

5-й Московский суперкомпьютерный форум (МСКФ 2014)

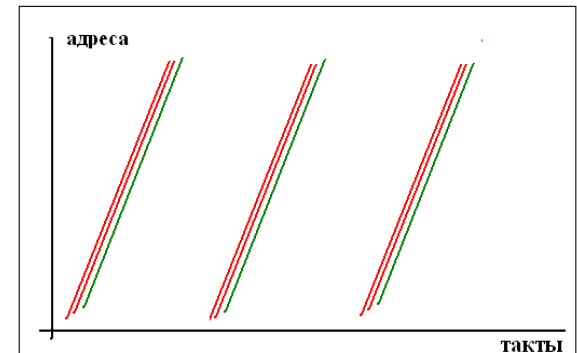
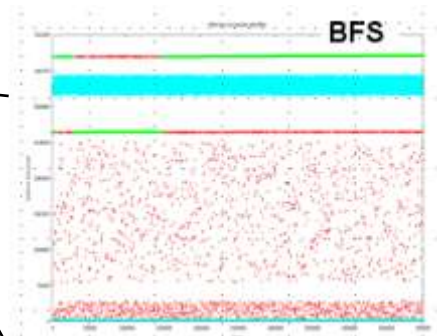
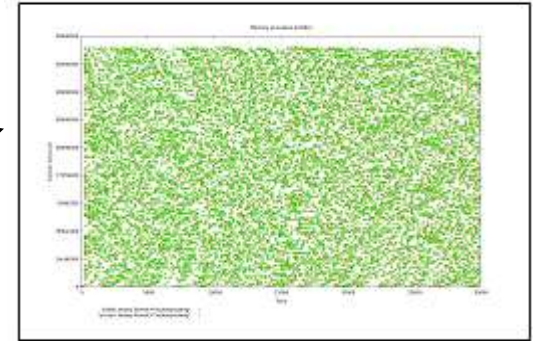
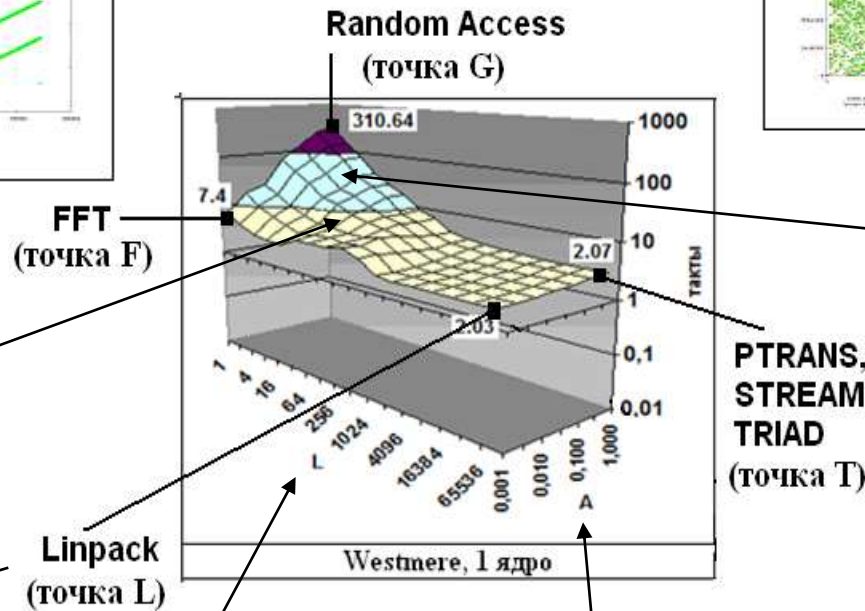
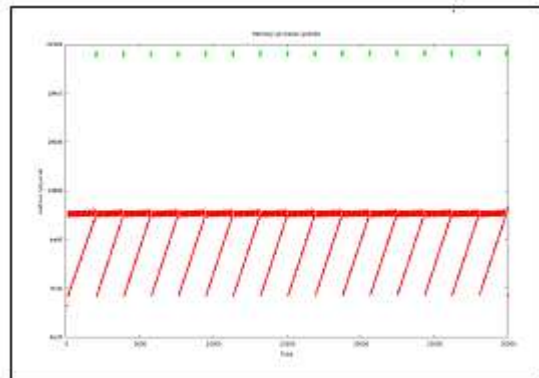
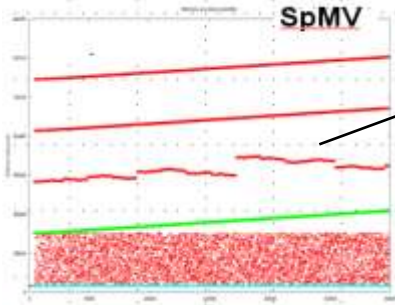
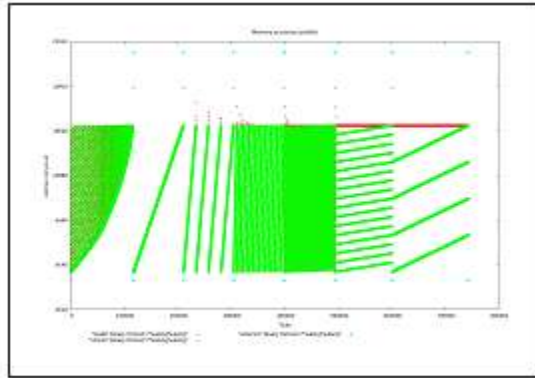
Инновационные суперкомпьютерные технологии и проблемы создания отечественной перспективной элементной базы

