

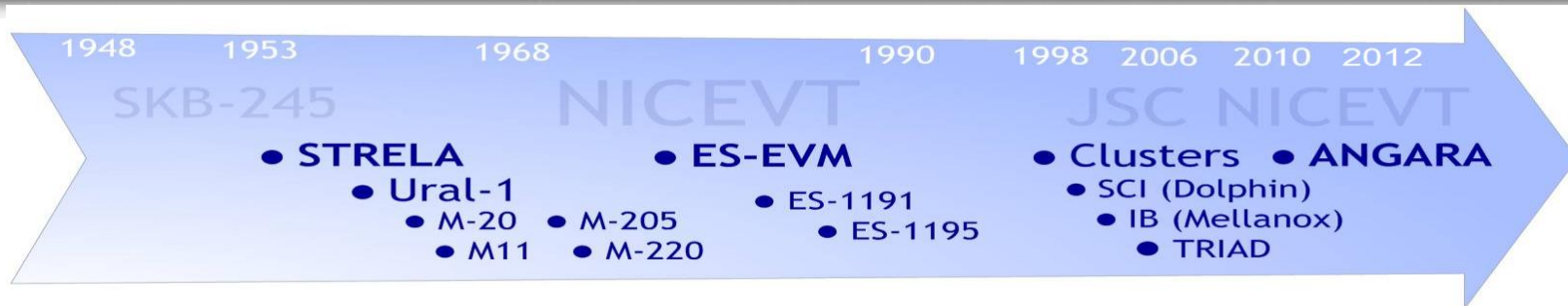


Ростех

Объединенная
приборостроительная
корпорация

Вычислительные системы с высокоскоростной сетью Ангара

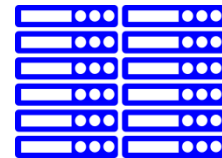
Симонов Алексей Сергеевич



1. Коммуникационная сеть Ангара. Общие сведения
2. Коммуникационная сеть Ангара. Достигнутые характеристики
3. Вычислительные системы с коммуникационной сетью Ангара
4. Перспективы коммуникационной сети Ангара

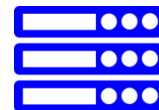
Назначение:

Коммуникационная сеть Ангара предназначена для осуществления передачи данных между узлами вычислительных систем с высокой скоростью и малой коммуникационной задержкой



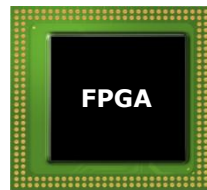
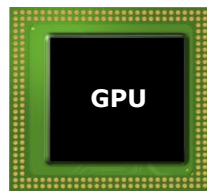
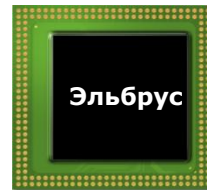
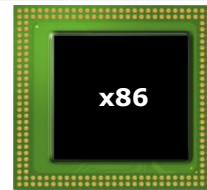
Области применения:

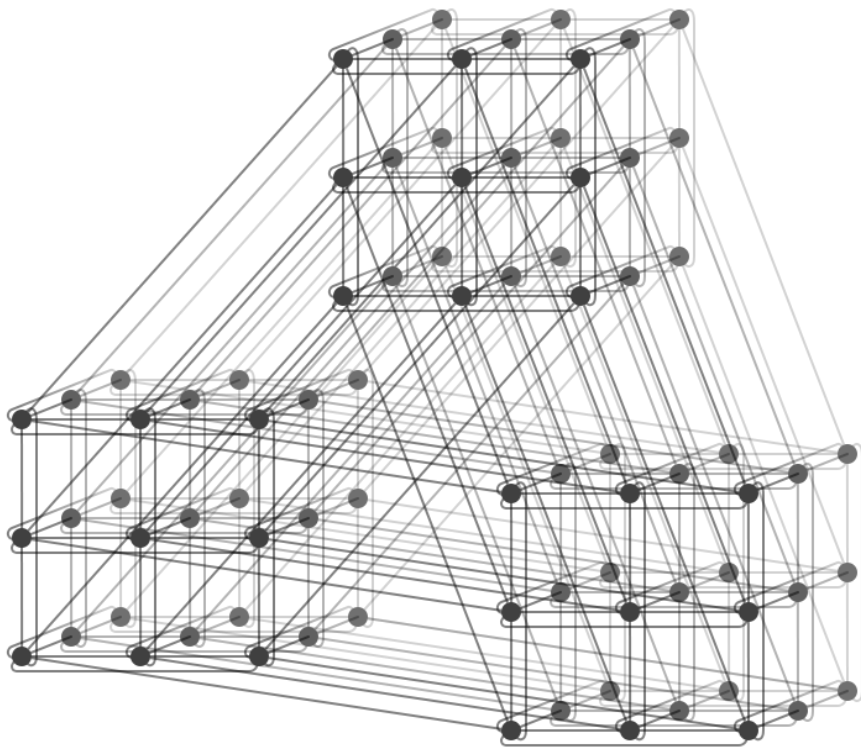
1. Вычислительные кластеры для проведения научных расчетов, решения задач инженерного анализа, обработки
2. Системы хранения и обработки больших данных
3. В качестве интерконнекта вычислительного поля в центрах обработки данных



Совместимость:

1. Системы на основе универсальных микропроцессоров x86, x86-64 (Intel, AMD), Эльбрус, GPU NVIDIA/AMD
2. Специализированные системы на основе FPGA
3. Проводятся работы по стыковке с микропроцессорами с архитектурой ARM



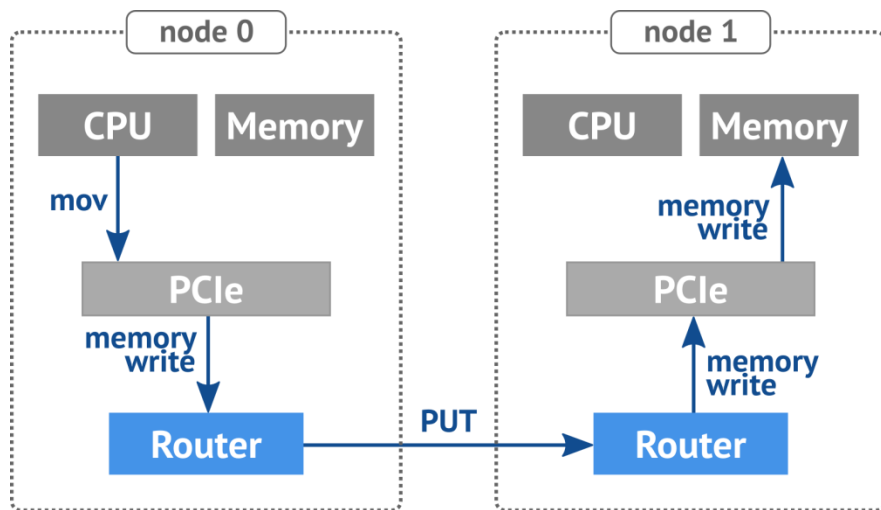


- Топология от «1D-mesh» до «4D-top»
- Односторонние коммуникации:
 - put (запись в удалённую память)
 - get (чтение из удалённой памяти)
 - атомарные операции add и xor
- Прямой доступ в память удаленного узла (RDMA)
- Коллективные операции:
 - broadcast
 - reduce
- Поддержка многоядерности
- Адаптивная передача пакетов
- Механизмы синхронизации

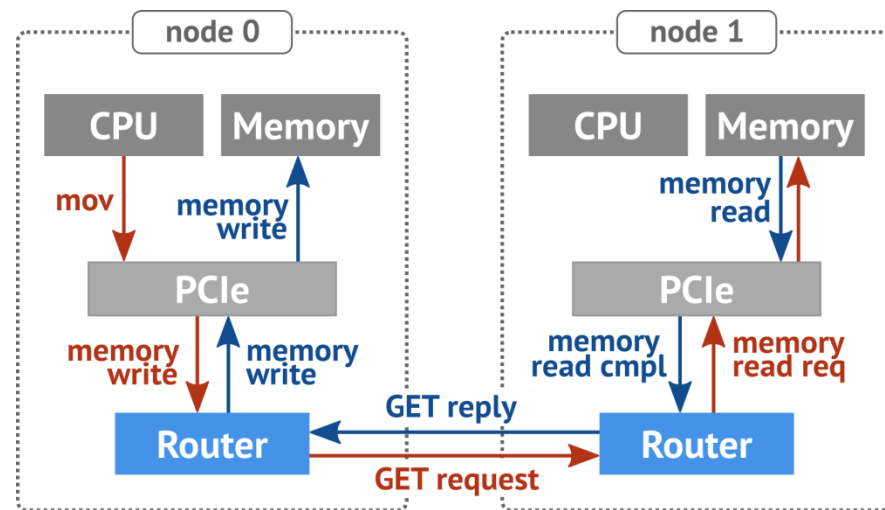
Операция	Кто инжектирует данные в сеть	Размер транзакции	Примечания
PUT vc0	Хост	1—256 байт	Эжекция мимо кроссбара (ответы)
PUT vc1	Хост	1—256 байт	Эжекция через кроссбар (запросы)
GET	Адаптер	1—256 байт	Гранулярность 8 байт
ADD/XOR	Хост	1—8 байт	Атомарность средствами адаптера
LWPUT	Адаптер	8— 2^{30} байт	Гранулярность 8 байт
LGET	Адаптер	8— 2^{30} байт	Гранулярность 8 байт

- При инжекции коротких пакетов (PUT) на хосте выполняются операции MOV в адресное пространство ресурса PCI Express адаптера (инжекционный буфер).
- При инжекции длинных пакетов (LWPUT) адаптер выполняет memory read из памяти хоста, после чего отправляет полученные из памяти данные в сеть.

