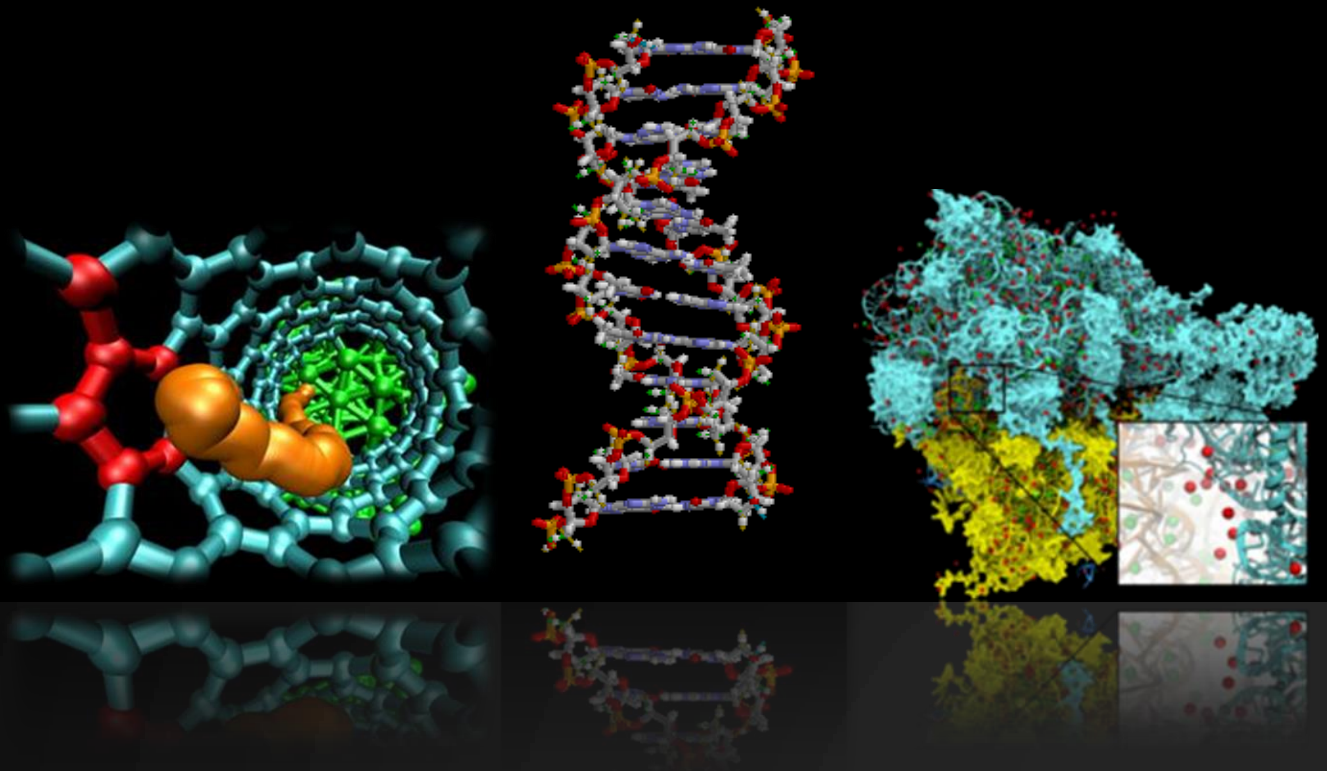


NVIDIA Computational Chemistry & Biology



Mark Berger
Senior Alliance Manager
Life and Material Sciences
mberger@nvidia.com

Updated: June 16, 2015

Overview of Life & Material Accelerated Apps



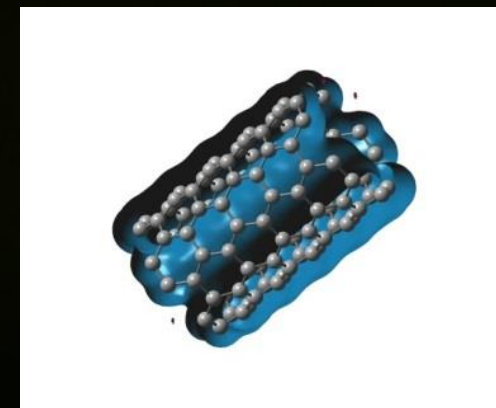
MD: All key codes are GPU-accelerated

- **ACEMD***, **AMBER (PMEMD)***, BAND, CHARMM, DESMOND, ESPResso, Folding@Home, GPUgrid.net, GROMACS, HALMD, **HOOMD-Blue***, LAMMPS, **Lattice Microbes***, mdcore, NAMD, OpenMM, **SOP-GPU***
- Great multi-GPU performance!
- Focus: on dense (up to 16) GPU nodes & large # of GPU nodes



QC: All key codes are ported or optimizing:

- GPU-accelerated and available today:
 - ABINIT, ACES III, ADF, BigDFT, CP2K, GAMESS, Quantum Espresso/PWscf, MOLCAS, MOPAC2012, NWChem, **OCTOPUS***, QUICK, Q-Chem, **TeraChem***
- Active GPU acceleration projects:
 - CASTEP, CPMD, GAMESS, **Gaussian**, NWChem, ONETEP, **Quantum Supercharger Library***, **VASP** & more
- Focus: on using GPU-accelerated math libraries, OpenACC directives



green* = application where all the workload is on GPU

3 Ways to Accelerate Applications



Applications

Libraries

“Drop-in”
Acceleration

OpenACC
Directives

Easily Accelerate
Applications

Programming
Languages

Maximum
Flexibility

GPU Accelerated Libraries

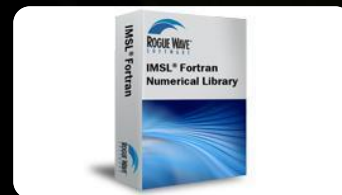
“Drop-in” Acceleration for your Applications



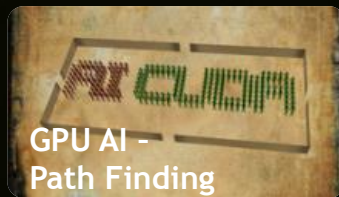
Linear Algebra
FFT, BLAS,
SPARSE, Matrix



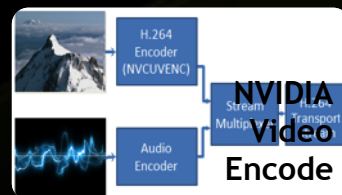
Numerical & Math
RAND, Statistics



Data Struct. & AI
Sort, Scan, Zero Sum

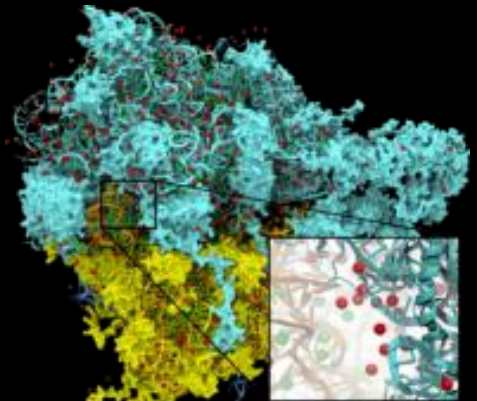
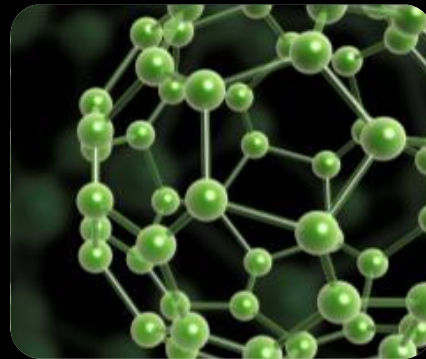


Visual Processing
Image & Video



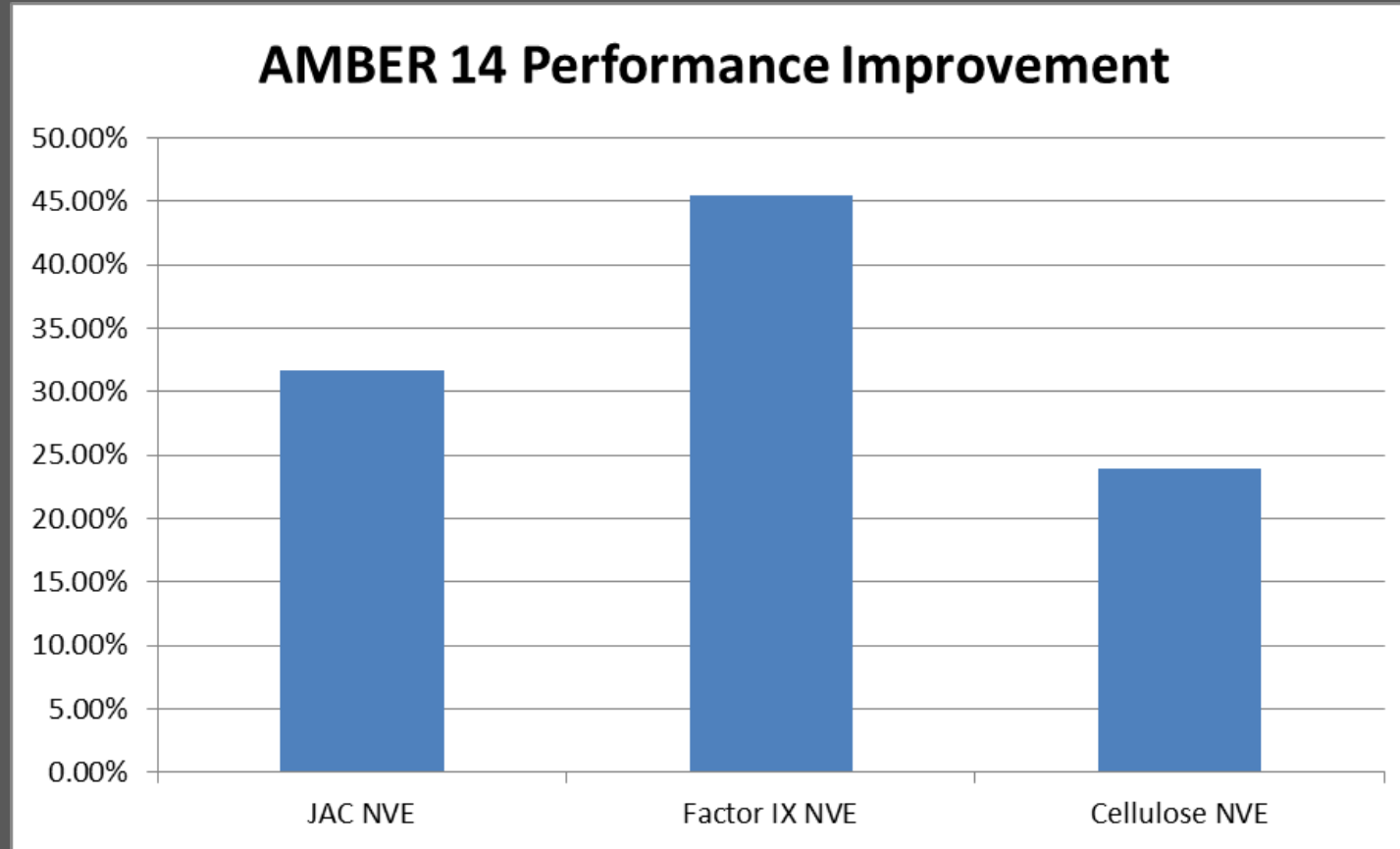
TESLA

Molecular Dynamics Module



AMBER 14

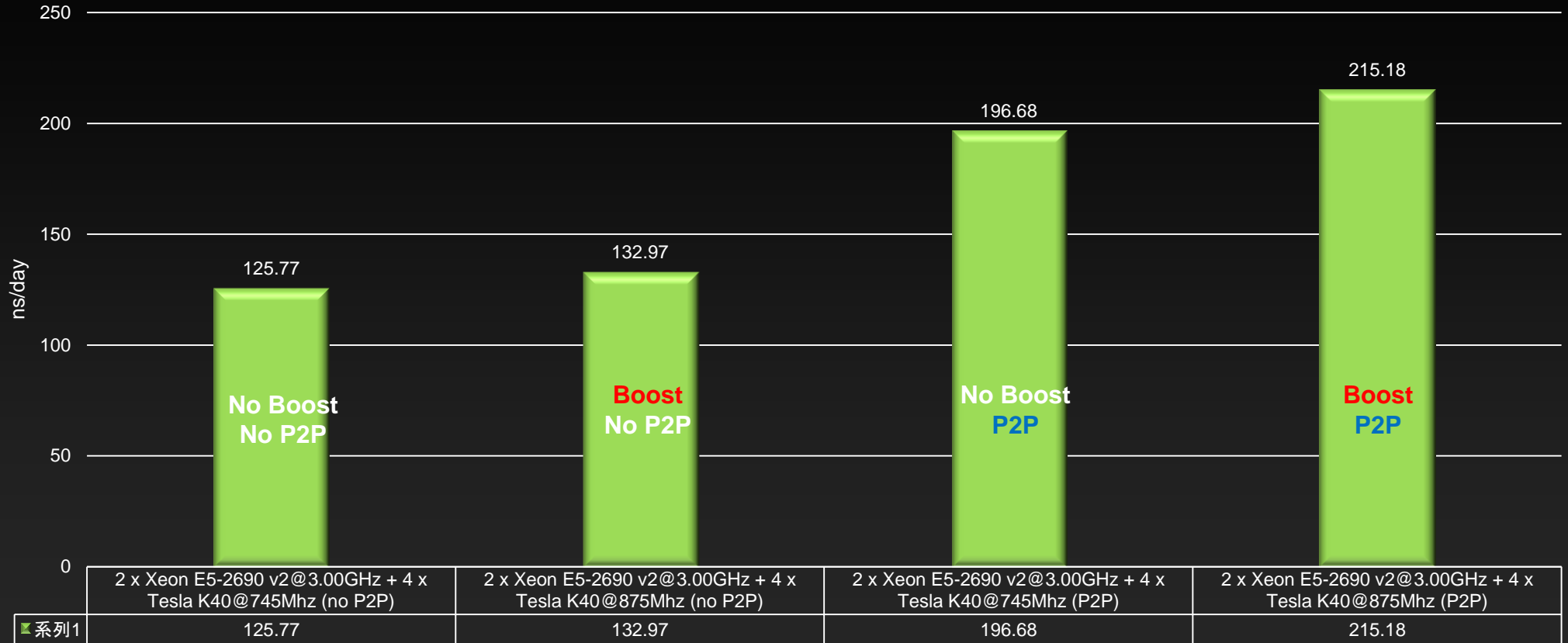
Compared to AMBER 12



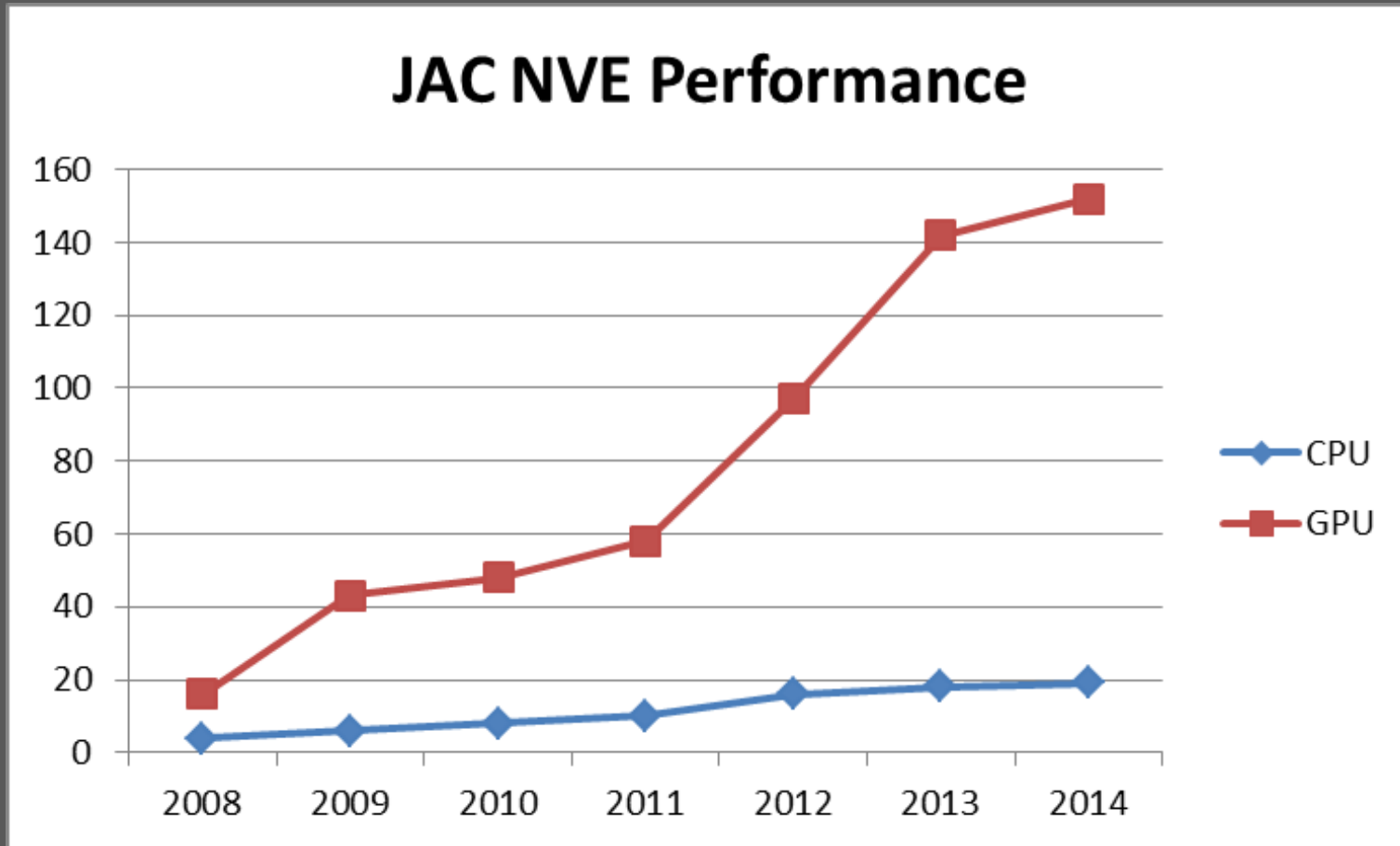
Courtesy of
Scott Le Grand
From GTC 2014
presentation

AMBER 14; large P2P and small Boost Clocks impacts

AMBER 14 (ns/day) on 4x K40; P2P and Boost Clocks Impact
DHFR NVE PME, 2fs Benchmark (CUDA 6.0, ECC off)



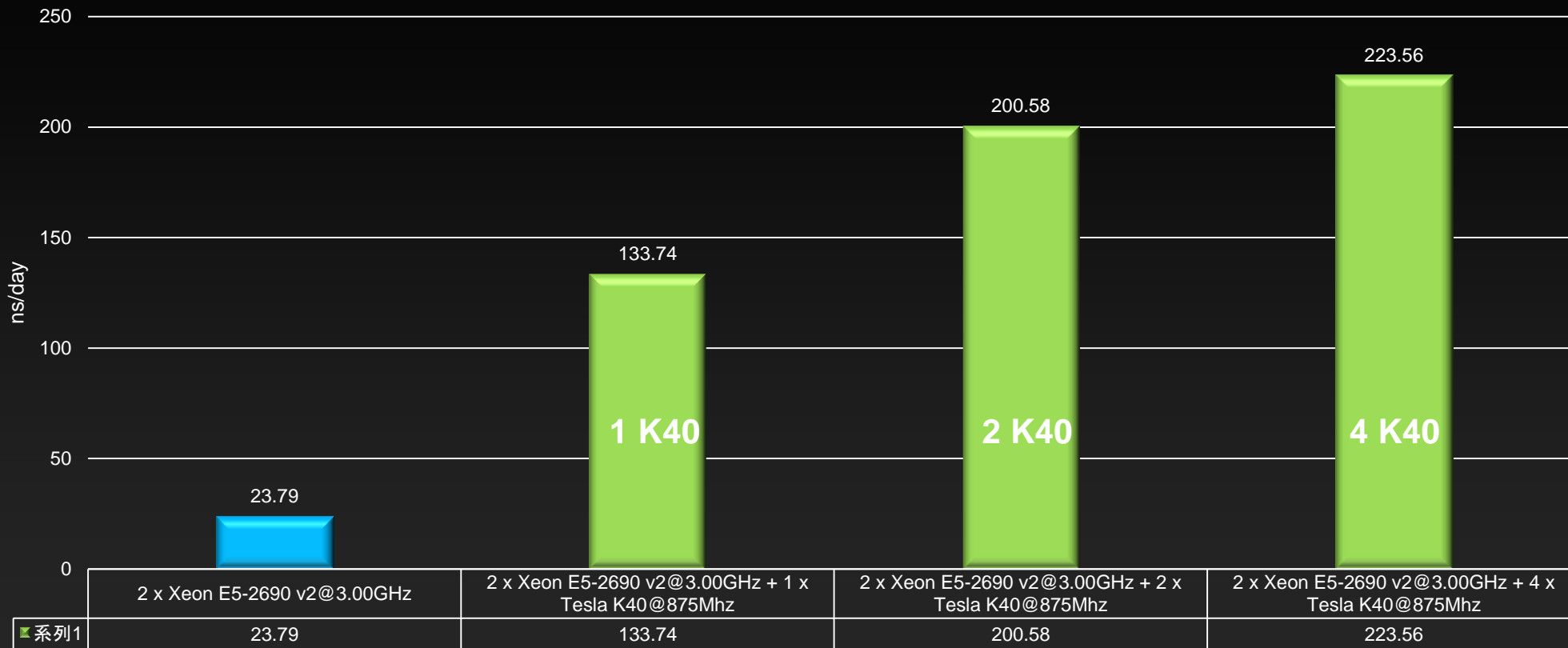
AMBER Performance



Courtesy of
Scott Le Grand
From GTC 2014
presentation

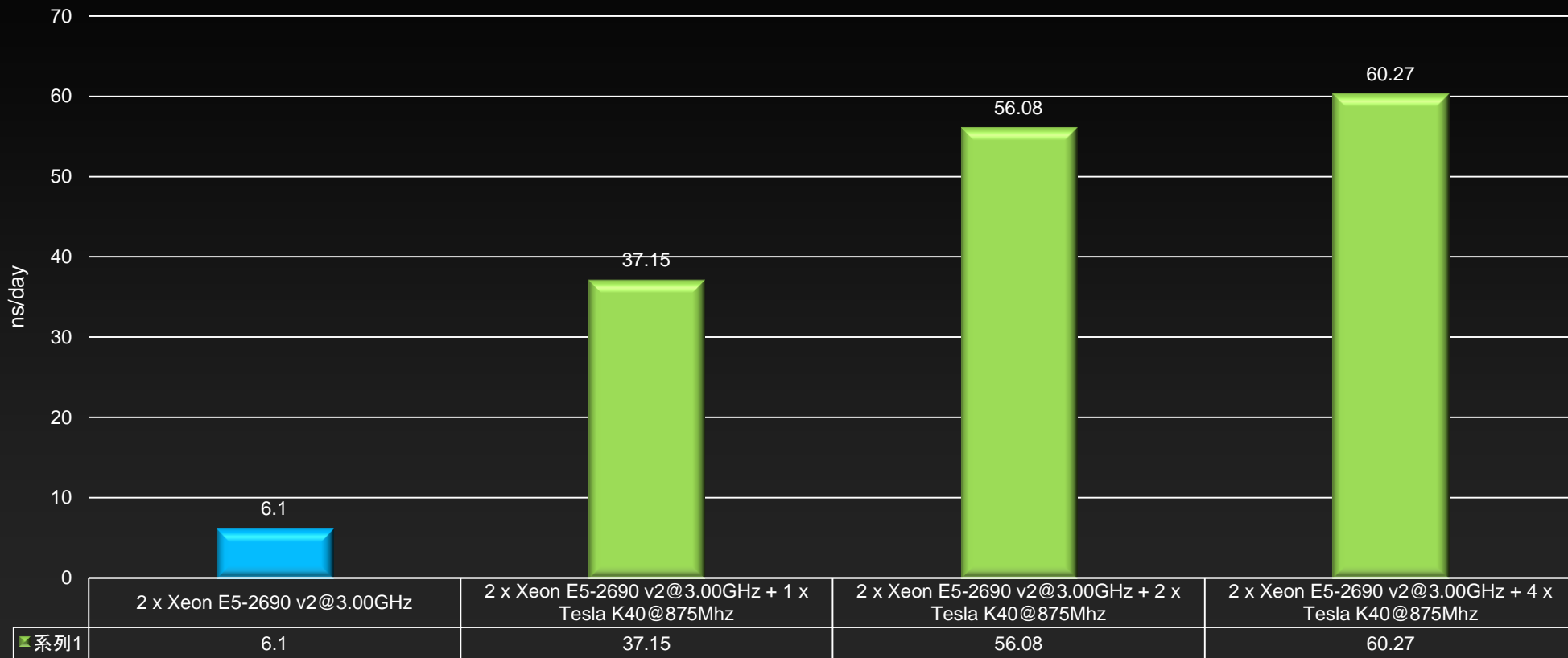
AMBER 14 with P2P, Higher Density Nodes

AMBER 14 (ns/day) on K40 with P2P and Boost Clocks
DHFR NVE PME, 2fs Benchmark (CUDA 5.5, ECC off)



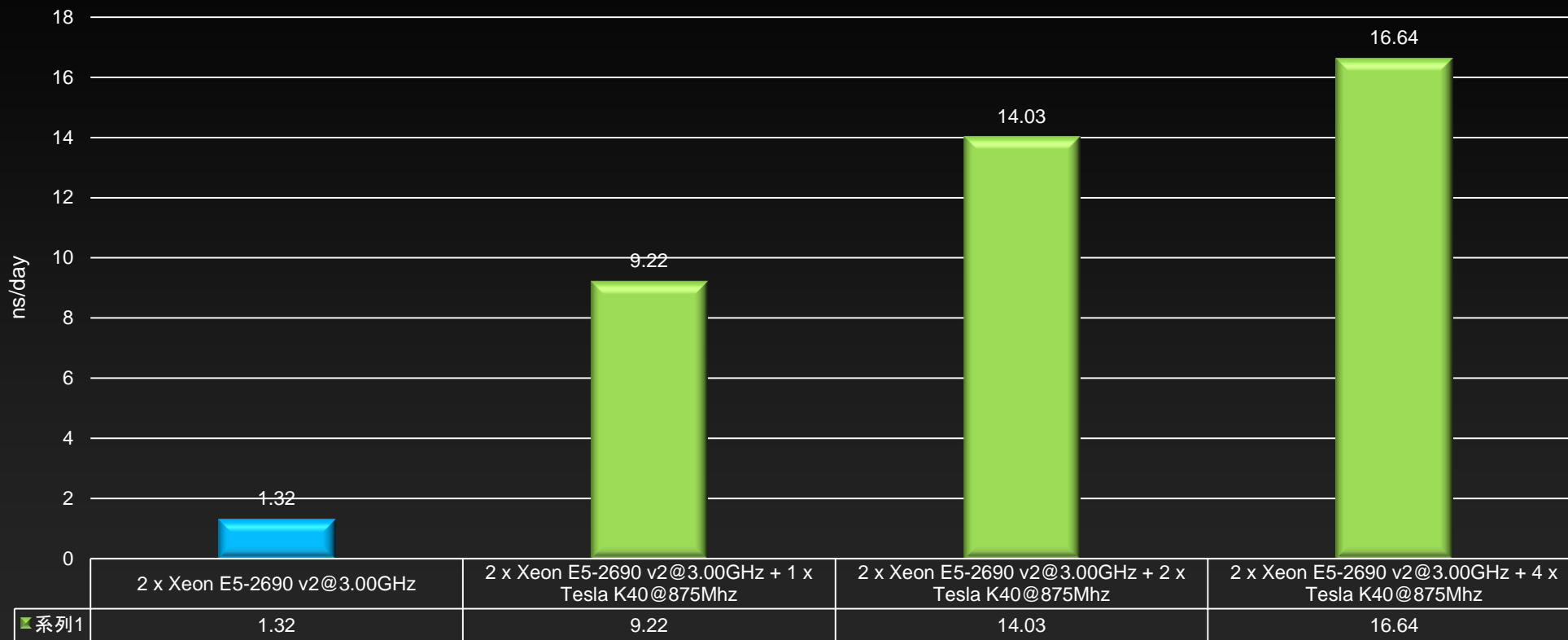
AMBER 14 and K40 with P2P, fastest GPU yet!

