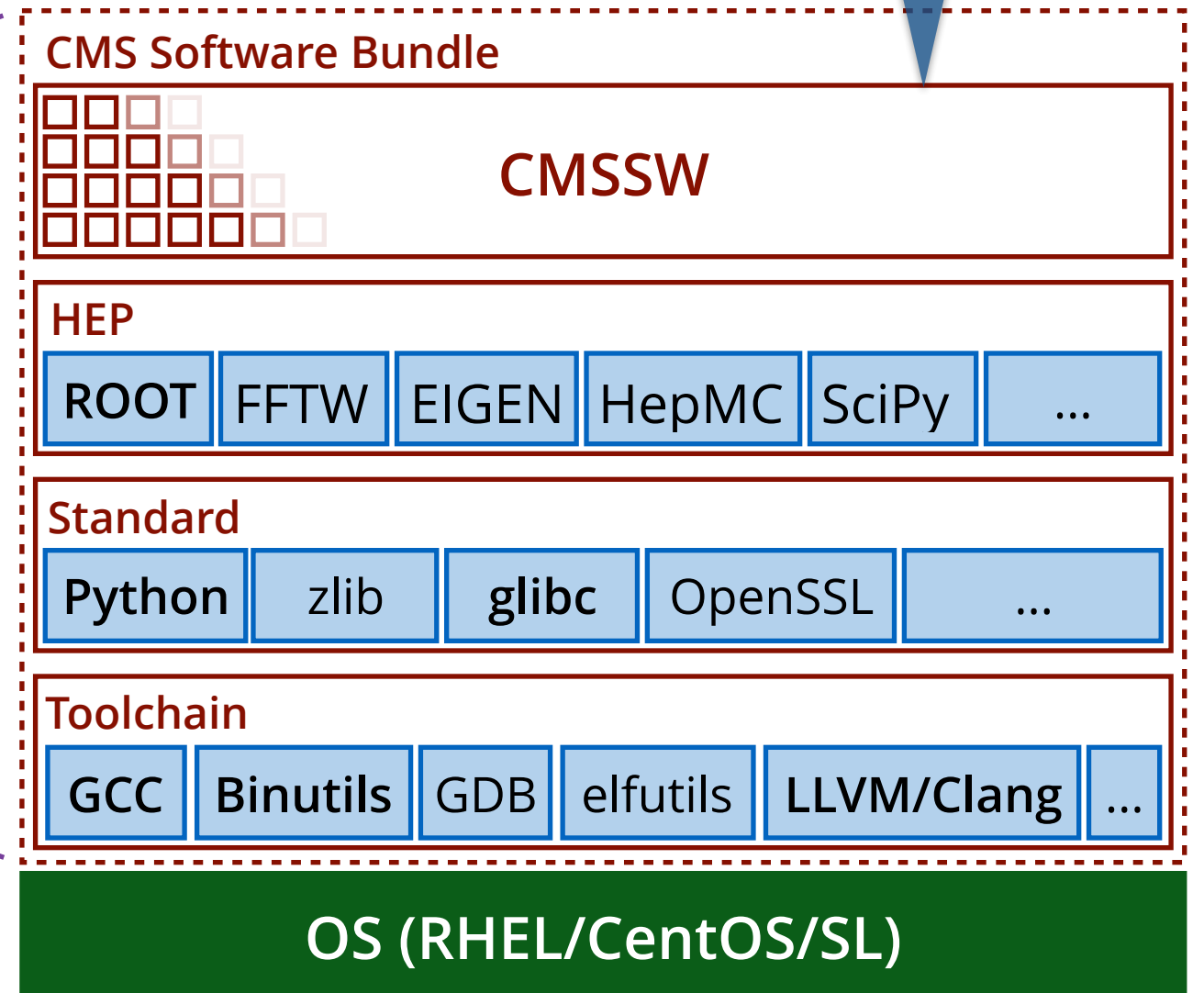
CMSSW is **open-source** and available at GitHub

Mostly written in **C++14**, **C**, **Python** and **Fortran**

CMSSW is like **Software Collection** package or **Linux Container** without actually being any of them

Quick comparison:

CVMFS



	CMSSW	Firefox	Other CERN developed software would increase SLOCs
SLOCs	6M	7M	
Initial Release	2005	2002	
Contributors	>1300	>1200	ROOT6 w/o Clang: 1.7M
Memory Footprint	~2GB	~0.3GB	GEANT4: 1.1M

# Porting to **ARMv8 (64-bit)**

CMSSW was originally ported to **ARMv7 (32-bit)** few years ago

High-end mobile SoC based development boards were used

ODROID-U2 (Exynos 4412 Prime), ODROID-XU2 (Exynos 541), Arndale Octa (Exynos 5420), Jetson TK1 (Tegra K1)

Resolved majority of porting issues and found numerous issues in CMSSW (even affecting x86\_64)

CMSSW for **ARMv8 (64-bit)** port was started early

**Step1:** ARM Foundation Model

**Step2:** QEMU + binfmt\_misc + user mode emulation

**Step3:** APM Mustang

**Step4:** HP Moonshot + m400

For ARMv8 we wanted to have CMSSW application software and GRID software (e.g., **HTCondor**) for software distribution, data transfers and job management

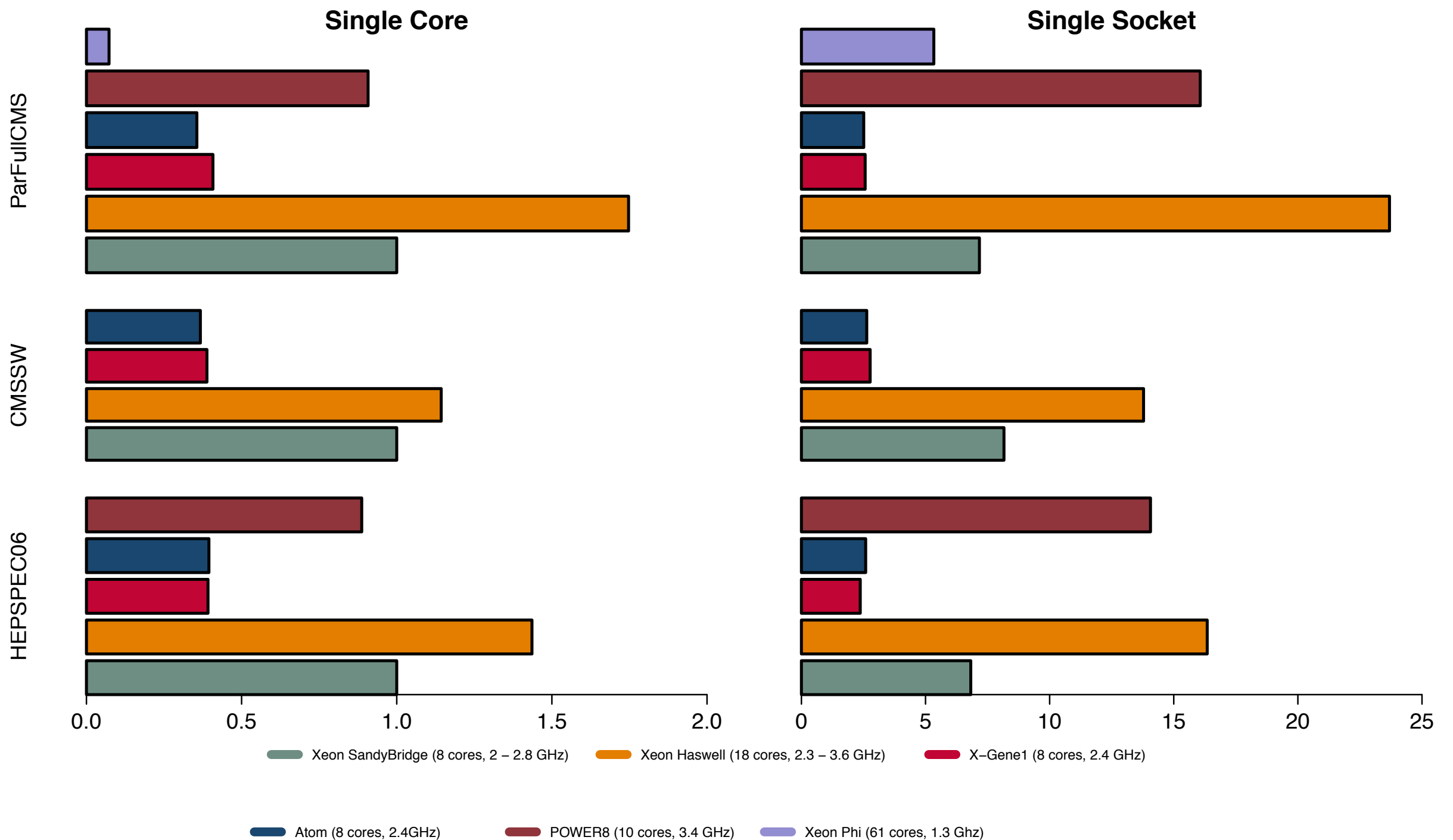


# CPU Specifications

	Vendor	Model	Year	Fab	Process
SandyBridge	Intel	E5-2650	Q1/12	Intel	32nm
Haswell	Intel	E5-2699	Q3/14	Intel	22nm
Atom	Intel	C2750	Q3/13	Intel	22nm
X-Gene 1	APM	883408	Q3/13	TSMC	40nm
POWER8	IBM	8247-22L	Late 13	IBM	22nm
Xeon Phi	Intel	KNC7100	Q2/14	Intel	22nm