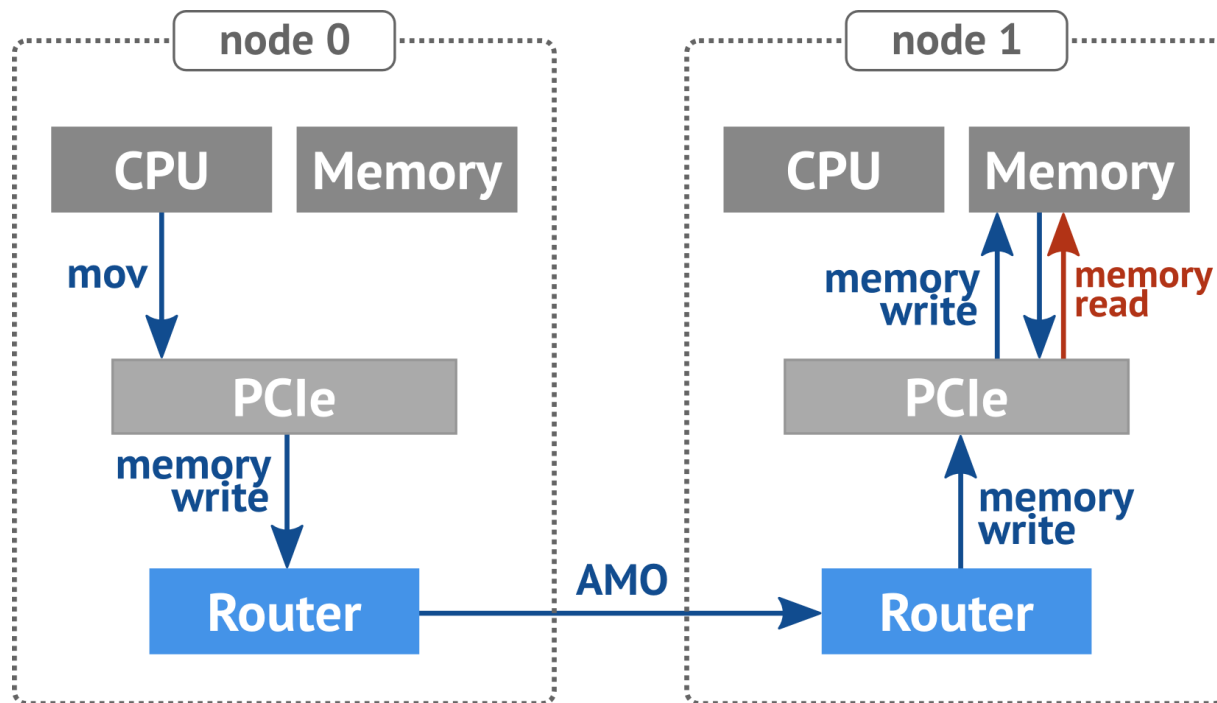
Запись



Чтение

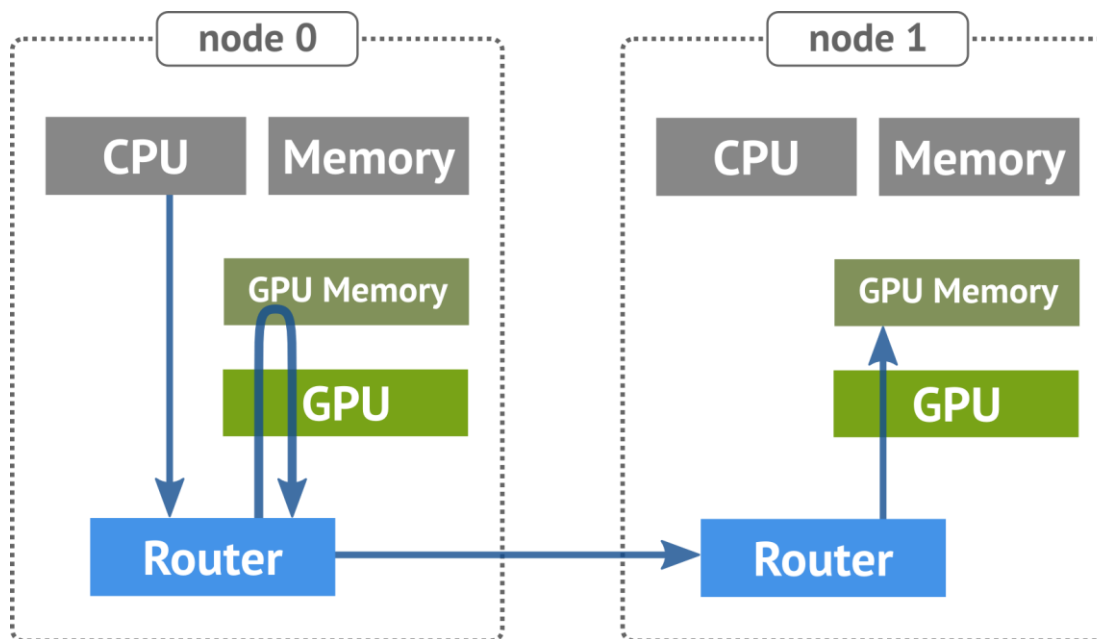


Удалённая атомарная операция (сложение, XOR)