**Л.К.Эйсымонт, В.С.Горбунов
(21 октября 2014 года)**

Общая картина в области СКТ (на примере инновационных проектов США)

- Внедрение результатов программы DARPA HPCS (2002-2010), коммерческие образцы и военные (заказные) суперЭВМ (2013-2017)**
- Выполнение программы DARPA UNPC (2010-2020) и программ DoE по экзамасштабным технологиям и суперЭВМ экза-уровня (2012-2023)**
- Выполнение программы DARPA STARNet (с 2013 года, на 10-15 лет) по оптимизации использования КМОП-технологий и разработки технологий пост-Муровской эры, зетта- и йотта-уровень производительности**

Главная проблема “стены памяти” проекта DARPA HPCS и аналогичных проектов Японии и Китая



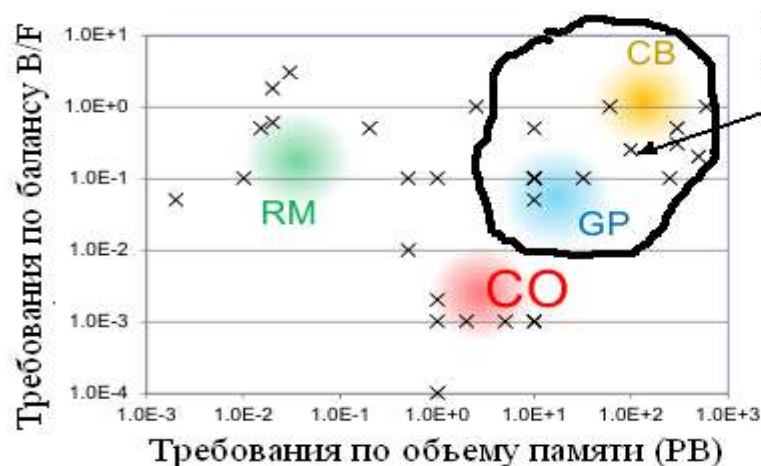
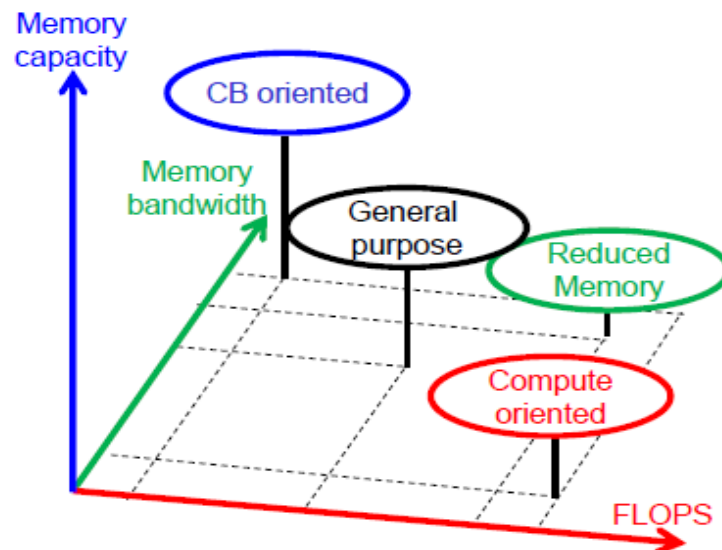
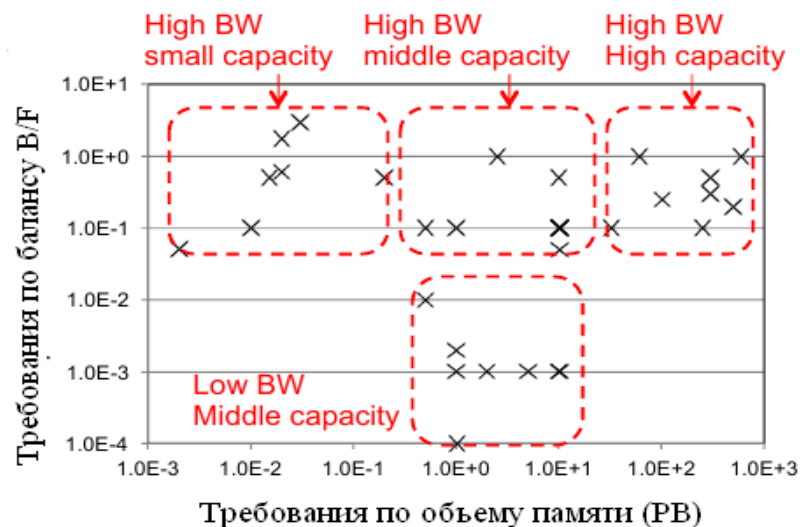
Пространственная локализация