	Frequency (GHz)	Cores	Threads/Core
SandyBridge	2.0 (2.8)	8	2
Haswell	2.3 (3.6)	18	2
Atom	2.4	8	1
X-Gene 1	2.4	8	1
POWER8	3.45	10	8
Xeon Phi	1.23	61	4

# Raw Performance

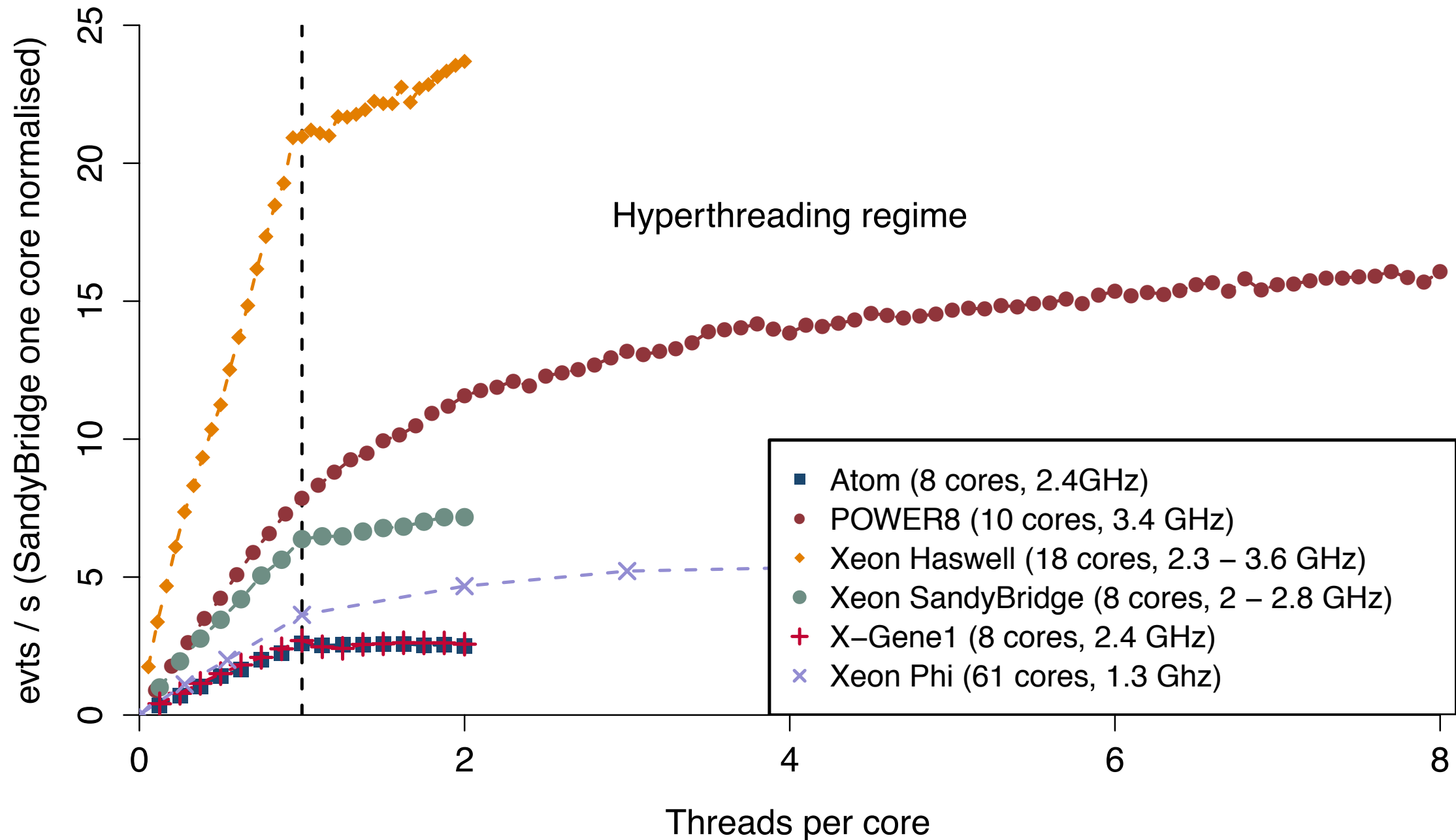


All numbers normalised to Xeon SandyBridge 1 core performance.



# Scalability #1

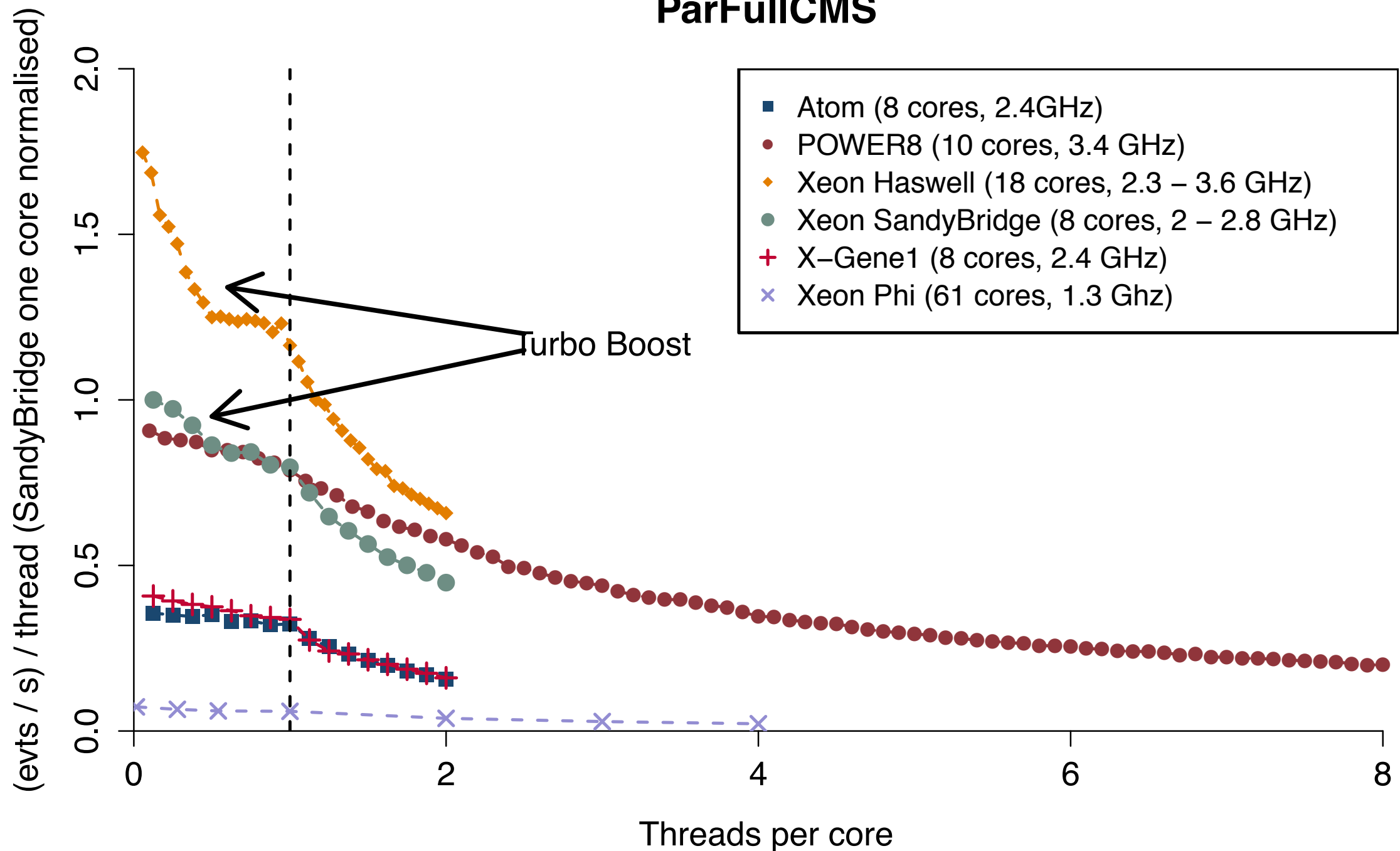
## ParFullCMS



All numbers normalised to Xeon SandyBridge 1 core performance.