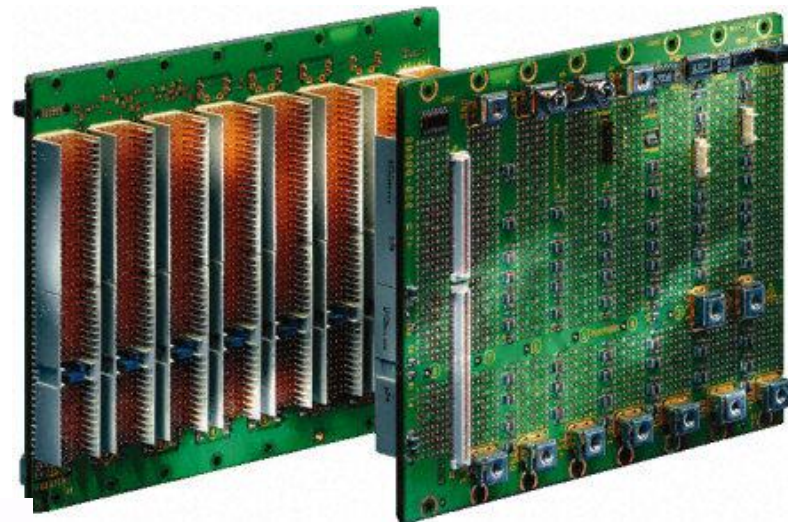


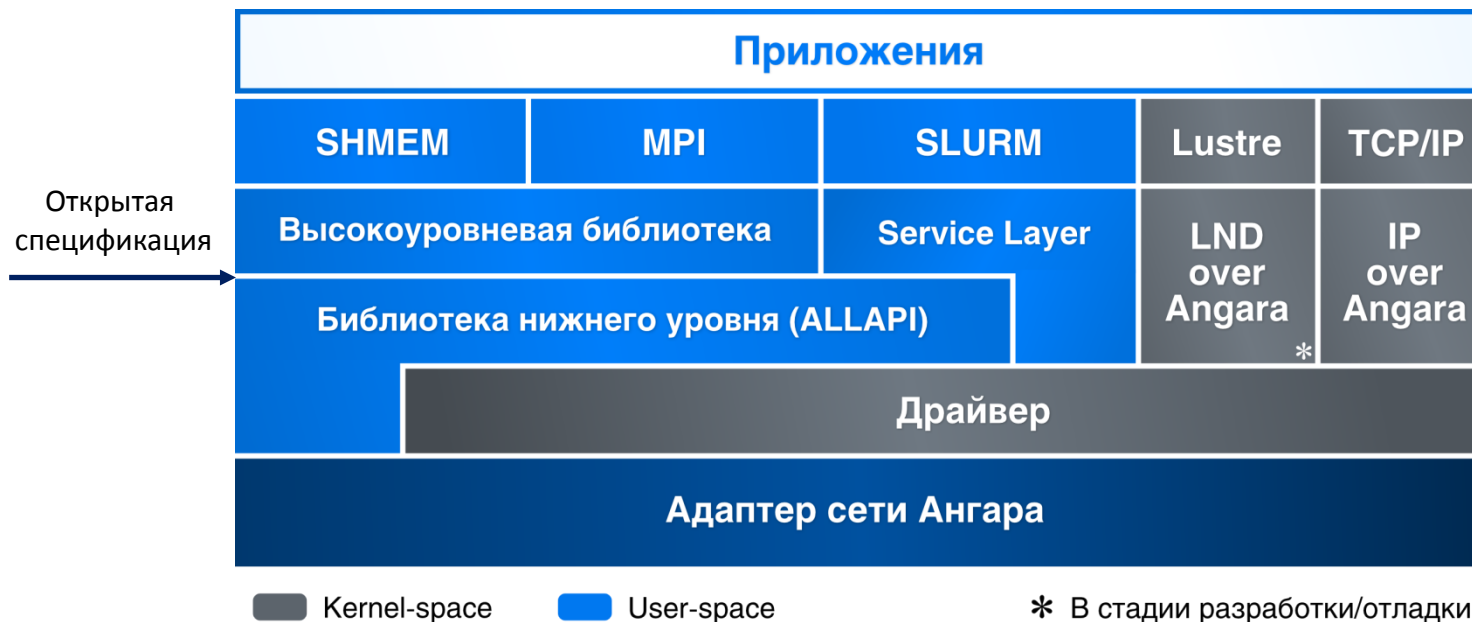
Реализация GPU Direct RDMA

Прямая запись из памяти локального GPU в память GPU на удаленном узле

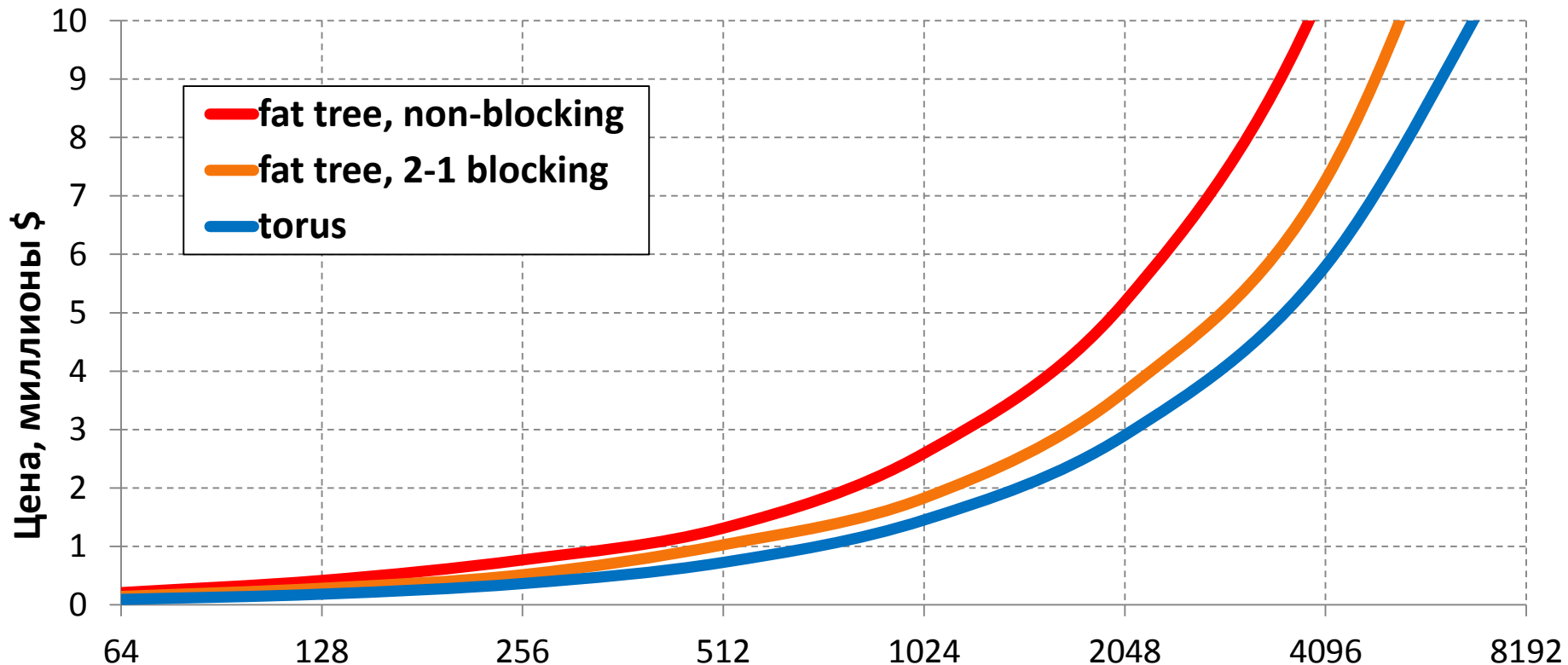


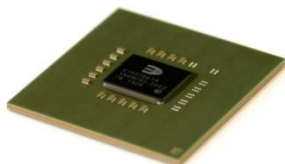
- Поддержка различных физических сред передачи данных (кабели медные и оптические, объединительные платы)
- Настройка и управление параметрами приемопередатчиков
- Коммуникационная задержка 670 нс





- Поддержка ОС : Astra Linux SE 1.3, ОС «Эльбрус» ,OpenSUSE/SLES 11SP3, CentOS 6.0-6.3, Версия ядра Linux от 2.6.21 до 3.16.0
- Поддержка компиляторов языков Fortran 77/90/95 (GNU, Intel), C/C++ (GNU, Intel)





Чип EC8430

FCBGA 1521
40×40 мм
35 Вт

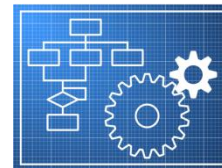


Заказная разработка
платы адаптера



Development kit

PCIe gen 2 x16
Full-height,
full length



Доработка
ПО адаптера
и адаптация
библиотек



Варианты исполнения Адаптера PCIe

PCIe gen 2 x16
- Full-height, full length
- Full-height, half length



Адаптация
и профилирование
прикладного ПО

1. Коммуникационная сеть Ангара. Общие сведения
2. Коммуникационная сеть Ангара. Достигнутые характеристики
3. Вычислительные системы с коммуникационной сетью Ангара
4. Перспективы коммуникационной сети Ангара

24 вычислительных узла

- Supermicro SuperServer 5017GR-TF
- 2 процессора Intel Xeon E5-2630 (LGA2011, 6 ядер, 2.3 ГГц)
- 64 ГБ

12 вычислительных узлов

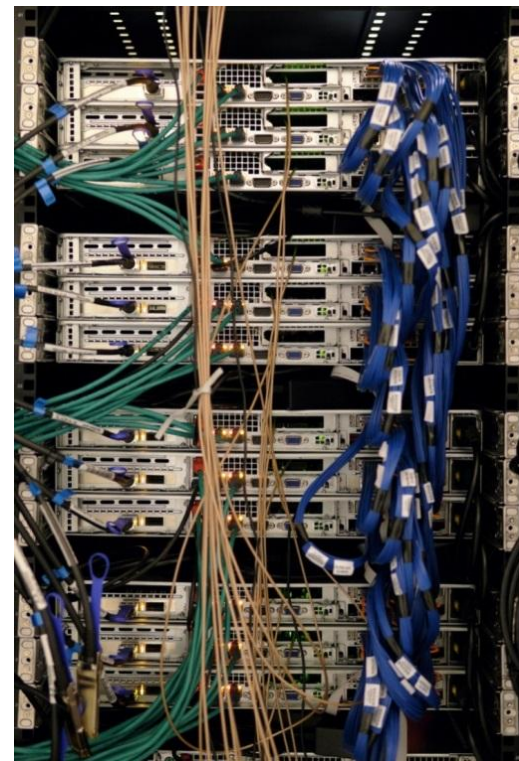
- Supermicro SuperServer 5017GR-TF
- процессор Intel Xeon E5-2660 (LGA2011, 8 ядер, 2.2 ГГц)
- 64 ГБ

Сеть «Ангара»