AMBER 14 (ns/day) on K40 with P2P and Boost Clocks
Factor IX NPT PME, 2fs Benchmark (CUDA 5.5, ECC off)



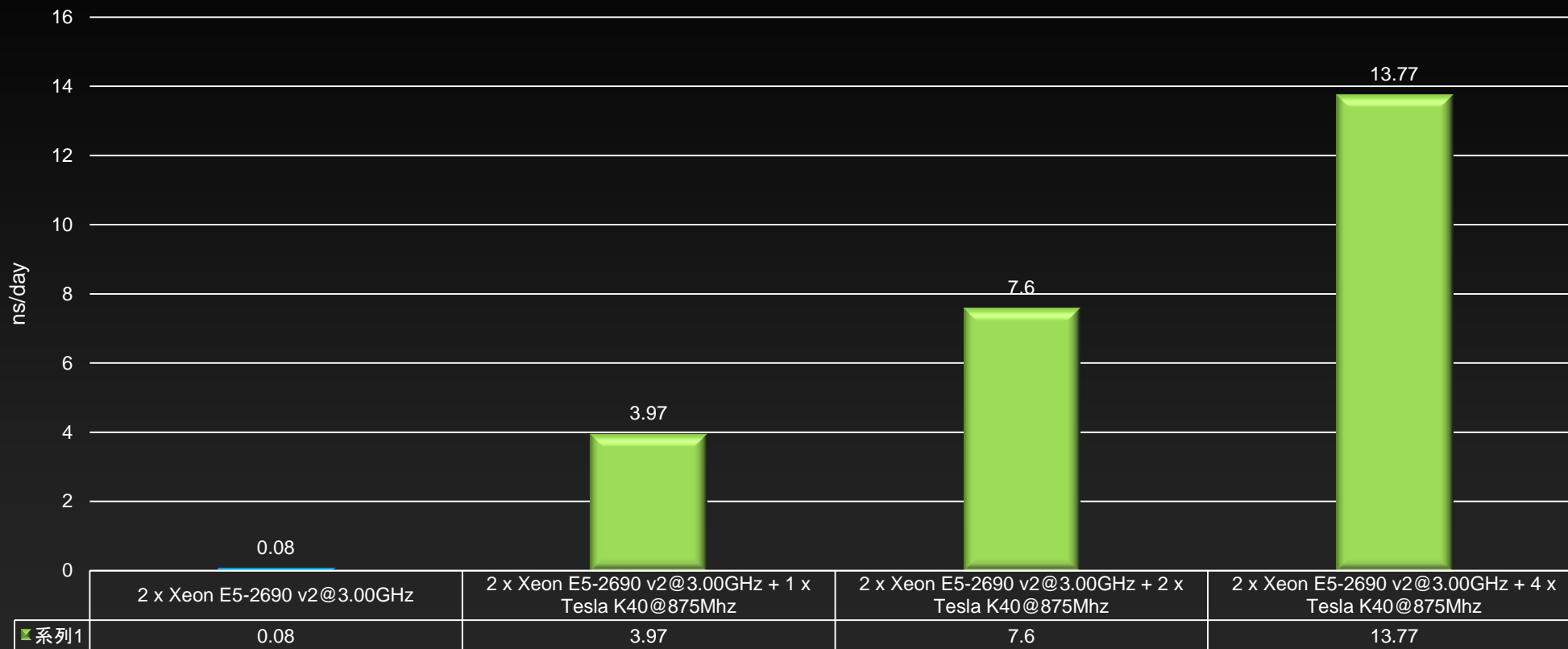
AMBER 14 and K40 with P2P, fastest GPU yet!

AMBER 14 (ns/day) on K40 with P2P and Boost Clocks
Cellulose NVE PME, 2fs Benchmark (CUDA 5.5, ECC off)



AMBER 14 and K40 with P2P, fastest GPU yet!

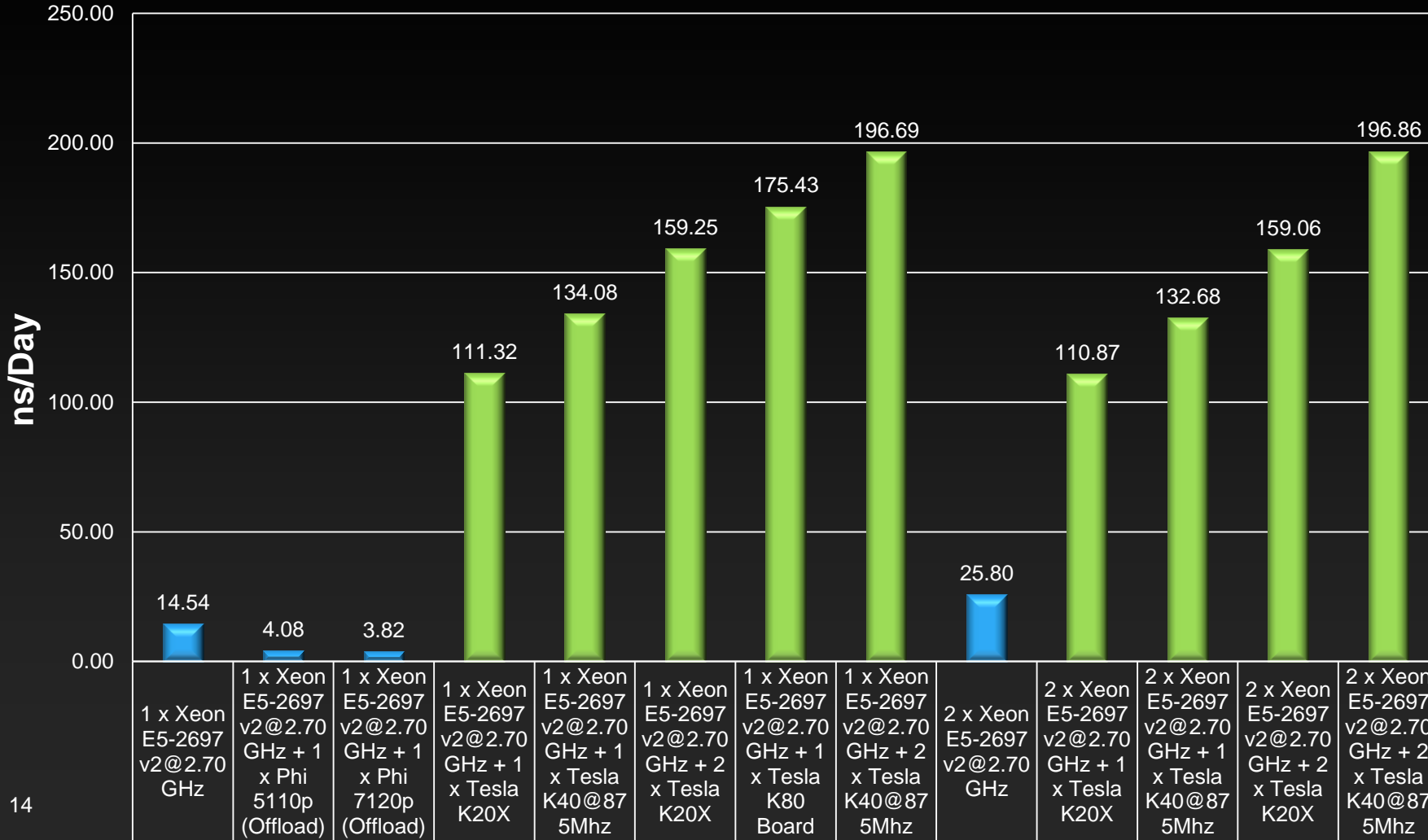
AMBER 14 (ns/day) on K40 with P2P and Boost Clocks
Nucleosome GB, 2fs Benchmark (CUDA 5.5, ECC off)



Kepler - Our Fastest Family of GPUs Yet



AMBER 14, SPFP-DHFR_production_NVE



Running AMBER 14

The blue node contains Dual E5-2697 CPUs (12 Cores per CPU).

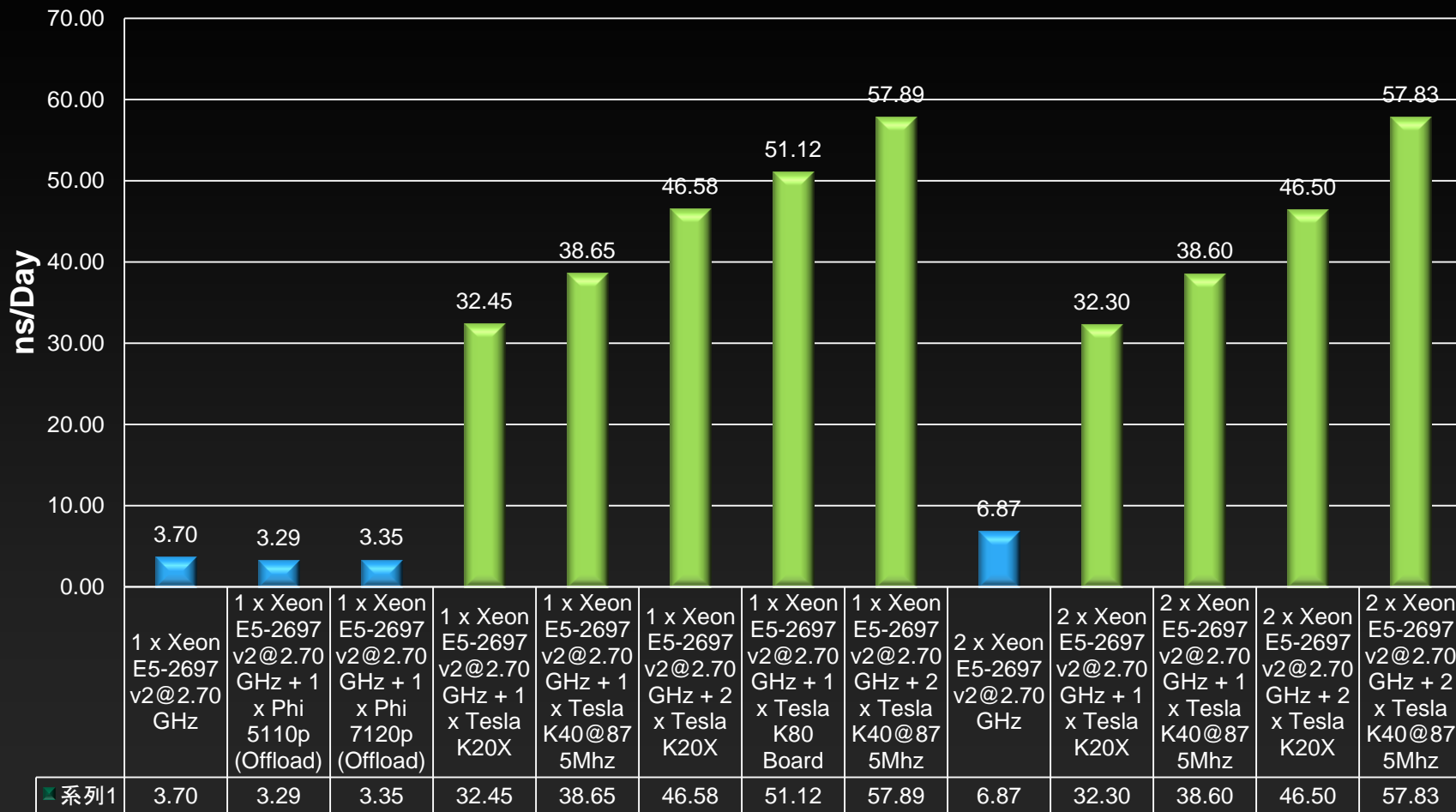
The green nodes contain Dual E5-2697 CPUs (12 Cores per CPU) and either 1x or 2x NVIDIA K20X, K40 or K80 for the GPU

DHFR (JAC)

Kepler - Our Fastest Family of GPUs Yet



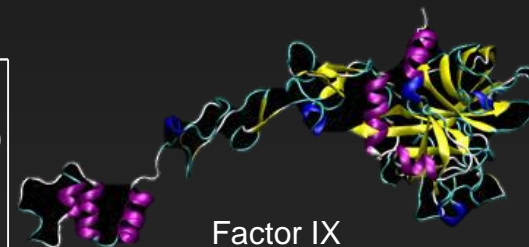
AMBER 14, SPFP-Factor_IX_Production_NVE



Running AMBER 14

The blue node contains Dual E5-2697 CPUs (12 Cores per CPU).

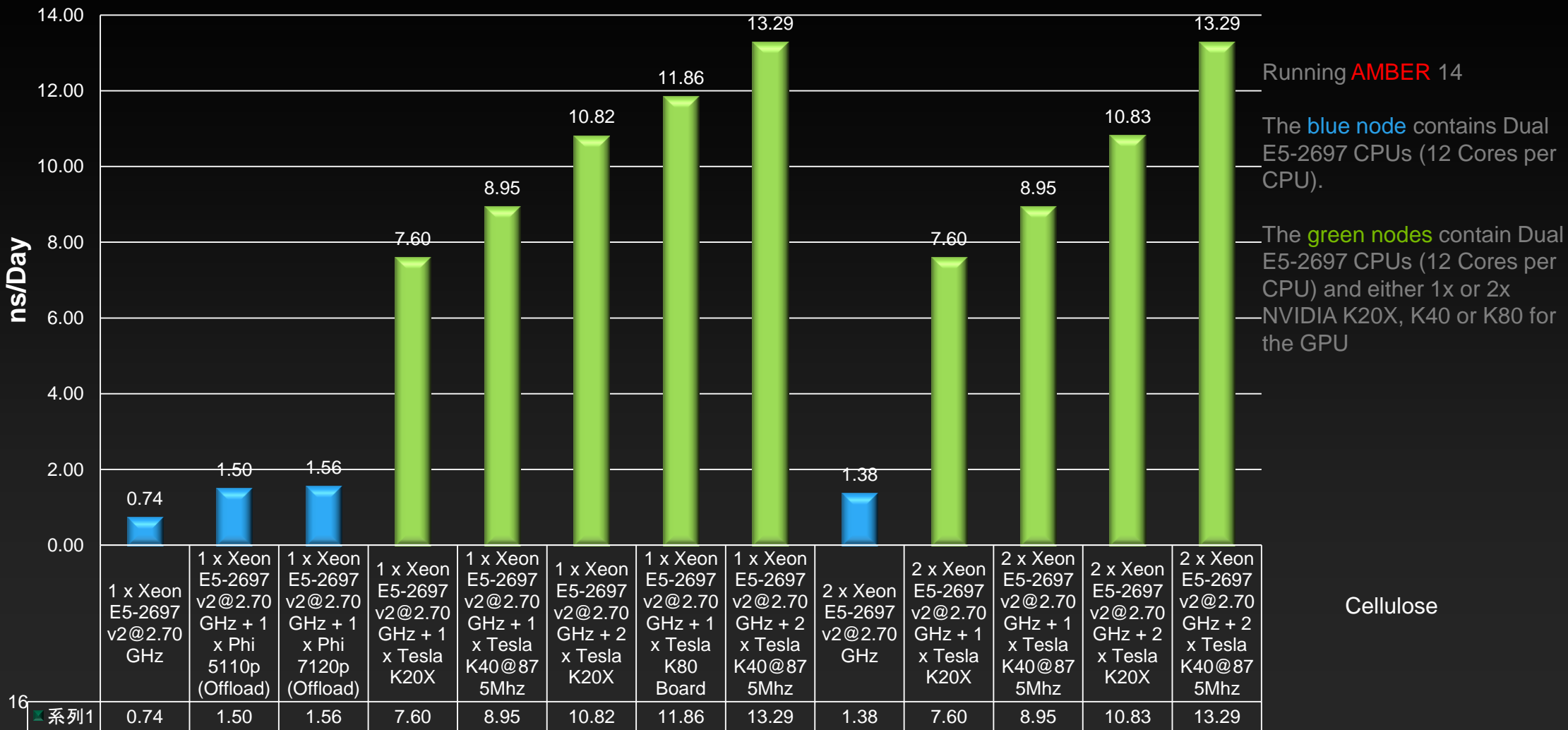
The green nodes contain Dual E5-2697 CPUs (12 Cores per CPU) and either 1x or 2x NVIDIA K20X, K40 or K80 for the GPU



Kepler - Our Fastest Family of GPUs Yet



AMBER 14, SPFP-Cellulose_Production_NVE



Replace 8 Nodes with 1 K20 GPU

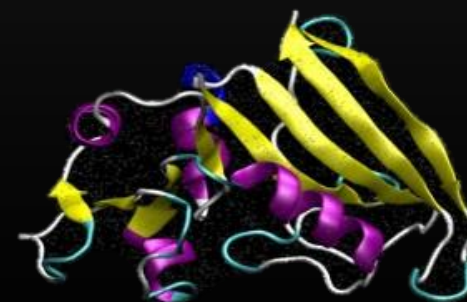


Running **AMBER** 12 GPU Support Revision 12.1 SPFP with CUDA 4.2.9 ECC Off

The **eight (8) blue nodes** each contain 2x Intel E5-2687W CPUs (8 Cores per CPU)

Each **green node** contains 2x Intel E5-2687W CPUs (8 Cores per CPU) plus 1x NVIDIA K20 GPU

Note: Typical CPU and GPU node pricing used. Pricing may vary depending on node configuration. Contact your preferred HW vendor for actual pricing.



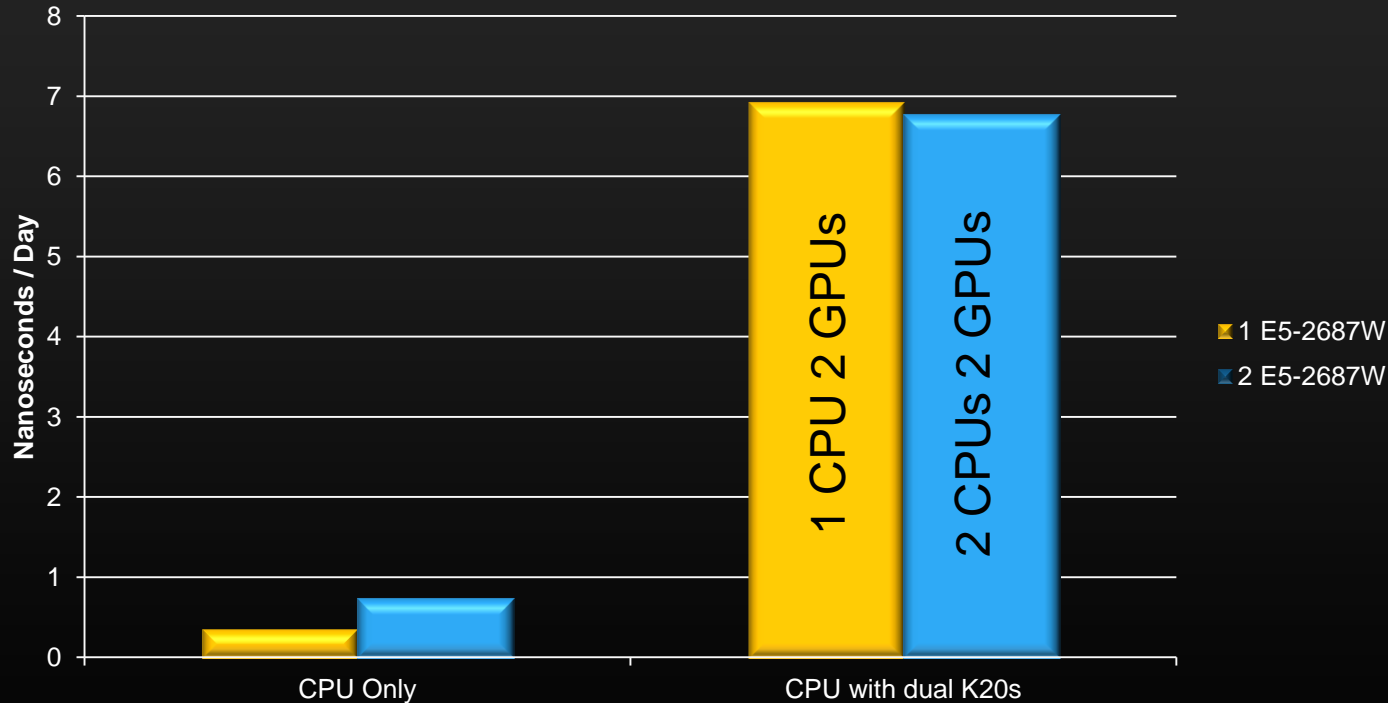
DHFR

Cut down simulation costs to $\frac{1}{4}$ and gain higher performance

Extra CPUs decrease Performance



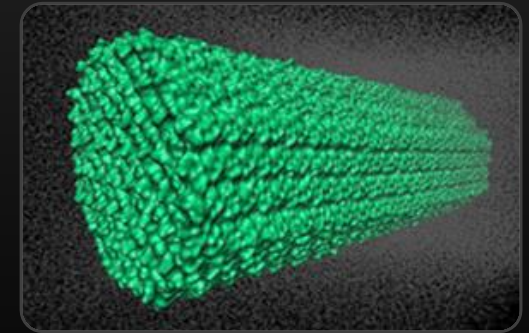
Cellulose NVE



Running **AMBER** 12 GPU Support Revision 12.1

The **orange bars** contains one E5-2687W CPUs (8 Cores per CPU).

The **blue bars** contain Dual E5-2687W CPUs (8 Cores per CPU)



Cellulose

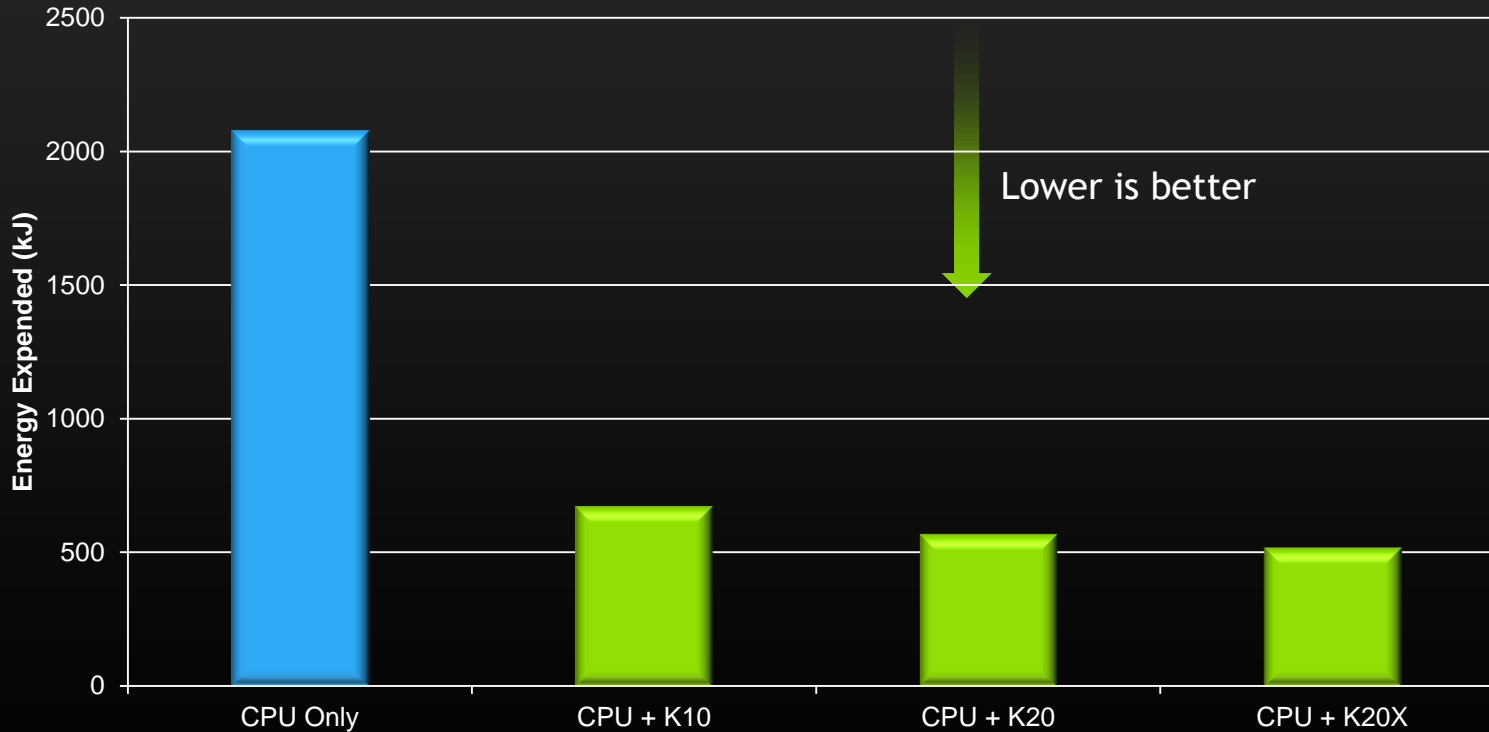
When used with GPUs, dual CPU sockets perform worse than single CPU sockets.

Kepler - Greener Science



Running **AMBER** 12 GPU Support Revision 12.1

Energy used in simulating 1 ns of DHFR JAC



The **blue node** contains Dual E5-2687W CPUs (150W each, 8 Cores per CPU).

The **green nodes** contain Dual E5-2687W CPUs (8 Cores per CPU) and 1x NVIDIA K10, K20, or K20X GPUs (235W each).

*Energy Expended
= Power x Time*

The GPU Accelerated systems use **65-75% less energy**

Recommended GPU Node Configuration for AMBER Computational Chemistry



Workstation or Single Node Configuration	
# of CPU sockets	2
Cores per CPU socket	6+ (1 CPU core drives 1 GPU)
CPU speed (Ghz)	2.66+
System memory per node (GB)	16
GPUs	Kepler K20, K40, K80
# of GPUs per CPU socket	1-4
GPU memory preference (GB)	6
GPU to CPU connection	PCIe 3.0 16x or higher
Server storage	2 TB
Network configuration	Infiniband QDR or better

Scale to multiple nodes with same single node configuration

Benefits of GPU AMBER Accelerated Computing



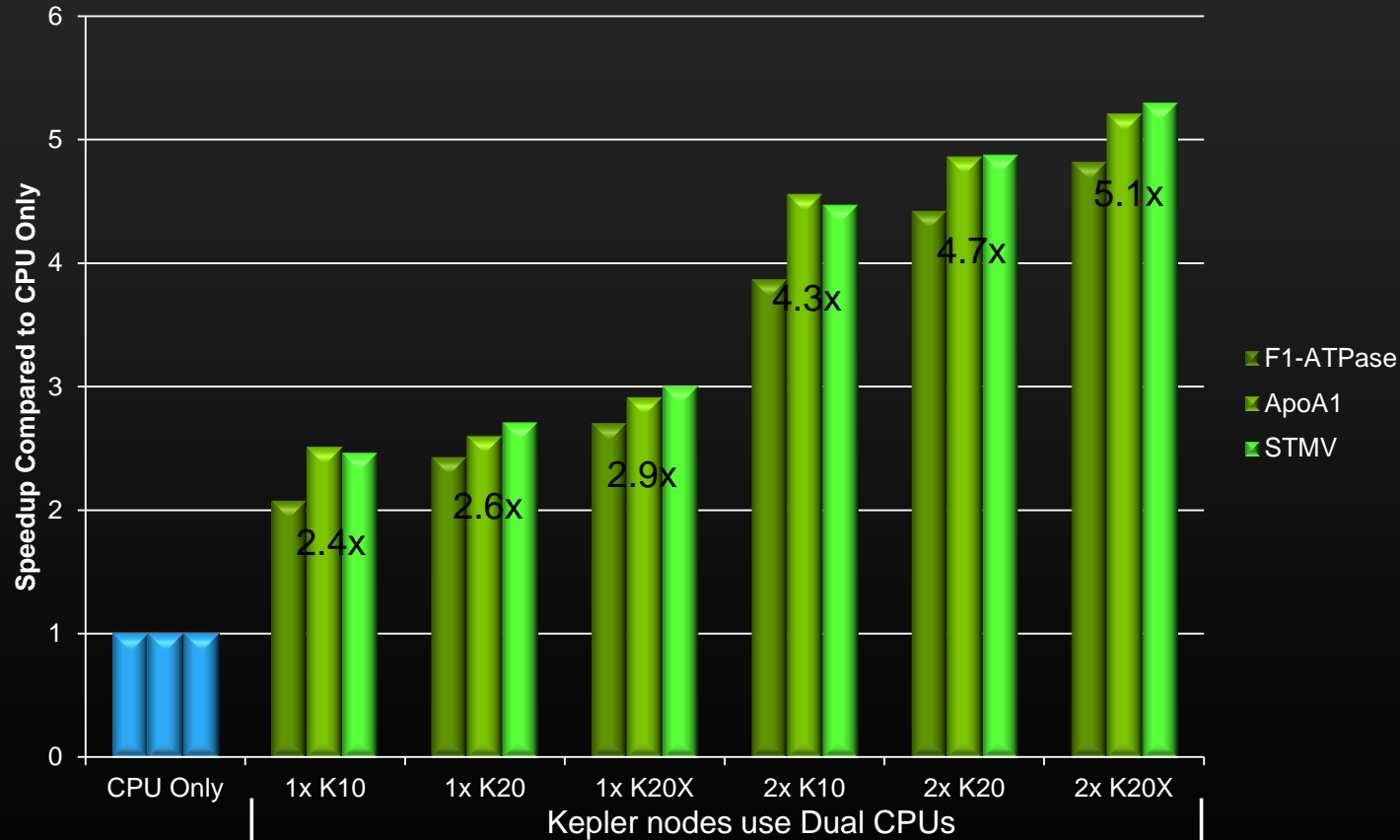
- Faster than CPU only systems in all tests
- Most major compute intensive aspects of classical MD ported
- Large performance boost with marginal price increase
- Energy usage cut by more than half
- GPUs scale well within a node and over multiple nodes
- K20 GPU is our fastest and lowest power high performance GPU yet

Try GPU accelerated AMBER for free – www.nvidia.com/GPUTestDrive



NAMD 2.10

Kepler - Universally Faster



Running **NAMD** version 2.9

The **CPU Only** node contains Dual E5-2687W CPUs (8 Cores per CPU).

The **Kepler nodes** contain Dual E5-2687W CPUs (8 Cores per CPU) and 1 or two NVIDIA K10, K20, or K20X GPUs.



