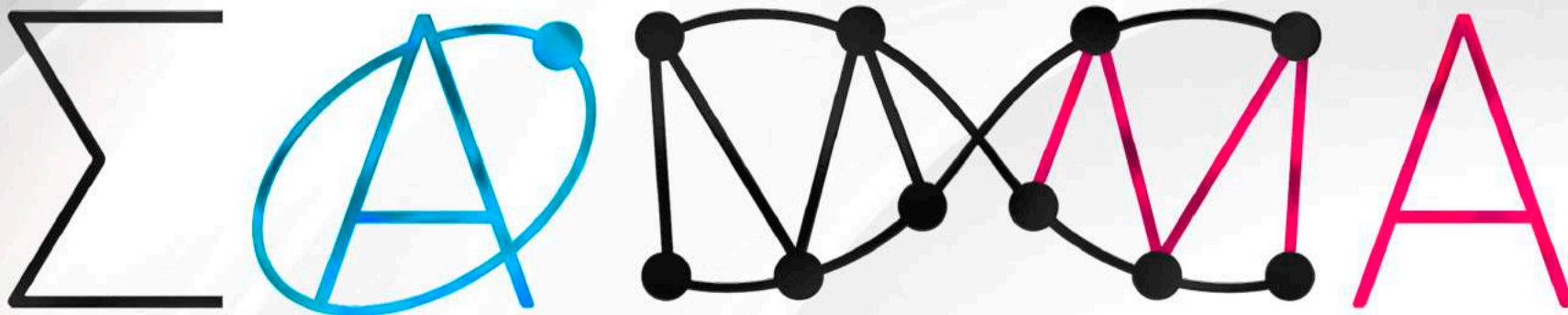# The Desmos supercomputer
# for computational materials science

Vladimir Stegailov, Nikolay Kondratyuk, Grigory Smirnov,
Alexei Timofeev

National Research University Higher School of Economics

Supercomputer Atomistic Modelling and Multi-scale Analysis

http://samma.hse.ru

# Outline of the talk

- Angara network and its developer JSC NICEVT

- Building process of the Desmos supercomputer

- Benchmarking of MD calculations

- Statistical data of Desmos deployment

- Comparison with a new hybrid supercomputer in HSE

- VASP benchmarks

- A next Angara-based HPC system
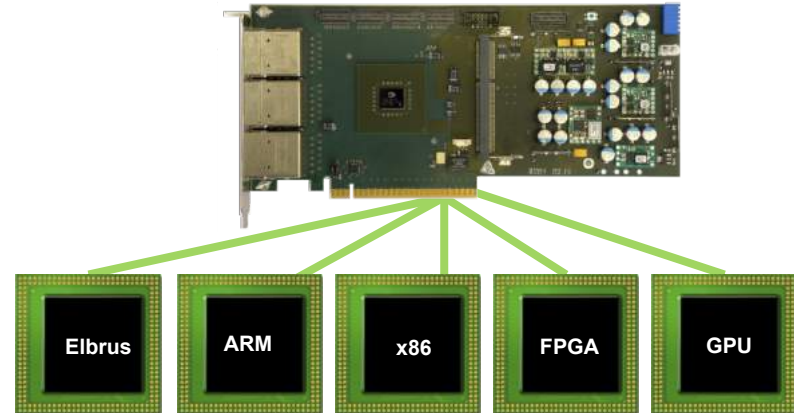
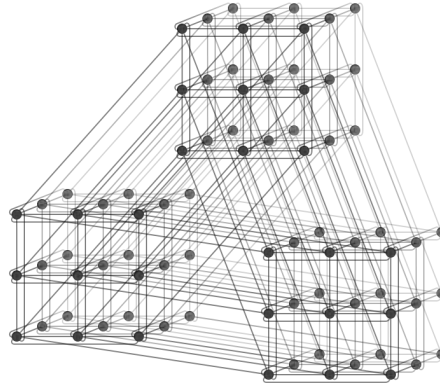# ANGARA NETWORK
# AND ITS DEVELOPER JSC NICEVT

# JSC NICEVT:
# from ES EVM computers to Angara network

# Angara interconnect

- Network topology: 1D..4D-torus
- ASIC-based network card
- Up to 8 communication channels
- Remote direct memory access (RDMA)
- Multi-core CPU support
- Adaptive packet transfer
- MPI ping-pong latency: 0,85 µs
- Single hop latency: 129 ns
- Scaling: up to 32K nodes
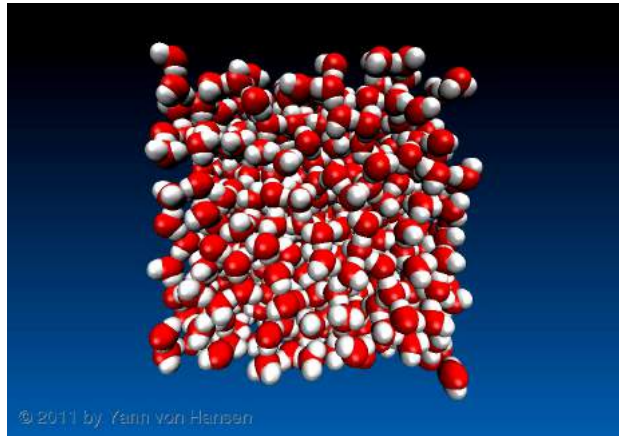- Power consumption: up to 20 W
- Various physical transmission media

| Elbrus | ARM | x86 | FPGA | GPU |

# Angara interconnect

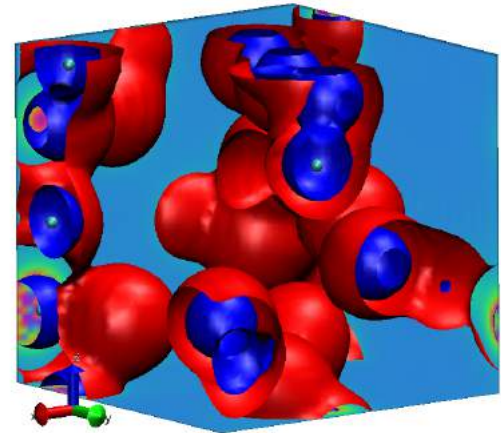| Interconnect | Mellanox IB FDR 4x | Angara | Cray Aries | Mellanox IB EDR 4x | Intel OmniPath |
|---|---|---|---|---|---|
| Year | 2011 | 2013 | 2012 | 2015 | 2015 |
| TOP500 | 13 | – | 5 | 20 | 6 |
| Topology | fat tree / kD-torus | 4D-torus | dragonfly | fat tree / kD-torus | fat tree |
| MPI latency, µs | 1 | 0.85 | 1.3 | 0.92 | 0.9-1.0 |
| Single hop latency, ns | – / 250 | 129 | 100 | n/a | n/a |

# BUILDING PROCESS
# OF THE SUPERCOMPUTER DESMOS

# The main goals of the project

- To make a supercomputer that is effective for classical molecular dynamics and usable of ab initio MD

- To use the Angara network and get the largest possible number of nodes within a limited budget limit