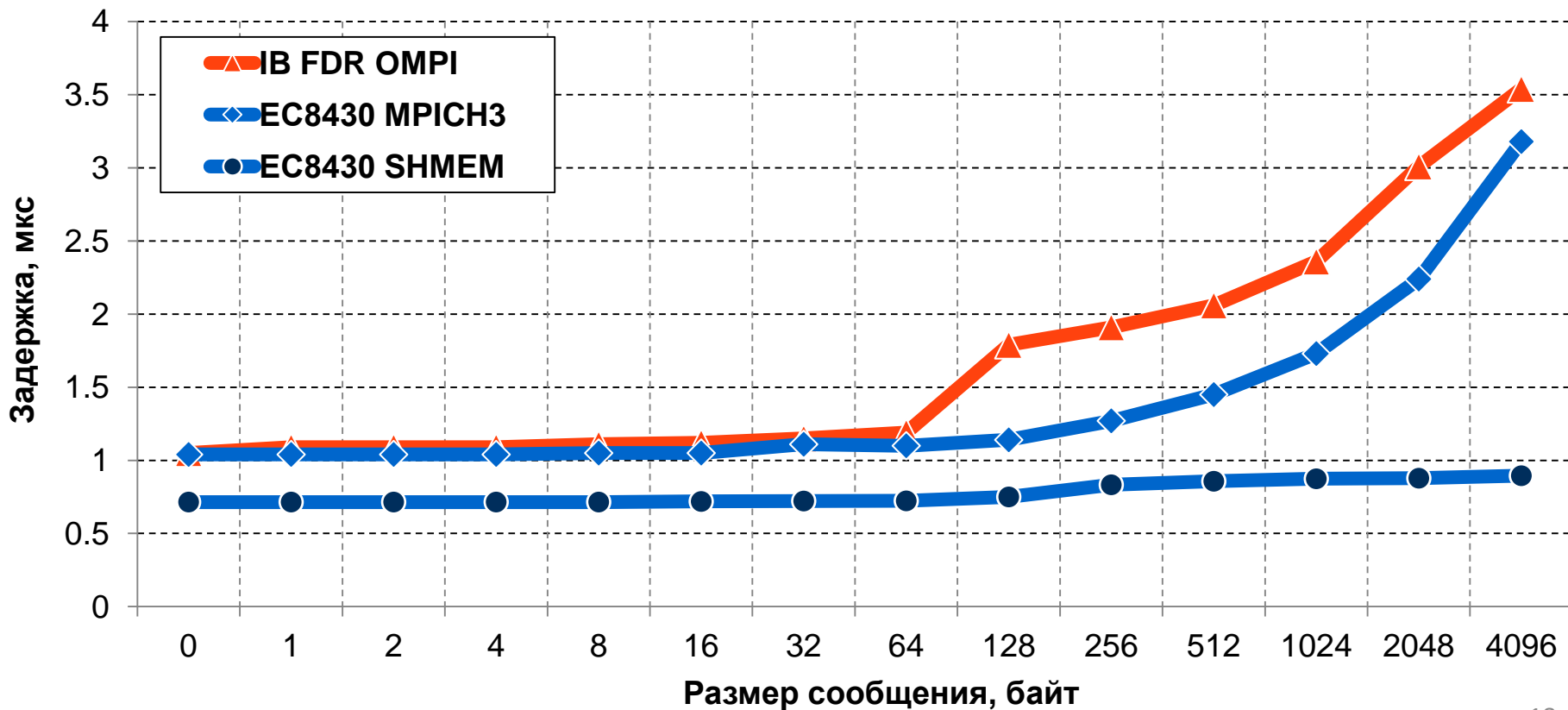
- Адаптер EC8430, топология 3D-тор 3x3x4
- Собственная реализация OpenSHMEM
- MPI: MPICH 3.0.4

Операционная система

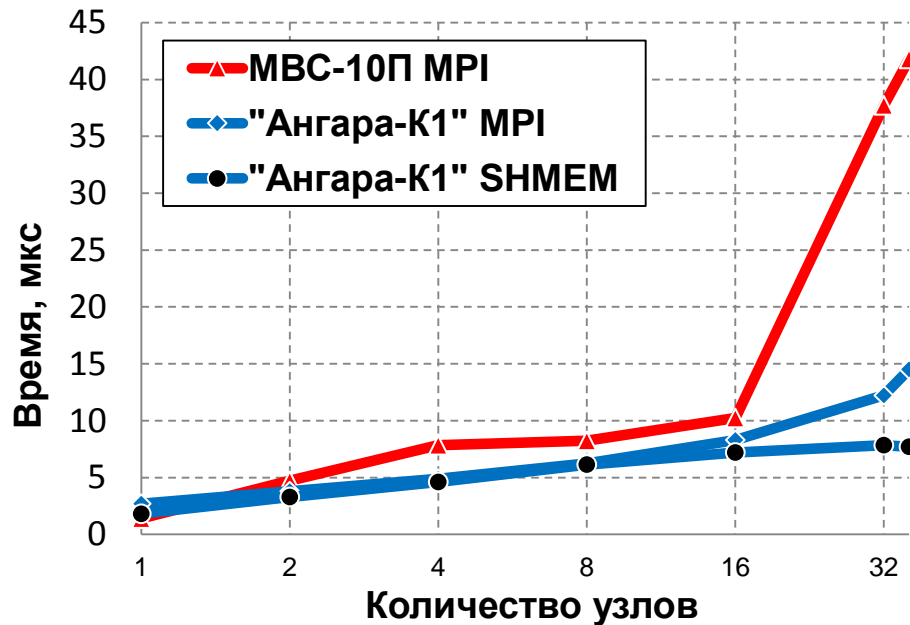
- SLES 11 SP2, Linux 3.0.13-0.27-default
- GCC 4.3.4 (revision 152973)



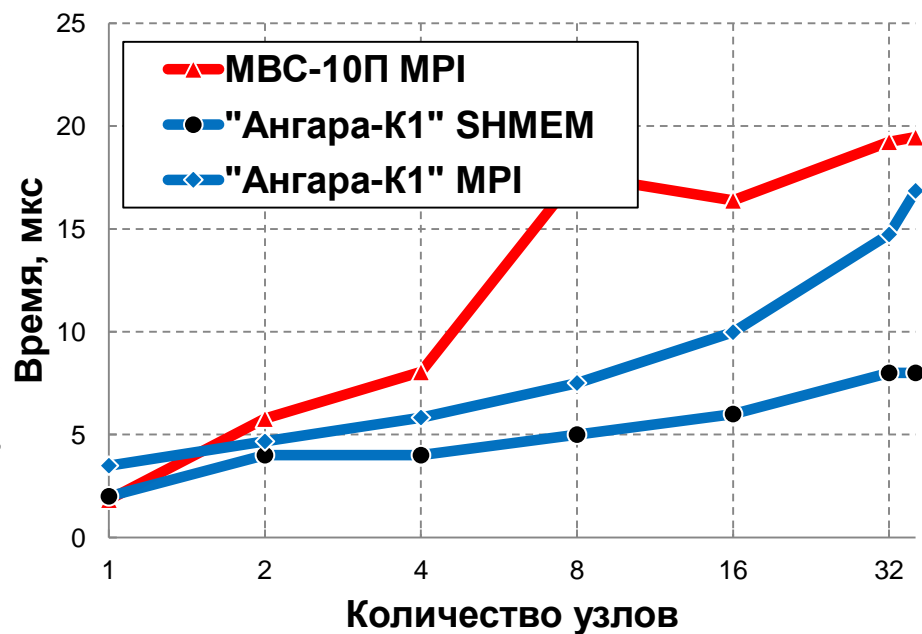
	Ангара-К1		МВС-10П
Узлы	А	2x Xeon E5-2630 по 6 ядер, 2.3 ГГц	2x Xeon E5-2690 по 8 ядер, 2.9 ГГц
	В	Xeon E5-2660 по 8 ядер, 2.2 ГГц	
Количество узлов	24*А+12*В = 36		207 (36)
Память узла	64 ГБ		64 ГБ
Сеть	Ангара 3D-топ 3x3x4		Infiniband 4xFDR Fat Tree