# Scalability #2

## ParFullCMS

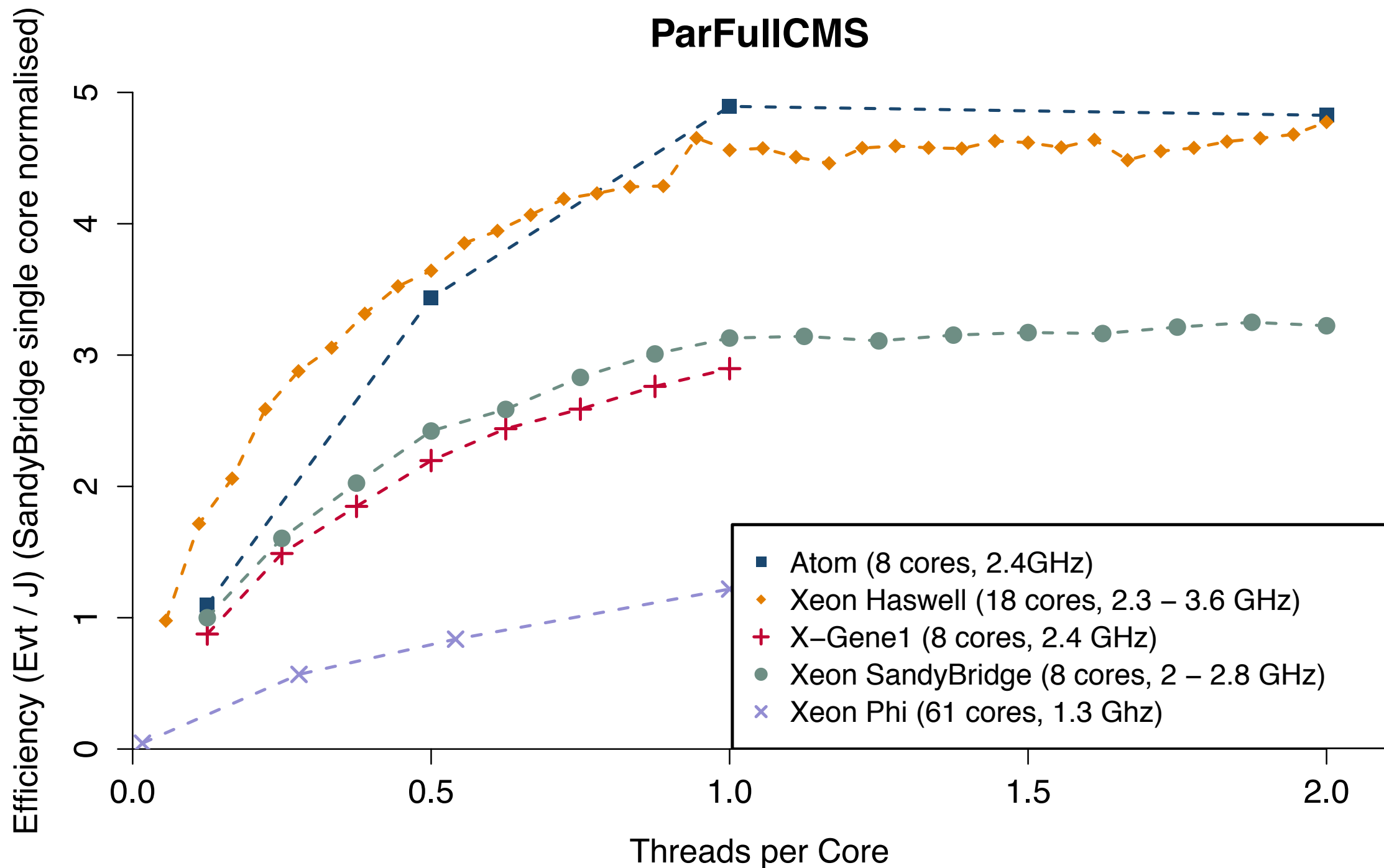


All numbers normalised to Xeon SandyBridge 1 core performance.





# Power Efficiency (1S) #2

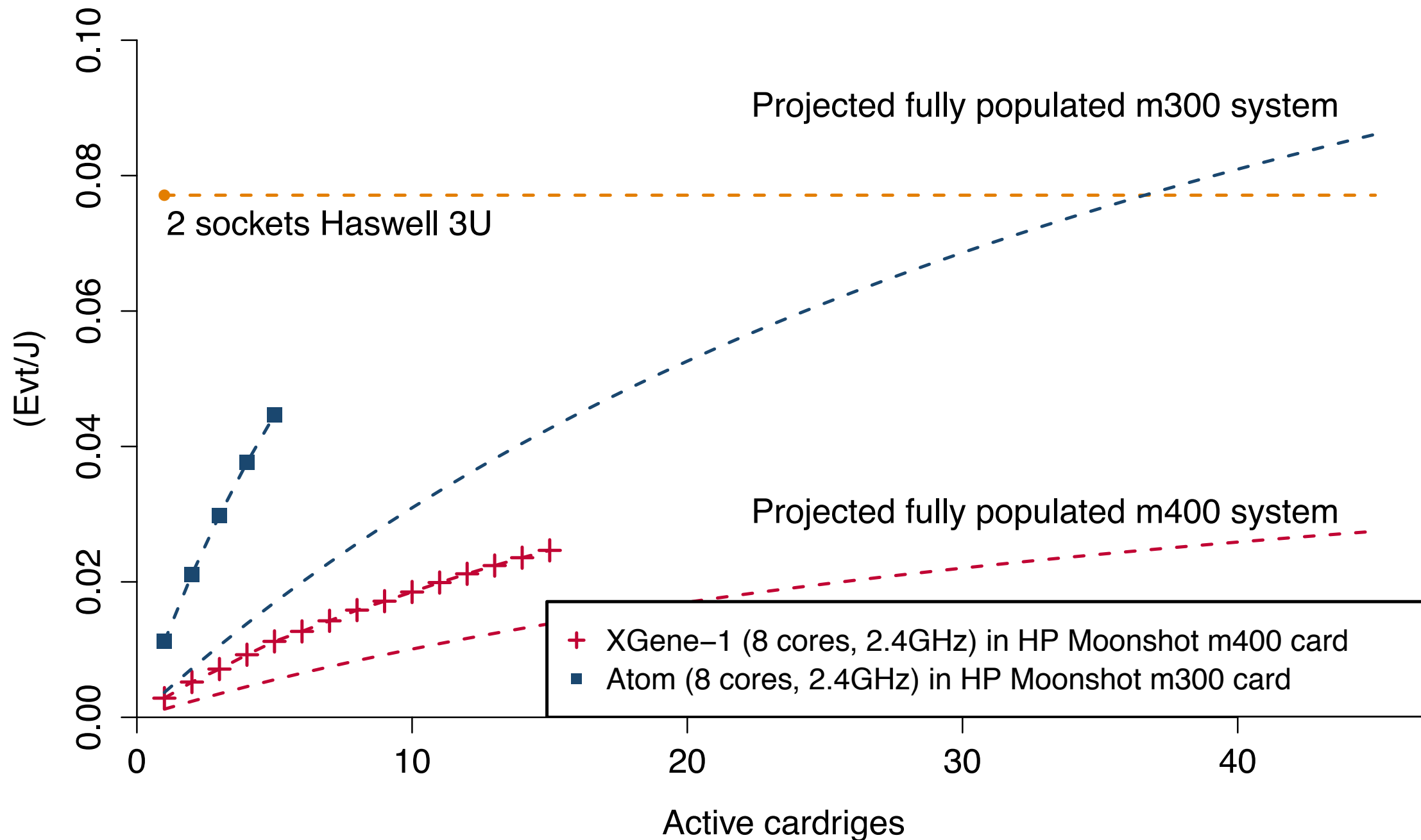


All numbers normalised to Xeon SandyBridge 1 core performance.