F1-ATPase

The Kepler GPUs **accelerate all simulations**, up to 5.1x
Average acceleration printed in bars

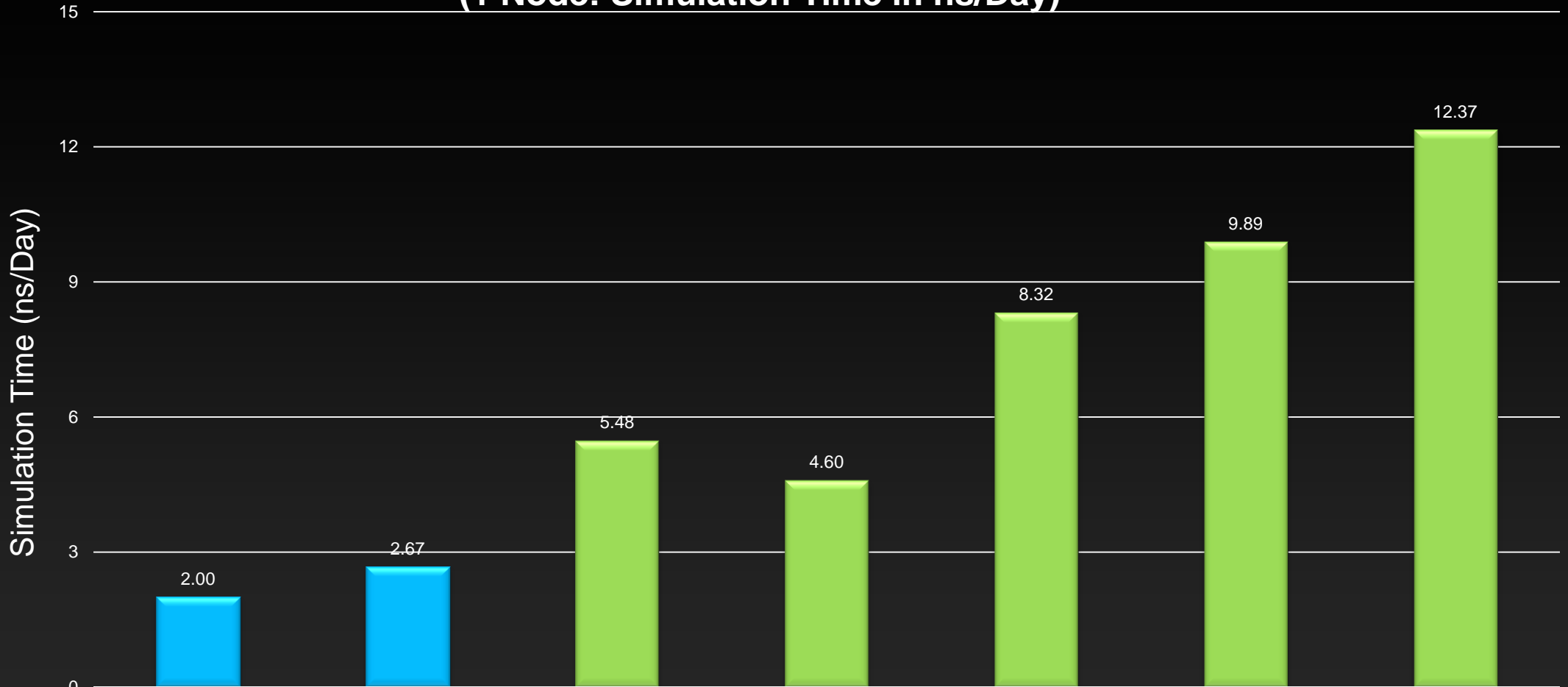


NAMD 2.10 (streaming patch)

NAMD 2.10; ApoA1 on Intel Phi, Tesla K40s and K80s & IVB CPUs



(1 Node: Simulation Time in ns/Day)



■ 系列1

2.00

2.67

5.48

4.60

8.32

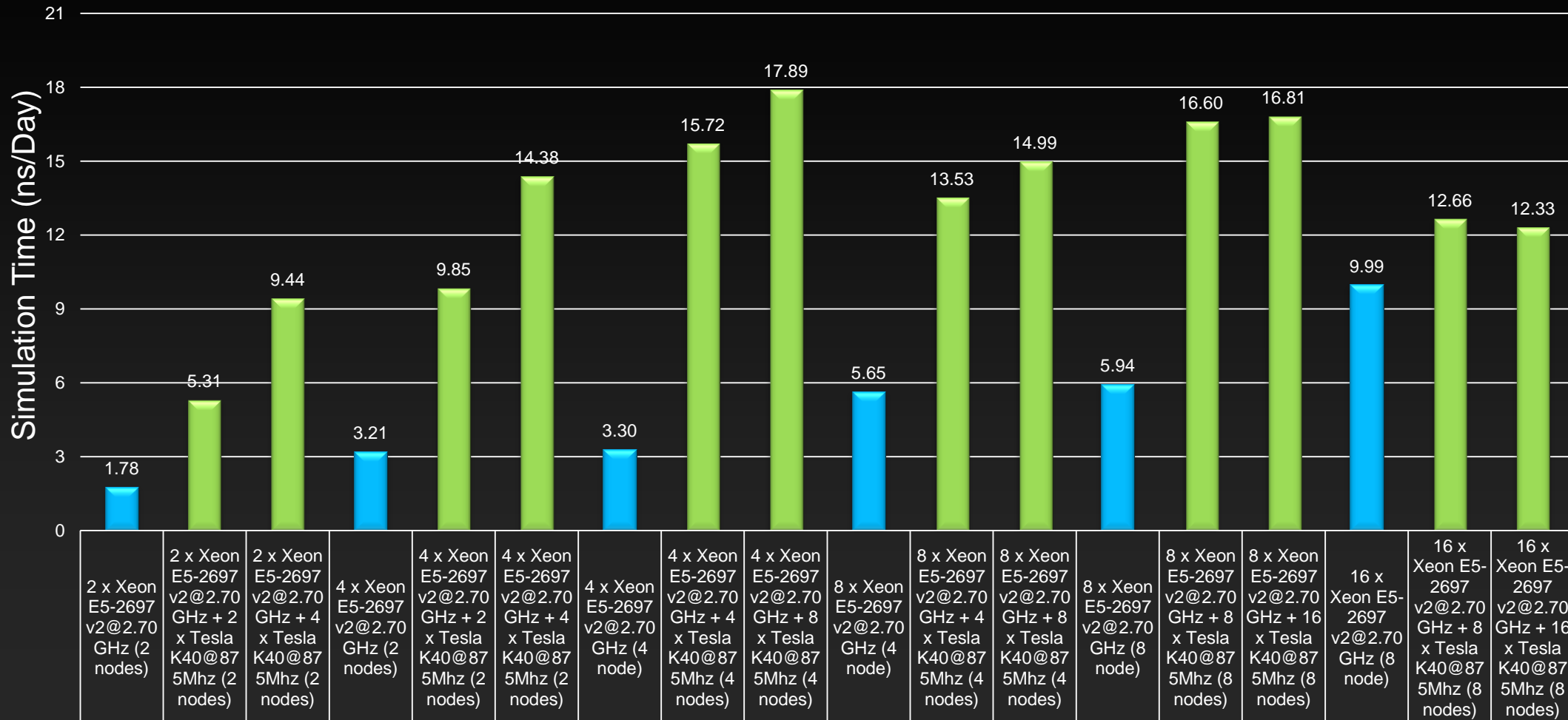
9.89

12.37

NAMD 2.10; ApoA1 on Tesla K40s & IVB CPUs



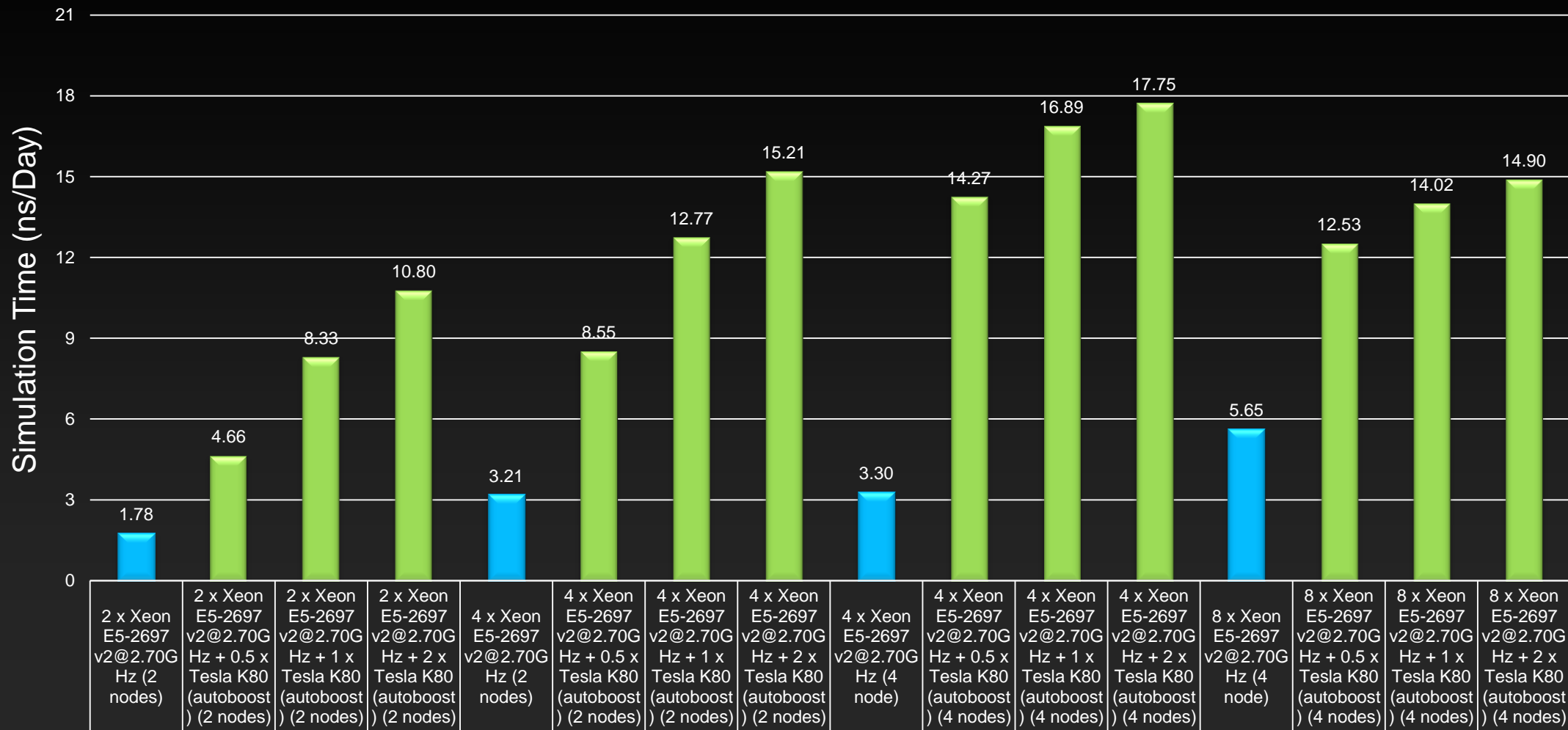
(2-8 Nodes: Simulation Time in ns/Day)



NAMD 2.10; ApoA1 on Tesla K80s & IVB CPUs



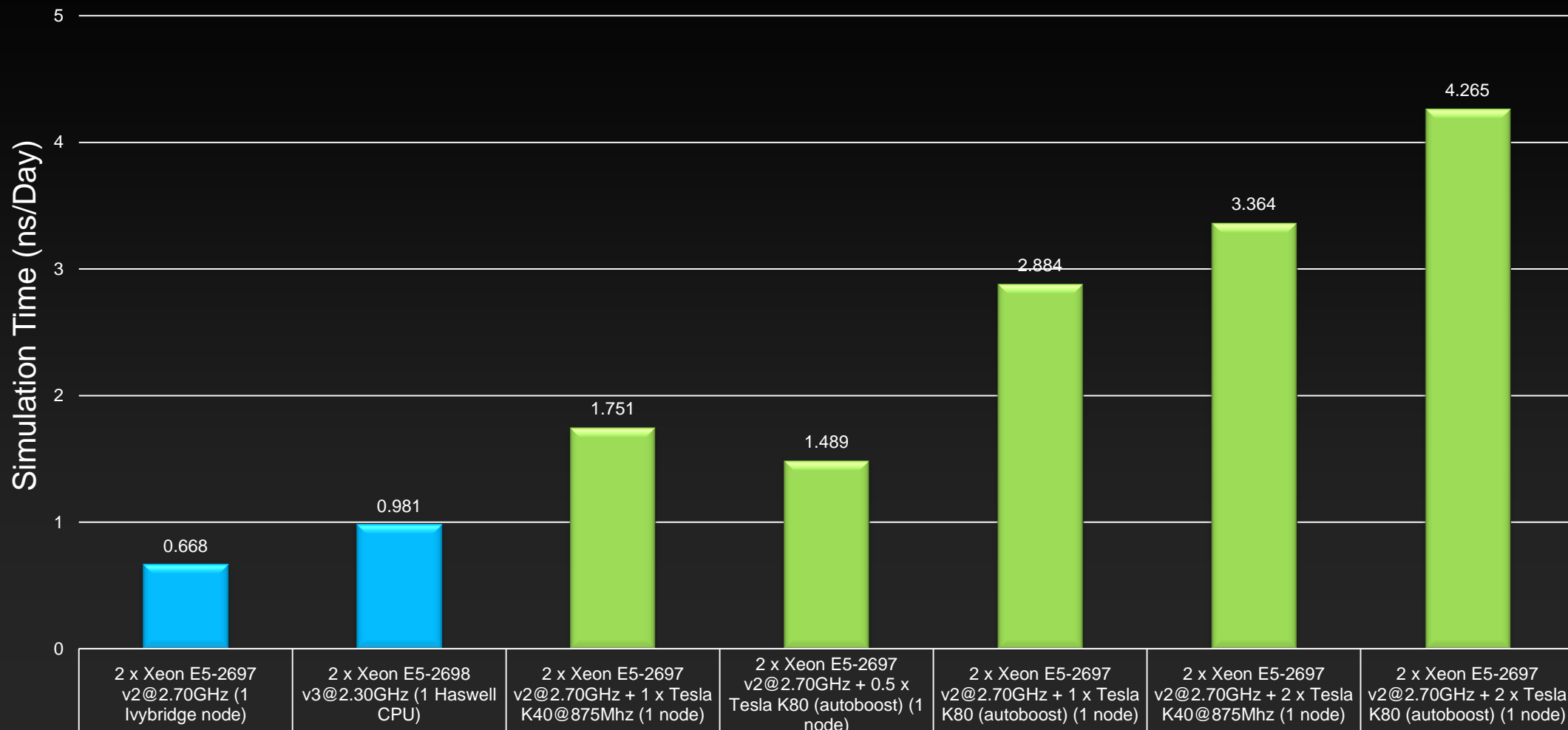
(2-4 Nodes: Simulation Time in ns/Day)



NAMD 2.10; F1-ATPase on Intel Phi, Tesla K40s and K80s & IVB CPUs



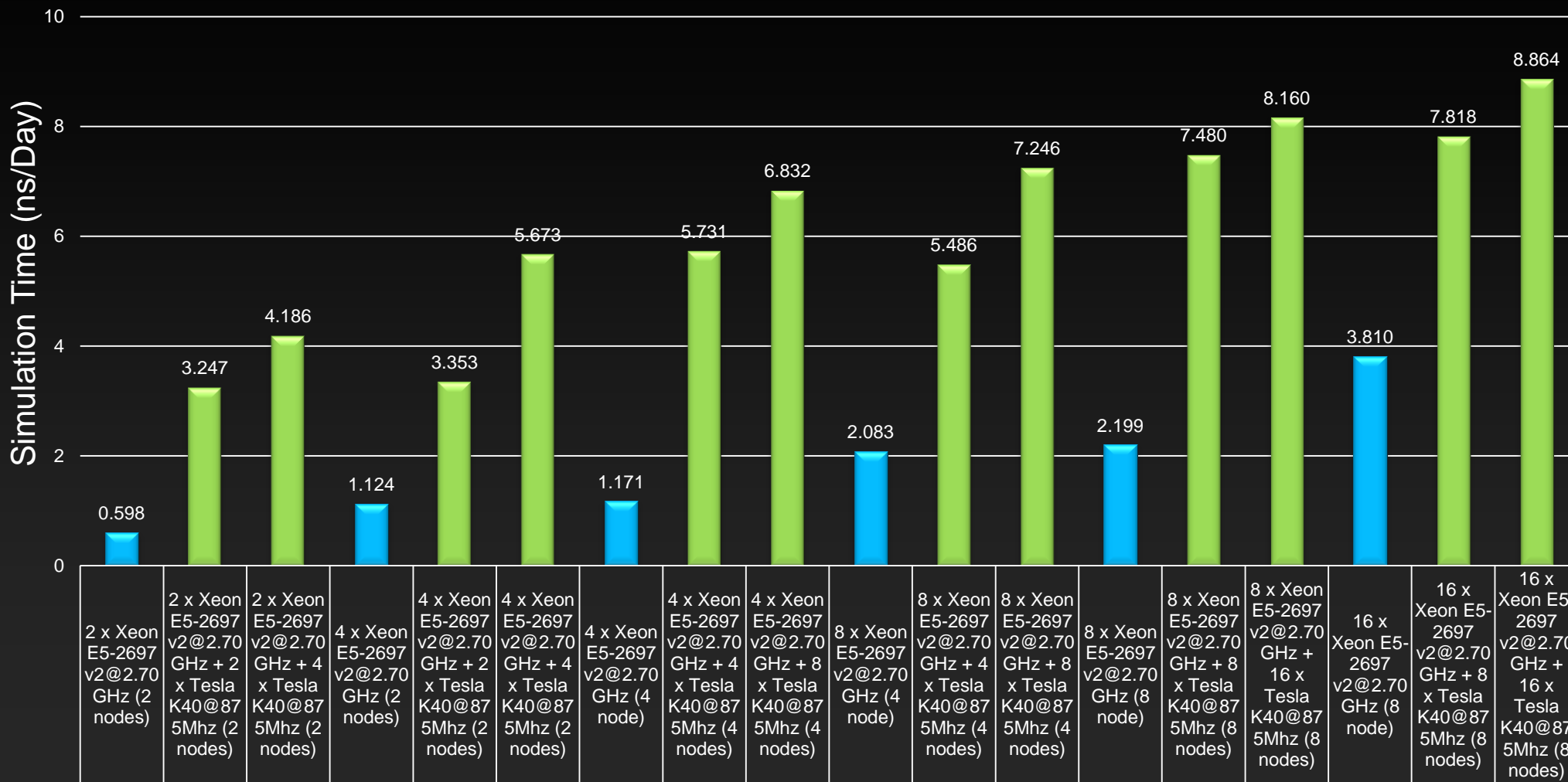
(1 Node: Simulation Time in ns/Day)



NAMD 2.10; F1-ATPase on Tesla K40s & IVB CPUs



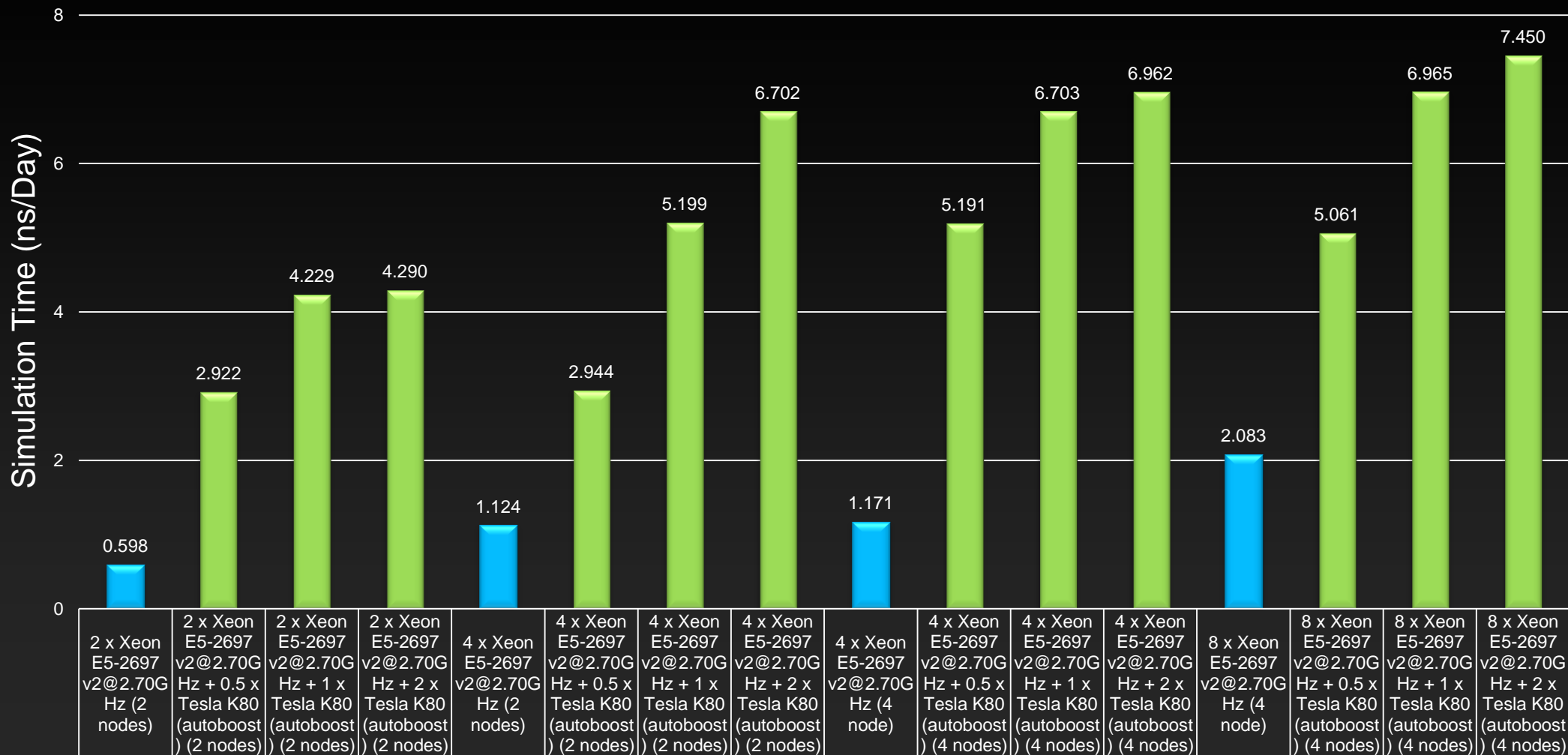
(2-8 Nodes: Simulation Time in ns/Day)



NAMD 2.10; F1-ATPase on Tesla K80s & IVB CPUs



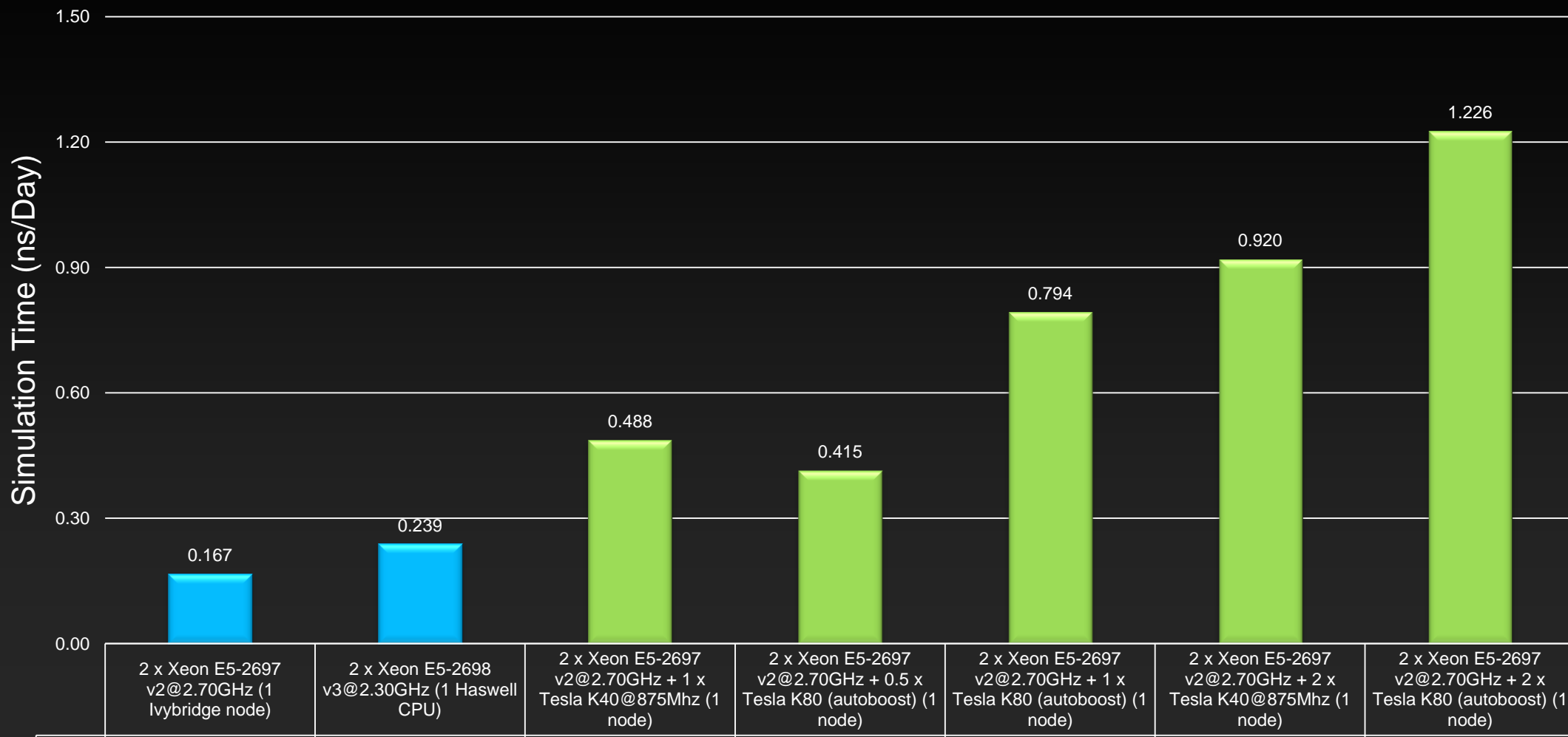
(2-4 Nodes: Simulation Time in ns/Day)



NAMD 2.10; STMV on Intel Phi, Tesla K40s and K80s & IVB CPUs



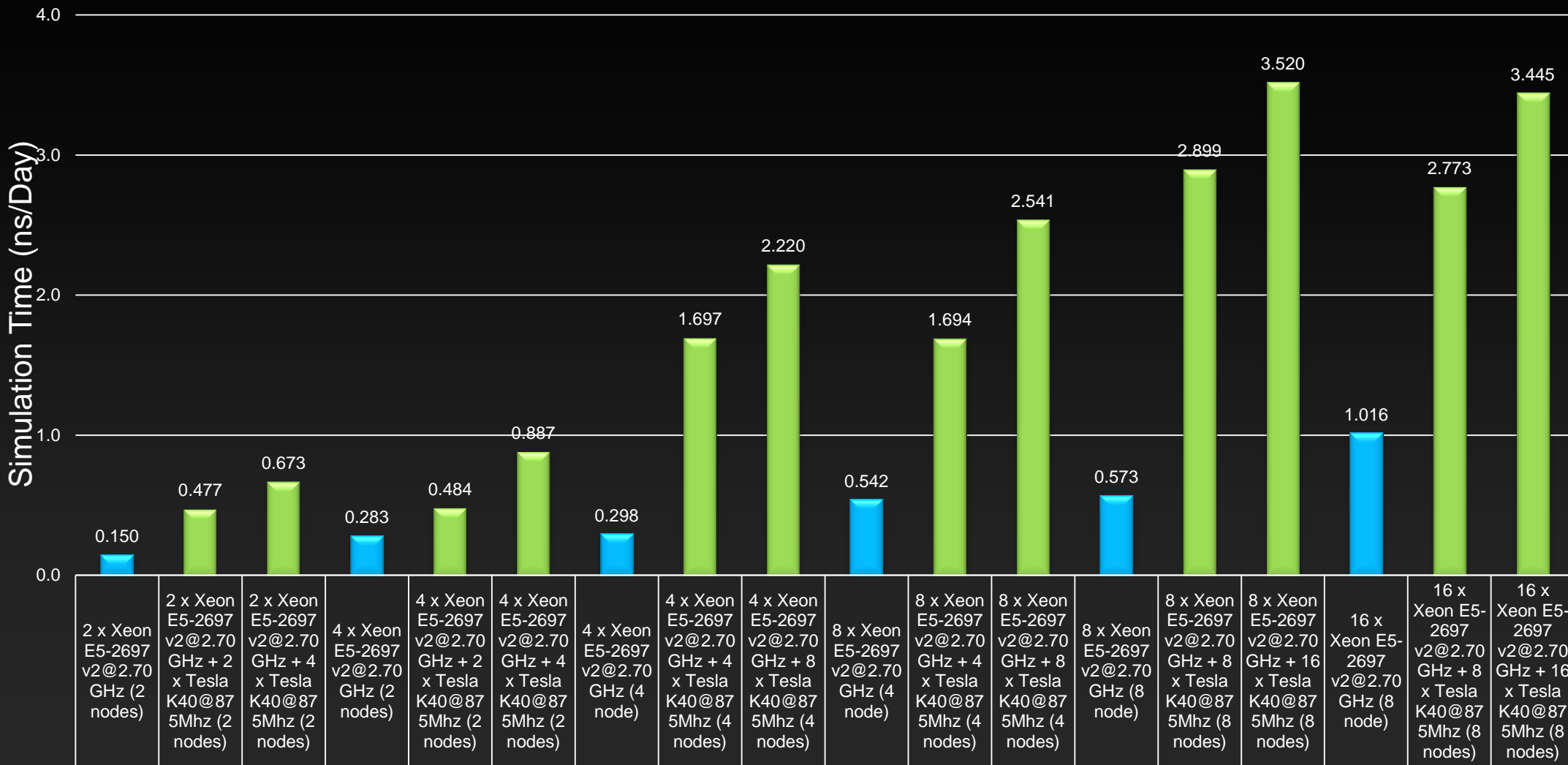
(1 Node: Simulation Time in ns/Day)