- **Access to novel hardware is very beneficial for HPC training**



LAMMPS, GROMACS



VASP, ABINIT, CP2K

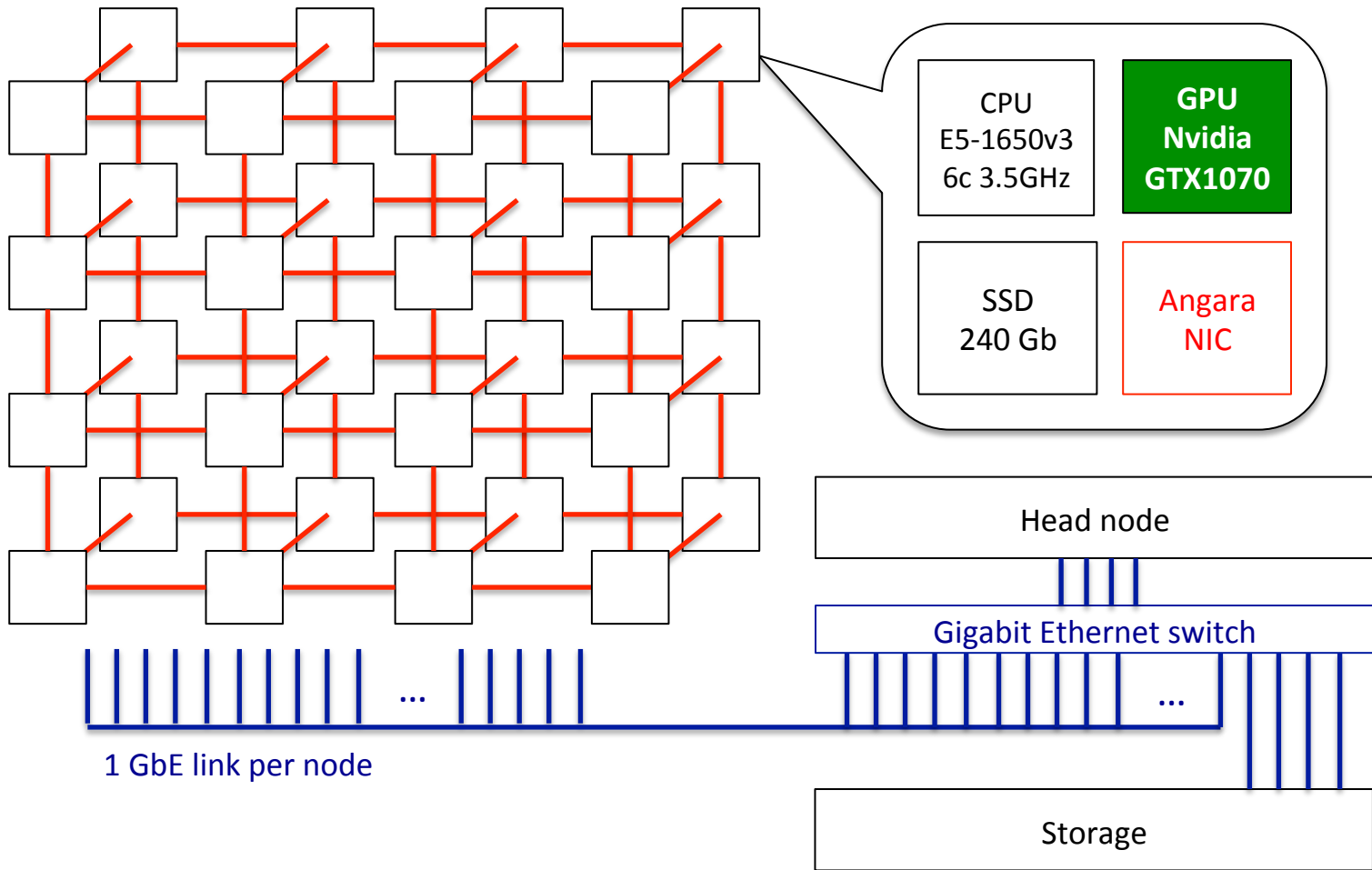Intel Xeon E5-1650v3



Nvidia GeForce 1070

**Desmos supercomputer**

CPU
E5-1650v3
6c 3.5GHz

**GPU Nvidia GTX1070**

SSD
240 Gb

Angara NIC

Head node

Gigabit Ethernet switch

...

...

1 GbE link per node

Storage

**Desmos supercomputer**



CPU
E5-1650v3
6c 3.5GHz

**GPU FirePro S9150**

SSD
240 Gb

Angara NIC

Head node

Gigabit Ethernet switch

...

...

Storage

1 GbE link per node
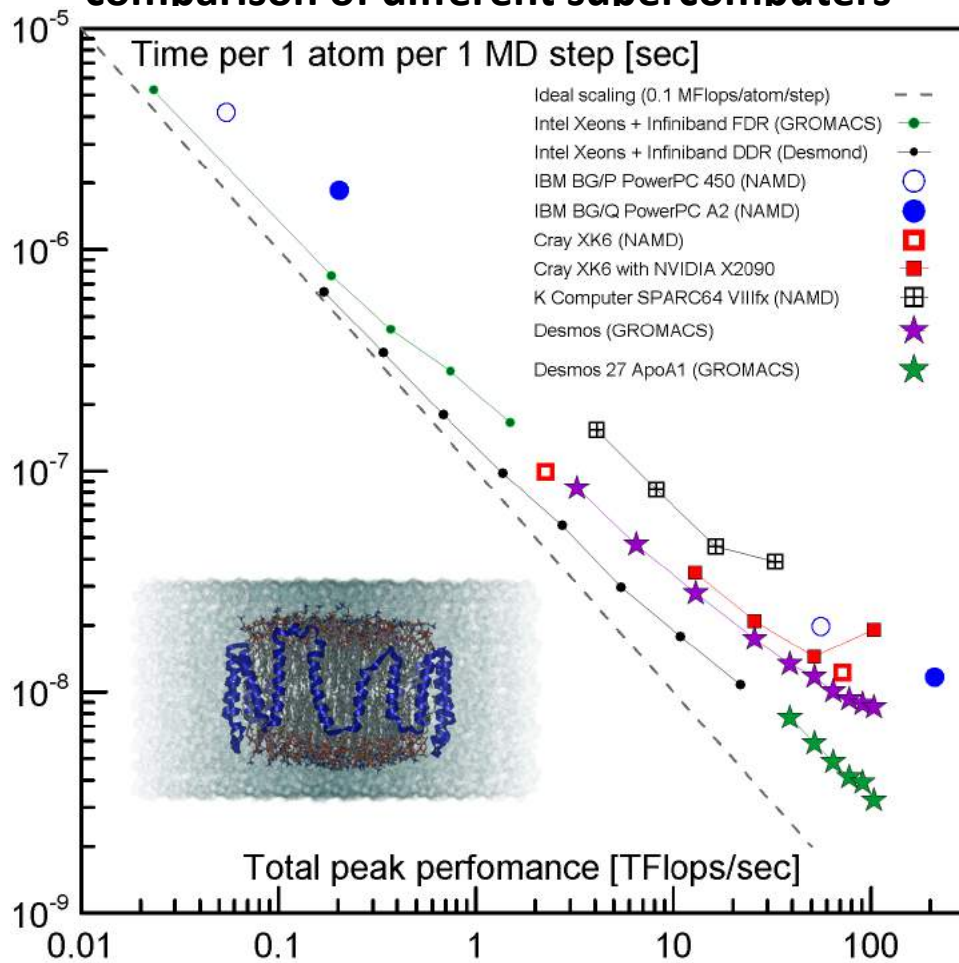
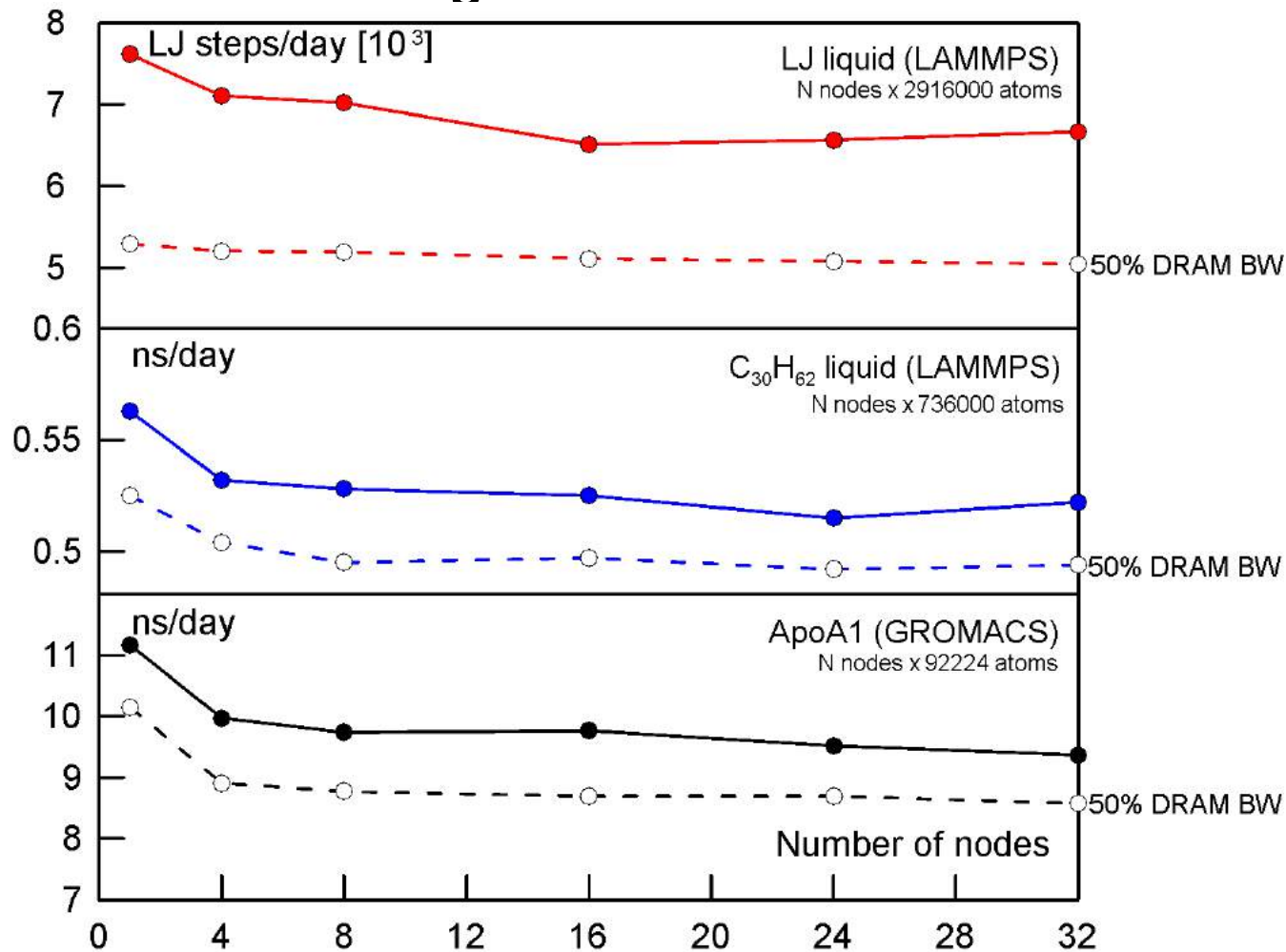# BENCHMARKING
# OF MD CALCULATIONS

## ApoA1 benchmark (protein in water, ~ 100 000 atoms): comparison of different supercomputers



Time per 1 atom per 1 MD step [sec]

Ideal scaling (0.1 MFlops/atom/step)
Intel Xeons + Infiniband FDR (GROMACS)
Intel Xeons + Infiniband DDR (Desmond)
IBM BG/P PowerPC 450 (NAMD)
IBM BG/Q PowerPC A2 (NAMD)
Cray XK6 (NAMD)
Cray XK6 with NVIDIA X2090
K Computer SPARC64 VIIIfx (NAMD)
Desmos (GROMACS)
Desmos 27 ApoA1 (GROMACS)