Временная локализация

Рейтинг на вычислительных задачах - тест HP Linpack в сравнении с HPCG

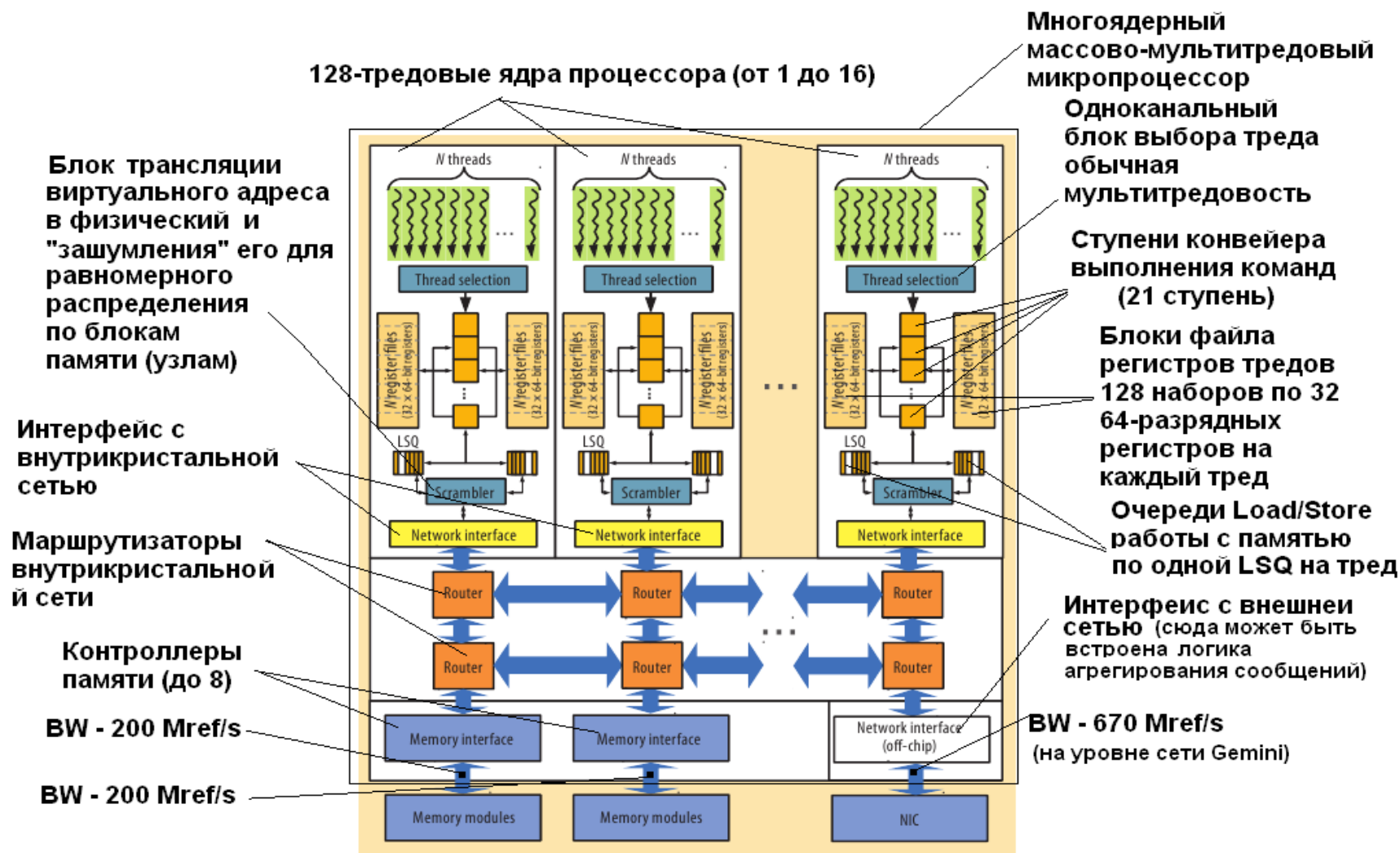
Site	Computer	Cores	HPL Rmax (Pflops)	HPL Rank	HPCG (Pflops)	HPCG/ HPL
NSCC / Guangzhou	Tianhe-2 NUDT, Xeon 12C 2.2GHz + Intel Xeon Phi 57C + Custom	3,120,000	33.9	1	.580	1.7%
RIKEN Advanced Inst for Comp Sci	K computer Fujitsu SPARC64 VIIIfx 8C + Custom	705,024	10.5	4	.427	4.1%
DOE/OS Oak Ridge Nat Lab	Titan, Cray XK7 AMD 16C + Nvidia Kepler GPU 14C + Custom	560,640	17.6	2	.322	1.8%
DOE/OS Argonne Nat Lab	Mira BlueGene/Q, Power BQC 16C 1.60GHz + Custom	786,432	8.59	5	.101 [#]	1.2%
Swiss CSCS	Piz Daint, Cray XC30, Xeon 8C + Nvidia Kepler 14C + Custom	115,984	6.27	6	.099	1.6%
Leibniz Rechenzentrum	SuperMUC, Intel 8C + IB	147,456	2.90	12	.0833	2.9%
CEA/TGCC-GENCI	Curie tine nodes Bullx B510 Intel Xeon 8C 2.7 GHz + IB	79,504	1.36	26	.0491	3.6%