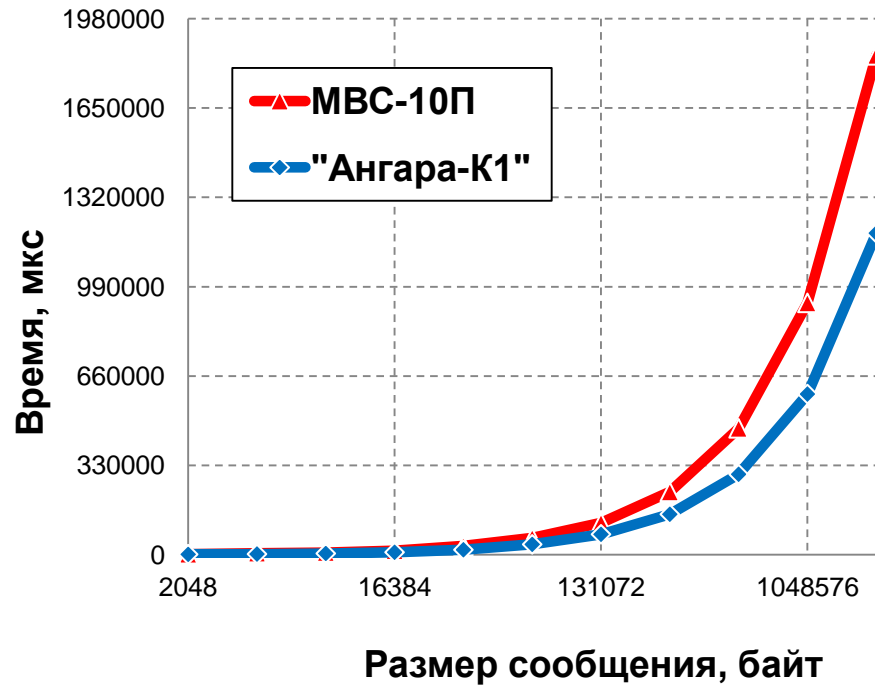
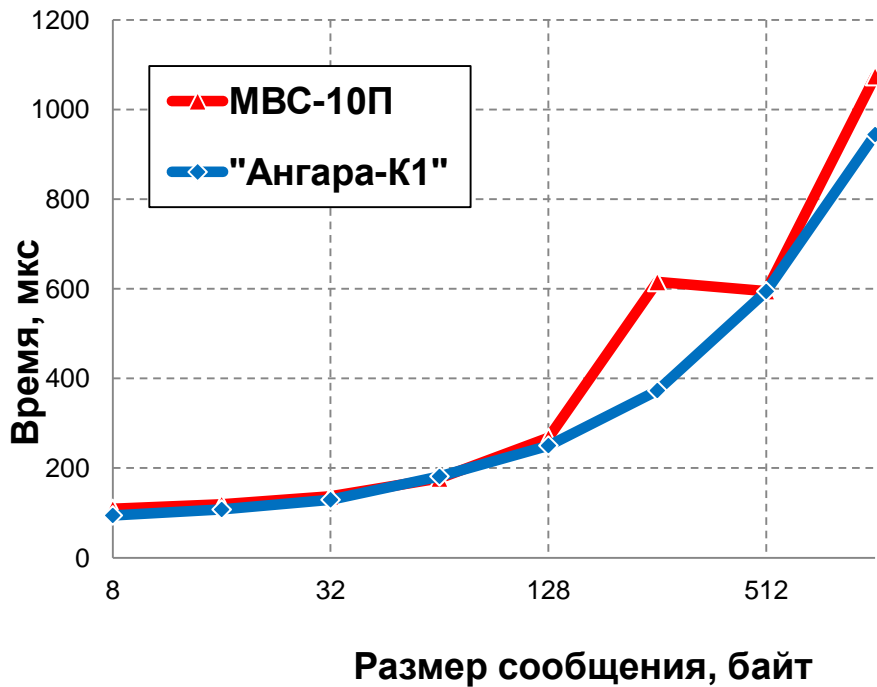
IMB Barrier



IMB Allreduce

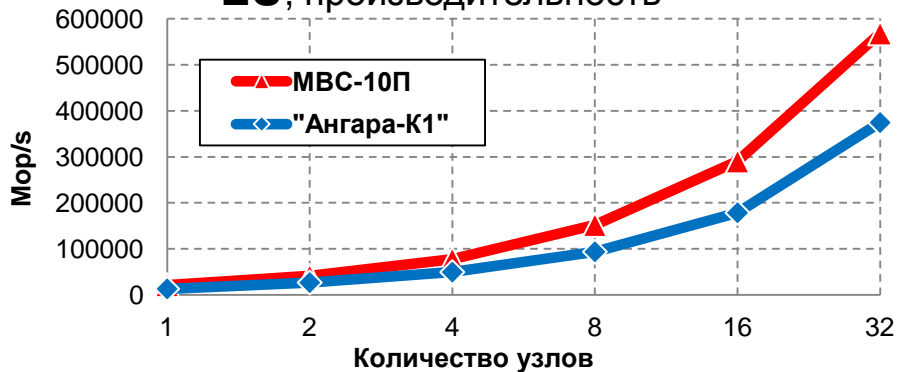


размер сообщения 8 байт

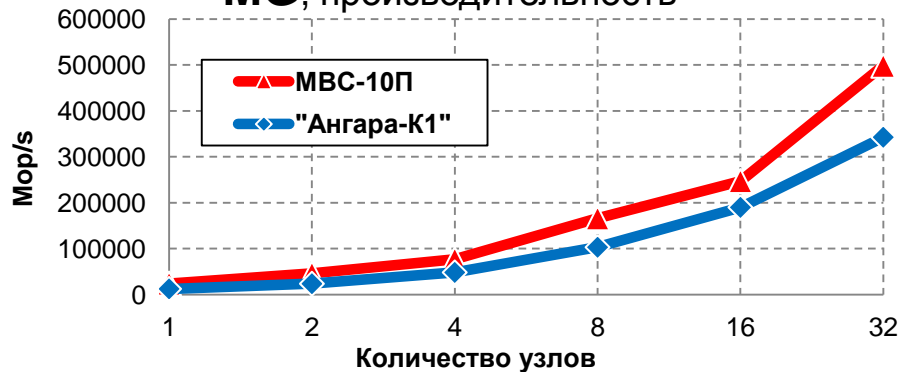


		Ангара	МВС-10П
HPL	Тфлопс	4.44	—
	% пиковой	85 %	—
HPCG MPI	Гфлопс	279	363
	% пиковой	5.3 %	5.4 %
HPCG SHMEM	Гфлопс	342	—
	% пиковой	6.5 %	—