Total peak perfomance [TFlops/sec]

# Weak scaling of different MD models

Journal of **COMPUTATIONAL CHEMISTRY**

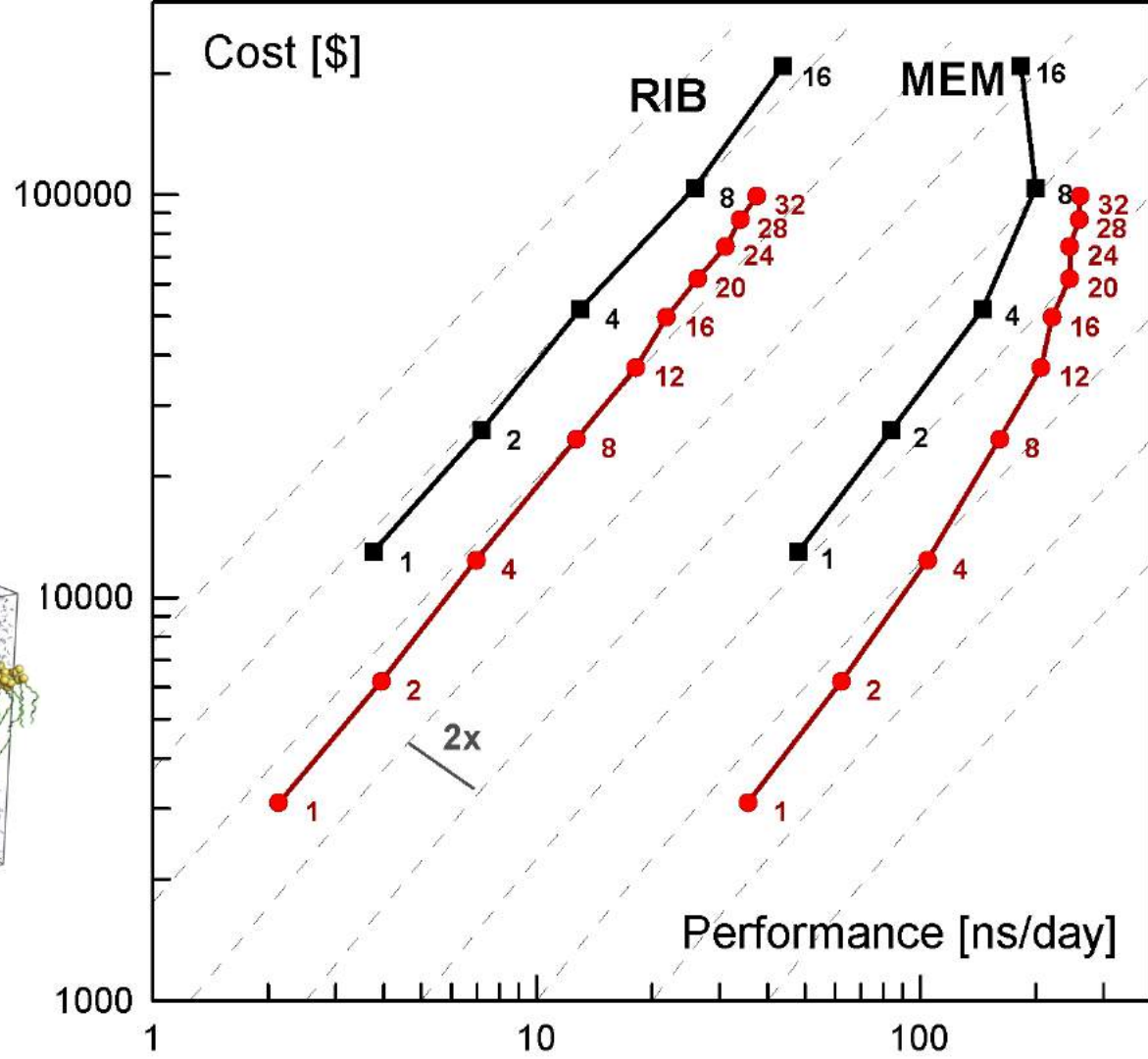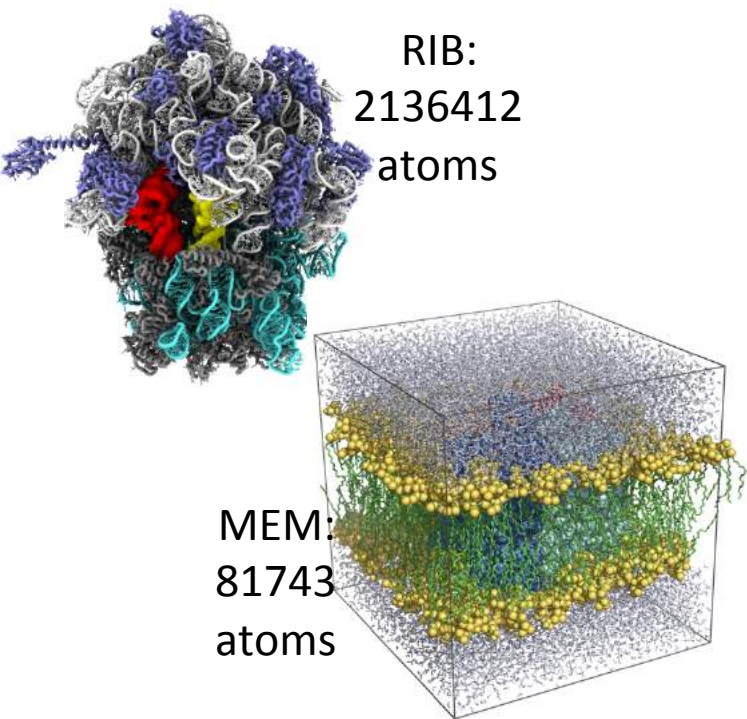# Best Bang for Your Buck: GPU Nodes for GROMACS Biomolecular Simulations

Carsten Kutzner,*[a] Szilárd Páll,[b] Martin Fechner,[a] Ansgar Esztermann,[a] Bert L. de Groot,[a] and Helmut Grubmüller[a]

The molecular dynamics simulation package GROMACS runs efficiently on a wide variety of hardware from commodity workstations to high performance computing clusters. Hardware features are well-exploited with a combination of single instruction multiple data, multithreading, and message passing interface (MPI)-based single program multiple data/multiple program multiple data parallelism while graphics processing units (GPUs) can be used as accelerators to compute interactions off-loaded from the CPU. Here, we evaluate which hardware produces trajectories with GROMACS 4.6 or 5.0 in the most economical way. We have assembled and benchmarked compute nodes with various CPU/GPU combinations to identify optimal compositions in terms of raw trajectory production rate, performance-to-price ratio, energy efficiency, and several other criteria. Although hardware prices are naturally subject to trends and fluctuations, general tendencies are clearly visible. Adding any type of GPU significantly boosts a node's simulation performance. For inex-