# Power Efficiency (box)

## ParFullCMS



# Heterogeneous Tier-3 Site

**Goal 1:** What is necessary for **AArch64-based** (or any other alternative architecture) production worker nodes to be a credible alternative to **x86\_64-based** nodes for use in **WLCG** computing sites (given the availability of application level software like **CMSSW**)?

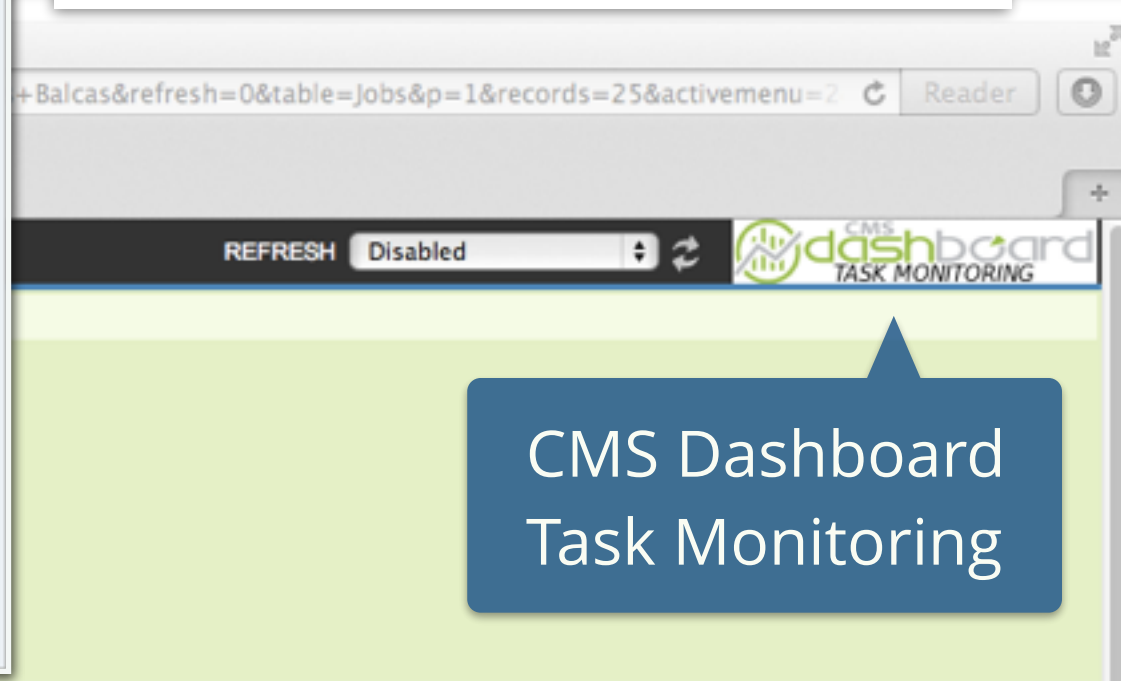
**Goal 2:** We wanted to demonstrate that such nodes can be added as a "drop-in" replacement for **x86\_64** nodes in **WLCG** and even mixed heterogeneously.

With above goals in mind, we created **US\_T3\_Princeton\_ARM** computing site using APM Mustang development board with Open Science Grid (OSG) infrastructure at Princeton University.





On mid-2015 CMS successfully executed CMSSW based job on AArch64 worker node via standard job injection pipeline



Start » [ Justas Balcas ] » Tasks » Jobs

Data Charts Show 25 entries Task: 150608\_200051:jbalkas\_crab\_ARM\_TEST\_2-output2 NJobTotal: 1000 Pending: 822 Running: 0 Unknown: 0 Cancelled: 0 Success: 168 Failed: 2 WNPostProc: 8 ToRetry: 0

Id	Status	AppExitCode	Site	Retries	Submitted	Started	Finished	Wall Time	Job Log	File Access	FTS File Status
1	finished	0	T3_US_Princeton_ARM	1	2015-06-08T20:01:22	2015-06-08T20:05:35	2015-06-08T20:15:16	00:09:41	Job Log, Job Log JSON, Post Job Log	File Info	N/A
2	finished	0	T3_US_Princeton_ARM	1	2015-06-08T20:01:22	2015-06-08T20:05:35	2015-06-08T20:15:17	00:09:38	Job Log, Job Log JSON, Post Job Log	File Info	N/A
3	finished	0	T3_US_Princeton_ARM	1	2015-06-08T20:01:22	2015-06-08T20:05:37	2015-06-08T20:15:25	00:09:48	Job Log, Job Log JSON, Post Job Log	File Info	N/A
4	finished	0	T3_US_Princeton_ARM	1	2015-06-08T20:01:22	2015-06-08T20:05:37	2015-06-08T20:15:32	00:09:55	Job Log, Job Log JSON, Post Job Log	File Info	N/A
5	finished	0	T3_US_Princeton_ARM	1	2015-06-08T20:01:22	2015-06-08T20:05:37	2015-06-08T20:15:34	00:09:57	Job Log, Job Log JSON, Post Job Log	File Info	N/A
6	finished	0	T3_US_Princeton_ARM	1	2015-06-08T20:01:22	2015-06-08T20:09:29	2015-06-08T20:16:00	00:06:31	Job Log, Job Log JSON, Post Job Log	File Info	N/A
7	finished	0	T3_US_Princeton_ARM	1	2015-06-08T20:01:22	2015-06-08T20:24:29	2015-06-08T20:28:52	00:04:23	Job Log, Job Log JSON, Post Job Log	File Info	N/A
8	finished	0	T3_US_Princeton_ARM	1	2015-06-08T20:01:22	2015-06-08T20:24:29	2015-06-08T20:29:03	00:04:34	Job Log, Job Log JSON, Post Job Log	File Info	N/A
9	finished	0	T3_US_Princeton_ARM	1	2015-06-08T20:01:22	2015-06-08T20:24:30	2015-06-08T20:29:32	00:05:02	Job Log, Job Log JSON, Post Job Log	File Info	N/A
10	finished	0	T3_US_Princeton_ARM	1	2015-06-08T20:01:22	2015-06-08T20:24:31	2015-06-08T20:28:10	00:03:39	Job Log, Job Log JSON, Post Job Log	File Info	N/A

The first AArch64 based WLCG site (demonstrator)



```

1 [|||||] 9 ] [|||||] 100.0%]
2 [|||||] 98.0%] 6 [|||||] 98.0 ]
3 [|||||] 100.0%] 7 [|||||] 98.0 ]
4 [|||||] 100.0%] 8 [|||||] 100.0%]
Mem[|||||] ] Tasks: 187, 19 thr; 11 running
Swp[|||||] ] Load average: 5.18 2.48
Uptime: 33 days, 02:07:36