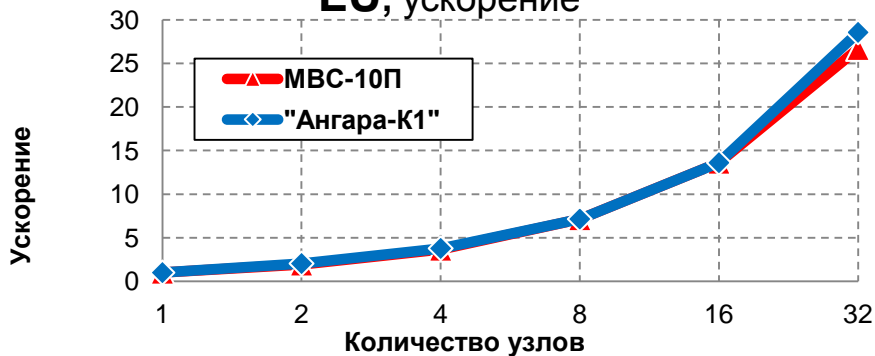
LU, производительность



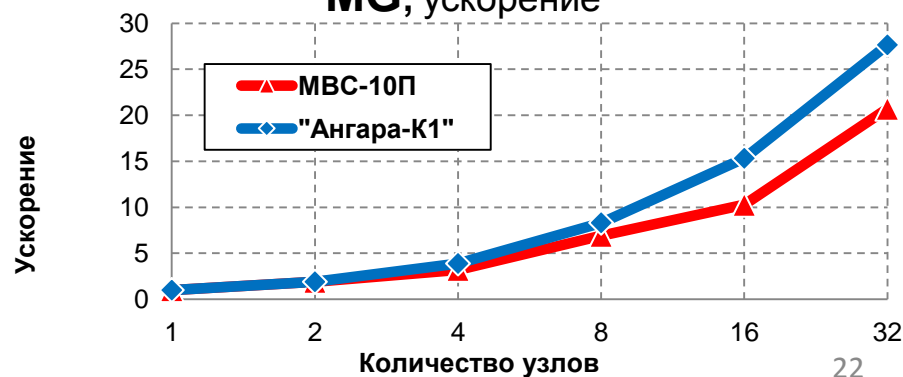
MG, производительность

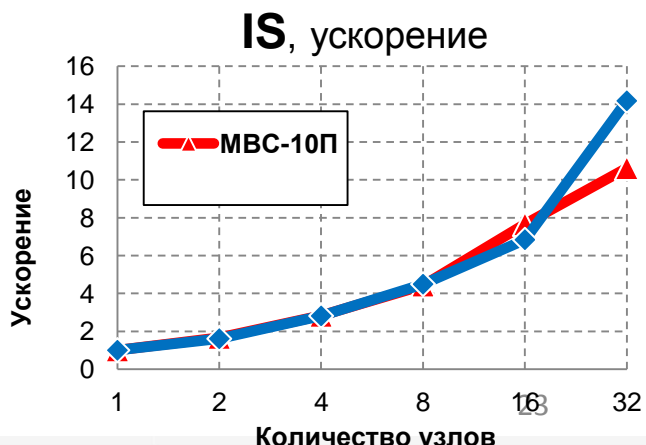
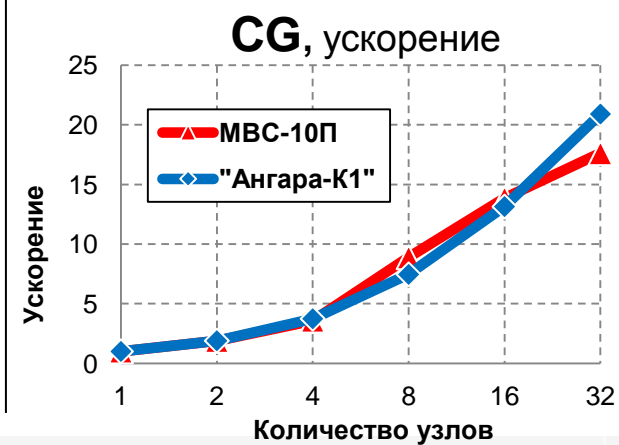
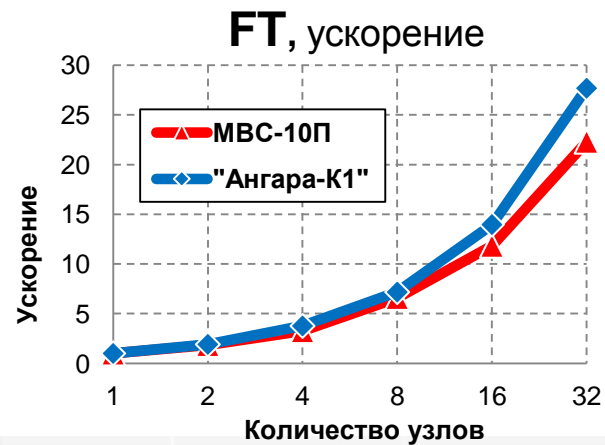
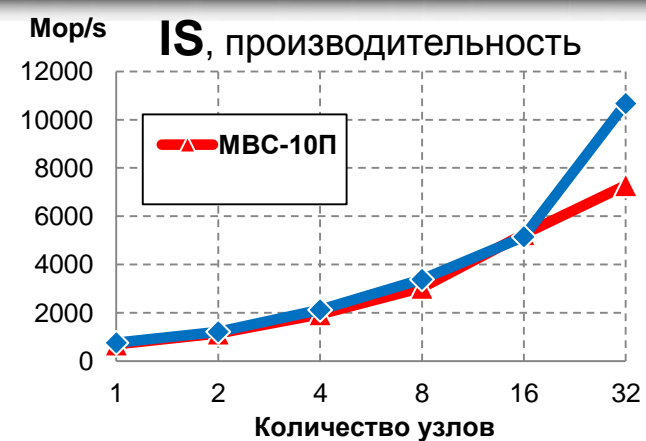
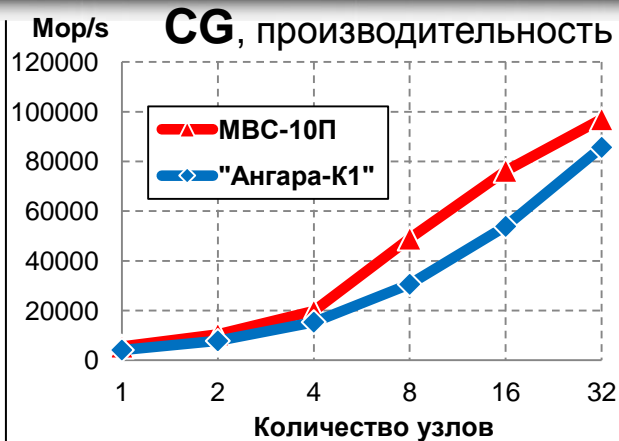
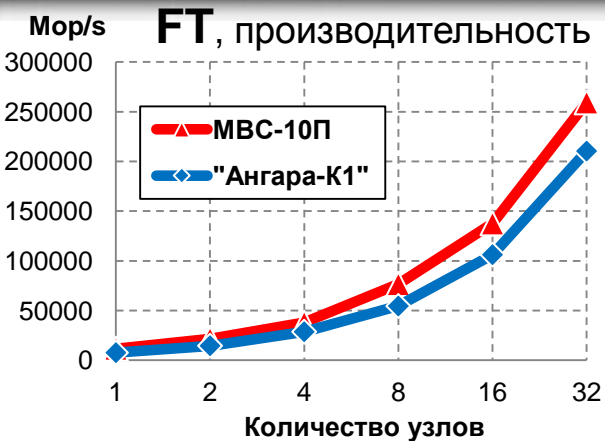