Exploration and Production Eni S.p.A.	HPC2, Intel Xeon 10C 2.8 GHz + Nvidia Kepler 14C + IB	62,640	3.00	11	.0489	1.6%
DOE/OS L Berkeley Nat Lab	Edison Cray XC30, Intel Xeon 12C 2.4GHz + Custom	132,840	1.65	18	.0439 [#]	2.7%
Texas Advanced Computing Center	Stampede, Dell Intel (8c) + Intel Xeon Phi (61c) + IB	78,848	.881*	7	.0161	1.8%
Meteo France	Beaufix Bullx B710 Intel Xeon 12C 2.7 GHz + IB	24,192	.469 (.467*)	79	.0110	2.4%
Meteo France	Prolix Bullx B710 Intel Xeon 2.7 GHz 12C + IB	23,760	.464 (.415*)	80	.00998	2.4%
U of Toulouse	CALMIP Bullx DLC Intel Xeon 10C 2.8 GHz + IB	12,240	.255	184	.00725	2.8%
Cambridge U	Wilkes, Intel Xeon 6C 2.6 GHz + Nvidia Kepler 14C + IB	3584	.240	201	.00385	1.6%
TiTech	TUSBAME-KFC Intel Xeon 6C 2.1 GHz + IB	2720	.150	436	.00370	2.5%

Классификация суперкомпьютеров (далее будем использовать)



Наиболее актуальна
проблема "стены памяти"