```

PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
6128		20	0	2042M	41224	1948	S	0.0	0.3	0:00.16	/usr/bin/cvmfs2 -o rw,fsname=cvmfs2,allow_other,grab_mountpoint,uid=997,gid=995 cms.cern.ch /cvmfs/cms.cern.ch
6127		20	0	2042M	41224	1948	S	0.0	0.3	0:00.17	/usr/bin/cvmfs2 -o rw,fsname=cvmfs2,allow_other,grab_mountpoint,uid=997,gid=995 cms.cern.ch /cvmfs/cms.cern.ch
6120		20	0	2042M	41224	1948	S	0.0	0.3	0:00.18	/usr/bin/cvmfs2 -o rw,fsname=cvmfs2,allow_other,grab_mountpoint,uid=997,gid=995 cms.cern.ch /cvmfs/cms.cern.ch
23248		20	0	16236	6420	5252	S	0.0	0.0	0:00.52	/usr/sbin/condor_master -f
23256		20	0	17728	8000	5688	S	2.0	0.0	5:04.76	condor_startd -f
30301		20	0	16688	6736	5492	S	0.0	0.0	0:00.09	condor_starter -f -a slot4 byggvir.Princeton.EDU
30305		30	10	3744	1848	1208	S	0.0	0.0	0:00.62	/bin/bash /var/lib/condor/execute/dir_30301/condor_exec.exe -v std -name gfactory_instance -entry CMS_T3_U
2478		30	10	3468	1548	1208	S	0.0	0.0	0:00.12	/bin/bash /var/lib/condor/execute/dir_30301/glide_NRPbun/main/condor_startup.sh glidein_config
3191		30	10	17884	8272	6320	S	0.0	0.1	0:00.16	/var/lib/condor/execute/dir_30301/glide_NRPbun/main/condor/sbin/condor_master -f -pidfile /var/lib/c
3194		30	10	18928	9140	6748	S	0.0	0.1	0:00.87	condor_startd -f
2898		30	10	17012	8324	6552	S	0.0	0.1	0:00.16	condor_starter -f vocns058.cern.ch
4428		30	10	3352	1456	1196	S	0.0	0.0	0:00.10	/bin/bash /var/lib/condor/execute/dir_30301/glide_NRPbun/execute/dir_2898/condor_exec.exe -
4585		30	10	3520	1520	1224	S	0.0	0.0	0:00.02	sh ./CMSRunAnalysis.sh -a sandbox.tar.gz --sourceURL=https://cmsweb.cern.ch/crabcache --
4631		30	10	23508	13492	1572	S	0.7	0.1	0:00.70	python CMSRunAnalysis.py -r /var/lib/condor/execute/dir_30301/glide_NRPbun/execute/di
5236		30	10	3624	1648	1160	S	0.0	0.0	0:00.01	/bin/bash /var/lib/condor/execute/dir_30301/glide_NRPbun/execute/dir_2898/cmsRun-m
5281	uscms01	30	10	921M	588M	115M	R	93.7	3.7	4:07.20	cmsRun -j FrameworkJobReport.xml PSet.py
3193		30	10	7024	4072	1100	S	0.0	0.0	0:00.71	condor_procd -A /var/lib/condor/execute/dir_30301/glide_NRPbun/log/procd_address -L /var/lib/cond
30119		20	0	16688	6724	5492	S	0.0	0.0	0:00.08	condor_starter -f -a slot1 byggvir.Princeton.EDU
30123		30	10	3744	1848	1208	S	0.0	0.0	0:00.62	/bin/bash /var/lib/condor/execute/dir_30119/condor_exec.exe -v std -name gfactory_instance -entry CMS_T3_U
2156		30	10	3472	1548	1208	S	0.0	0.0	0:00.12	/bin/bash /var/lib/condor/execute/dir_30119/glide_LreWcj/main/condor_startup.sh glidein_config
2871		30	10	17884	8272	6320	S	0.0	0.1	0:00.16	/var/lib/condor/execute/dir_30119/glide_LreWcj/main/condor/sbin/condor_master -f -pidfile /var/lib/c
2874		30	10	18952	9168	6748	S	0.0	0.1	0:00.87	condor_startd -f
2892		30	10	17416	8676	6568	S	0.0	0.1	0:00.16	condor_starter -f vocns058.cern.ch
3431		30	10	3352	1456	1196	S	0.0	0.0	0:00.10	/bin/bash /var/lib/condor/execute/dir_30119/glide_LreWcj/execute/dir_2892/condor_exec.exe -
3638		30	10	3520	1516	1224	S	0.0	0.0	0:00.02	sh ./CMSRunAnalysis.sh -a sandbox.tar.gz --sourceURL=https://cmsweb.cern.ch/crabcache --
3692		30	10	23508	13256	1340	S	0.0	0.1	0:00.70	python CMSRunAnalysis.py -r /var/lib/condor/execute/dir_30119/glide_LreWcj/execute/di
4965		30	10	3624	1648	1160	S	0.0	0.0	0:00.01	/bin/bash /var/lib/condor/execute/dir_30119/glide_LreWcj/execute/dir_2892/cmsRun-m
5104		30	10	917M	566M	98616	R	97.6	3.5	4:07.37	cmsRun -j FrameworkJobReport.xml PSet.py
2873		30	10	6924	3412	1100	S	0.0	0.0	0:00.63	condor_procd -A /var/lib/condor/execute/dir_30119/glide_LreWcj/log/procd_address -L /var/lib/cond
24914		20	0	16688	6740	5492	S	1.3	0.0	0:00.09	condor_starter -f -a slot7 byggvir.Princeton.EDU
24918		30	10	3744	1848	1208	S	0.0	0.0	0:00.61	/bin/bash /var/lib/condor/execute/dir_24914/condor_exec.exe -v std -name gfactory_instance -entry CMS_T3_U
29404		30	10	3472	1548	1208	S	0.0	0.0	0:00.12	/bin/bash /var/lib/condor/execute/dir_24914/glide_iEheSD/main/condor_startup.sh glidein_config
30115		30	10	17884	8272	6320	S	0.0	0.1	0:00.16	/var/lib/condor/execute/dir_24914/glide_iEheSD/main/condor/sbin/condor_master -f -pidfile /var/lib/c
30118		30	10	18928	9140	6748	S	0.0	0.1	0:00.88	condor_startd -f
2894		30	10	17012	8336	6568	S	0.0	0.1	0:00.16	condor_starter -f vocns058.cern.ch
3697		30	10	3352	1456	1196	S	0.0	0.0	0:00.10	/bin/bash /var/lib/condor/execute/dir_24914/glide_iEheSD/execute/dir_2894/condor_exec.exe -
3823		30	10	3520	1520	1224	S	0.0	0.0	0:00.02	sh ./CMSRunAnalysis.sh -a sandbox.tar.gz --sourceURL=https://cmsweb.cern.ch/crabcache --
3852		30	10	23508	13228	1312	R	0.0	0.1	0:00.71	python CMSRunAnalysis.py -r /var/lib/condor/execute/dir_24914/glide_iEheSD/execute/di
5049		30	10	3624	1648	1160	S	0.0	0.0	0:00.01	/bin/bash /var/lib/condor/execute/dir_24914/glide_iEheSD/execute/dir_2894/cmsRun-m
5152		30	10	919M	567M	98404	R	98.9	3.5	4:07.56	cmsRun -j FrameworkJobReport.xml PSet.py
30117		30	10	7048	4000	1100	S	0.0	0.0	0:00.73	condor_procd -A /var/lib/condor/execute/dir_24914/glide_iEheSD/log/procd_address -L /var/lib/cond

Heterogeneous computing: batch job submitted from x86\_64 machine at CERN to AArch64 worker node at Princeton University