NAMD 2.10; STMV on Tesla K40s & IVB CPUs



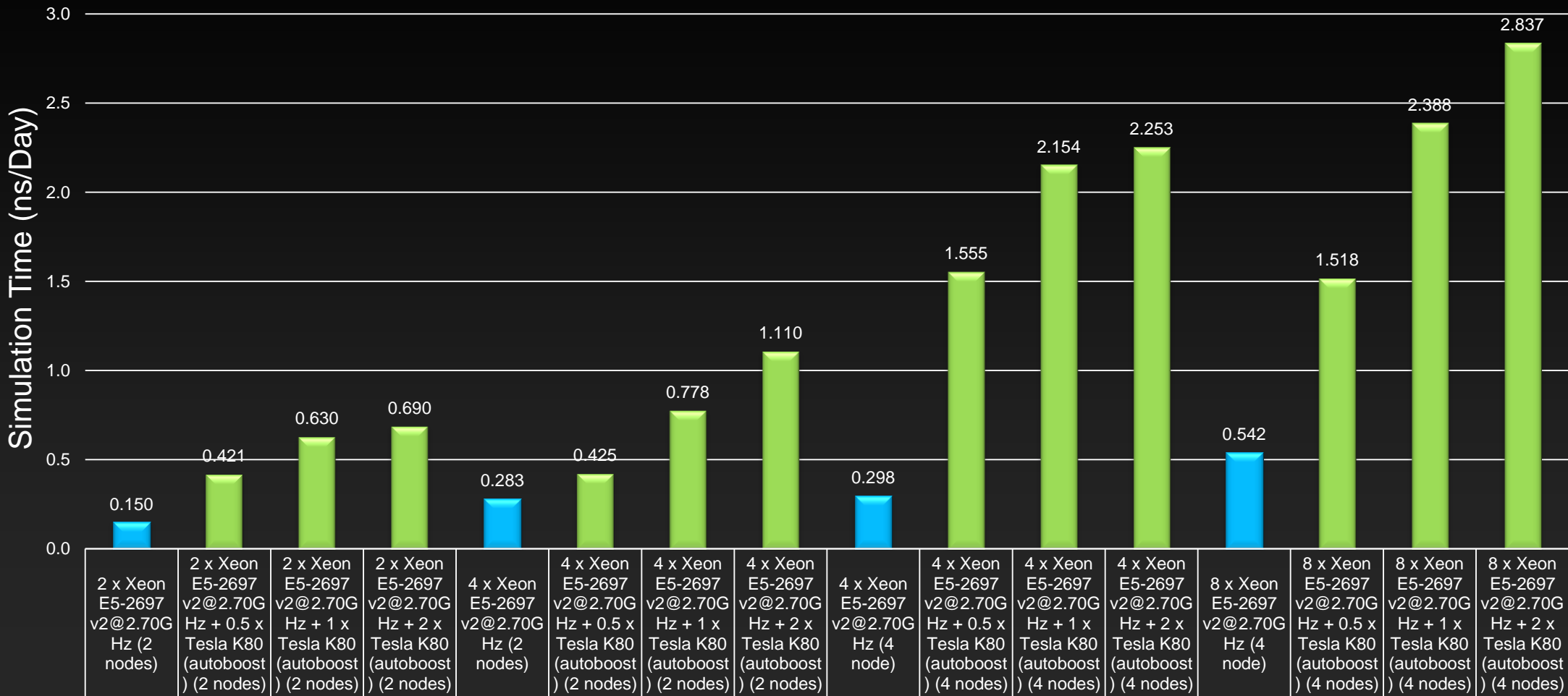
(2-8 Nodes: Simulation Time in ns/Day)



NAMD 2.10; STMV on Tesla K80s & IVB CPUs



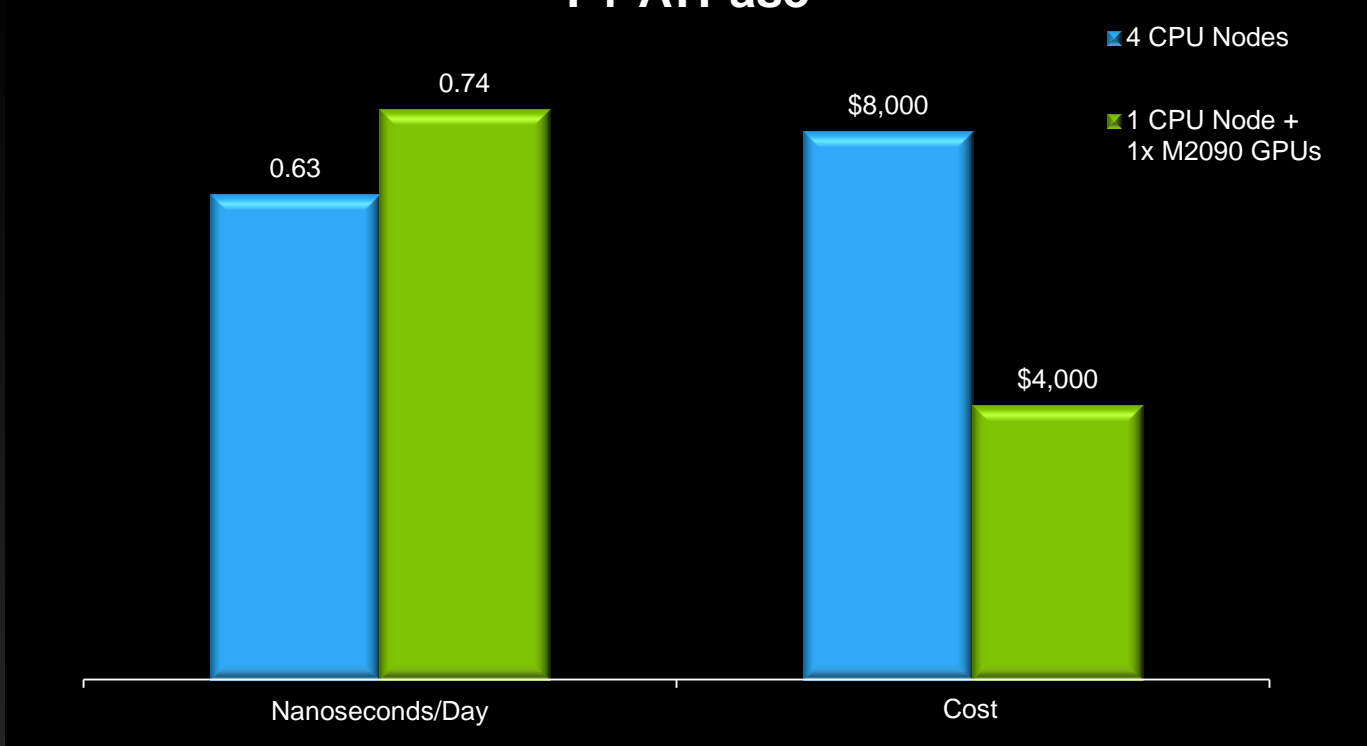
(2-4 Nodes: Simulation Time in ns/Day)





Replace 3 Nodes with 1 2090 GPU

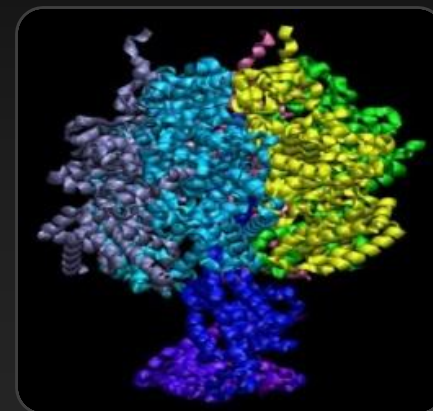
F1-ATPase



Running **NAMD** version 2.9
Each **blue node** contains 2x Intel Xeon X5550 CPUs (4 Cores per CPU).

The **green node** contains 2x Intel Xeon X5550 CPUs (4 Cores per CPU) and 1x NVIDIA M2090 GPU

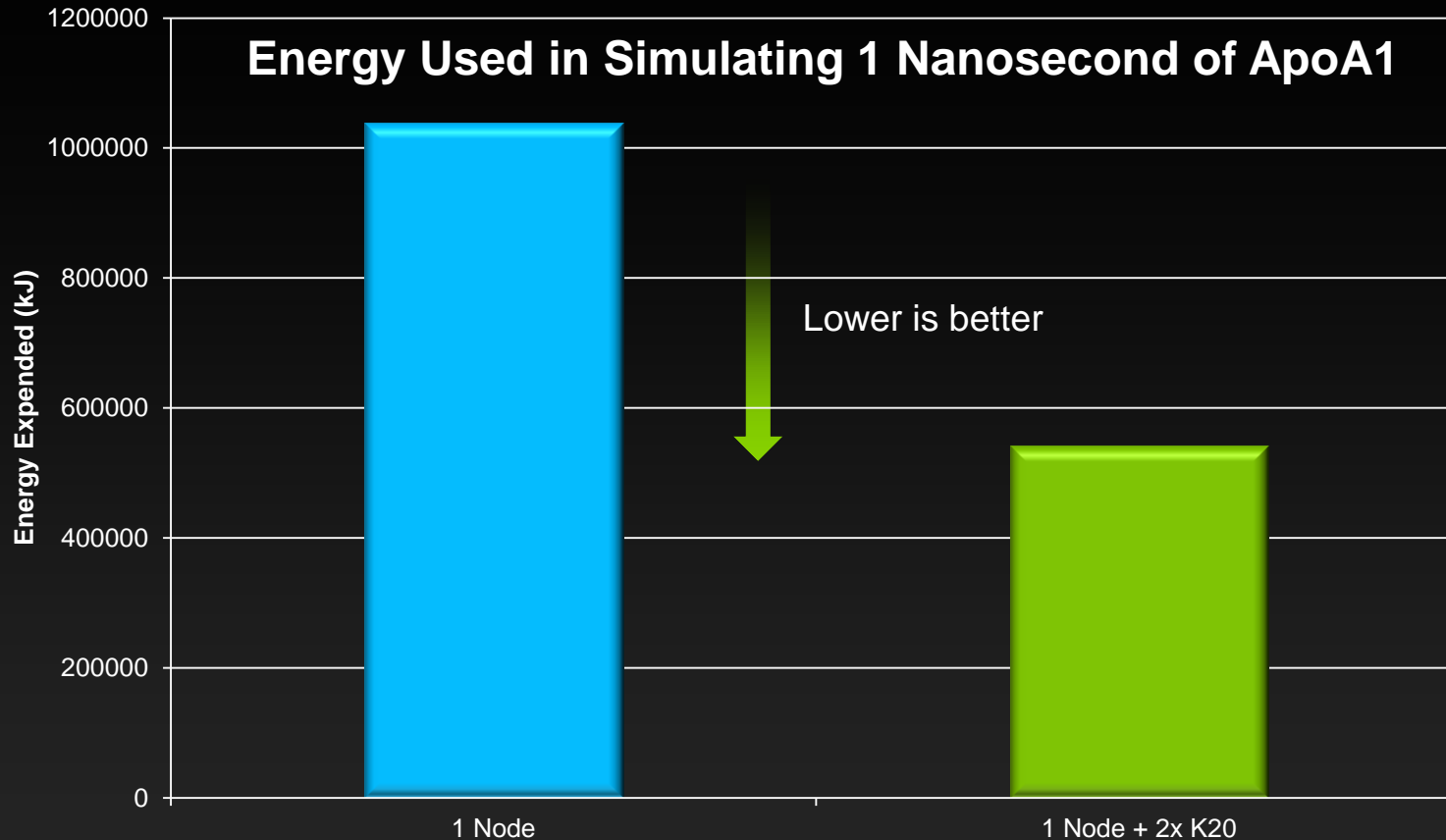
Note: Typical CPU and GPU node pricing used. Pricing may vary depending on node configuration. Contact your preferred HW vendor for actual pricing.



F1-ATPase

Speedup of **1.2x** for **50%** the cost

K20 - Greener: Twice The Science Per Watt



Energy Used in Simulating 1 Nanosecond of ApoA1

Lower is better

Cut down energy usage by $\frac{1}{2}$ with GPUs

Running **NAMD** version 2.9
Each **blue node** contains Dual E5-2687W CPUs (95W, 4 Cores per CPU).

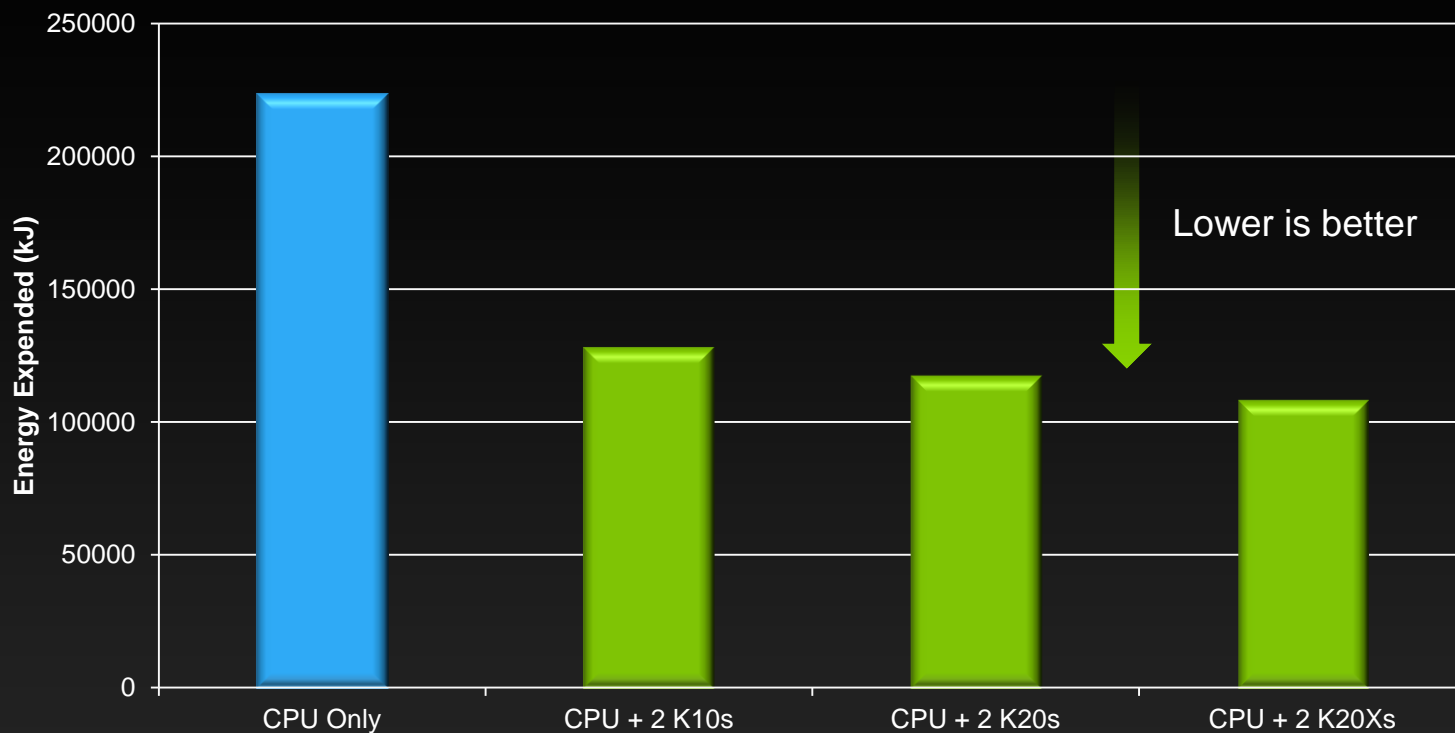
Each **green node** contains 2x Intel Xeon X5550 CPUs (95W, 4 Cores per CPU) and 2x NVIDIA K20 GPUs (225W per GPU)

$$\text{Energy Expended} = \text{Power} \times \text{Time}$$

Kepler - Greener: Twice The Science/Joule



Energy used in simulating 1 ns of SMTV

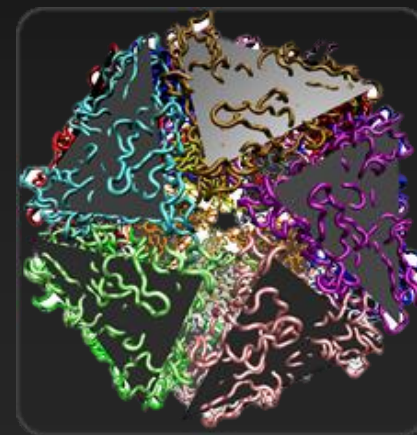


Running **NAMD** version 2.9

The **blue node** contains Dual E5-2687W CPUs (150W each, 8 Cores per CPU).

The **green nodes** contain Dual E5-2687W CPUs (8 Cores per CPU) and 2x NVIDIA K10, K20, or K20X GPUs (235W each).

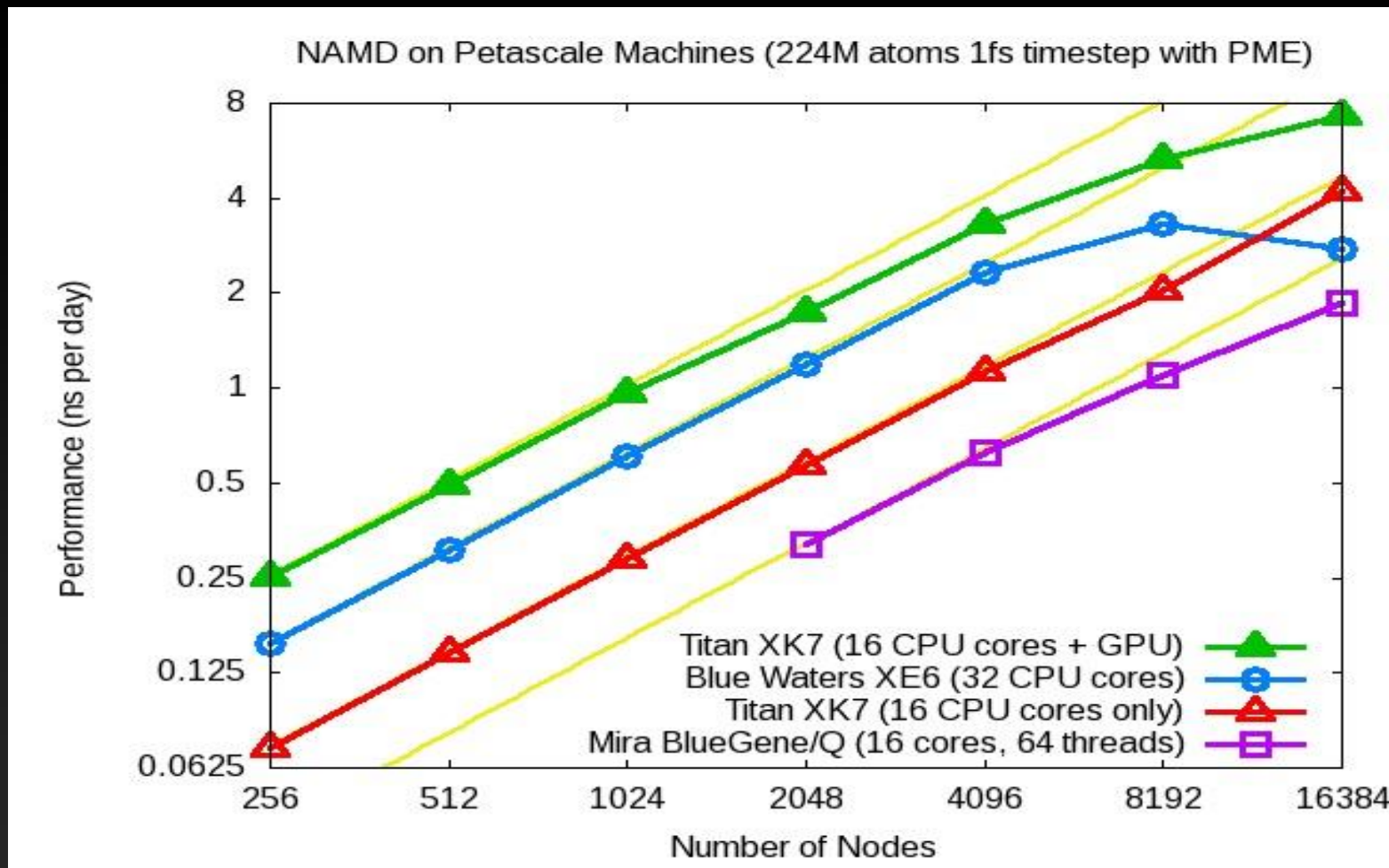
*Energy Expended
= Power x Time*



Satellite Tobacco Mosaic Virus

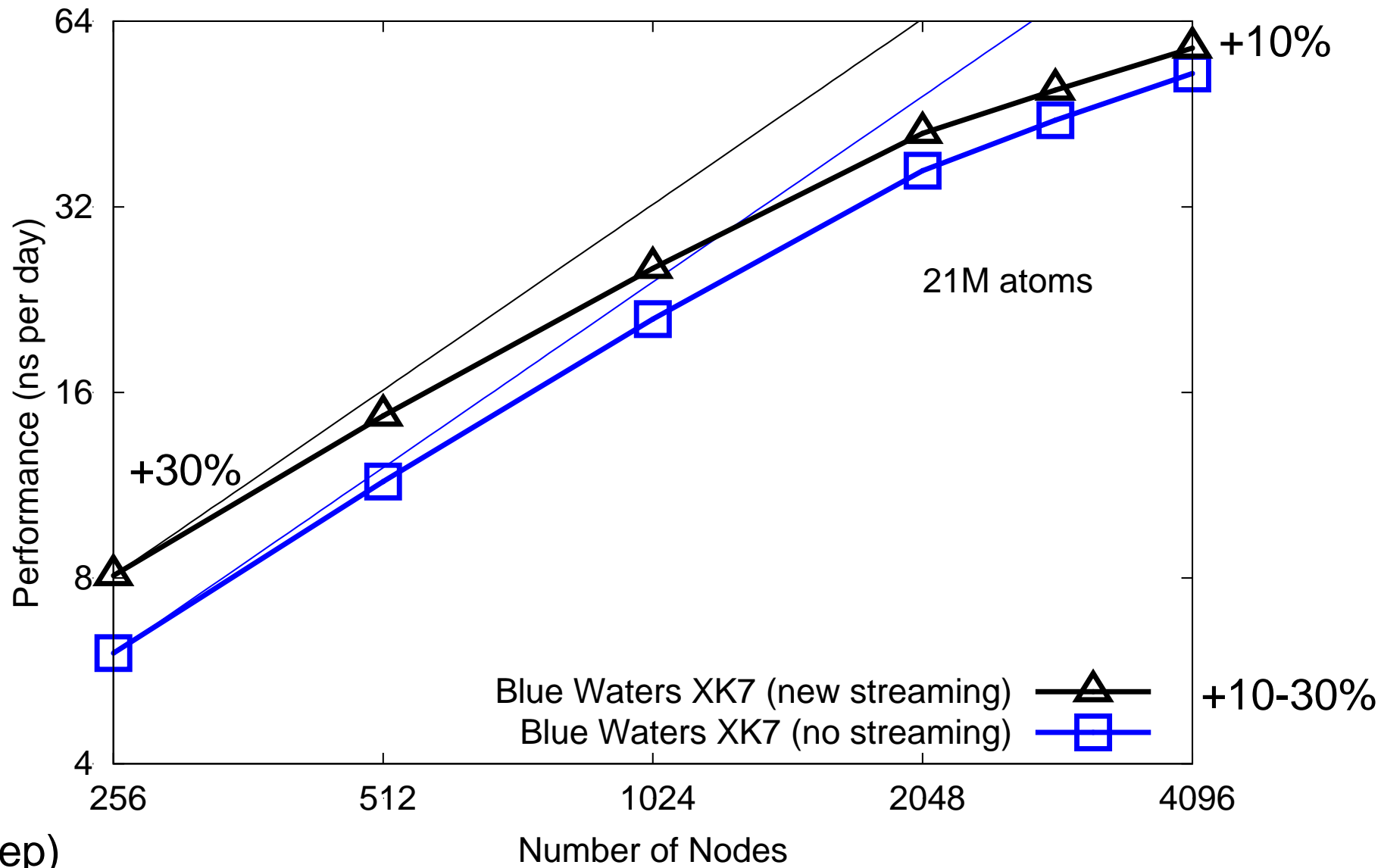
Cut down energy usage by $\frac{1}{2}$ with GPUs