pensive consumer-class GPUs this improvement equally reflects in the performance-to-price ratio. Although memory issues in consumer-class GPUs could pass unnoticed as these cards do not support error checking and correction memory, unreliable GPUs can be sorted out with memory checking tools. Apart from the obvious determinants for cost-efficiency like hardware expenses and raw performance, the energy consumption of a node is a major cost factor. Over the typical hardware lifetime until replacement of a few years, the costs for electrical power and cooling can become larger than the costs of the hardware itself. Taking that into account, nodes with a well-balanced ratio of CPU and consumer-class GPU resources produce the maximum amount of GROMACS trajectory over their lifetime. © 2015 The Authors. Journal of Computational Chemistry Published by Wiley Periodicals, Inc.
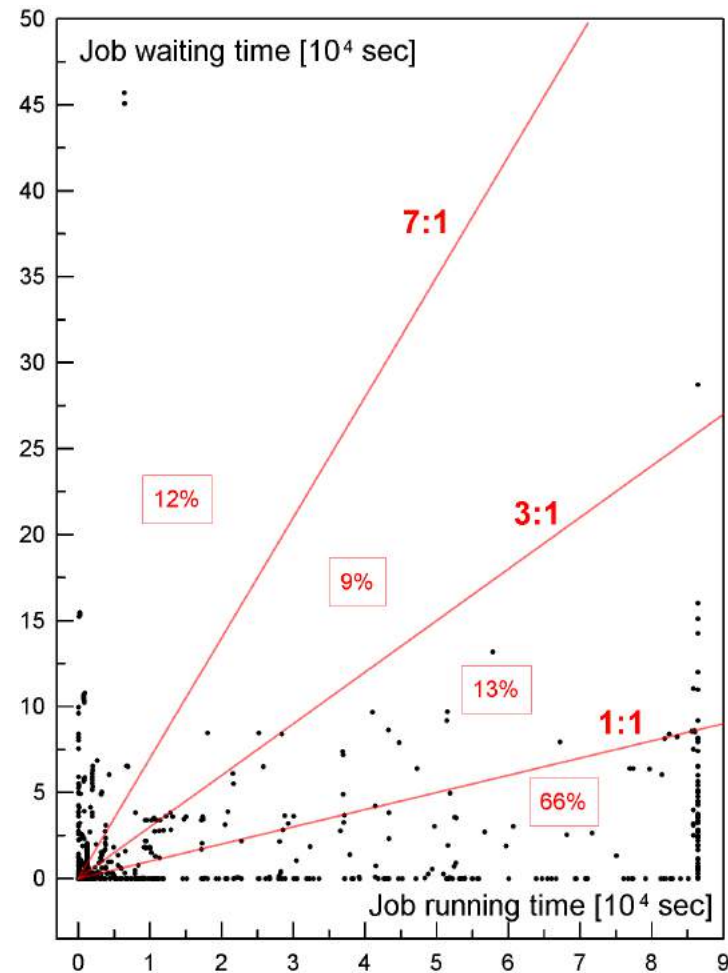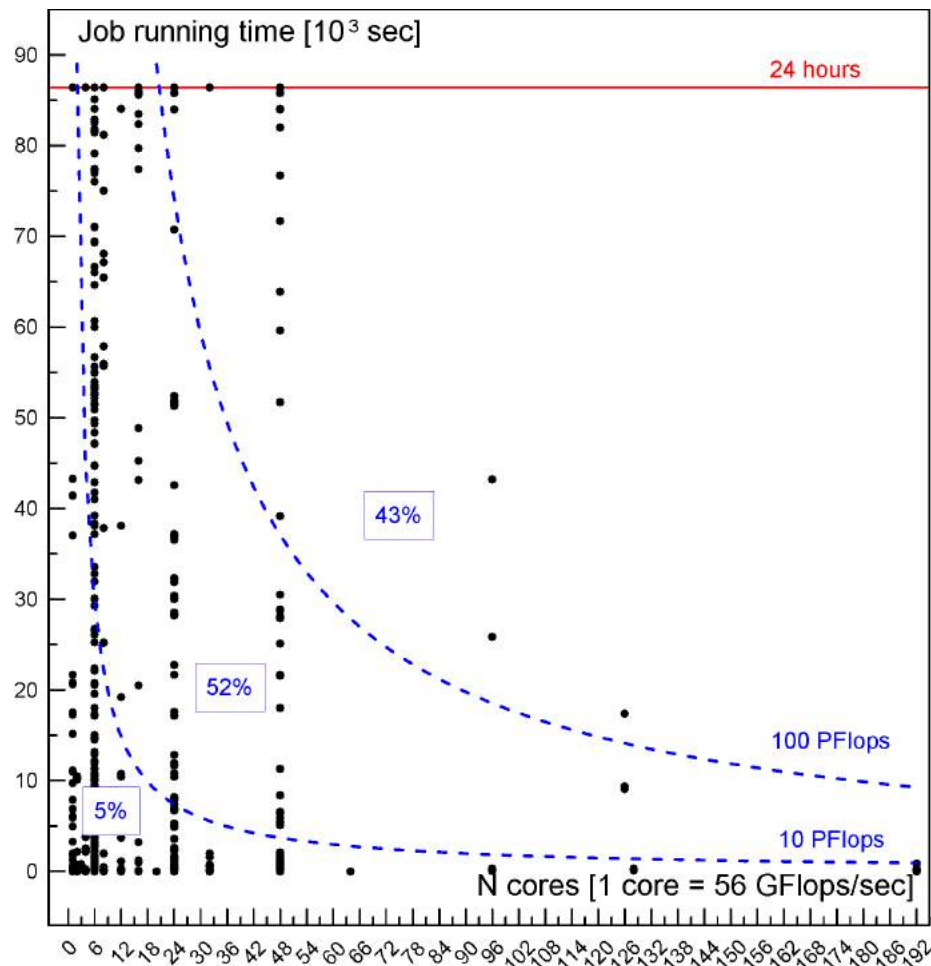
Comparison with
the IRUS17 supercomputer
2 x Intel Xeon E5-2699 v4 per node
with Intel Omni-Path

RIB:
2136412
atoms

MEM:
81743
atoms

Prices from thinkmate.com
(November 2017)

# STATISTICAL DATA
# OF  DESMOS DEPLOYMENT

# Angara interconnect makes GPU-based Desmos supercomputer an efficient tool for molecular dynamics calculations
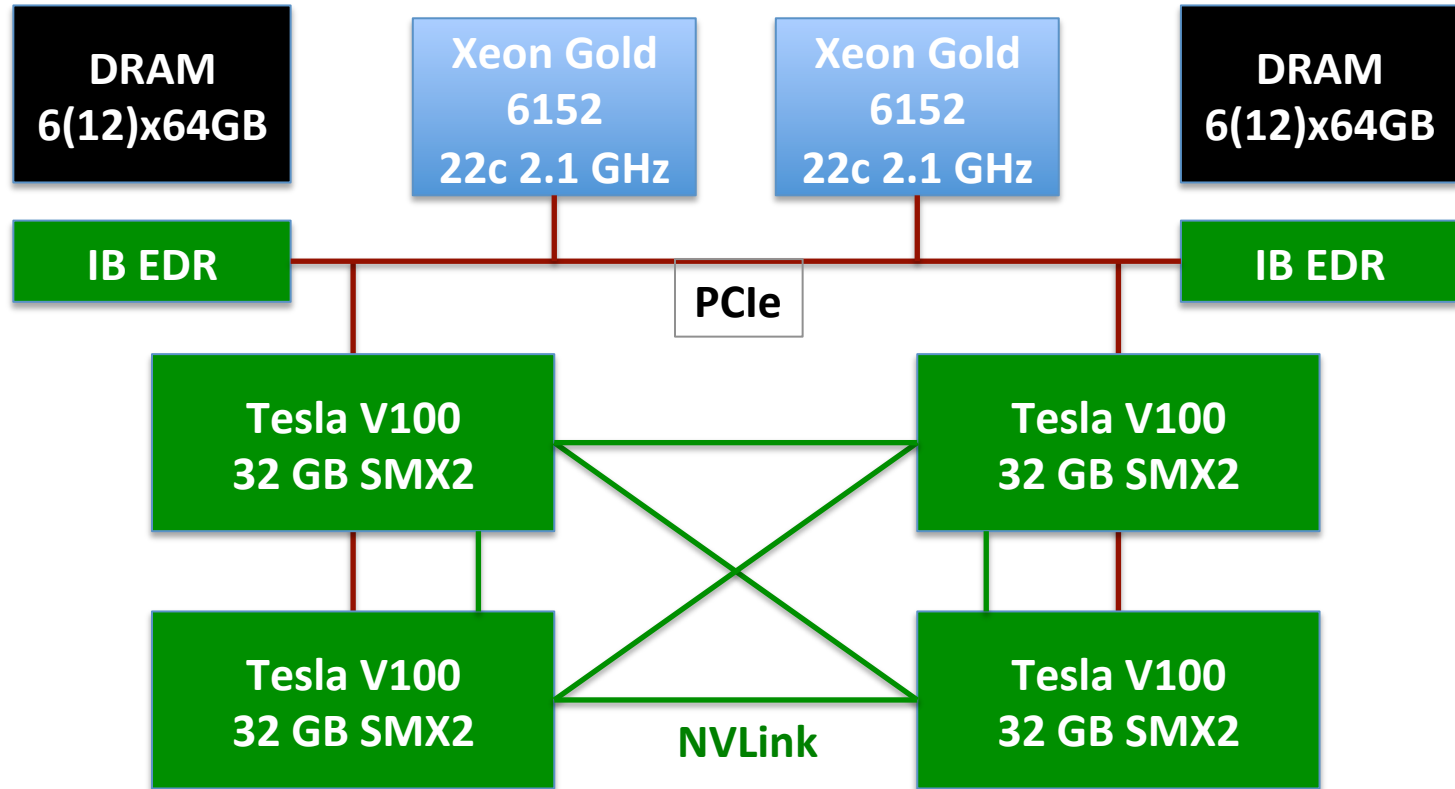
Vladimir Stegailov[1,2] ⑩, Ekaterina Dlinnova[2],
Timur Ismagilov[3], Mikhail Khalilov[2], Nikolay Kondratyuk[1,2],
Dmitry Makagon[3], Alexander Semenov[2,3], Alexei Simonov[3],
Grigory Smirnov[1,2] and Alexey Timofeev[1,2]

## Abstract

In this article, we describe the Desmos supercomputer that consists of 32 hybrid nodes connected by a low-latency high-bandwidth Angara interconnect with torus topology. This supercomputer is aimed at cost-effective classical molecular dynamics calculations. Desmos serves as a test bed for the Angara interconnect that supports 3-D and 4-D torus network topologies and verifies its ability to unite massively parallel programming systems speeding-up effectively message-passing interface (MPI)-based applications. We describe the Angara interconnect presenting typical MPI benchmarks. Desmos benchmarks results for GROMACS, LAMMPS, VASP and CP2K are compared with the data for other high-performance computing (HPC) systems. Also, we consider the job scheduling statistics for several months of Desmos deployment.
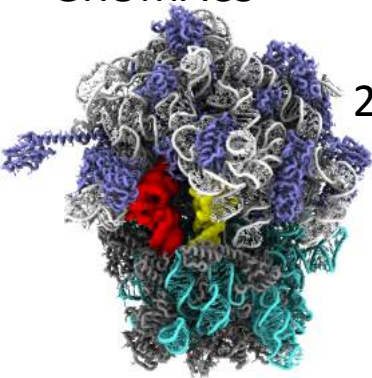
# COMPARISON WITH A NEW HYBRID SUPERCOMPUTER IN HSE