Showcased on Fedora 19 on APM Mustang development board

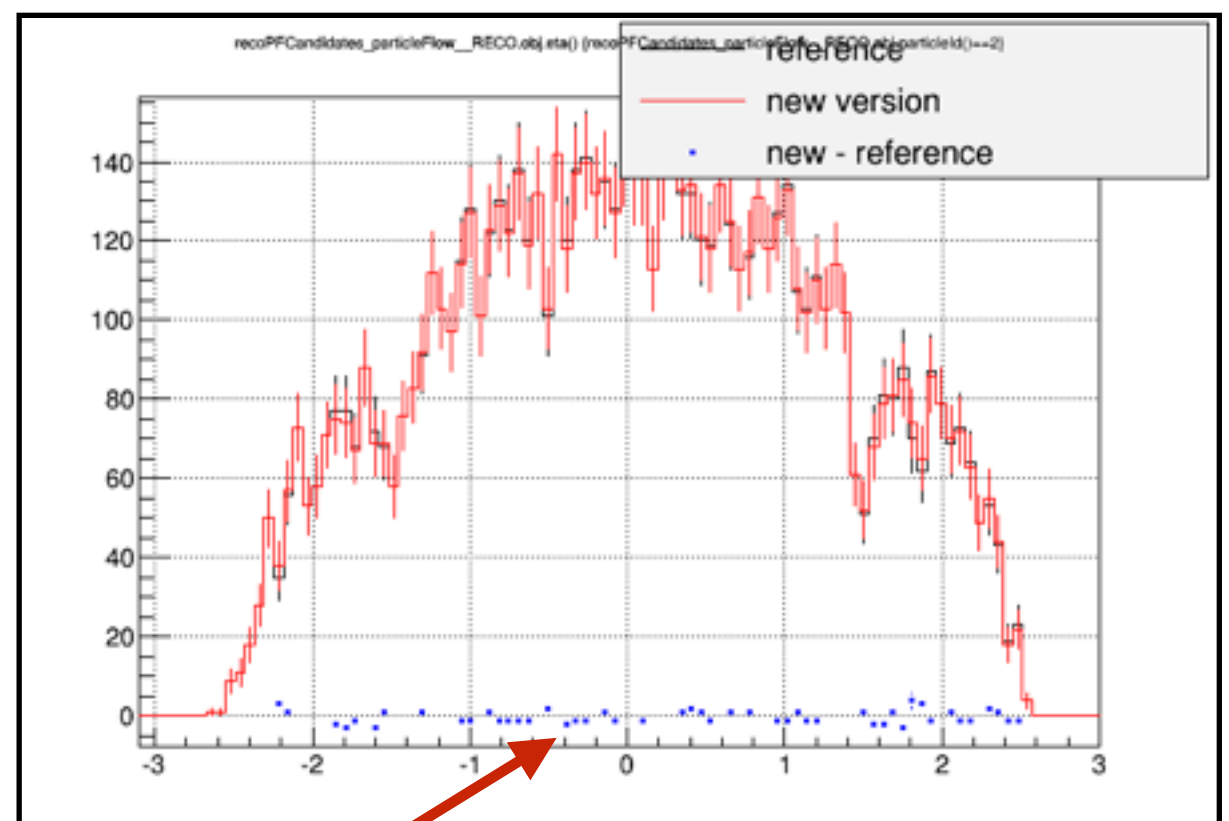
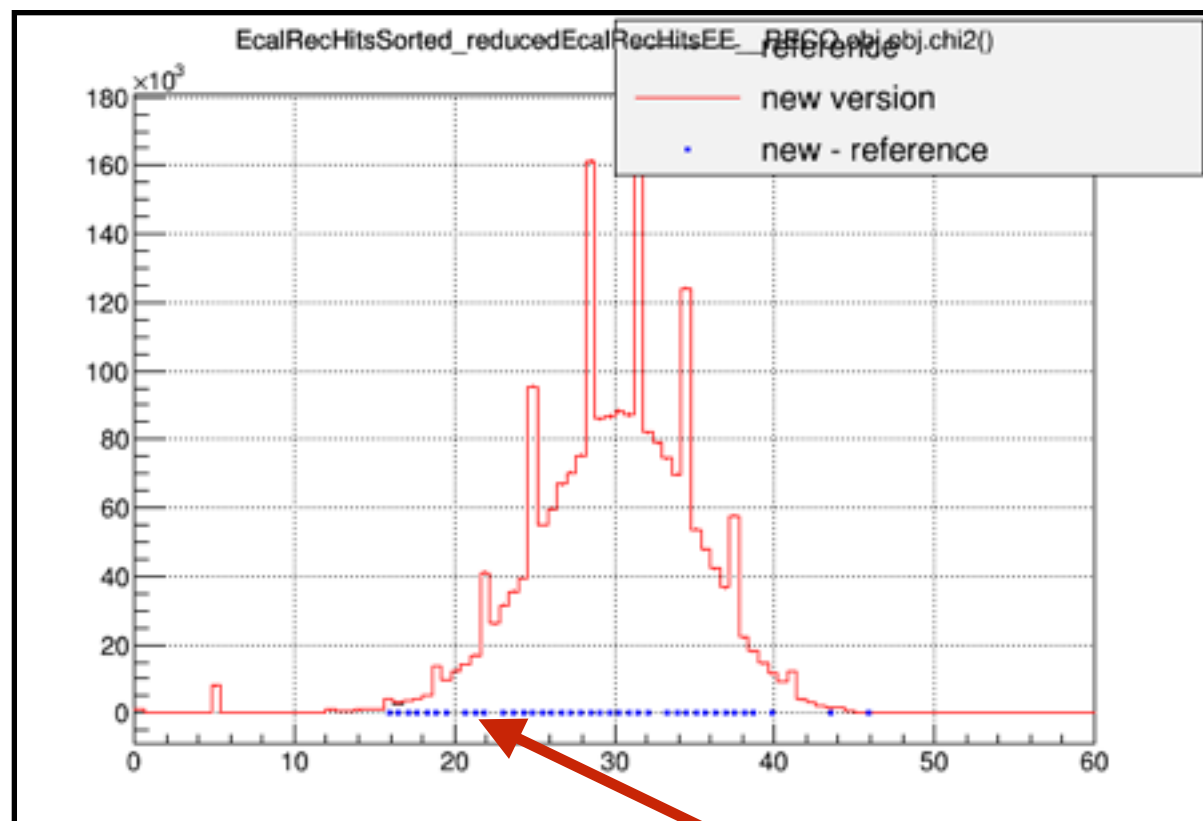
Moving to CentOS 7.2 on HP Moonshot + 6 x m400 (production-level system)



# Numerical Validation (ARMv8)

We used CMSSW\_7\_2\_0 pre-release for x86\_64 and AArch64 reconstruction comparison, where input were generated on x86\_64.

~950 differences were detected, but majority were minimal, i.e. non-significant; Examples:



Blue dots mark the difference

# What's new? #1

Official CMSSW Integration Builds (IBs) now include **aarch64** and **ppc64le**!

	Architectures	Builds	Unit Tests	RelVals	Other Tests	FWLite	Q/A
<b>CMSSW_8_1_X_2016-02-21-0000</b> IB Tag Static Analyzer Modules to thread unsafe statics Modules to thread unsafe EventSetup products HLT Validation Valgrind DQM Tests	slc6_amd64_gcc493	See Details	4 Tests Failing	Pass: 841 Fail: 12	See Details	See Details	Q
	Full Build						
	<u>slc7_aarch64_gcc530</u>	26 Warnings	Unknown	<b>Missing tests!</b>			Q
	Full Build						
	slc7_amd64_gcc530	19 Warnings	5 Tests Failing	Pass: 1193 Fail: 18	See Details		Q
	Full Build						
slc6_amd64_gcc530	19 Warnings	3 Tests Failing	Pass: 1233 Fail: 3	See Details		Q	
Full Build							
<u>fc22_ppc64le_gcc530</u>	47 Warnings	32 Tests Failing	Pass: 1158 Fail: 30	See Details		Q	
Full Build							

- No new pull requests since CMSSW\_8\_1\_X\_2016-02-19-2300

**Note:** more work is needed to make everything stable

```
$ file /cvmfs/cms.cern.ch/{slc7_aarch64_gcc530,fc22_ppc64le_gcc530}/cms/cmssw/CMSSW_8_0_0
/cvmfs/cms.cern.ch/slc7_aarch64_gcc530/cms/cmssw/CMSSW_8_0_0: directory
/cvmfs/cms.cern.ch/fc22_ppc64le_gcc530/cms/cmssw/CMSSW_8_0_0: directory
```

Available for all sites now!

# What's new? #2

- (1) CVMFS 2.3.0 (dev)** is building on **CentOS 7.2/aarch64** since **Feb 14th** nightly build
  - See "Technology Previews" under "CernVM-FS Downloads" page
  - We are running CVMFS client under aarch64 since before ACAT '14 (September 2014) and have not observed issues
  - Server to be tested once OverlayFS issues are solved (currently one needs aufs)
- (2) Static PRoot/QEMU + CentOS/Fedora rootfs** setup for doing non-native installations
- (3) CMSSW port to ppc64le (POWER8)** discovered 2 issues in LLVM (all resolved upstream)
- (4) Attempt to for CMSSW ppc64 (big-endian)** port revealed issues
  - Bundled LLVM inside ROOT 6.06 is broken (waiting for move to 3.8.0)
  - pyroot is not endian safe (patch WIP)
- (5) Preparations for Open Science Grid (OSG)** full (hopefully) repository rebuild for **aarch64** (will also require some EPEL packages to be built)



# POWER8 Very Early Comparison

CMSSW reconstruction, Run II-like

2xIBM 8247-22L 2xHaswell E5-2699

No impact in performance

## # Physical core comparison (8 vs 2 threads/proc)

Single thread (performance)	0.156907 ev/s	0.200261 ev/s
Multi threaded (performance)	0.155383 ev/s	0.198463 ev/s
Single thread (peak RSS)	15'190.5 MB	3'341.89 MB
Multi threaded (peak RSS)	3'145.62 MB	1'859.4 MB