NAMD CPU and GPU Scaling

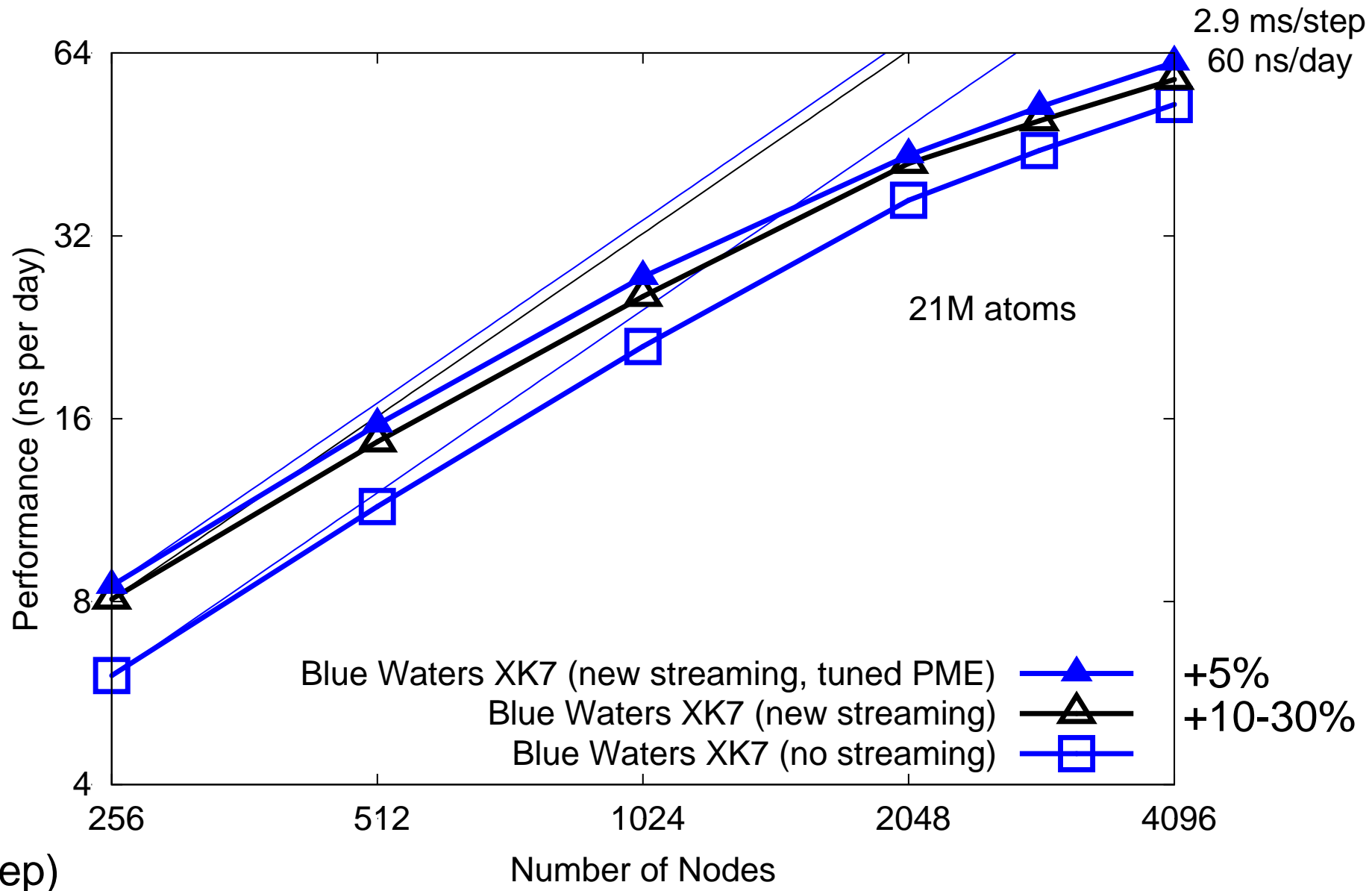


The next 8 slides are courtesy of Jim Phillips @ Beckman Institute UIUC

New Streaming Kernel Performance

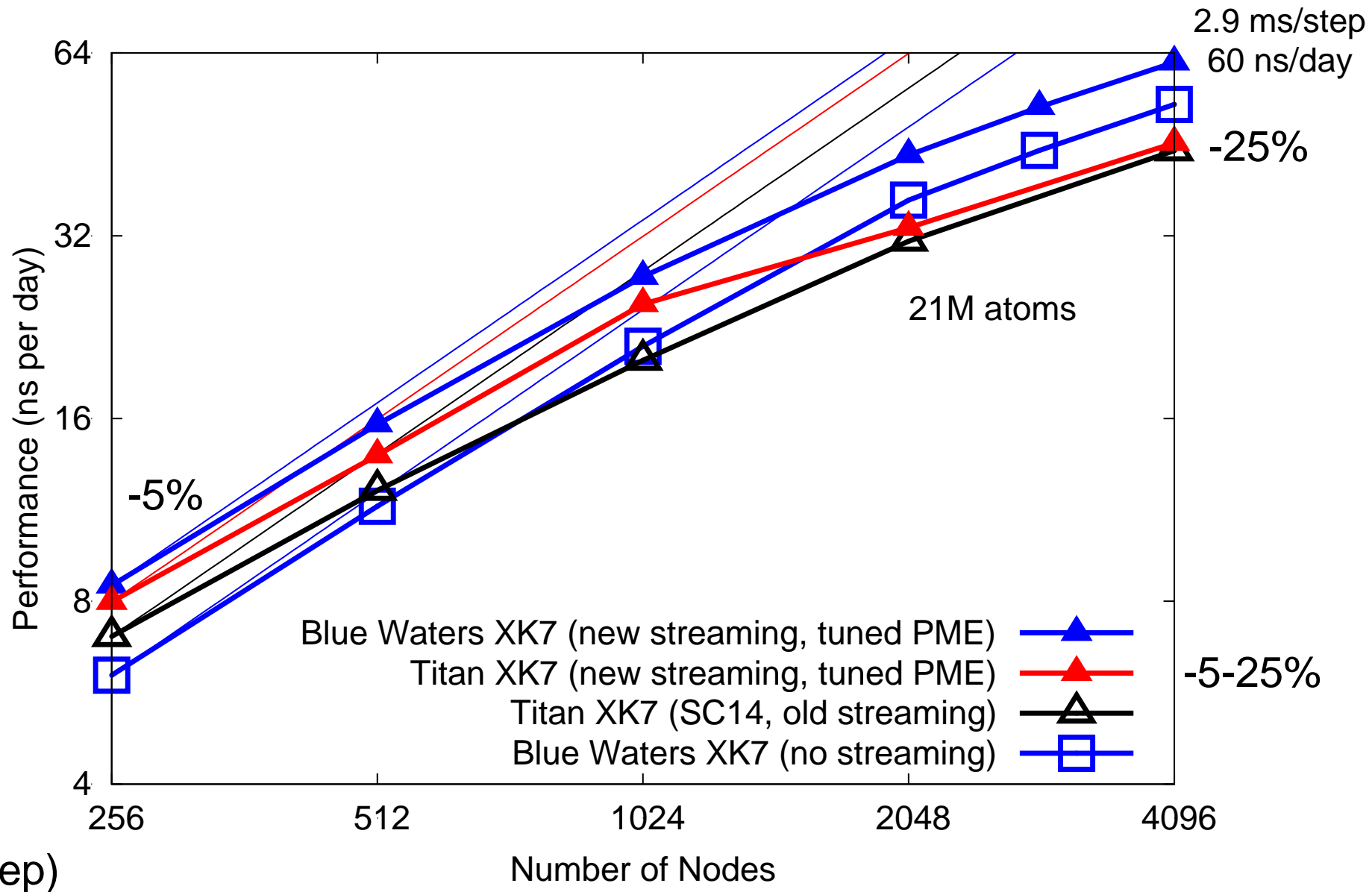


Parallelize PME Within Node



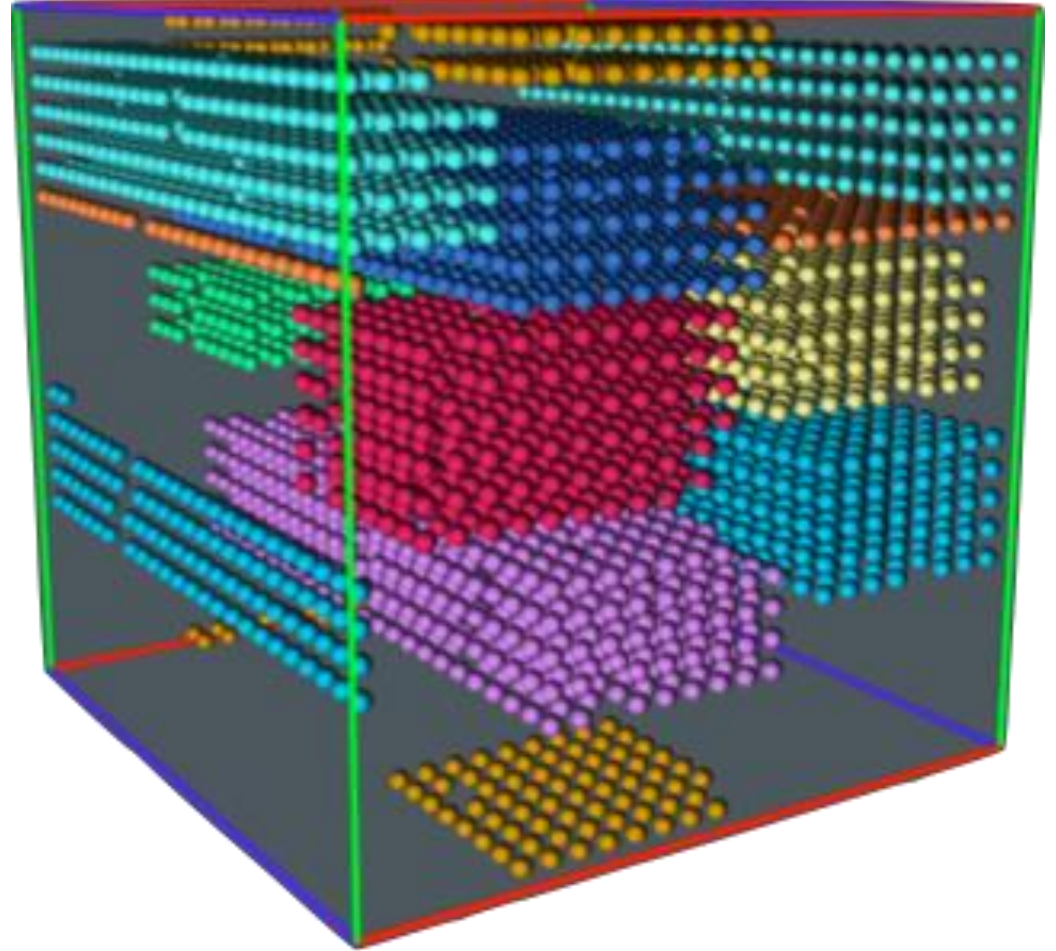
(2fs timestep)

Blue Waters vs Titan

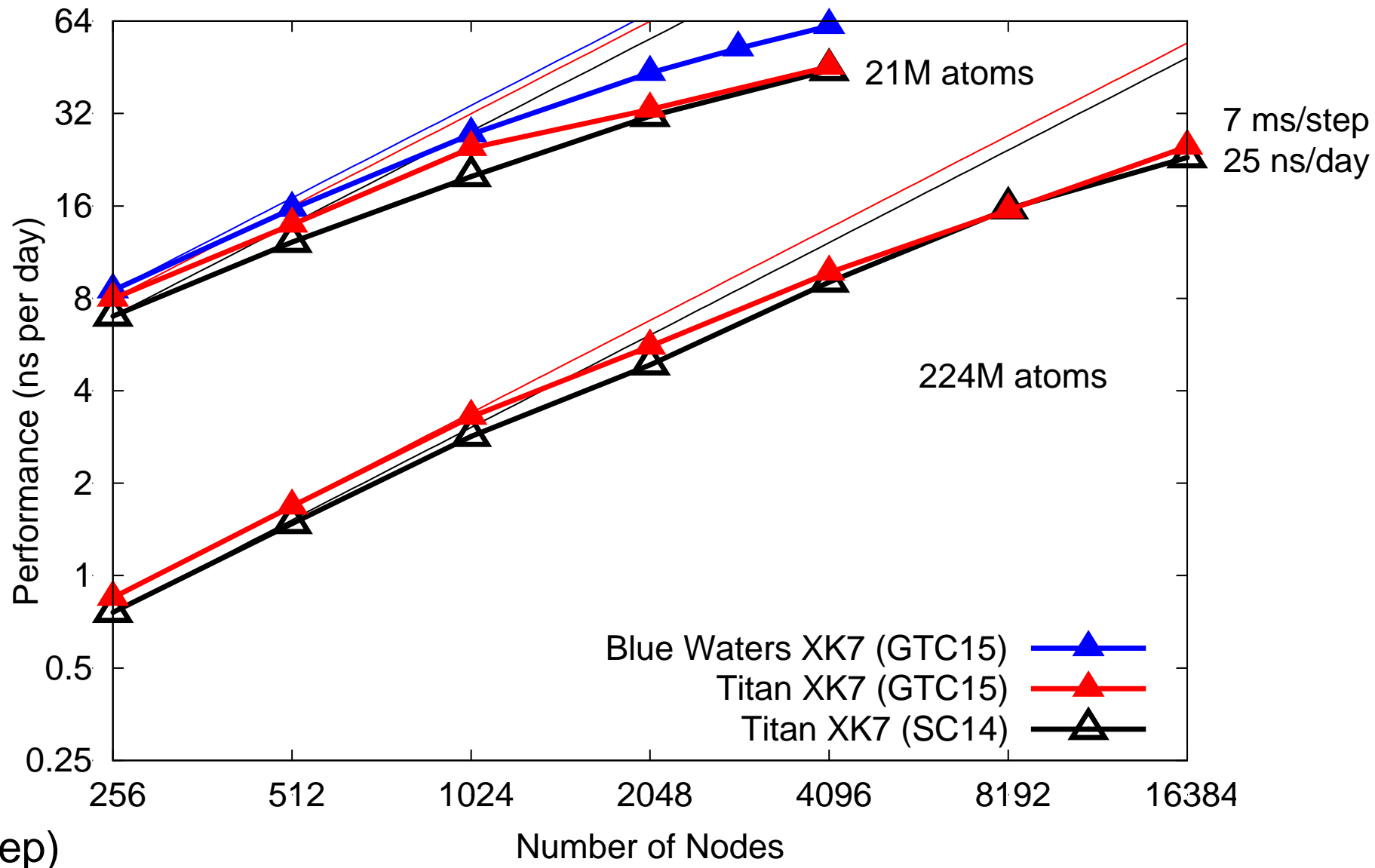


Topology-Aware Scheduling on Blue Waters

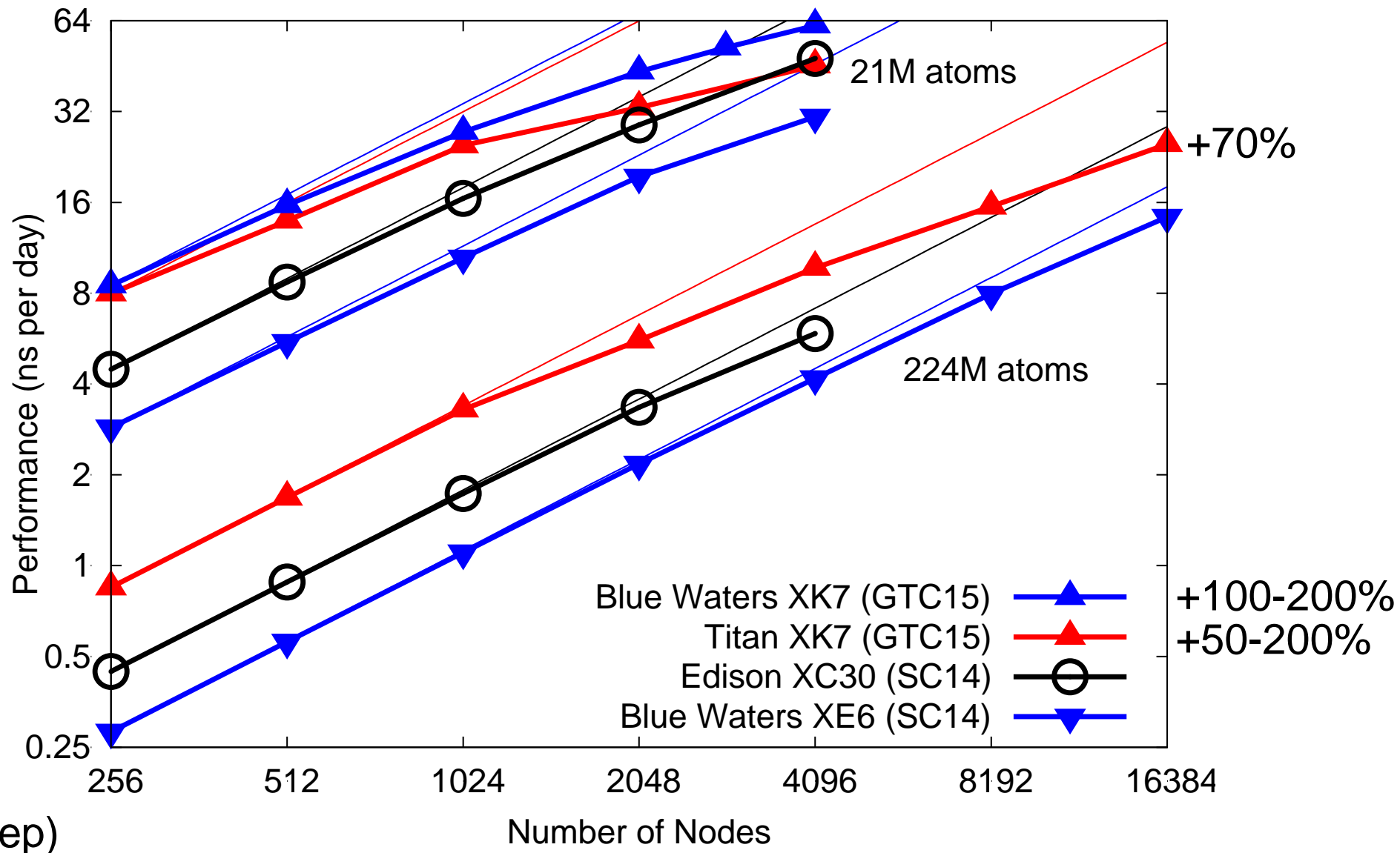
- Map jobs to convex sets to avoid network interference
- NCSA, Cray, Adaptive
- Just enabled January 13
- Most likely explanation for Blue Waters performance advantage over Titan
- See Enos *et al.*, CUG 2014



Blue Waters vs Titan



Comparison with CPU-only Machines



Recommended GPU Node Configuration for NAMD Computational Chemistry



Workstation or Single Node Configuration	
# of CPU sockets	2
Cores per CPU socket	6+
CPU speed (Ghz)	2.66+
System memory per socket (GB)	32
GPUs	Kepler K20, K40, K80
# of GPUs per CPU socket	1-2
GPU memory preference (GB)	6-12
GPU to CPU connection	PCIe 3.0 or higher
Server storage	500 GB or higher
Network configuration	Gemini, InfiniBand

Scale to multiple nodes with same single node configuration

Summary/Conclusions

Benefits of GPU Accelerated Computing

- Faster than CPU only systems in all tests
- Large performance boost with small marginal price increase
- Energy usage cut in half
- GPUs scale very well within a node and over multiple nodes
- Tesla K20 GPU is our fastest and lowest power high performance GPU to date

Try GPU accelerated NAMD for free – www.nvidia.com/GPUTestDrive

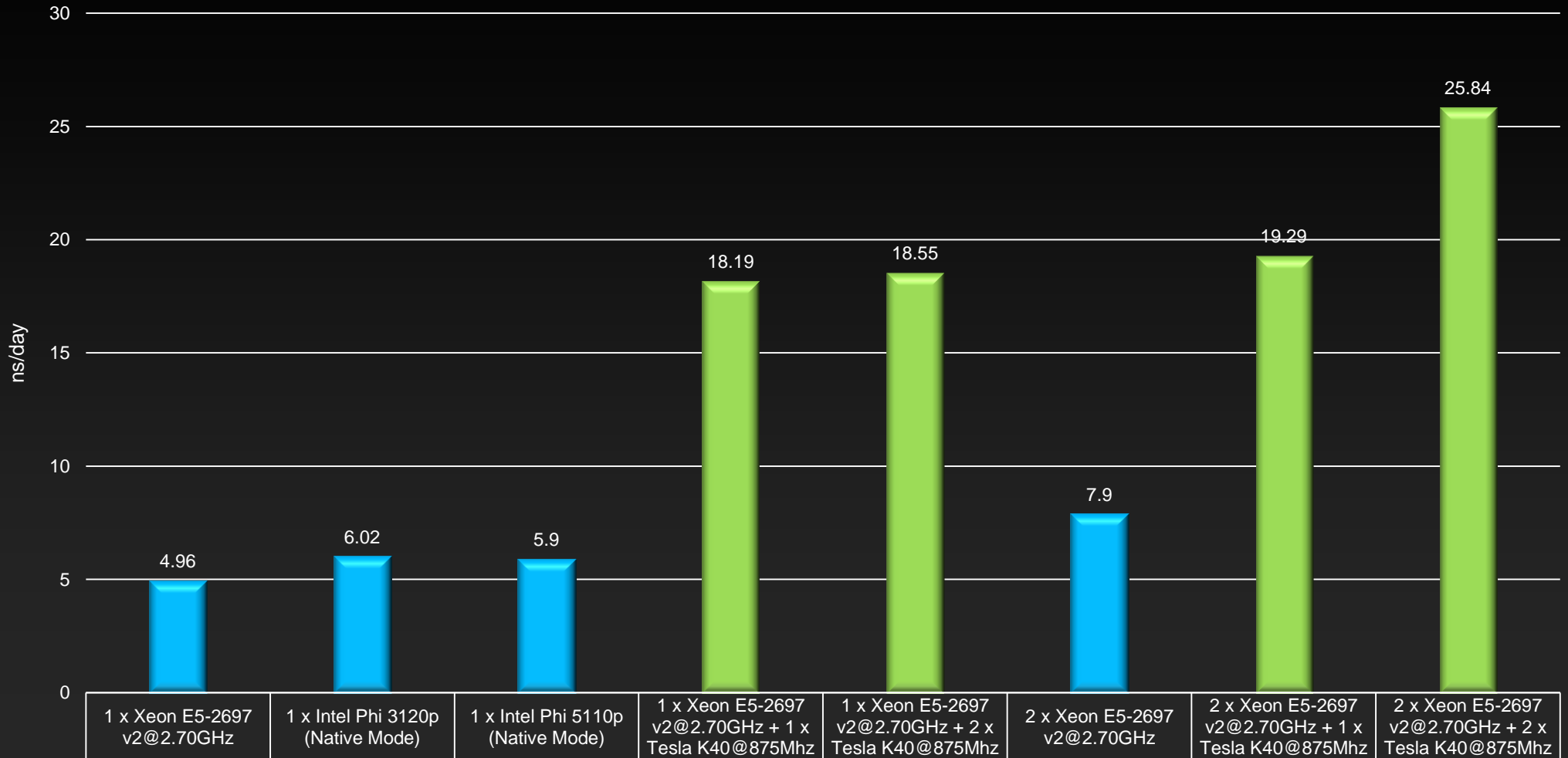


GROMACS 5.0

GROMACS 5.0: Phi vs. Kepler K40 fastest GPU!



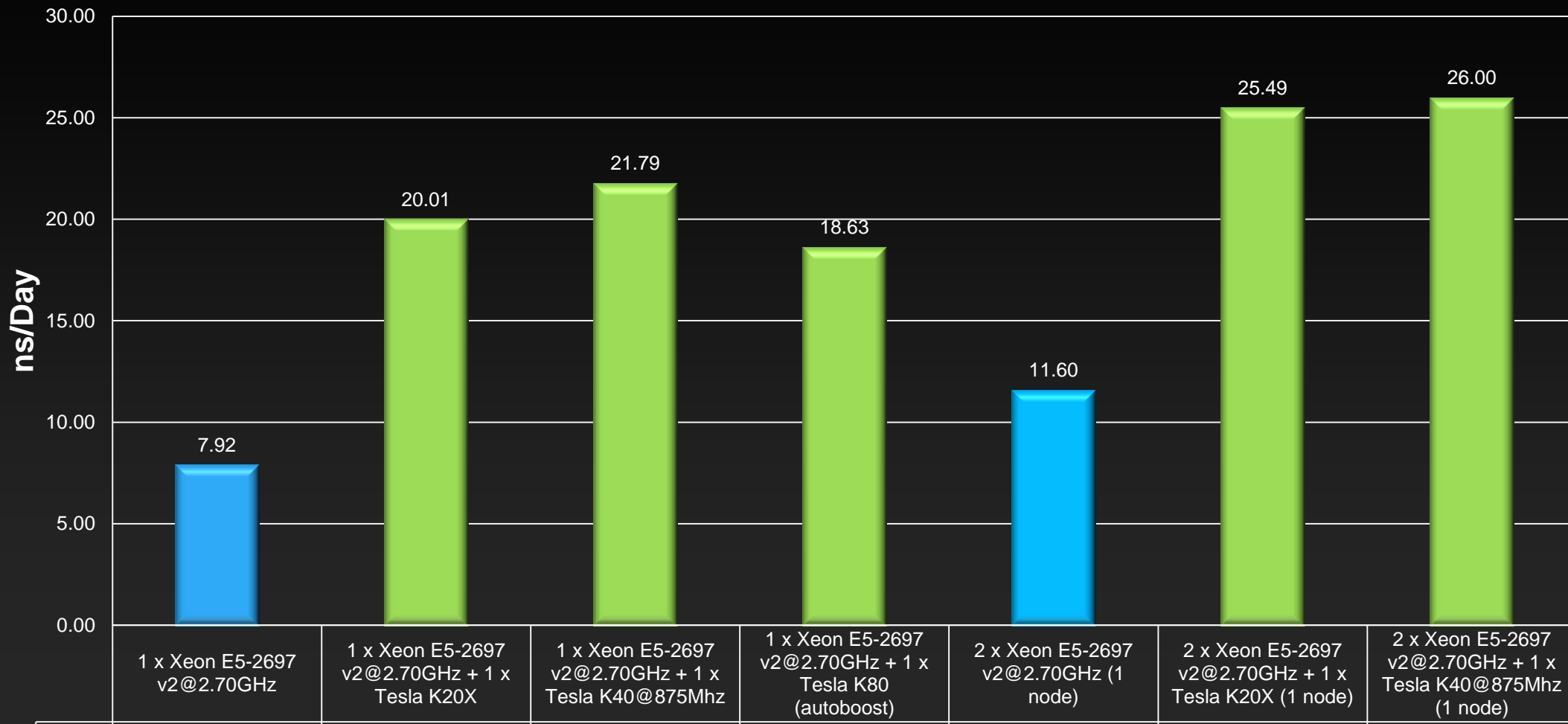
GROMACS 5.0 RC1 (ns/day) on K40 with Boost Clocks and Intel Phi
192K Waters Benchmark (CUDA 6.0)



GROMACS 5.0 & Fastest Kepler GPUs yet!



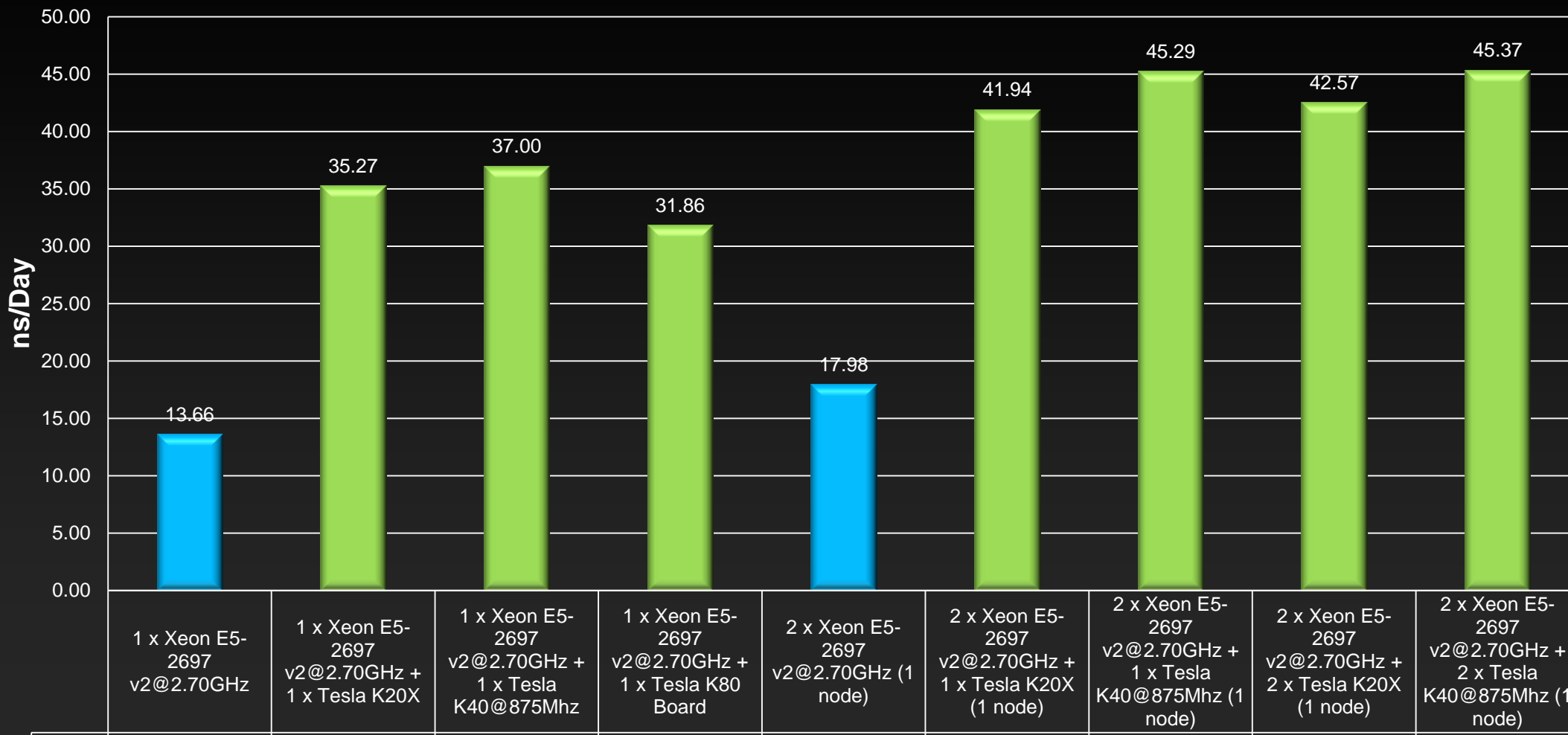
GROMACS 5.0, cresta_ion_channel
Single Node with & without Kepler GPUs



GROMACS 5.0 & Fastest Kepler GPUs yet!



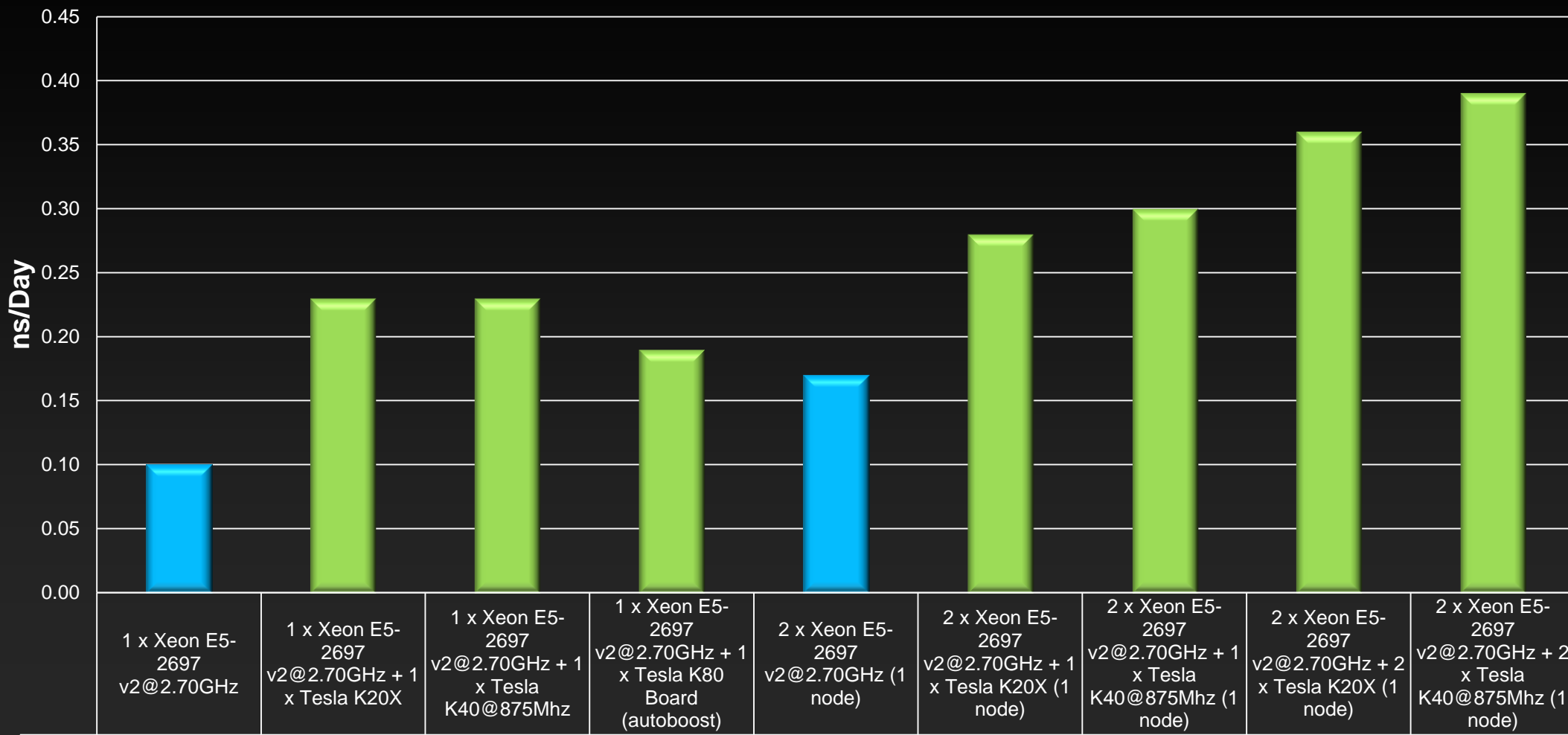
GROMACS 5.0, cresta_ion_channel_vsites
Single node with & without Kepler GPUs



GROMACS 5.0 & Fastest Kepler GPUs yet!