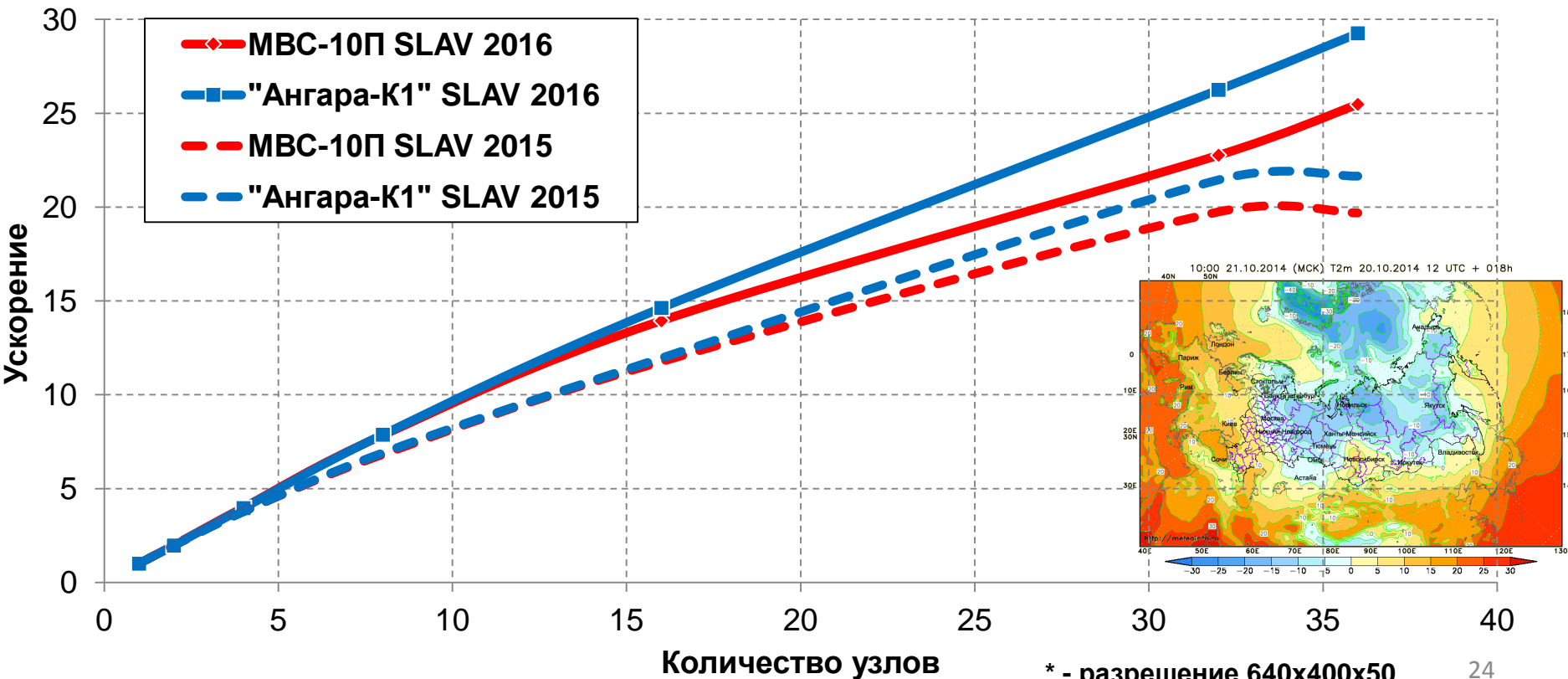
LU, ускорение

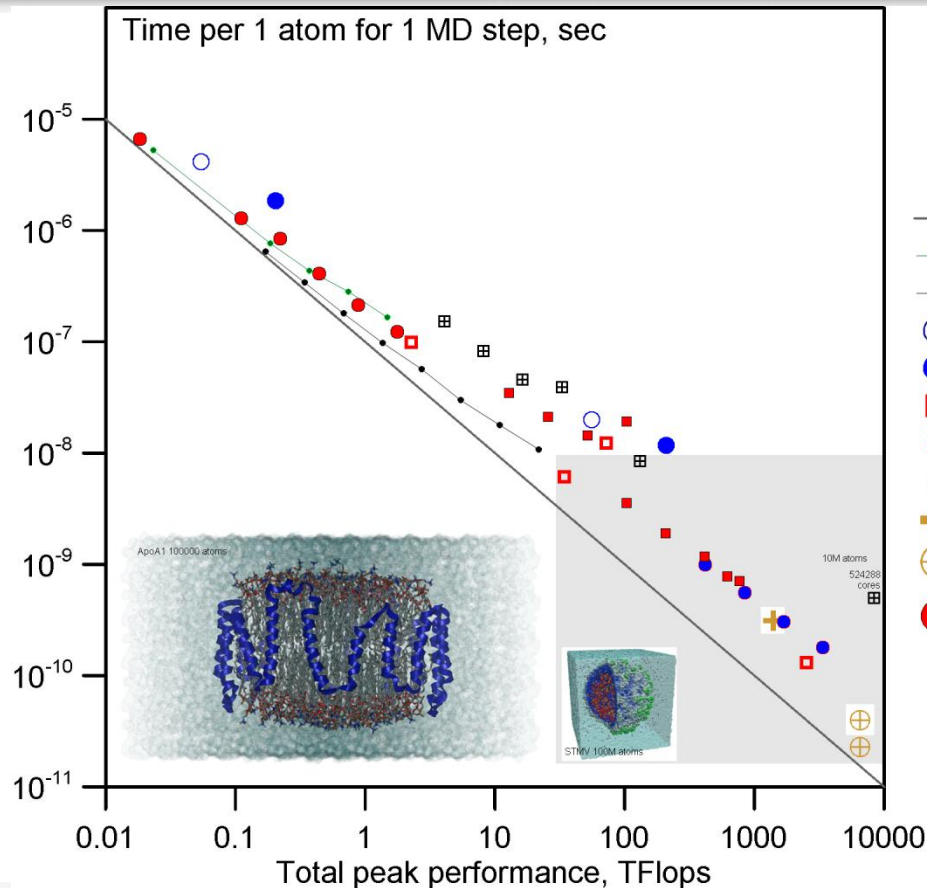


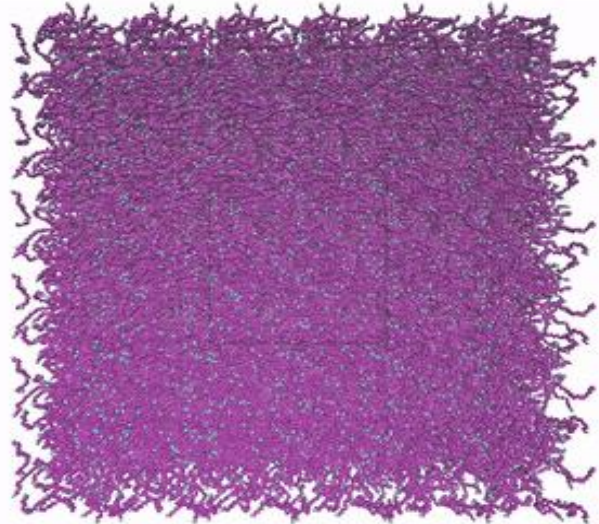
MG, ускорение



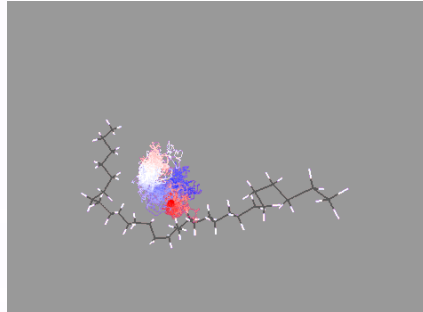








н-триактановая жидкость
 $T = 350 \div 490 \text{ K}$; $P = 1 \text{ атм}$
Количество молекул $\sim 4\ 000$



Траектория 1-й молекулы
в исследуемой жидкости

Диффузия, вязкость жидких углеводородов,
т.к. они входят в состав

трансформаторных масел,
топлив и смазочных материалов

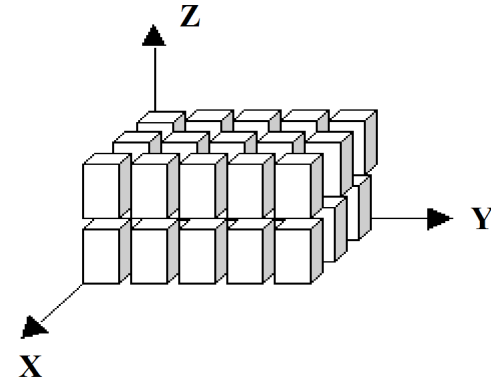
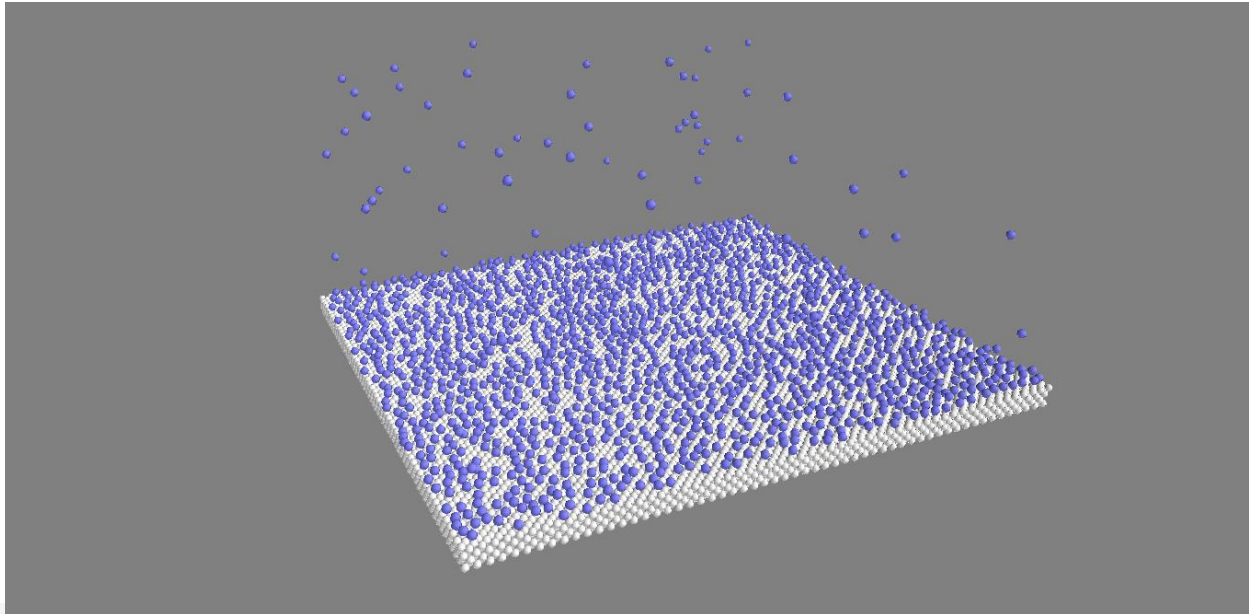
Молекулярная динамика -> **макроскопические свойства**

LAMMPS, 30 July 2016

Расчет по взаимодействию азота со стенками никелевого микроканала

Число частиц: $8\,128\,512 + 423\,840 = 8\,552\,352$,Температура термостатов: $T_{Ni} = 273.15\text{ K}$, $T_{N_2} = 273.15\text{ K}$

Число шагов по времени: 2 000 000 шагов, 1 шаг = 2 фс

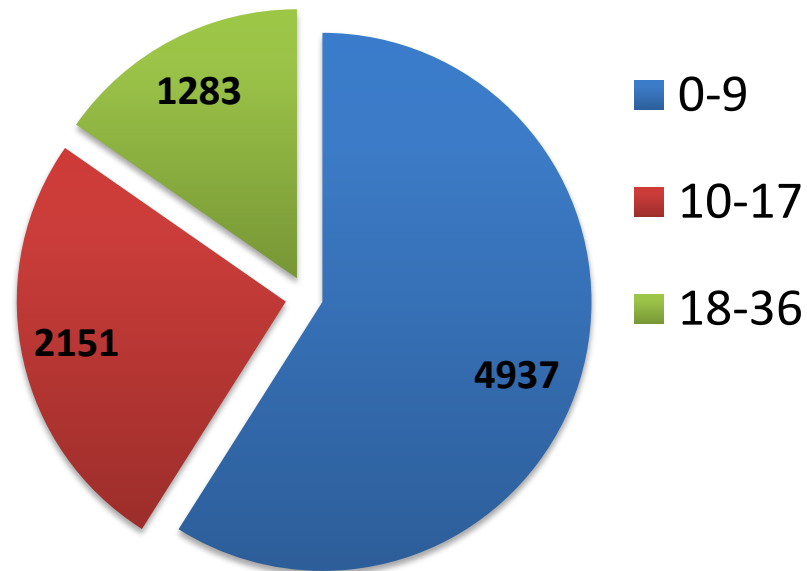
Размер системы: $102 \times 102 \times 1534\text{ нм}^3$ 

MPI+OpenMP

Фрагмент распределения
молекул азота (область
 $20 \times 20\text{ нм}$) на поверхности
никелевой пластины, в
момент времени 2.3 нс

- В режиме внешнего доступа работает с 1.11.2015 г.
- Число пользователей: 28
- Всего использовано 183778 процессоро-часов

Часы работы по количеству
используемых узлов кластера
(на 19.09.2016)



1. Коммуникационная сеть Ангара. Общие сведения
2. Коммуникационная сеть Ангара. Достигнутые характеристики
3. Вычислительные системы с коммуникационной сетью Ангара
4. Перспективы коммуникационной сети Ангара