## # Full machine comparison (160 vs 72 threads) or (40 vs 18 4-threaded jobs)

MT memory savings

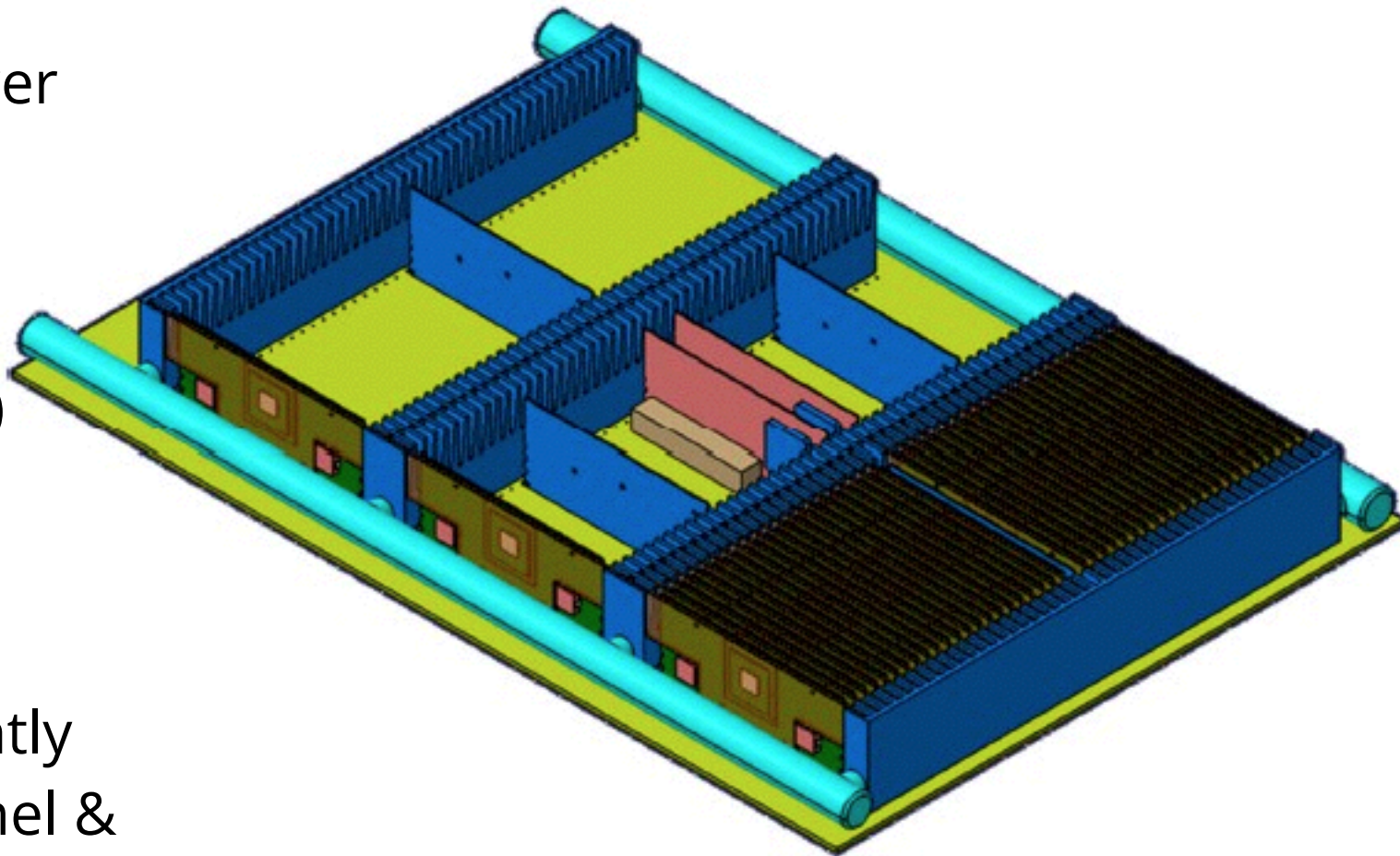
Multi threaded (performance)	2.78965 ev/s	3.65784 ev/s
Multi threaded (peak RSS)	97'844 MB	38'824.2 MB

Intel Xeon Haswell (E5-2699) provided 1.31x more events/s compared to IBM POWER8 (8247-22L)

# Datacenter-in-a-box

IBM/ASTRON (in Zurich) DOME 64-bit  $\mu$ Server for SKA big data challenge

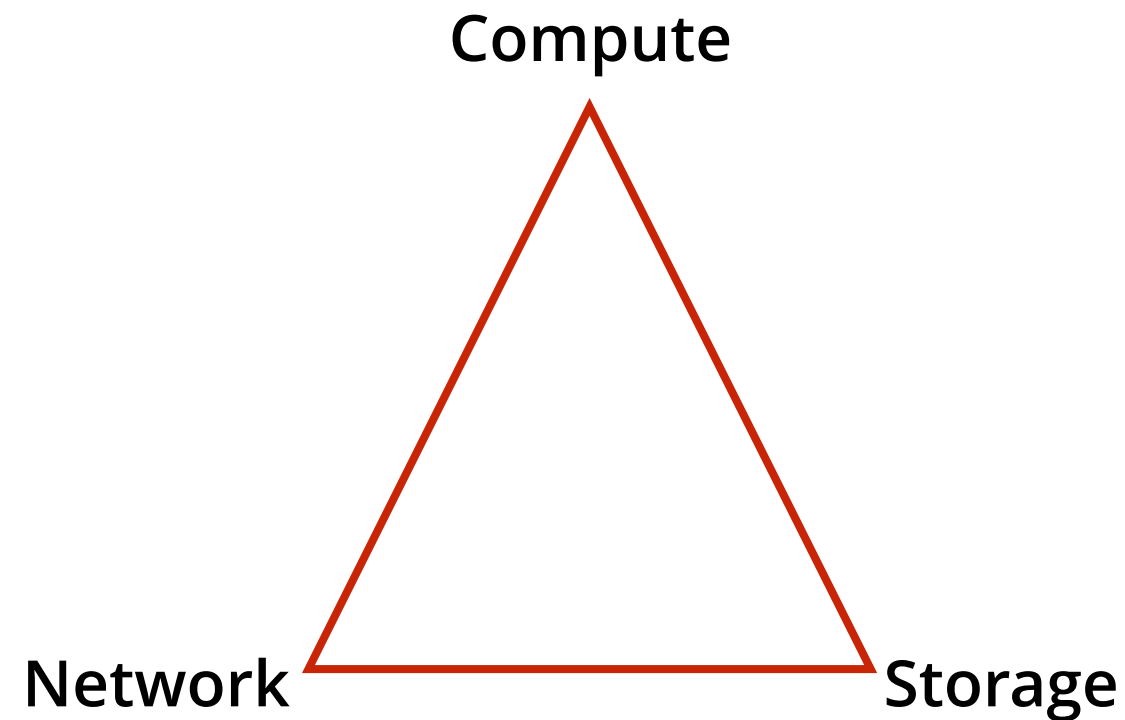
- ▶ 19" 2U w/ combined cooling & power
- ▶ 128 compute nodes
- ▶ **1536 ppc64** cores / 3074 threads
- ▶ 6 TB DRAM
- ▶ 1.28 Tbps Ethernet (@40Gbps x 32)
- ▶ Expected total power is ~6kW
- ▶ **Hot-water cooled for efficiency and density**
- ▶ Upstream support in future, currently runs Fedora rootfs + Freescale kernel & uboot
- ▶ **Memory bandwidth density:**
  - DOME 128 nodes 2U: 159GB/s/Liter (peak)
  - POWER8 S822L (2S) 13.9GB/s/Liter (peak)



It's all about SoC and packaging!  
There will be **aarch64** version!