# The new supercomputer in Higher School of Economics: a hybrid system based on 26 DELL C4140 nodes
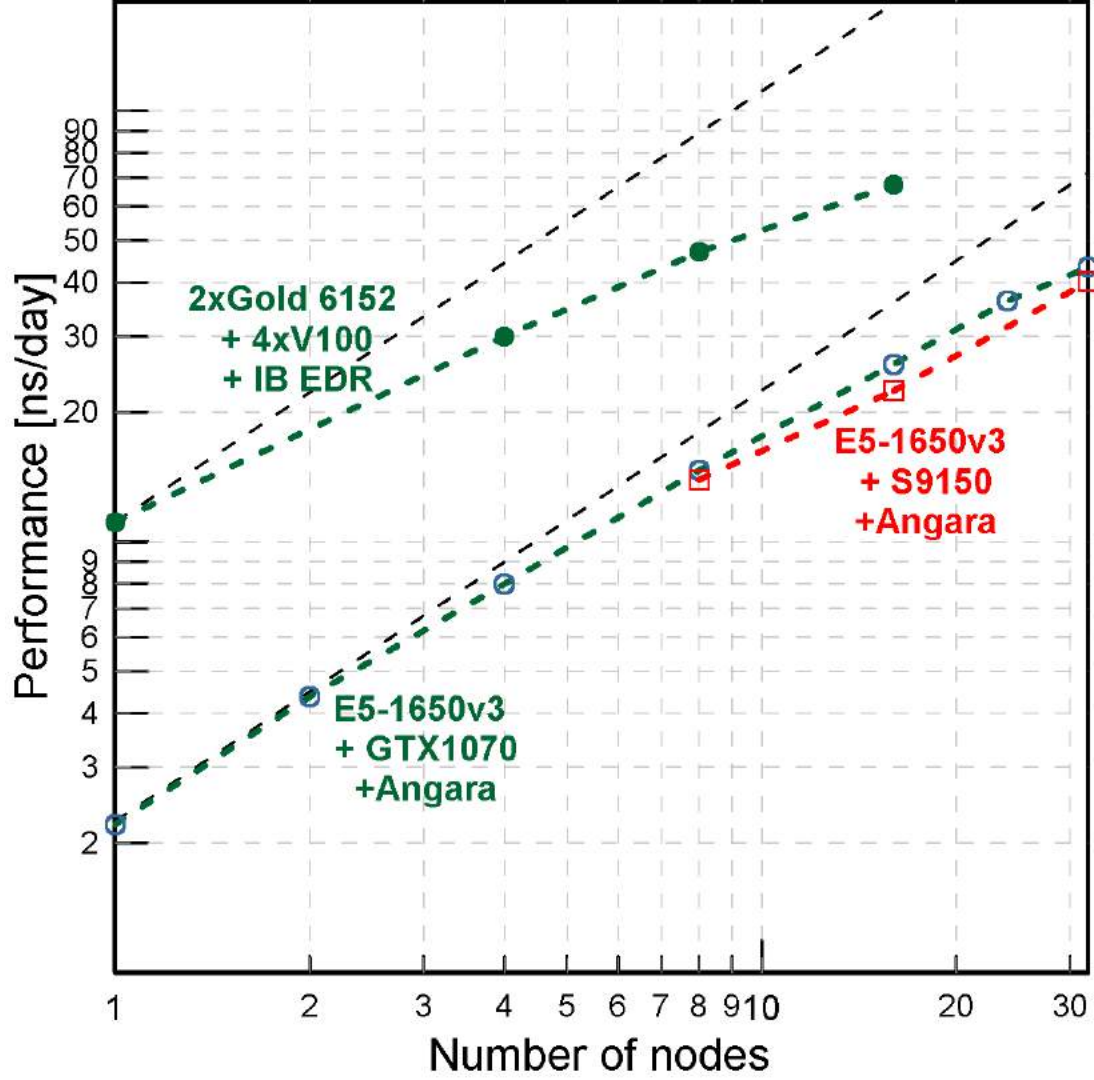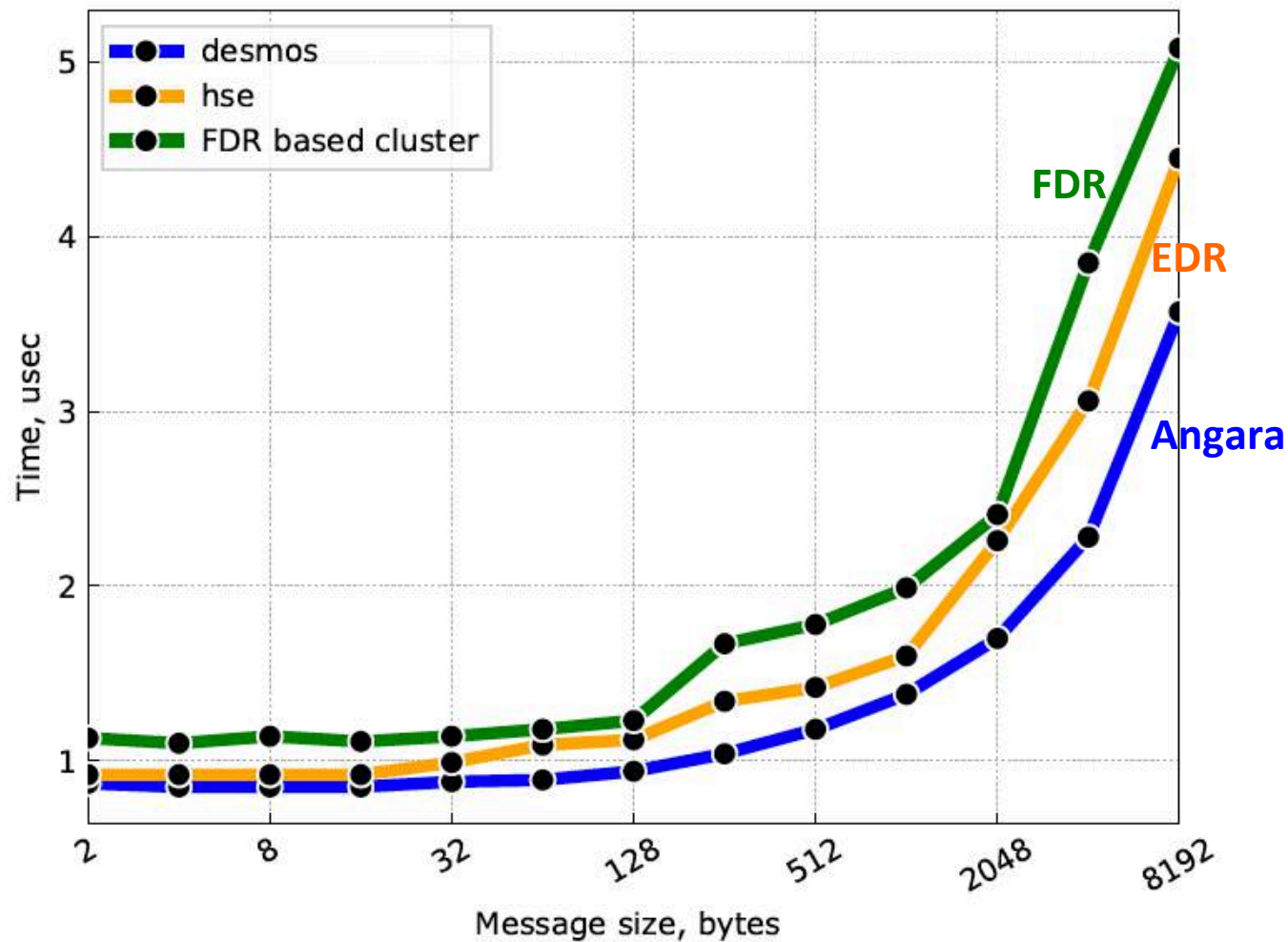