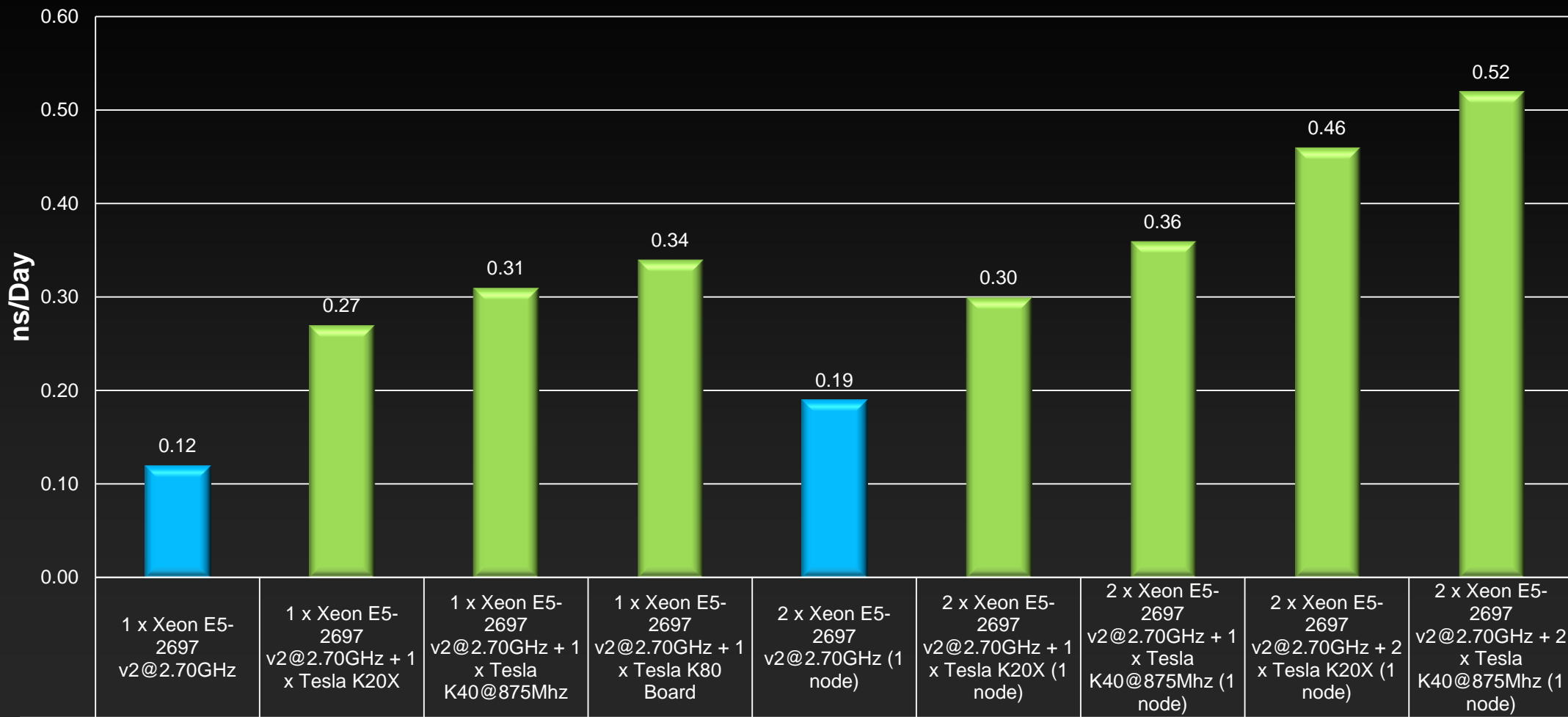
GROMACS 5.0, cresta_methanol
Single node with & without Kepler GPUs



GROMACS 5.0 & Fastest Kepler GPUs yet!



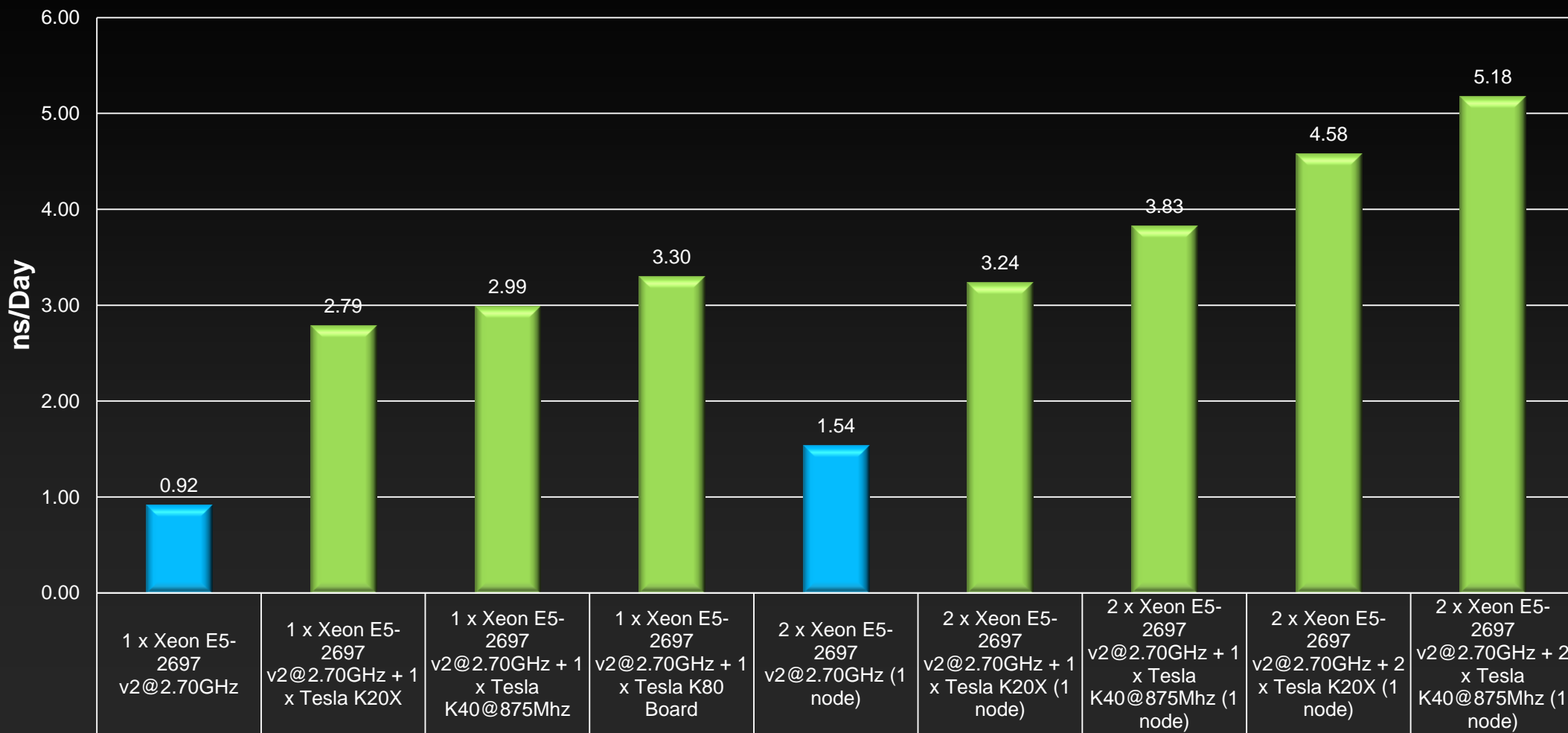
GROMACS 5.0, cresta_methanol_rf
Single Node with & without Kepler GPUs



GROMACS 5.0 & Fastest Kepler GPUs yet!



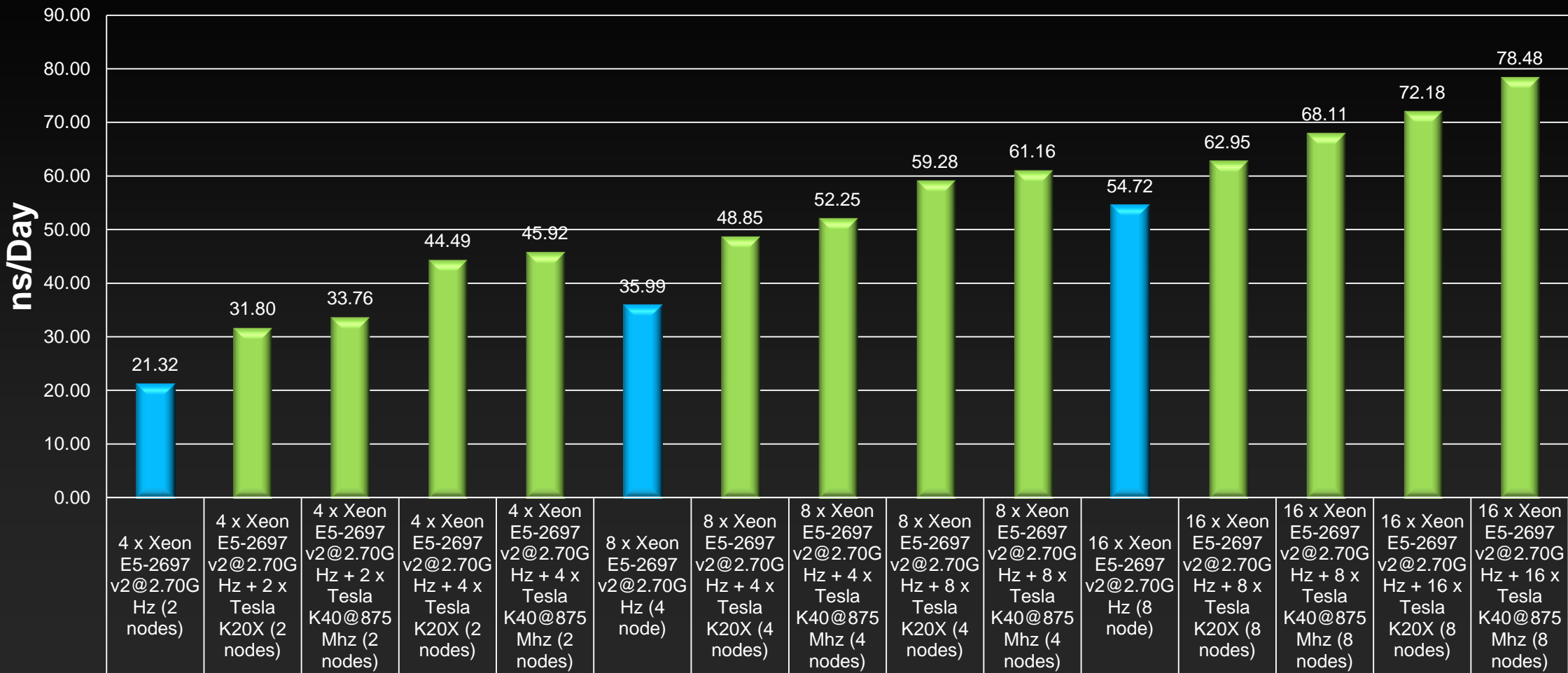
GROMACS 5.0, cresta_virus_capsid
Single Node with & without Kepler GPUs



GROMACS 5.0 & Fastest Kepler GPUs yet!



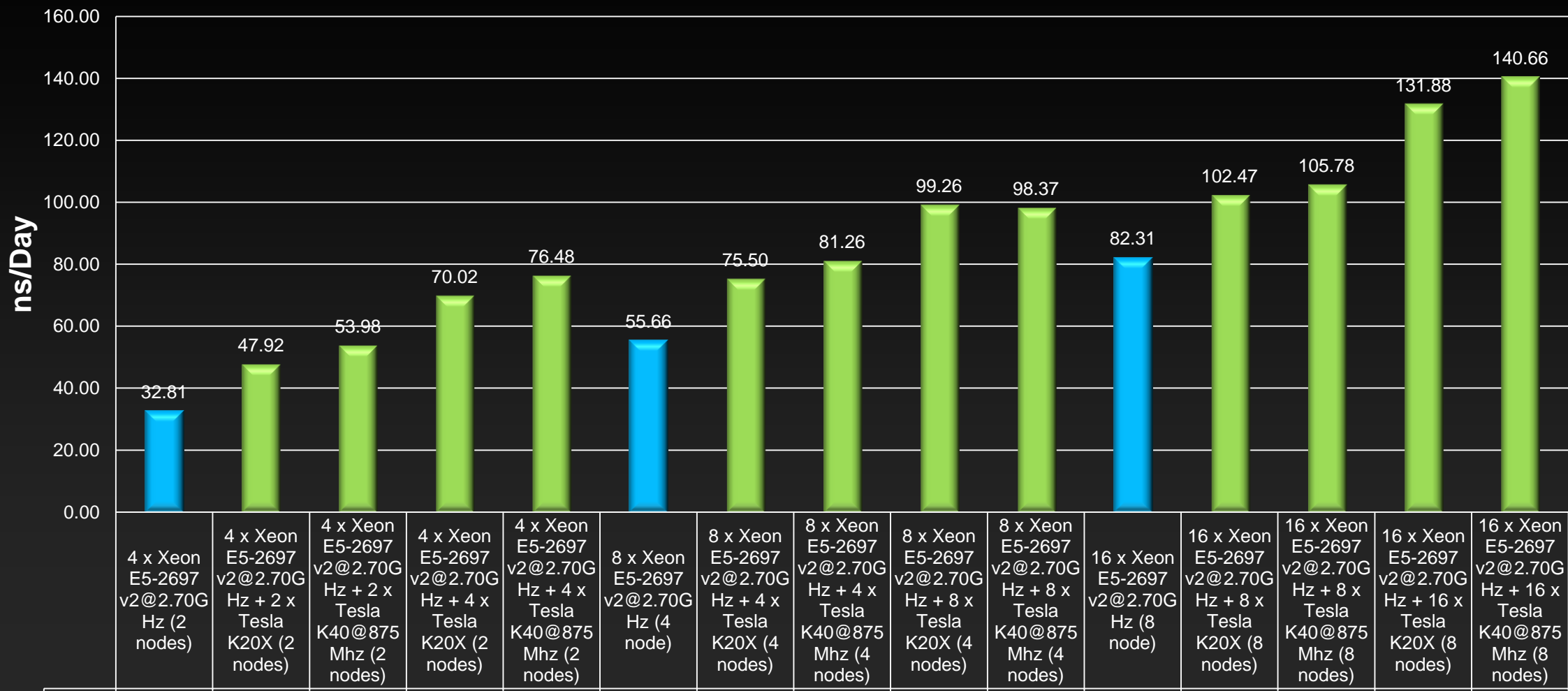
GROMACS 5.0, cresta_ion_channel
2 to 8 Nodes, with & without Kepler GPUs



GROMACS 5.0 & Fastest Kepler GPUs yet!



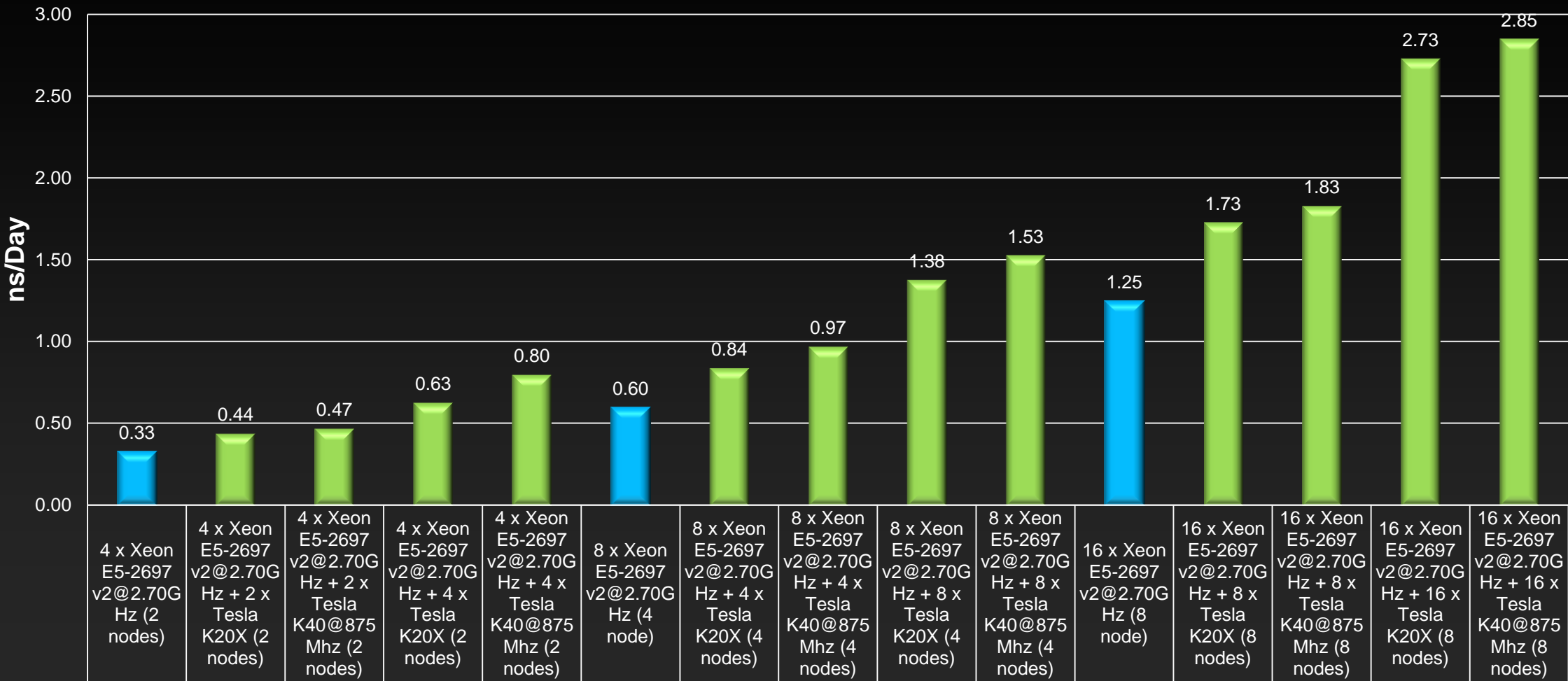
GROMACS 5.0, cresta_ion_channel_vsites
2 to 8 Nodes, with & without Kepler GPUs



GROMACS 5.0 & Fastest Kepler GPUs yet!



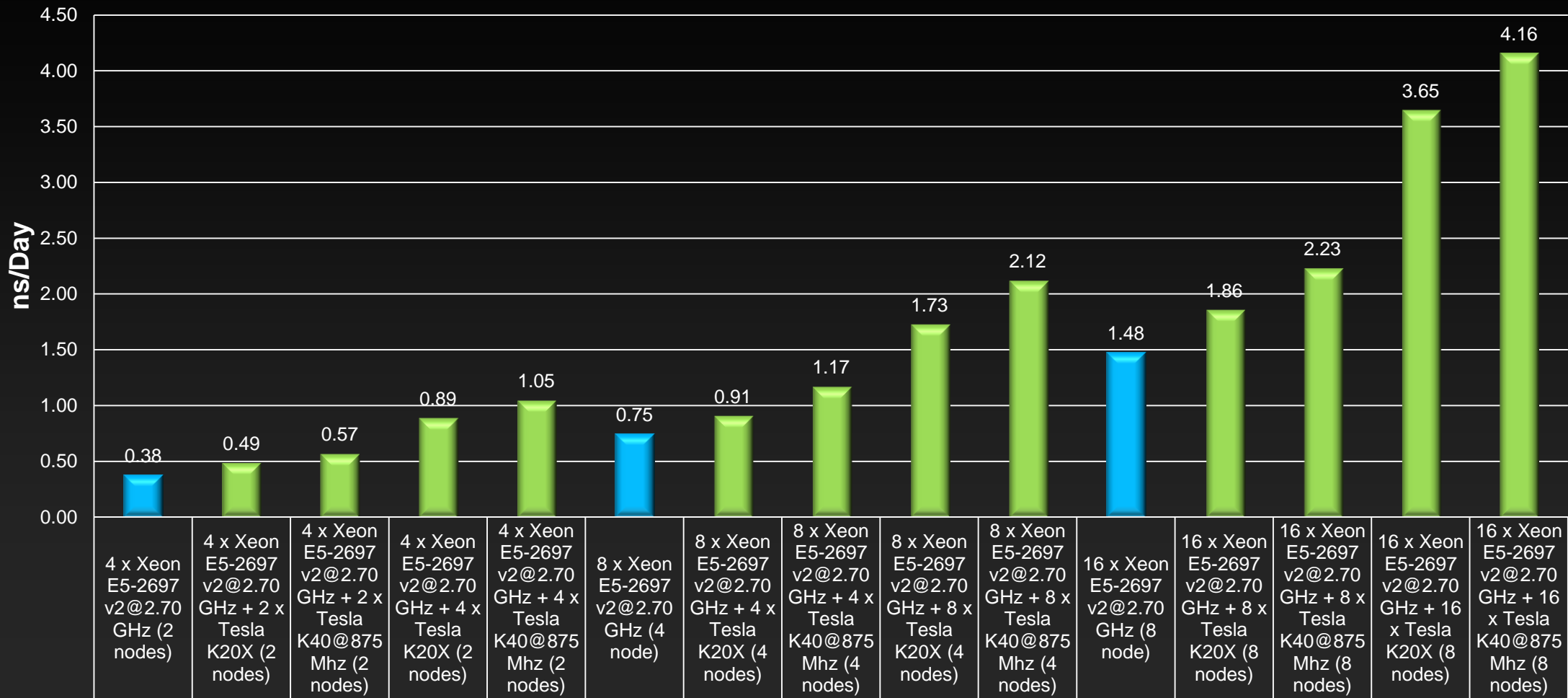
GROMACS 5.0, cresta_methanol
2 to 8 Nodes, with & without Kepler GPUs



GROMACS 5.0 & Fastest Kepler GPUs yet!



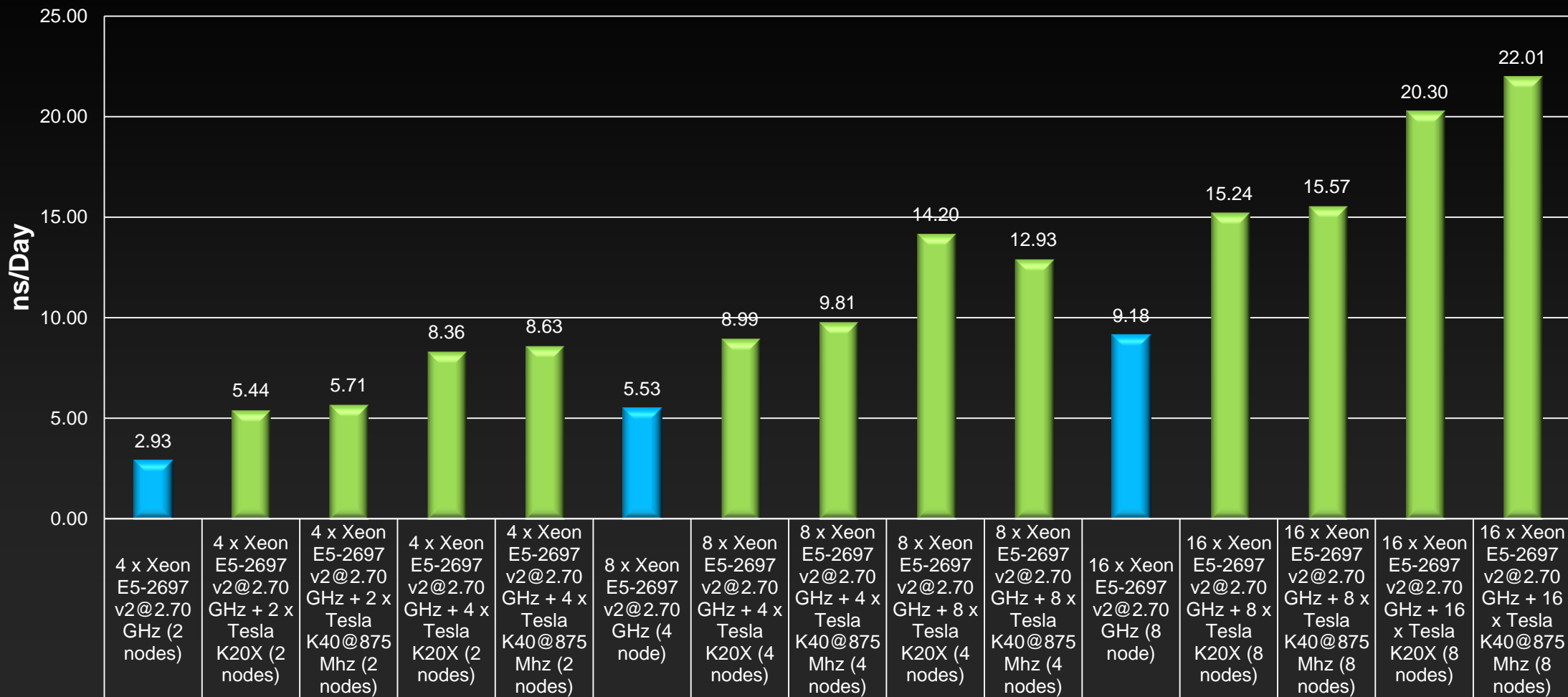
GROMACS 5.0, cresta_methanol_rf
2 to 8 Nodes, with & without Kepler GPUs



GROMACS 5.0 & Fastest Kepler GPUs yet!



GROMACS 5.0, cresta_virus_capsid
2 to 8 Nodes, with & without Kepler GPUs



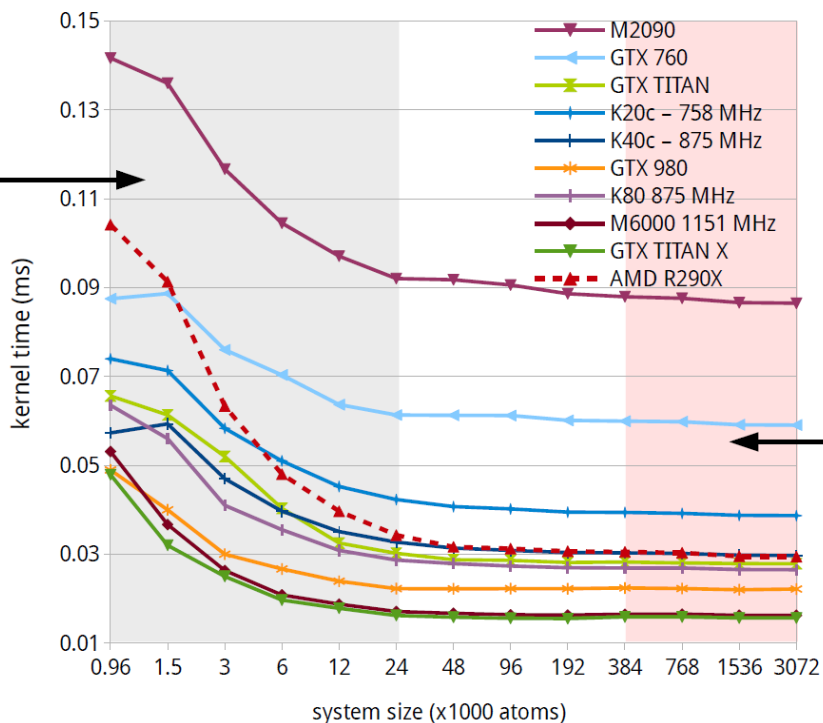
Slides - courtesy of GROMACS Dev Team



Kernel performance and scaling

Strong scaling regime:

This is where most of our efforts go!