Motivation for porting **CMSSW** to **ppc64**

# "Bring Balance to the Force"



**Changing one (e.g. Compute) might disturb existing balance in the Force**



# Summary

- ▶ Power constraints and market evolution may drive change in the kinds of **processors we use**
- ▶ Application diversity could drive heterogeneity to aid in {**performance, power, cost**} optimizations
- ▶ **The race is heating up, and Intel/platform vendors are not sitting idle**
- ▶ We have been exploring **alternative general purpose architectures** to the current **x86\_64** cores, incl. ARMv7 32-bit, ARMv8 64-bit, PowerPC (LE and BE), Xeon Phi
- ▶ We have demonstrated both application software (**CMSSW**) as well as job submission using **CRAB** (CMS Remote Analysis Builder) to **aarch64** nodes using a demonstrator cluster, and we will keep improving it
- ▶ We showed that **heterogeneity** by submitting jobs **from x86\_64** machine and landing them **on aarch64** worker nodes
- ▶ We are involved with open source communities and industry partners
- ▶ We need to continue investigating new SoCs/CPU's and platforms (e.g. Xeon D, new ARMv8.{0,1} SoCs and platforms using better processes (i.e. 14/16nm FinFET))



# Q&A

CONTACT



DAVIDLT AT CERN DOT CH

# BACKUP



# CPU Evolution

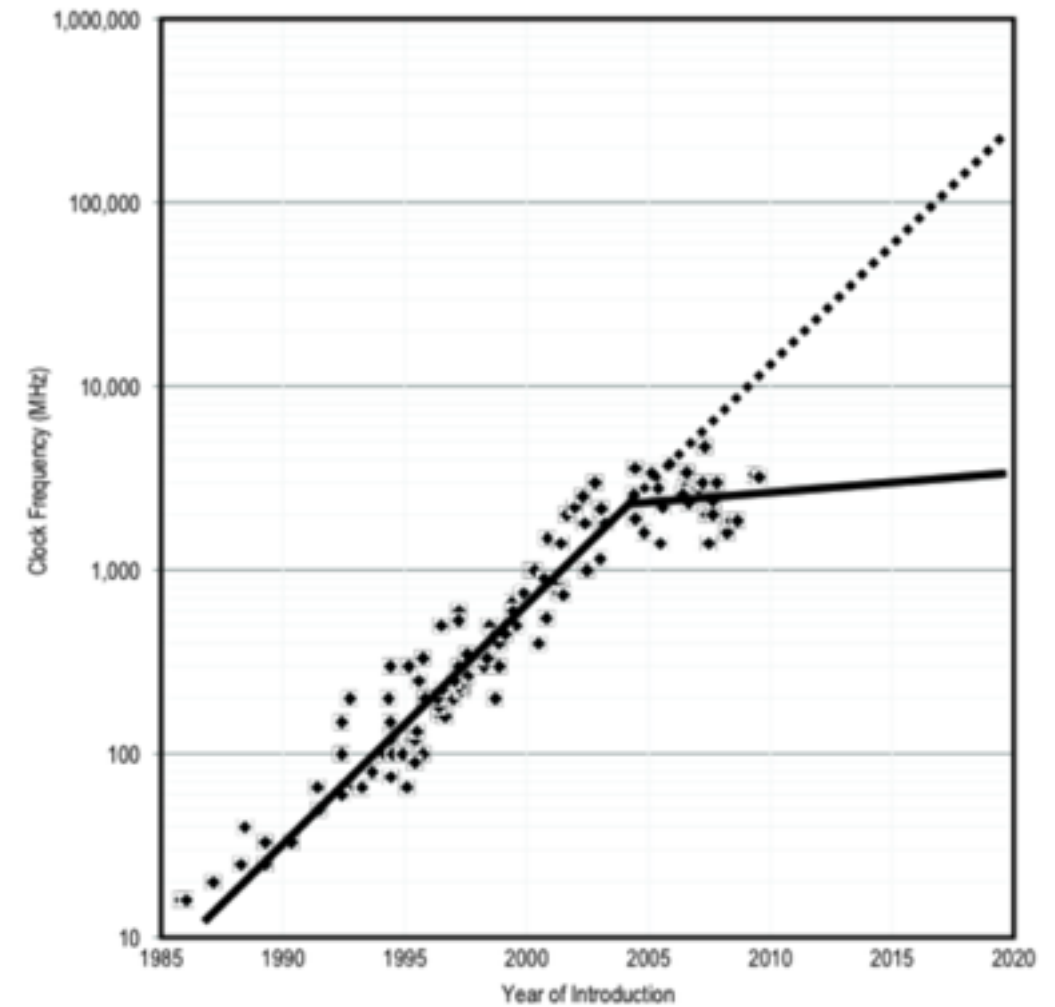
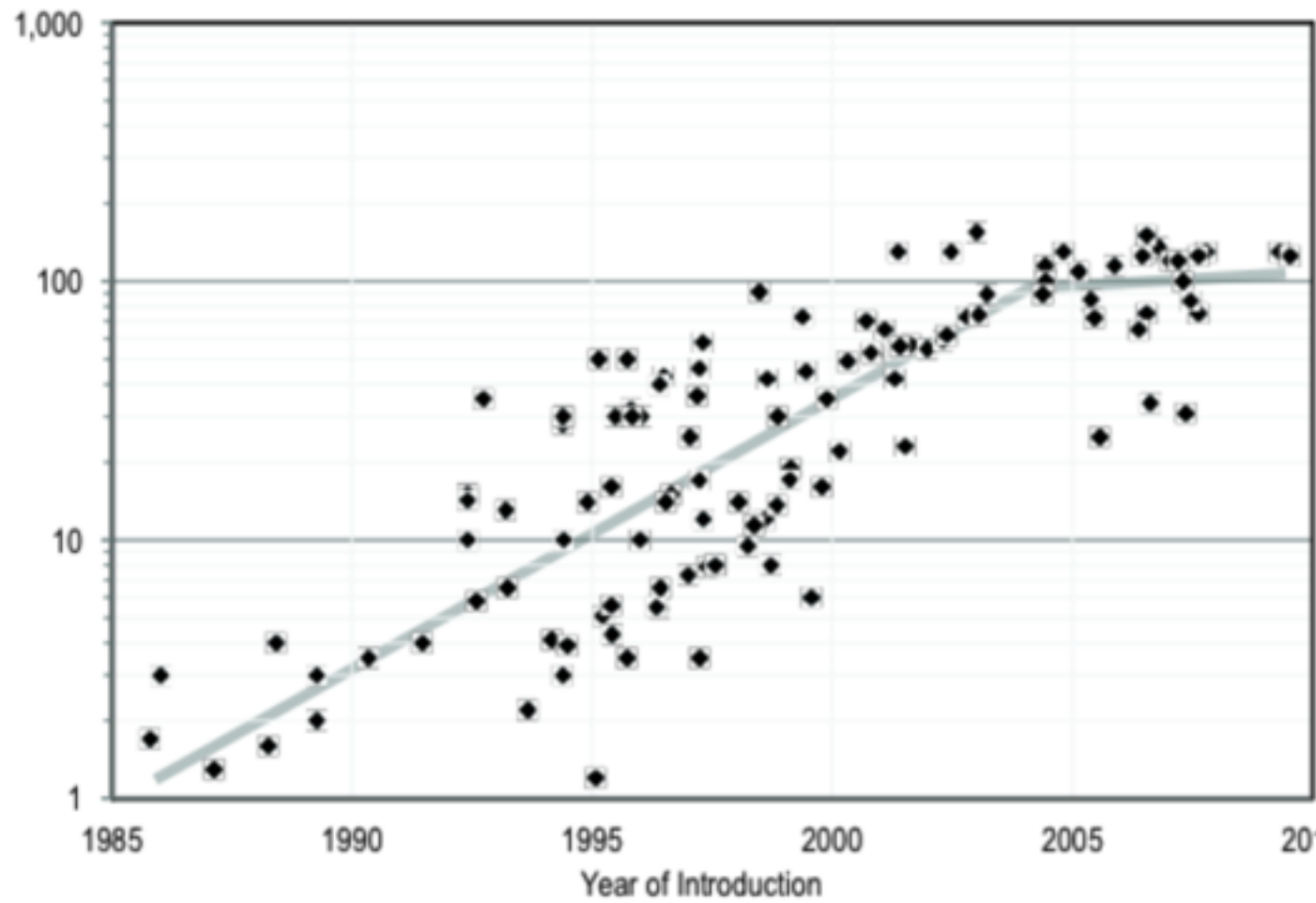


FIGURE A.4 Microprocessor power dissipation (watts) over time (1985-2010).

Source: "The Future of Computing Performance: Game Over or Next Level?"