Comparison Desmos
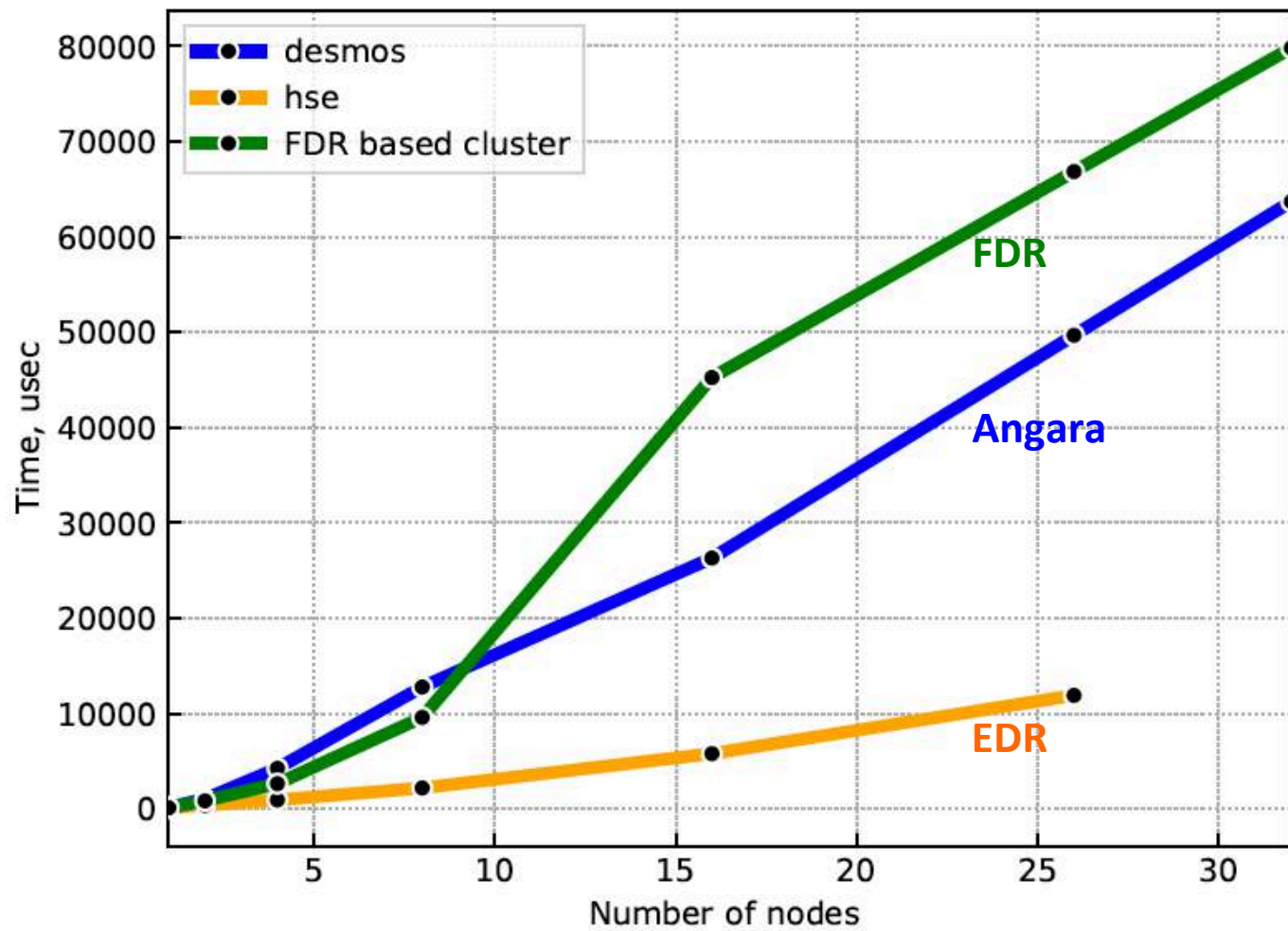with the new HSE supercomputer
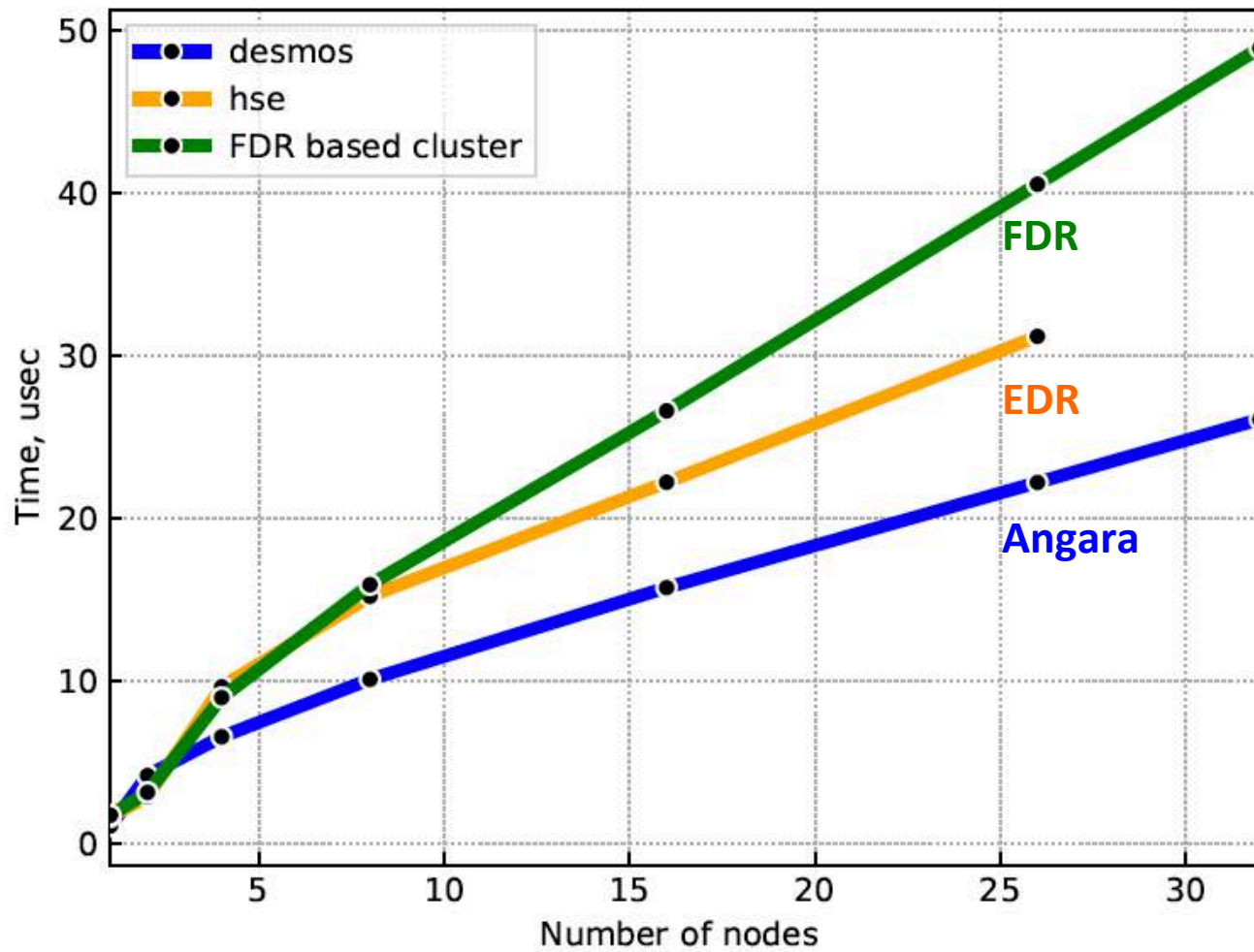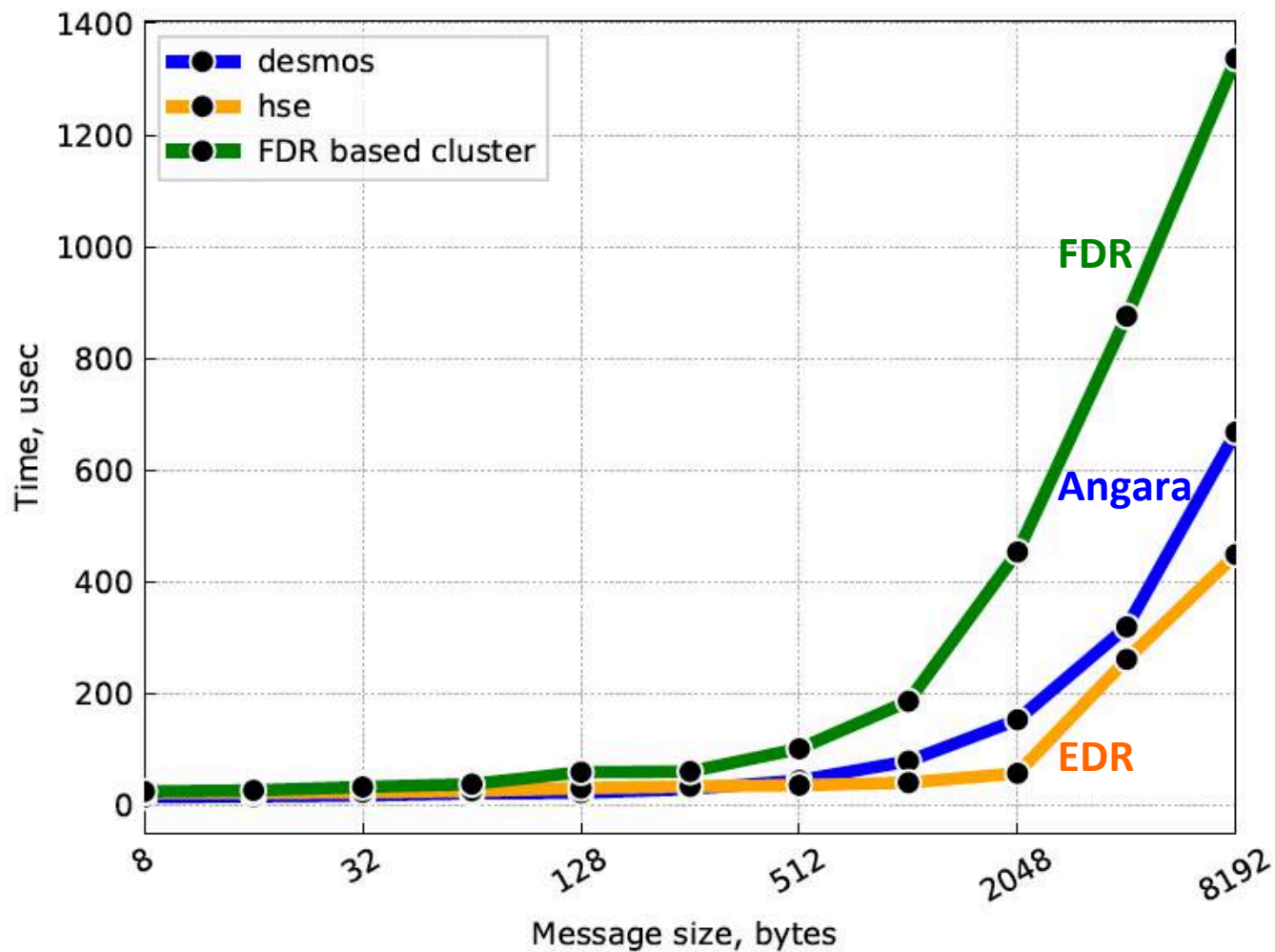
GROMACS