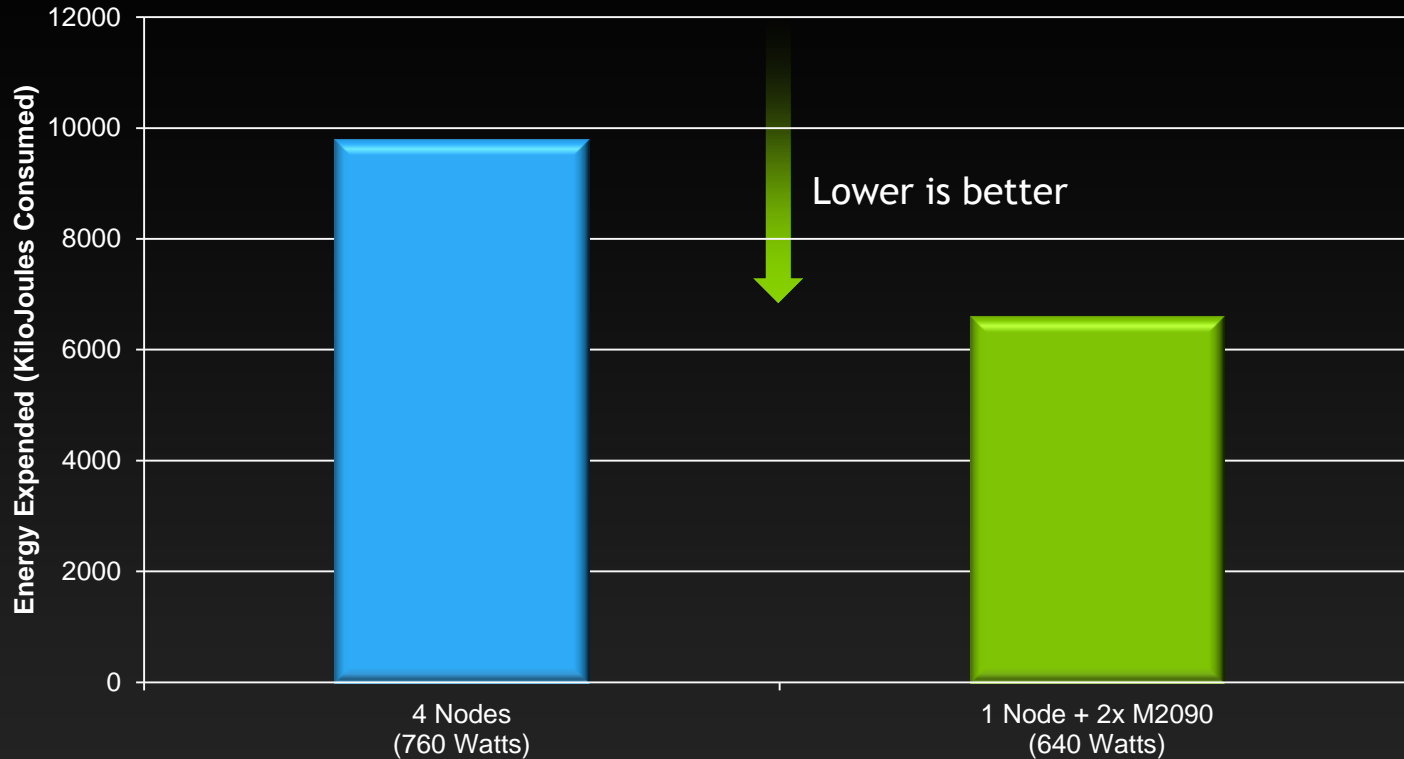
Benchmark "show-off" regime:

This is where the "free lunch" from new hardware comes in full effect

Greener Science



ADH in Water (134K Atoms)



Running **GROMACS** 4.6 with CUDA 4.1

The **blue nodes** contain 2x Intel X5550 CPUs (95W TDP, 4 Cores per CPU)

The **green node** contains 2x Intel X5550 CPUs, 4 Cores per CPU) and 2x NVIDIA M2090s GPUs (225W TDP per GPU)

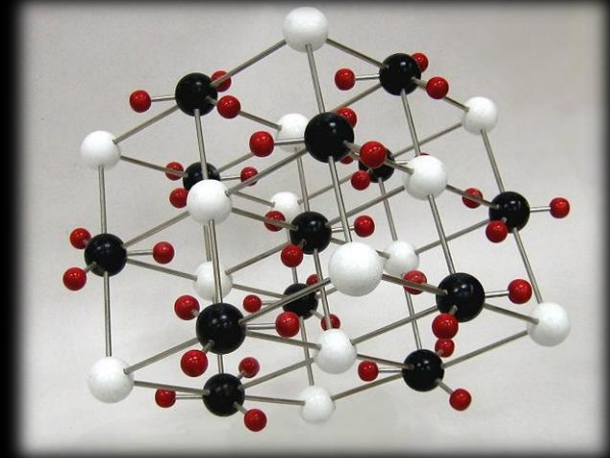
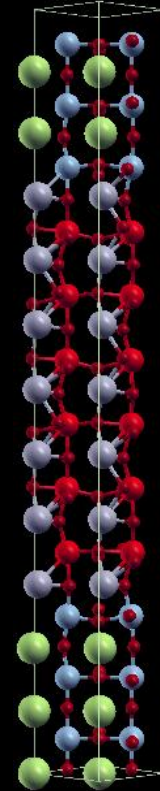
$$\text{Energy Expended} = \text{Power} \times \text{Time}$$

In simulating each nanosecond, the GPU-accelerated system uses **33% less energy**

Recommended GPU Node Configuration for GROMACS Computational Chemistry

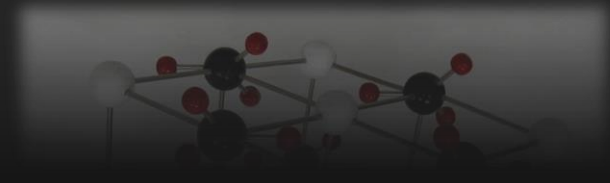


Workstation or Single Node Configuration	
# of CPU sockets	2
Cores per CPU socket	6+
CPU speed (Ghz)	2.66+
System memory per socket (GB)	32
GPUs	Kepler K20, K40, K80
# of GPUs per CPU socket	1x Kepler GPUs: need fast Sandy Bridge or Ivy Bridge, or high-end AMD Opterons
GPU memory preference (GB)	6
GPU to CPU connection	PCIe 3.0 or higher
Server storage	500 GB or higher
Network configuration	Gemini, InfiniBand



TESLA

Quantum Chemistry Module





GAUSSIAN

Gaussian

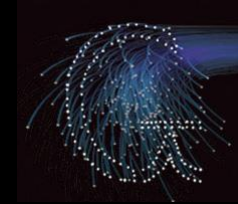


- ACS Fall 2011 press release
 - Joint collaboration between Gaussian, NVDA and PGI for GPU acceleration:
http://www.gaussian.com/g_press/nvidia_press.htm
 - **No such press release exists for Intel MIC or AMD GPUs**
 - Mike Frisch quote from press release:
 - *“Calculations using Gaussian are limited primarily by the available computing resources,” said Dr. Michael Frisch, president of Gaussian, Inc. “By coordinating the development of hardware, compiler technology and application software among the three companies, the new application will bring the speed and cost-effectiveness of GPUs to the challenging problems and applications that Gaussian’s customers need to address.”*

Select Slides from “Enabling Gaussian 09 on GPGPUs” at GTC March 2014



- In 2011 Gaussian, Inc., NVIDIA Corp. and PGI started a long-term project to enable all the performance critical paths of Gaussian on GPGPUs.
 - Ultimate goal is to show significant performance improvement by using accelerators in conjunction with CPUs
 - Initial efforts are directed towards creating an infrastructure that will leverage the current CPU code base and at the same time minimize the additional maintenance effort associated with running on GPUs.
- Current status of this work for Direct Hartree-Fock and triples-correction calculations as applied in for example Coupled Cluster calculations that uses mostly the directives based OpenACC framework.
- Slides & Audio: <http://on-demand.gputechconf.com/gtc/2014/video/S4613-enabling-gaussian-09-gpgpus.mp4>



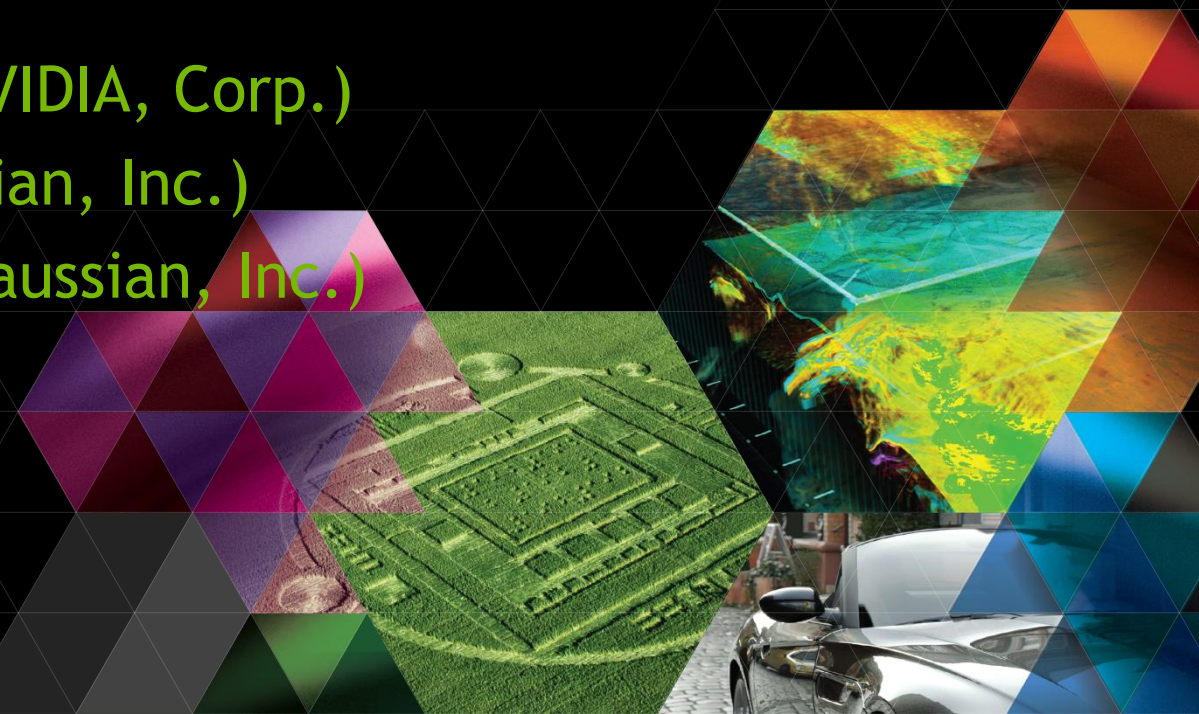
CURRENT STATUS OF THE PROJECT TO ENABLE GAUSSIAN 09 ON GPGPUS

Roberto Gomperts (NVIDIA, Corp.)

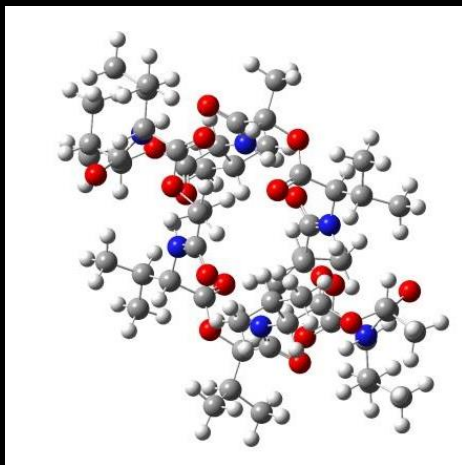
Michael Frisch (Gaussian, Inc.)

Giovanni Scalmani (Gaussian, Inc.)

Brent Leback (PGI)

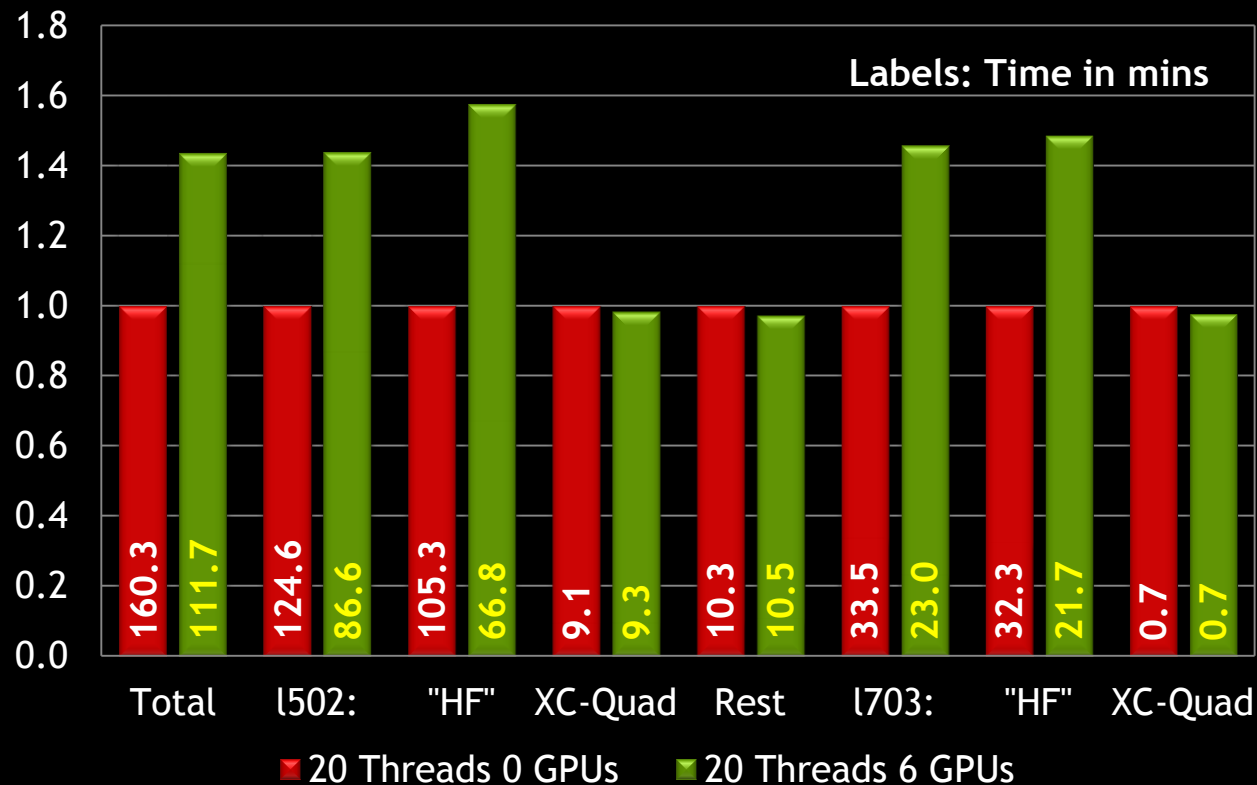


EARLY PERFORMANCE RESULTS (DIRECT SCF)



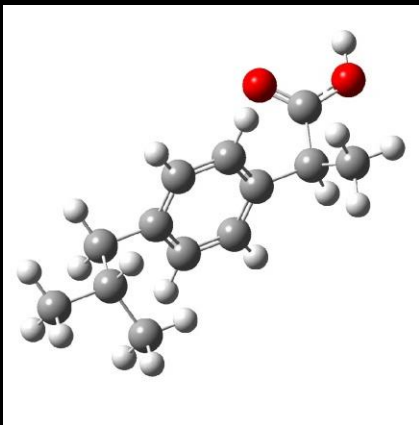
Method	rB3LYP
No. of Atoms	168
Basis Set	6-31G(3df,3p)
No. of Basis Funcs	3 642
No. of Cycles	17

Valinomycin Force Calculation
Speed Ups Relative to CPU-Only Full Node



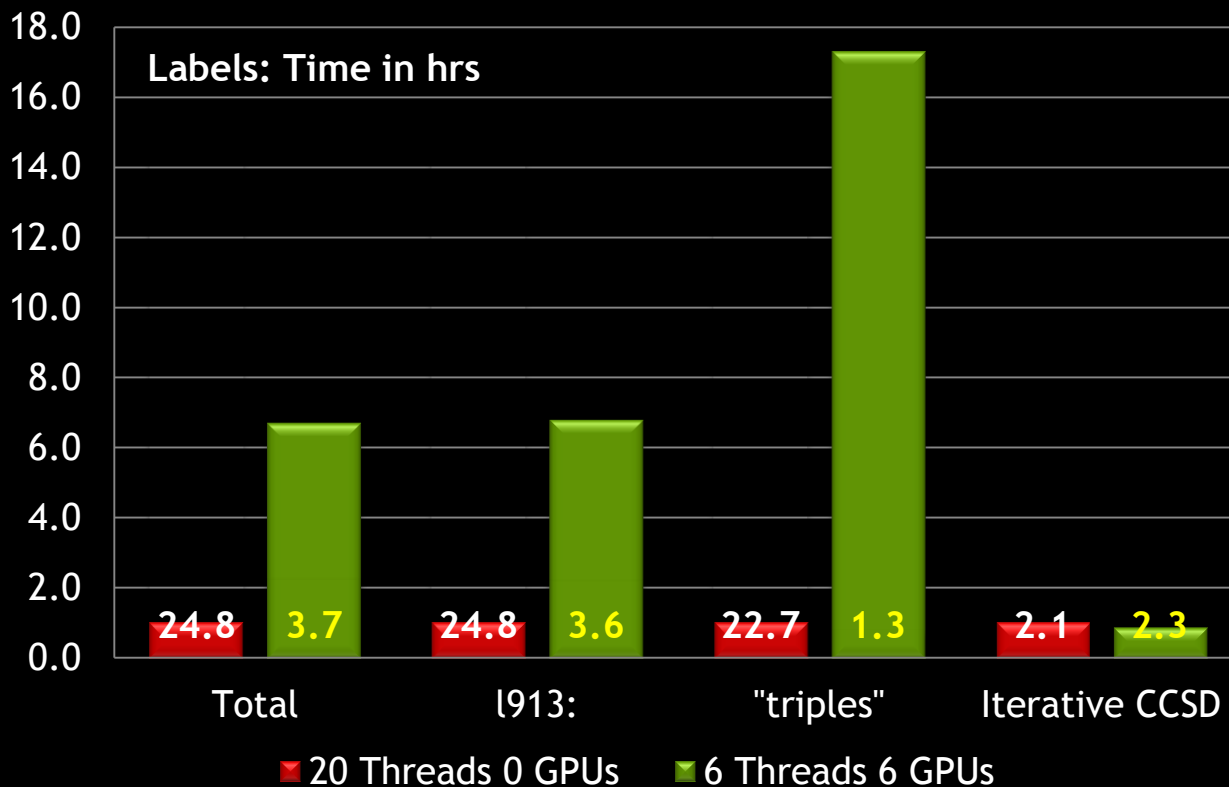
System: 2 Sockets E5-2690 V2 (2x 10 Cores @ 3.0 GHz); 128 GB RAM (DD3-1600); Used 108 GB
 GPUs: 6 Tesla K40m (15 SMPs @ 875 MHz); 12 GB Global Memory

EARLY PERFORMANCE RESULTS (CCSD(T))



Method	CCSD(t)
No. of Atoms	33
Basis Set	6-31G(d,p)
No. of Basis Funcs	315
No. Occ Orbitals	41
No. Virt Orbitals	259
No. of Cycles	15
No. CCSD iters	16

Ibuprofen CCSD(t) Calculation
Speed Ups Relative to CPU-Only Full Node



System: 2 Sockets E5-2690 V2 (2x 10 Cores @ 3.0 GHz); 128 GB RAM (DD3-1600); Used 108 GB
 GPUs: 6 Tesla K40m (15 SMPs @ 875 MHz); 12 GB Global Memory

CLOSING REMARKS

- Significant progress has been made in creating a framework that keeps an unified code structure for GPU enabled Gaussian
- There is room for performance improvement in the Direct SCF work
- The (t) correction performance looks promising
- Further work: Continue working towards a “product” quality version of Gaussian to be released to customers
 - Continue unification of the code base
 - Tackle non-default paths of the currently enabled code
 - Expand enabling of other Gaussian functionality (2nd Derivatives, XC-quadrature, TDDFT, MP2, etc.)
 - Performance tuning



GAMESS

GAMESS Partnership Overview

- **Mark Gordon and Andrey Asadchev, key developers of GAMESS, in collaboration with NVIDIA. Mark Gordon is a recipient of a NVIDIA Professor Partnership Award.**
- **Quantum Chemistry one of major consumers of CPU cycles at national supercomputer centers**
- **NVIDIA developer resources fully allocated to GAMESS code**

We like to push the envelope as much as we can in the direction of highly scalable efficient codes. GPU technology seems like a good way to achieve this goal. Also, since we are associated with a DOE Laboratory, energy efficiency is important, and this is another reason to explore quantum chemistry on GPUs.

”

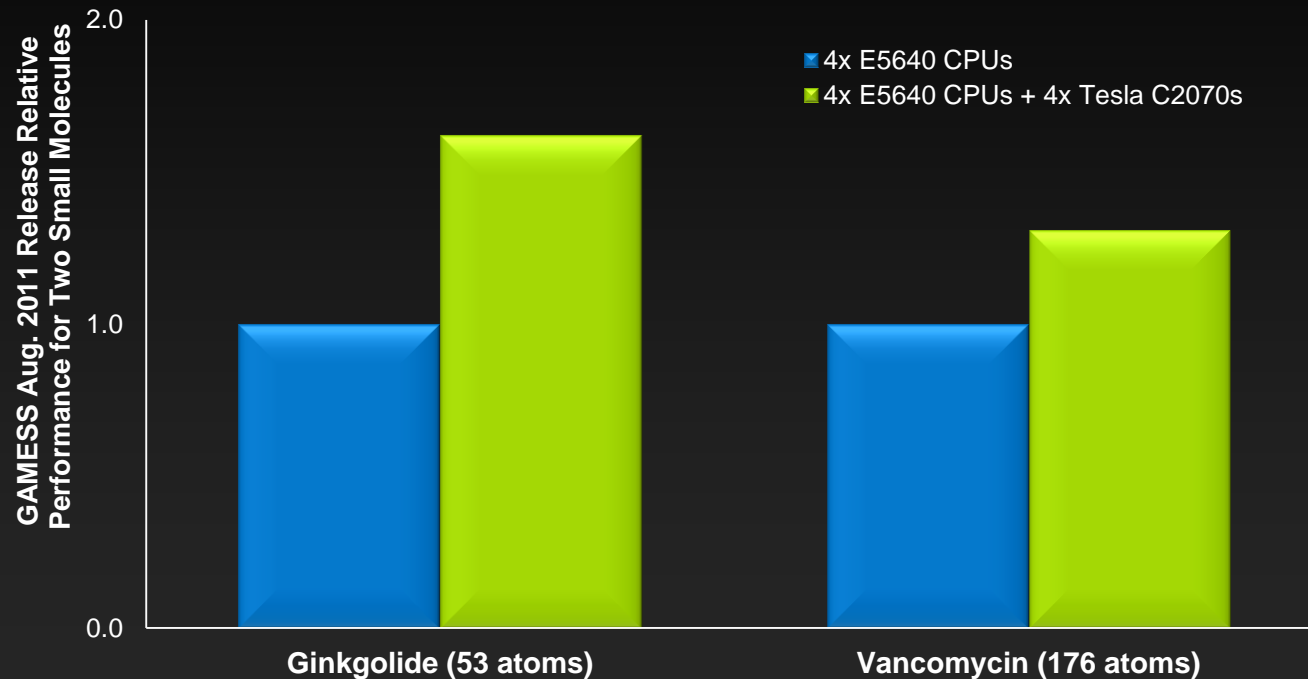
Prof. Mark Gordon

*Distinguished Professor, Department of Chemistry, Iowa State University and
Director, Applied Mathematical Sciences Program, AMES Laboratory*

GAMESS August 2011 GPU Performance



- First GPU supported GAMESS release via "libqc", a library for fast quantum chemistry on multiple NVIDIA GPUs in multiple nodes, with CUDA software
- 2e- AO integrals and their assembly into a closed shell Fock matrix

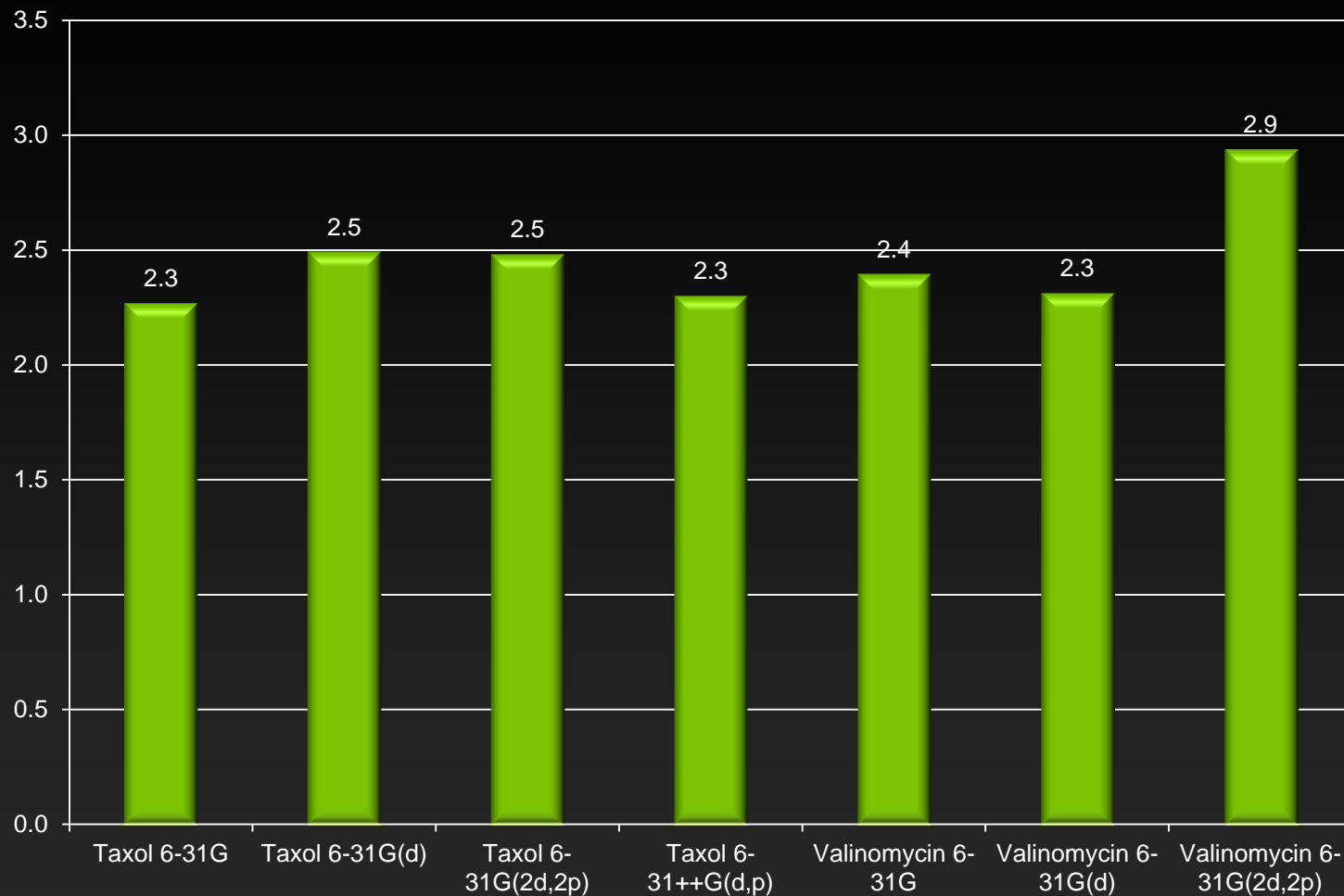


Upcoming GAMESS Q4 2013 Release

- **Multi-nodes with multi-GPUs supported**
- **Rys Quadrature**
- **Hartree-Fock**
 - 8 CPU cores: 8 CPU cores + M2070 yields 2.3-2.9x speedup. See 2012 publication
- **Møller–Plesset perturbation theory (MP2):**
 - Paper published
- **Coupled Cluster SD(T): CCSD code completed, (T) in progress**