1. Внедрение в коммерческий сегмент (вычислительные кластеры для науки, образовательных учреждений, промышленных предприятий, создаваемые на основе коммерчески доступных комплектующих зарубежного производства) осуществляется в партнерстве с компаниями - интеграторами путем поставки им сетевого оборудования Ангара в стандартных форм-факторах и формирования для них гибкой ценовой политики

2. Внедрение в сегмент заказных разработок осуществляется путем работы напрямую с Заказчиками и создания под них специализированных версий сетевого оборудования Ангара или вычислительных систем на его основе



Кластер на базе ПК



До 16 узлов

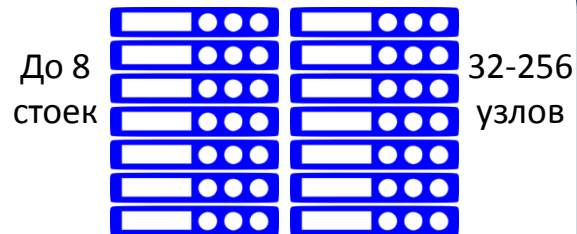
CPU	Хеон 1600 v3/4
HDD	500 ГБ
ОЗУ	8-64 ГБ
GPU*	NVidia GTX™/ AMD Radeon™
Сеть	Ангара 2D-тор

Кластер: 1U сервера



CPU	Хеон 1600 v3/4
HDD/SSD*	1 ТБ
ОЗУ	16-256 ГБ
GPU*	NVidia Tesla™/ AMD FirePro™
Сеть	Ангара 3D-тор

Кластер: 1U сервера



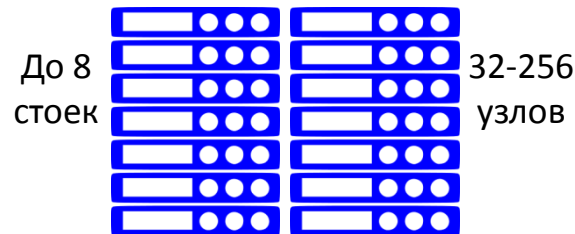
CPU	Хеон 1600 v3/4
HDD/SSD*	1 ТБ
ОЗУ	64-256 ГБ
GPU*	NVidia Tesla™/ AMD FirePro™
Сеть	Ангара 4D-тор

Кластер: 1U сервера



CPU	2x Xeon 2600 v3/4
HDD/SSD*	1 ТБ
ОЗУ	16-256 ГБ
GPU*	NVidia Tesla™/ AMD FirePro™
Сеть	Ангара 3D/4D-тор

Кластер: 1U сервера



CPU	2x Xeon 2600 v3/4
HDD/SSD*	1 ТБ
ОЗУ	64-256 ГБ
GPU*	NVidia Tesla™/ AMD FirePro™
Сеть	Ангара 4D-тор



24 вычислительных узла

- Supermicro SuperServer 5017GR-TF
- 2 процессора Intel Xeon E5-2630 (LGA2011, 6 ядер, 2.3 ГГц)
- 64 ГБ

12 вычислительных узлов

- Supermicro SuperServer 5017GR-TF
- процессор Intel Xeon E5-2660 (LGA2011, 8 ядер, 2.2 ГГц)
- 64 ГБ

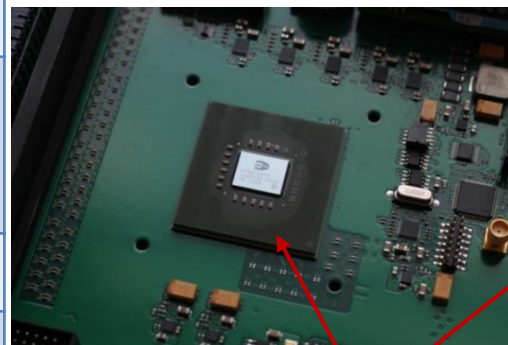
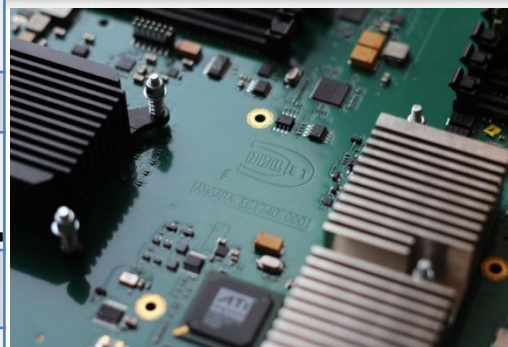
Сеть «Ангара»

- Адаптер EC8430, топология 3D-тор 3x3x4
- Собственная реализация OpenSHMEM
- MPI: MPICH 3.0.4

Операционная система

- SLES 11 SP2, Linux 3.0.13-0.27-default
- GCC 4.3.4 (revision 152973)

Характеристика	EC1740.0003 (EATX)
Сокеты x86/ядра x86	2/32
Тип ОЗУ	4-х канальная DDR3 1866/1600/1333/1066 МГц
Объем ОЗУ	128 Гбайт DDR3
Модули ОЗУ	8 шт.
Интерфейсы	Ангара 2x10G Ethernet 4xSATA RAID
Накопителей	4 шт.
Потребляемая мощность	400 Вт
Габаритные размеры	321 x 367 x 30 мм



СБИС ВКС Ангара



10G Ethernet

Оптические линки
ВКС Ангара

Испытательный стенд**Узлы**На основе процессоров
Эльбрус-4С,
частота 750 МГц, и КПИ1**Количество
узлов**

4

**Память
узла**

24

СетьАнгара
1D-тор**HPL***

75.09 Гфлопс (78.2%)



* - предварительный результат

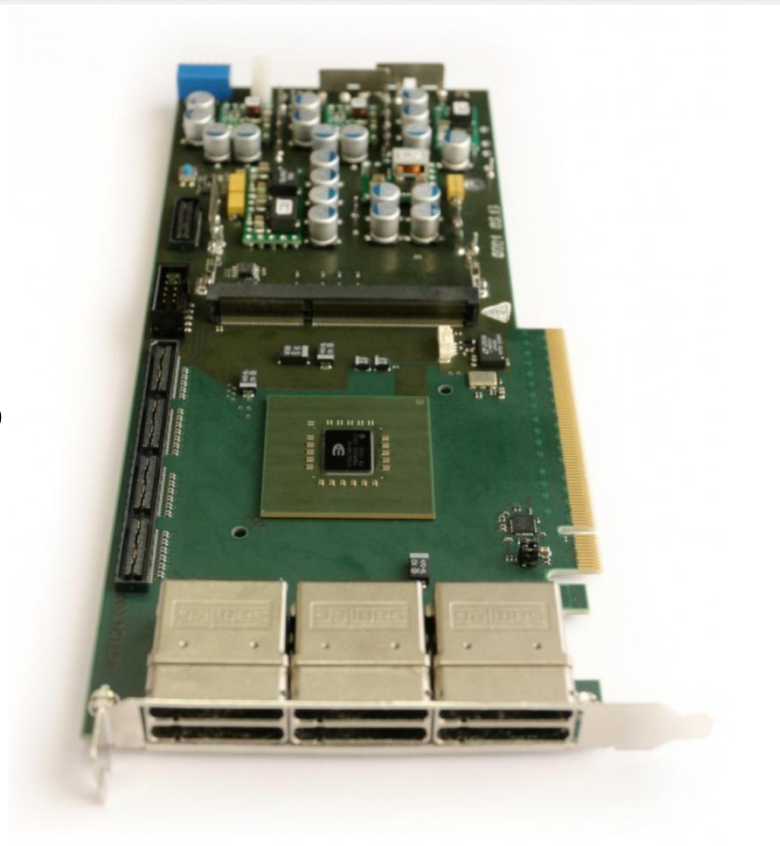
1. Коммуникационная сеть Ангара. Общие сведения
2. Коммуникационная сеть Ангара. Достигнутые характеристики
3. Вычислительные системы с коммуникационной сетью Ангара
4. Перспективы коммуникационной сети Ангара

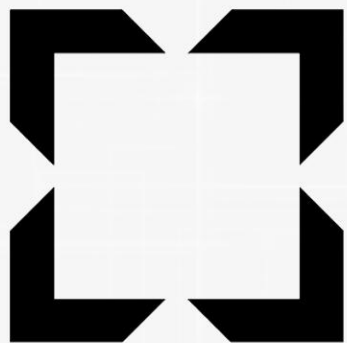
- Конференция GraphHPC: 2014, 2015, 2016, 2017 годы
- Конкурс с автоматической системой проведения
- Доступ на кластер с сетью Ангара
- Сайт конференции graphhpc.dislab.org
- Открытое ПО

Контакты:

117587, Москва, Варшавское ш, 125

angara@nicevt.ru





Ростех

*Объединенная
приборостроительная
корпорация*