RIB:

2136412

atoms

osu_latency

Legend: desmos, hse, FDR based cluster

FDR, EDR, Angara

Y-axis: Time, usec
X-axis: Message size, bytes

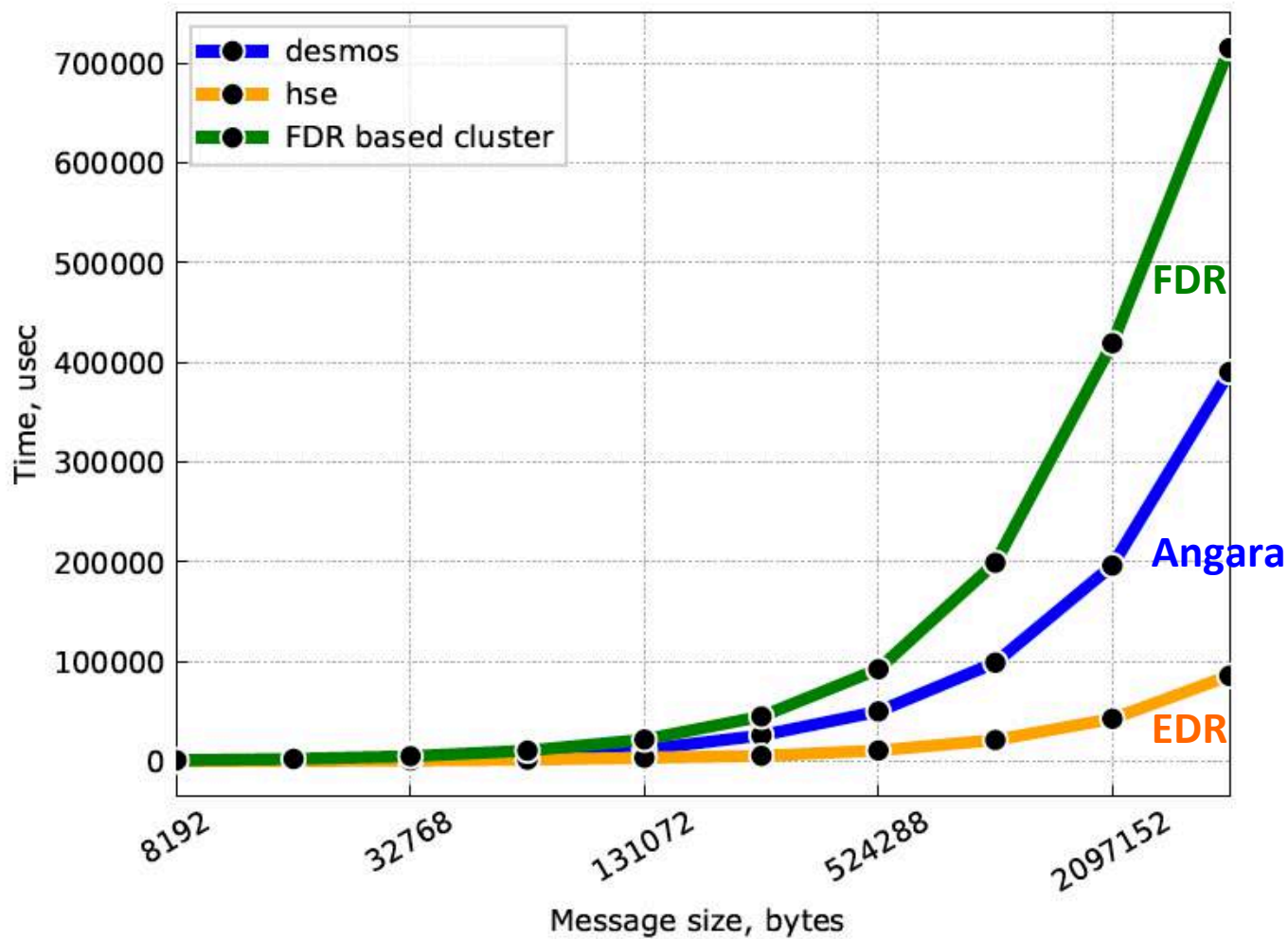alltoall, ppn=4, msg_size=262144

alltoall, ppn=4, msg_size=16

alltoall, ppn=4, nnodes=16

alltoall, ppn=4, nnodes=16

# VASP BENCHMARKS

WILEY

**SPECIAL ISSUE PAPER**

# VASP hits the memory wall: Processors efficiency comparison

**Vladimir Stegailov**[1,2] | **Grigory Smirnov**[1,2] | **Vyacheslav Vecher**[1,2]

[1] Joint Institute for High Temperatures of the Russian Academy of Sciences, Moscow, Russia
[2] Moscow Institute of Physics and Technology (State University), Dolgoprudny, Russia

**Correspondence**
Vladimir Stegailov, Moscow Institute of Physics and Technology (State University), Institutskiy Pereulok 9, Dolgoprudny 141701, Russia.
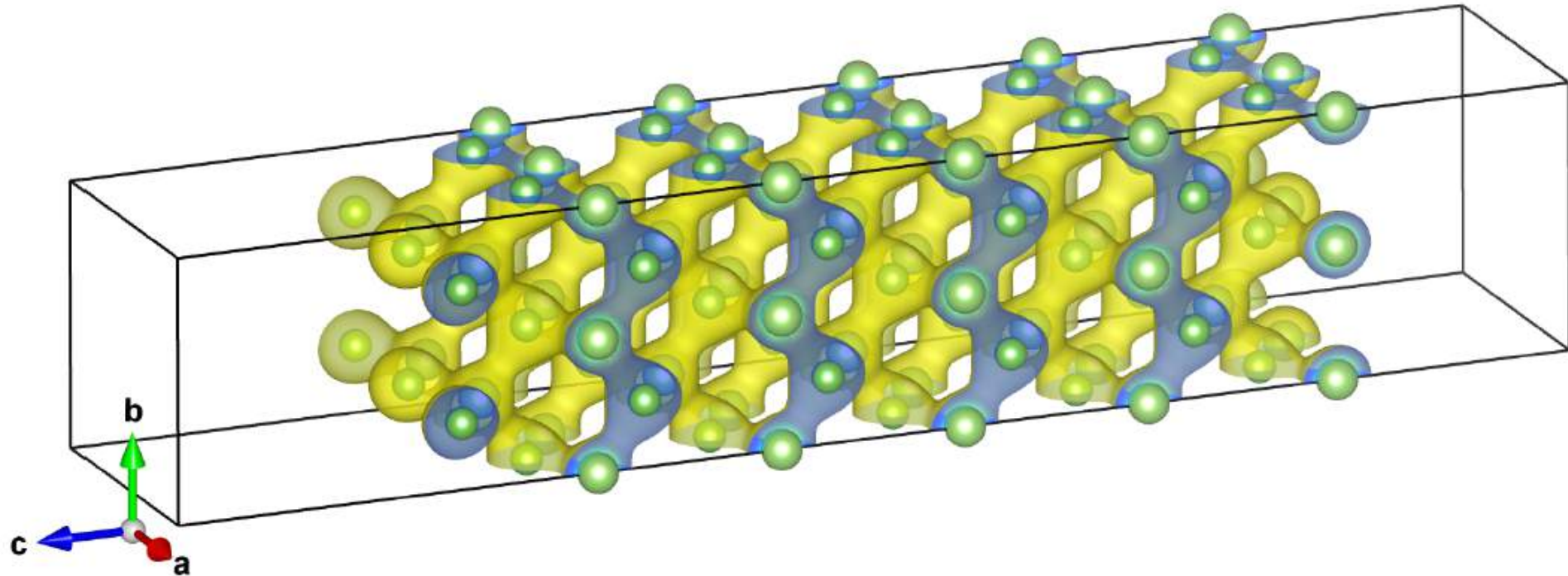Email: stegailov.vv@mipt.ru

**Summary**

First-principles calculations of electronic structure have been one of the most important classes of supercomputer applications for a long time. In this paper, we consider VASP as a *de facto* standard tool for density functional theory calculations widely used in materials science, condensed matter research, and other related fields. The choice of hardware for the efficient VASP calculations is not easy because of the large number of processor types available. We use the benchmark metric that is based on the balance of the peak floating point performance and the memory bandwidth. This metric gives us the possibility to compare different types of processors. We consider time-to-solution and energy-to-solution criteria and compare different Intel, AMD, and ARM 64-bit CPUs and hybrid CPU-GPU systems based on Nvidia Tesla P100.