GAMESS - New Multithreaded Hybrid CPU/GPU Approach to H-F

Hartree-Fock GPU Speedups*



Adding 1x 2070 GPU speeds up computations by 2.3x to 2.9x

■ Speedup

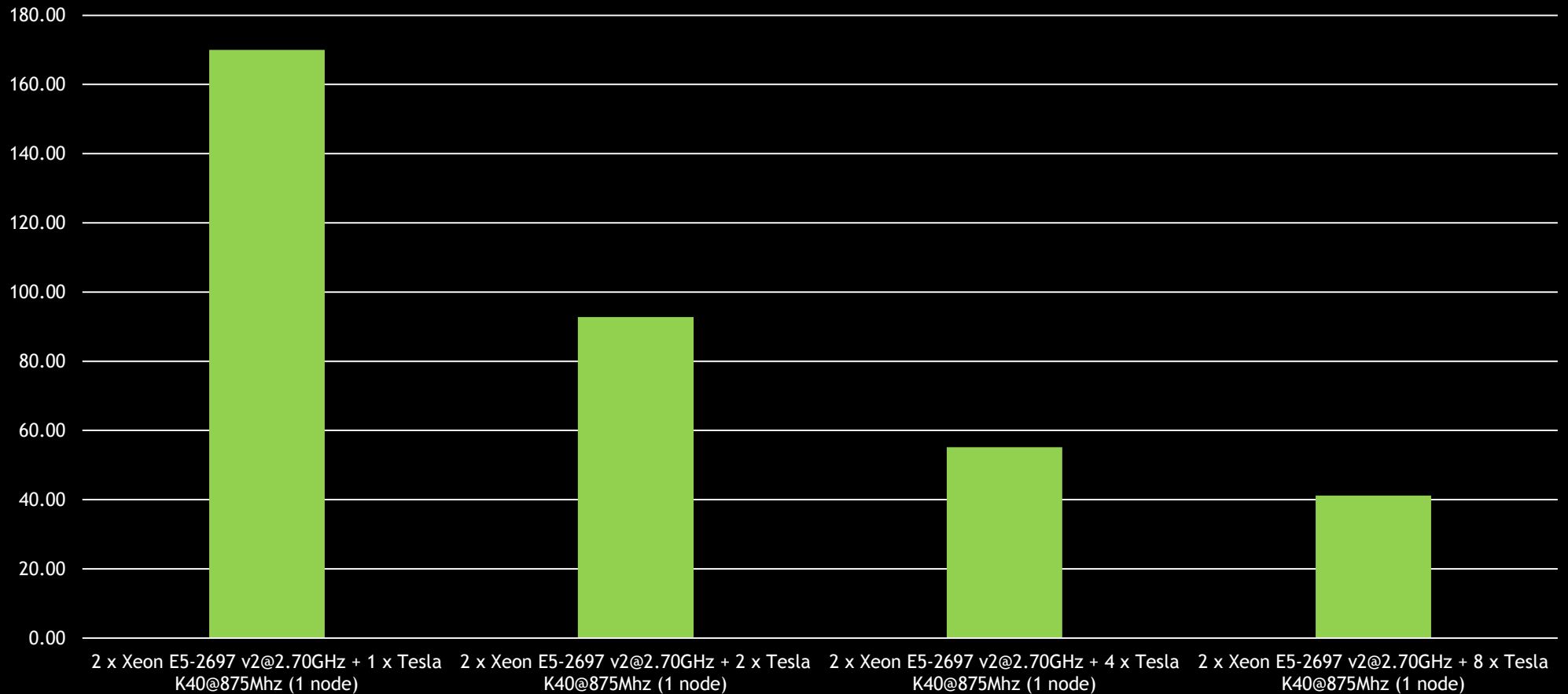
* A. Asadchev, M.S. Gordon, "New Multithreaded Hybrid CPU/GPU Approach to Hartree-Fock," Journal of Chemical Theory and Computation (2012)



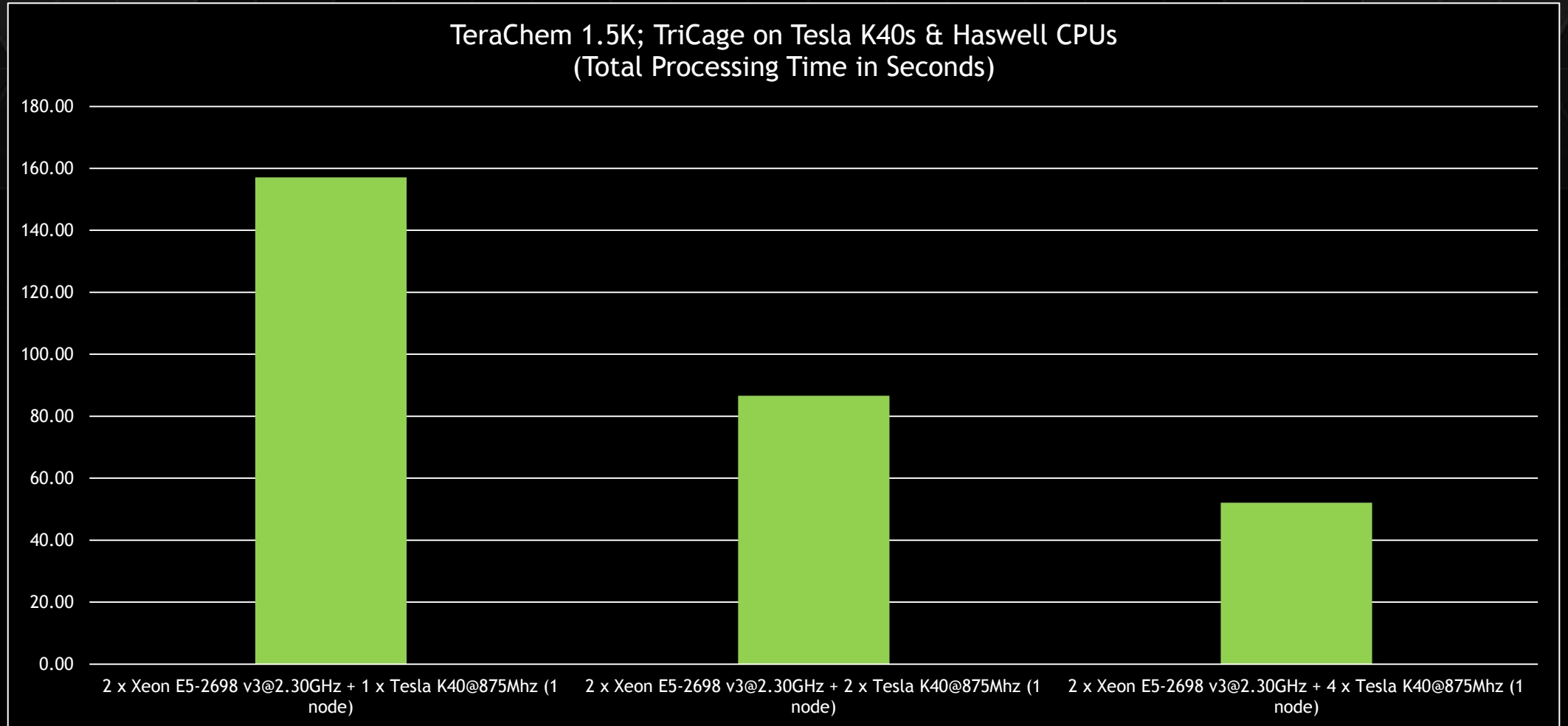
TeraChem

TERACHEM 1.5K; TRPCAGE ON TESLA K40S

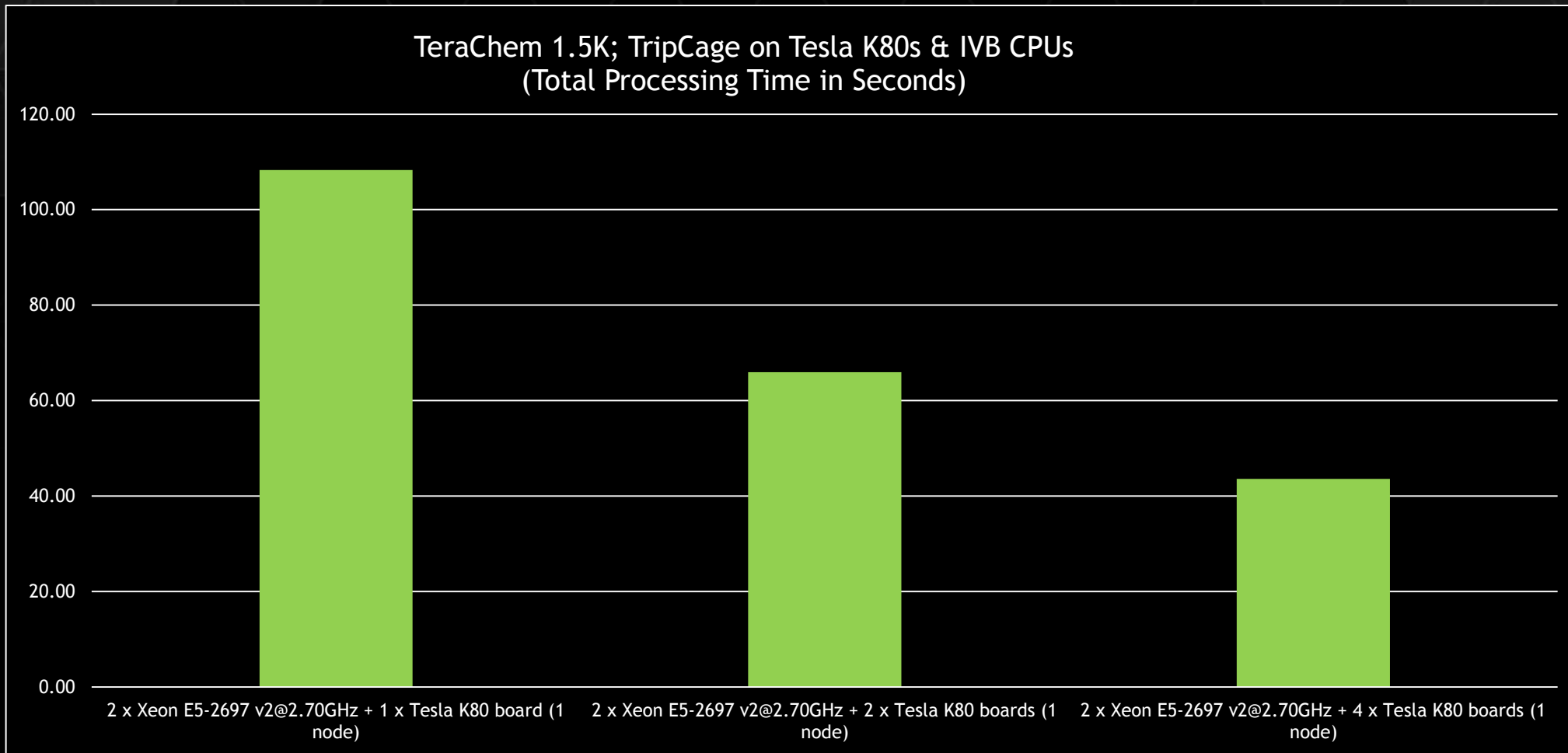
TeraChem 1.5K; TriCage on Tesla K40s & IVB CPUs
(Total Processing Time in Seconds)



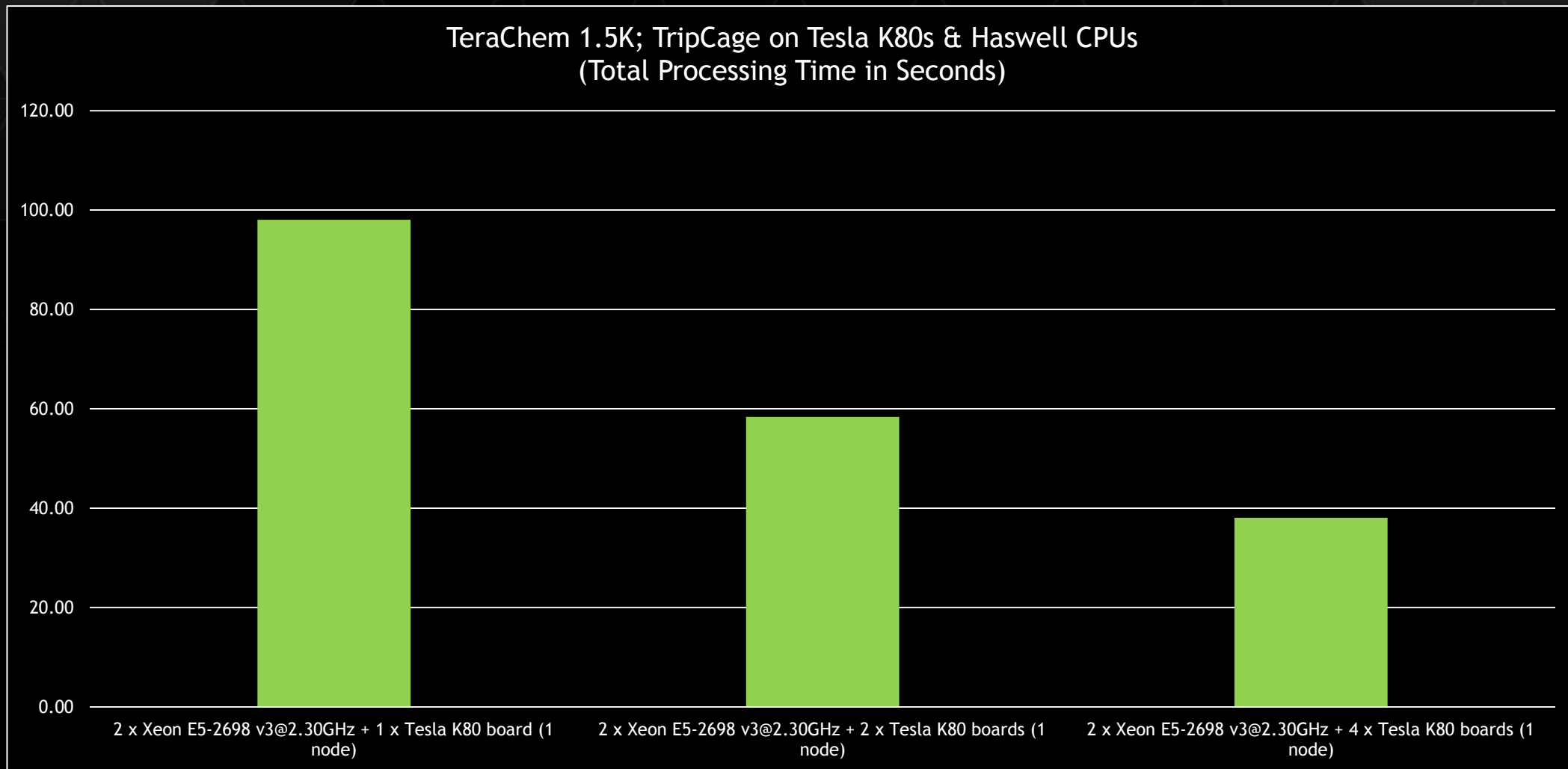
TERACHEM 1.5K; TRPCAGE ON TESLA K40S & HASWELL CPUS



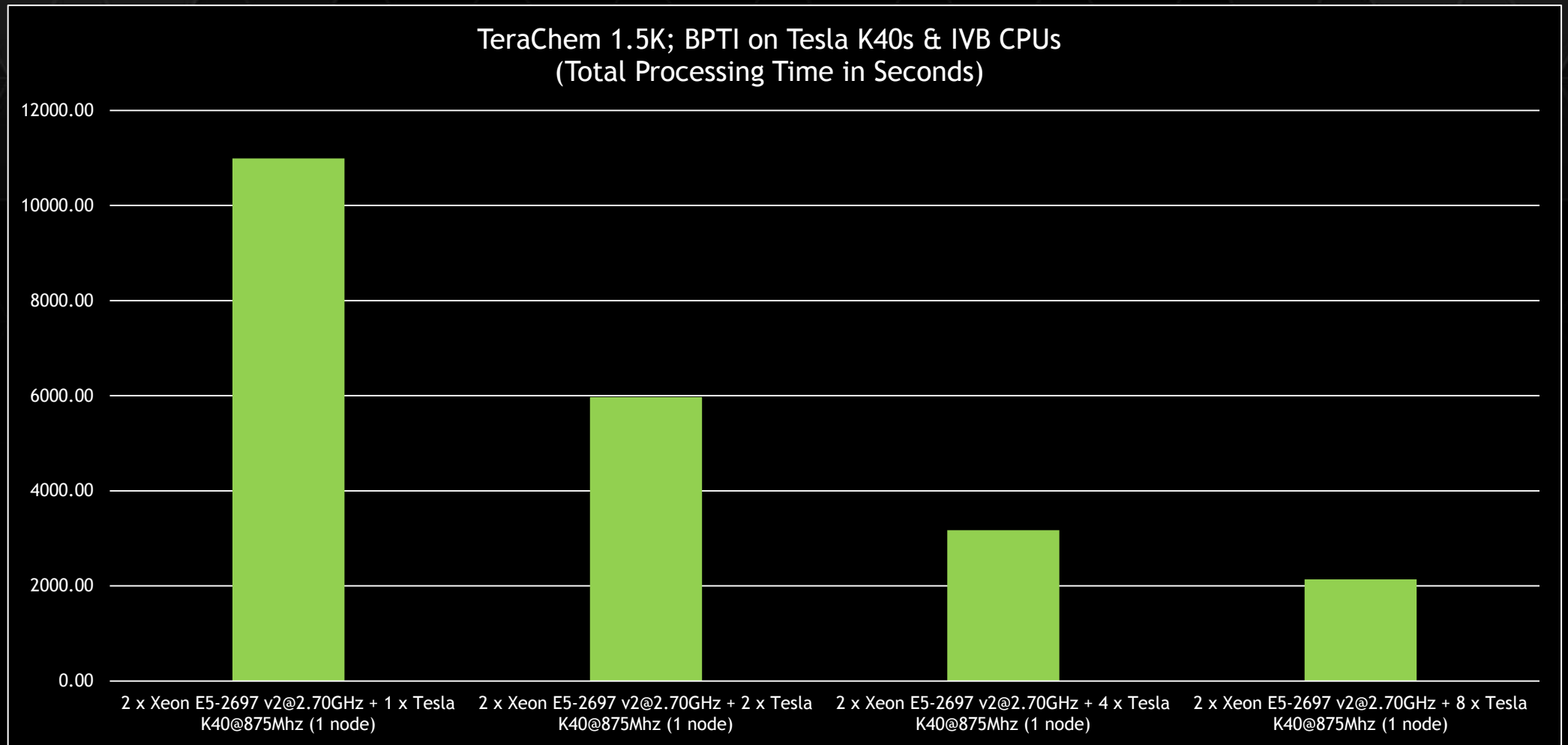
TERACHEM 1.5K; TRPCAGE ON TESLA K80S & IVB CPUS



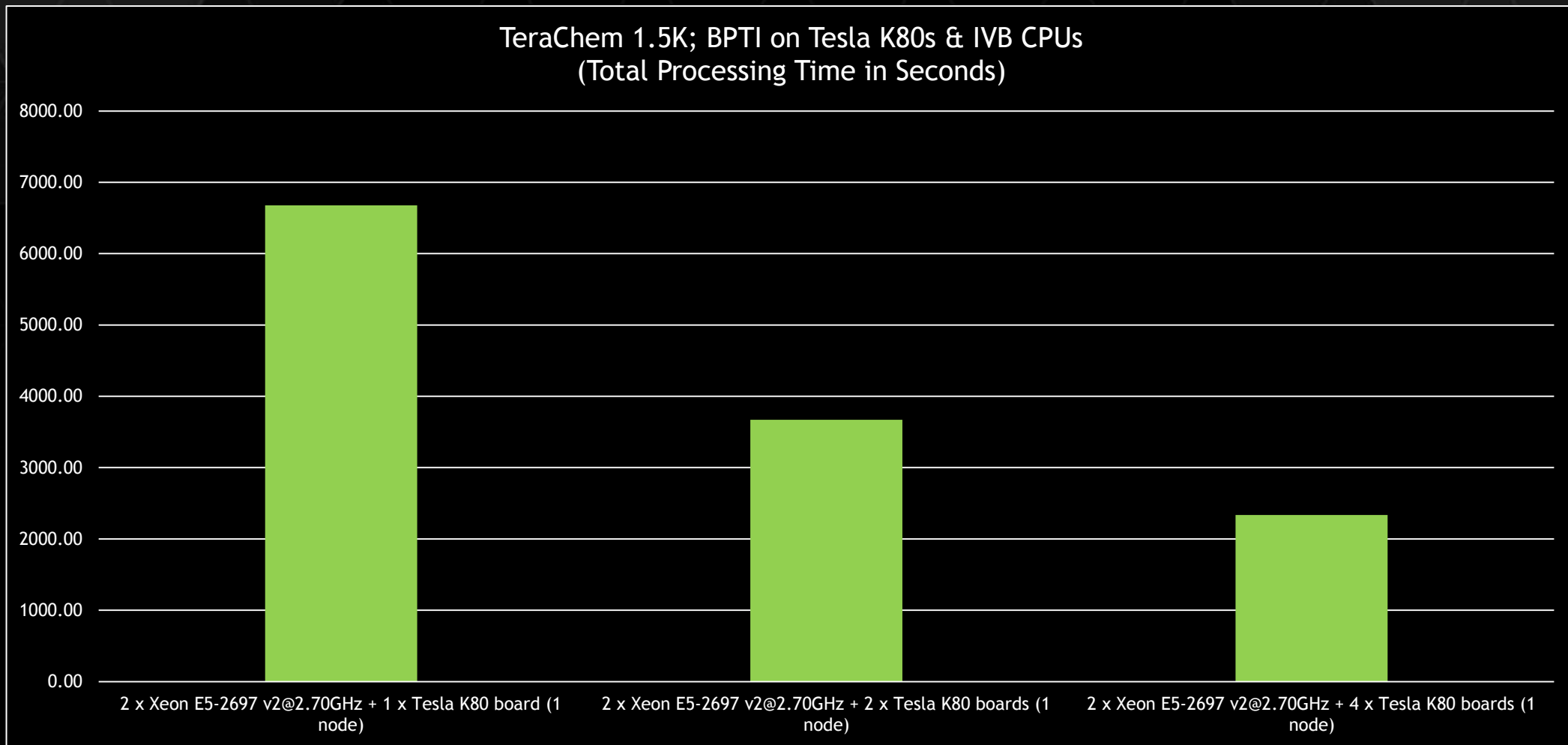
TERACHEM 1.5K; TRIPCAGE ON TESLA K80S & HASWELL CPUS



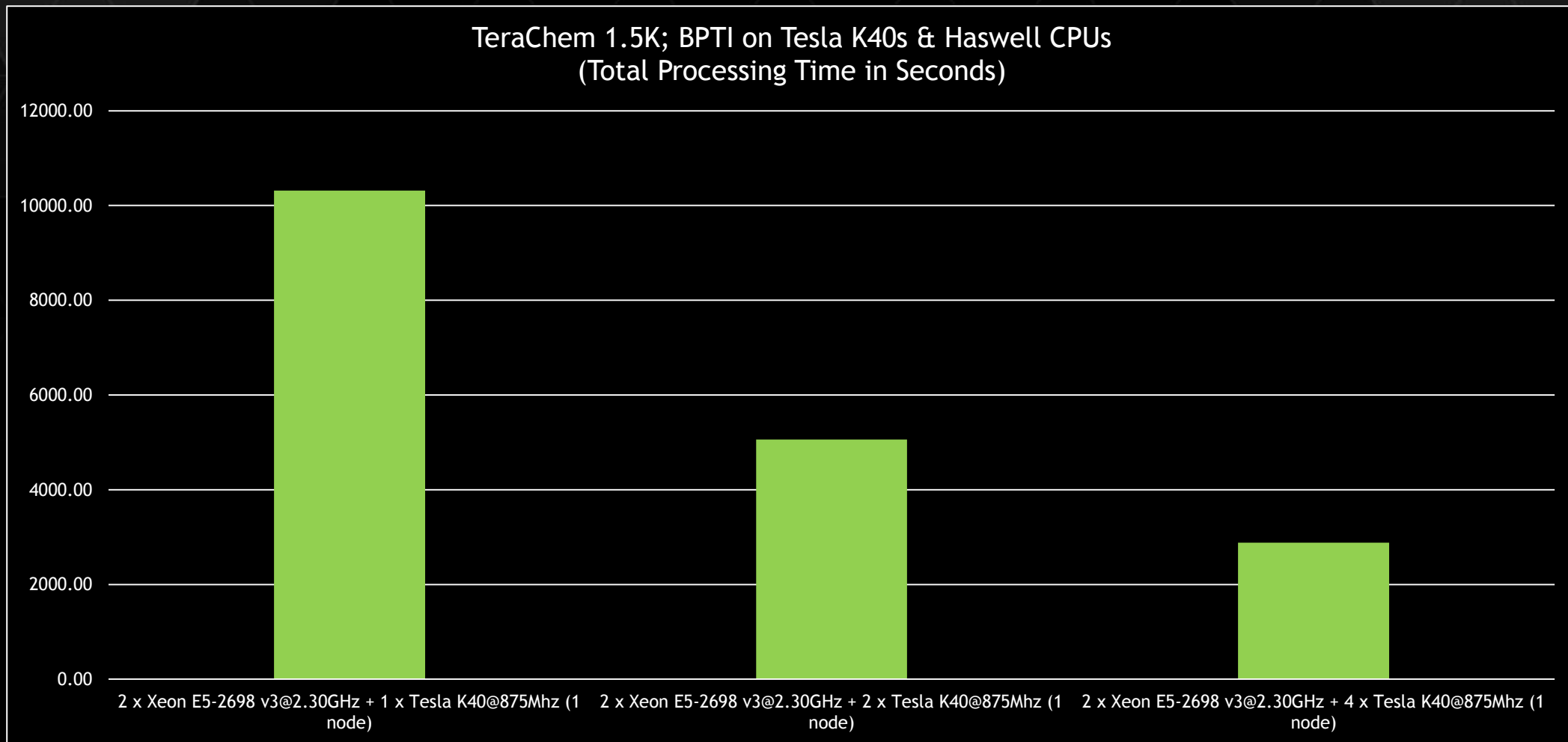
TERACHEM 1.5K; BPTI ON TESLA K40S & IVB CPUS



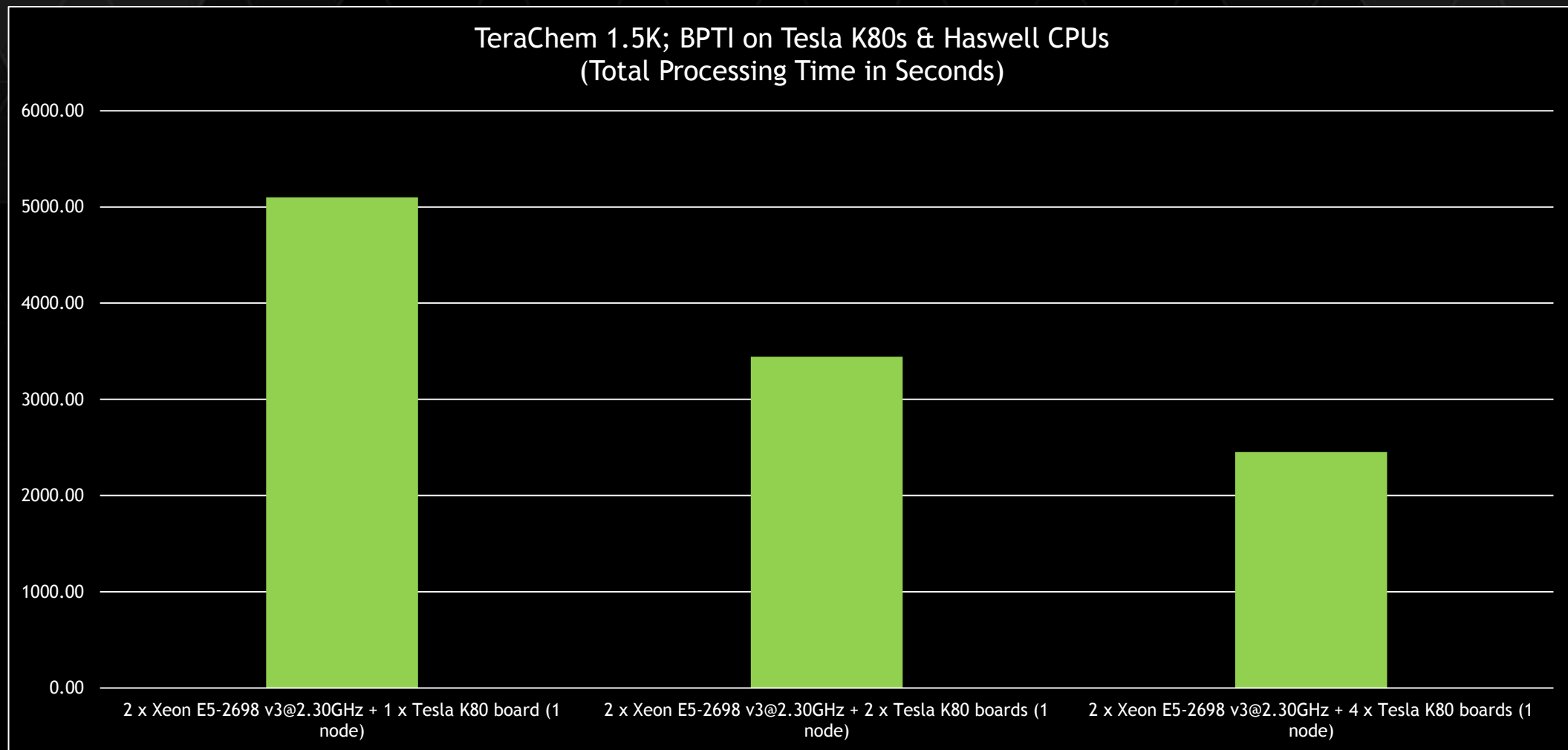
TERACHEM 1.5K; BPTI ON TESLA K80S & IVB CPUS



TERACHEM 1.5K; BPTI ON TESLA K40S & HASWELL CPUS



TERACHEM 1.5K; BPTI ON TESLA K80S & HASWELL CPUS



Benefits of GPU Accelerated Computing



- Faster than CPU only systems in all tests
- Large performance boost with marginal price increase
- Energy usage cut by more than half
- GPUs scale well within a node and over multiple nodes
- K20 GPU is our fastest and lowest power high performance GPU yet

Try GPU accelerated TeraChem for free – www.nvidia.com/GPUTestDrive

Test Drive K80 GPUs!

Experience The Acceleration

www.nvidia.com/GPUTestDrive



Run Computational Chemistry Apps on Tesla K80 GPUs today



Try Preconfigured Apps:
AMBER, NAMD, GROMACS, LAMMPS, Quantum Espresso, TeraChem
Or Load Your Own



Sign up for FREE GPU Test Drive on remotely hosted clusters