**KEYWORDS**

ab initio calculations, energy-to-solution, memory bandwidth, peak performance, time-to-solution
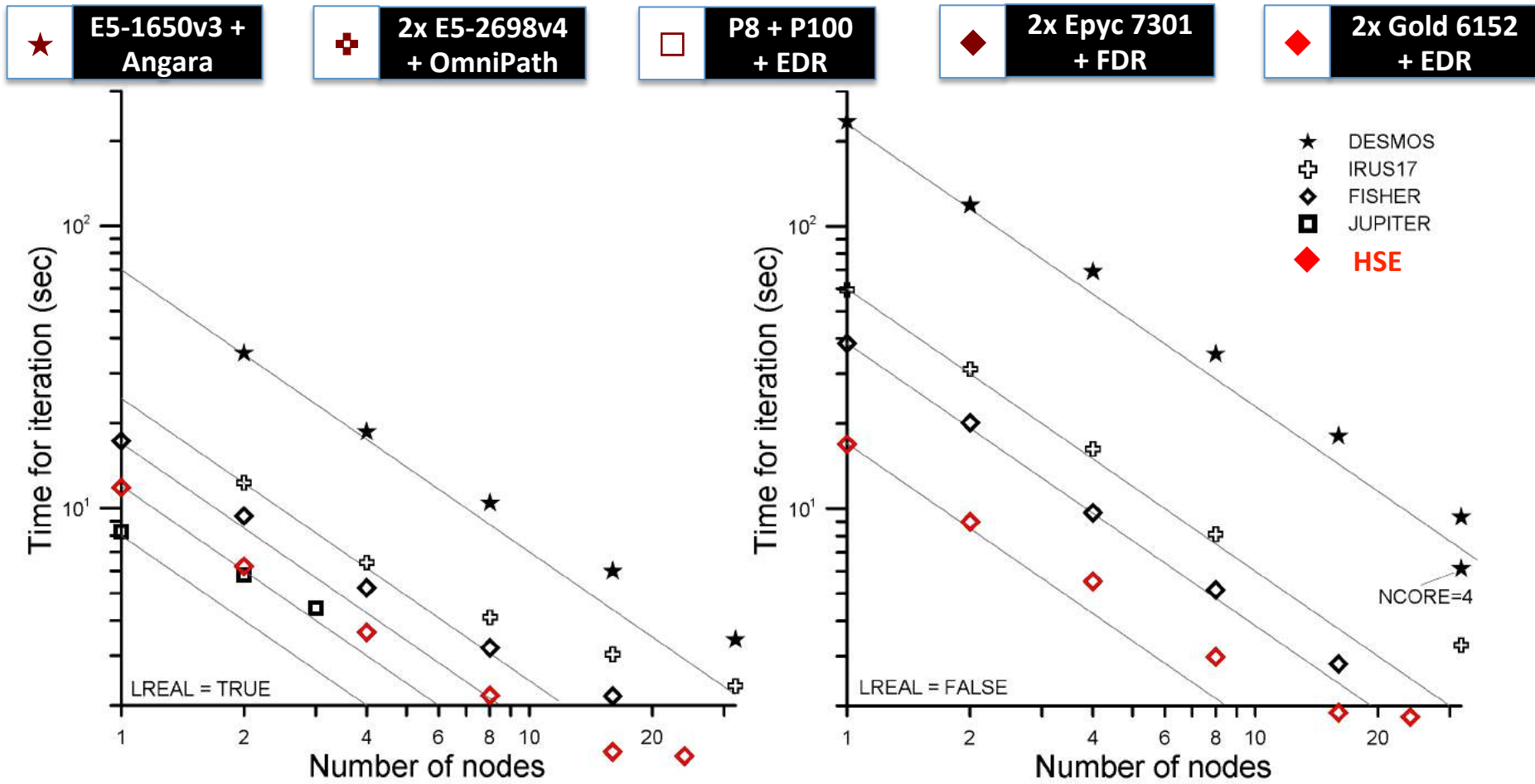
# VASP: GaAs benchmark with 80 atoms

# VASP: GaAs benchmark with 80 atoms

**TABLE 2** The selected best time-to-solution values (NCORE=1, KPAR=4). The number of MPI processes per socket is given together with the number of OpenMP threads per one MPI process

| PE type | LREAL=TRUE $\tau_{iter}$ (sec) | Options | LREAL=FALSE $\tau_{iter}$ (sec) | Options |
|---|---|---|---|---|
| 2 x Xeon E5-1650v3 | 36 | 6 MPI | 66 | 2 MPI x 3 OpenMP |
| 2 x Xeon E5-2660v4 | 25 | 14 MPI | 66 | 4 MPI x 3 OpenMP |
| 2 x Xeon E5-2698v4 | 24 | 20 MPI | 59 | 4 MPI x 4 OpenMP |
| | | | 59 | 8 MPI x 2 OpenMP |
| 2 x Epyc 7301 | 20 | 16 MPI | 39 | 8 MPI x 2 OpenMP |
| 2 x Epyc 7551 | 13 | 32 MPI | 33 | 8 MPI x 4 OpenMP |
| 1 x Tesla P100 (JUPITER) | 12 | NSIM=8 | n/a | |
| 2 x Tesla P100 (JUPITER) | 7.5 | NSIM=8 | n/a | |
| 2 x Tesla P100 (DGX-1) | 7.5 | NSIM=8 | n/a | |
| 4 x Tesla P100 (DGX-1) | 5.5 | NSIM=8 | n/a | |

# VASP: GaAs benchmark with 80 atoms

# VASP: GaAs benchmark with 80 atoms

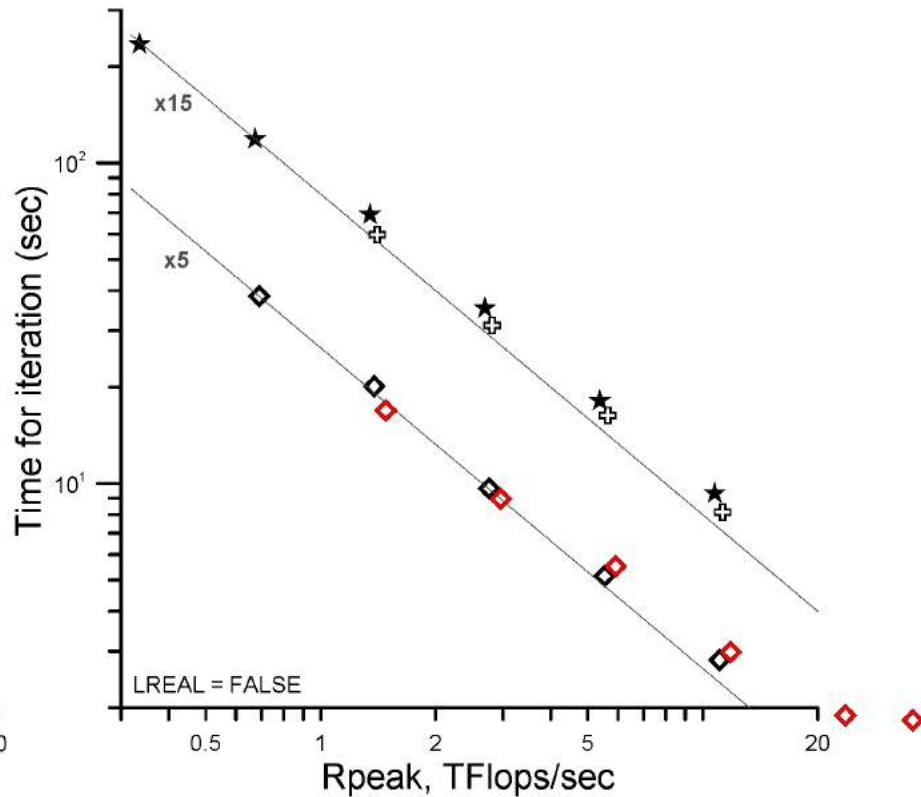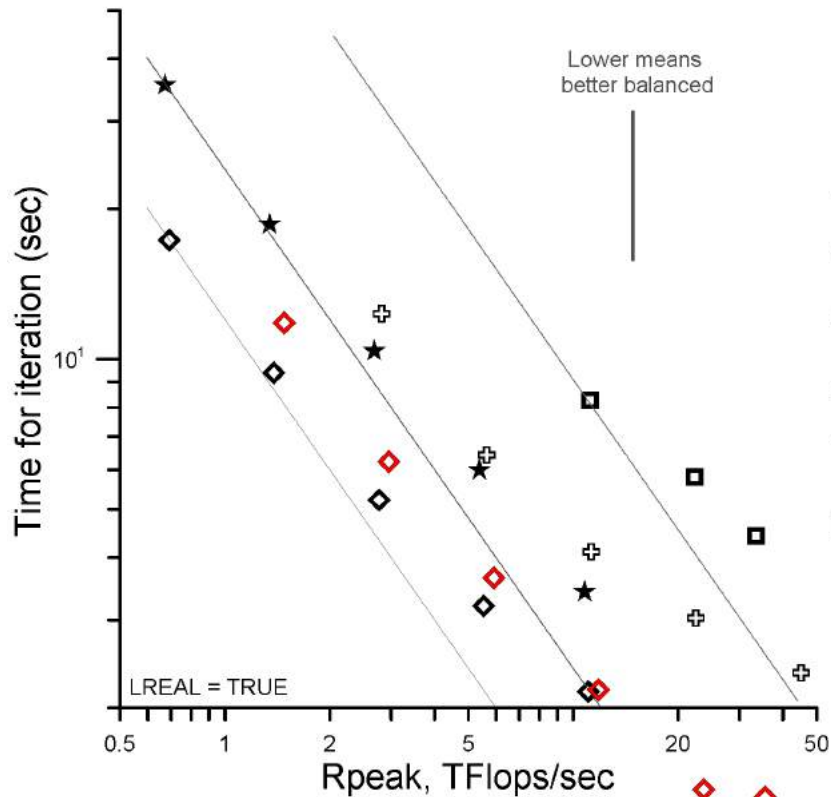| | E5-1650v3 + Angara | | 2x E5-2698v4 + OmniPath | | P8 + P100 + EDR | | 2x Epyc 7301 + FDR | | 2x Gold 6152 + EDR |

# A NEXT ANGARA-BASED HPC SYSTEM

# In march 2019 a next HPC system in JIHT RAS has been installed based on the Angara switch and one-port low-profile PCIe cards

# Conclusions

- The scaling tests for the classical MD and electronic structure calculations show the high efficiency of the MPI-exchanges over the Angara network.

- GPU-accelerated classical MD with Gromacs runs faster and is more cost-effective on supercomputers similar to Desmos than on wide-spread supercomputers based on expensive Intel Xeon multi-core CPUs.

- The job accounting statistic of the Desmos supercomputer has been reviewed. Two methods of quantitative efficiency monitoring have been proposed.

- All-to-all performance of the Angara network confirms its efficiency for electronic structure codes. The detailed analysis of modern CPUs efficiency for typical VASP model calculations has helped to select the best option for the next system.