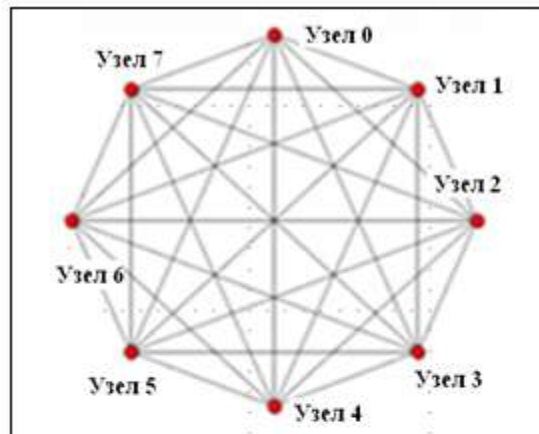
Предполагаемый базовый микропроцессор заказных суперкомпьютеров СВ-класса - развитие Threadstorm (Cray XMT)



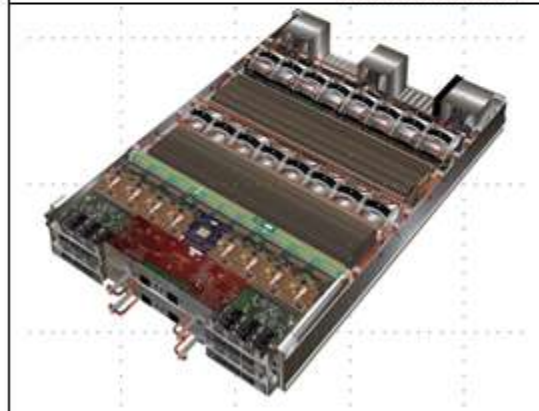
Петафлопс и траспетафлопс

(США, Япония, Китай – инновационные СКТ
Россия – эволюционные СКТ)

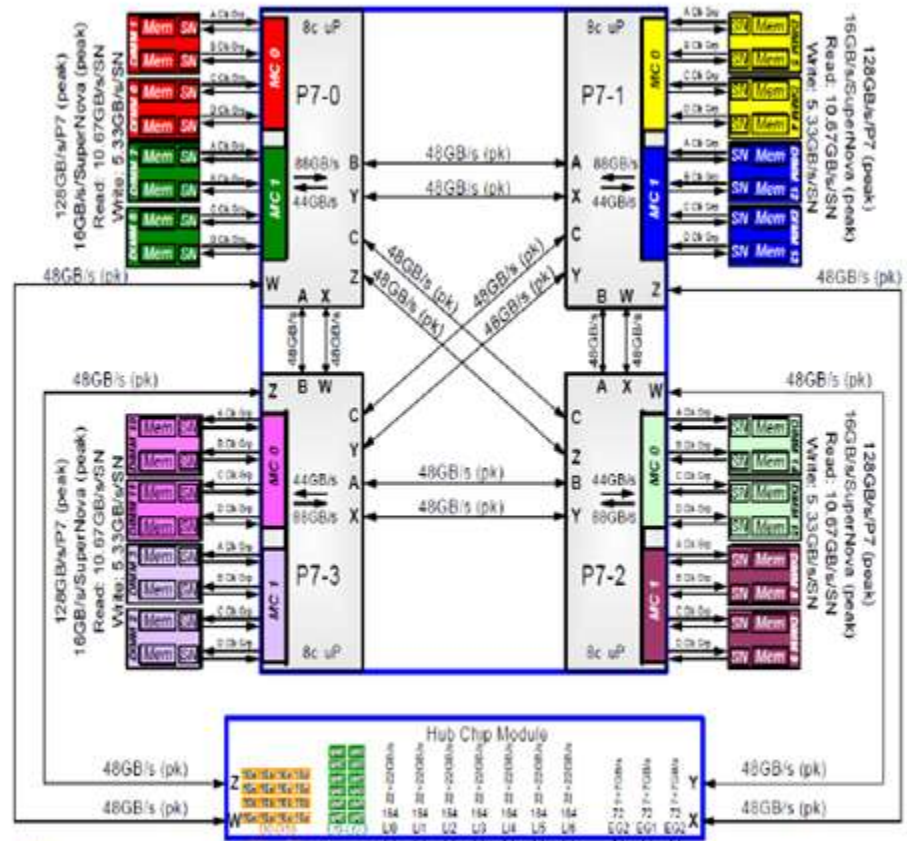
Узел и макроузел IBM Power 775



Внутренняя сеть макроузла

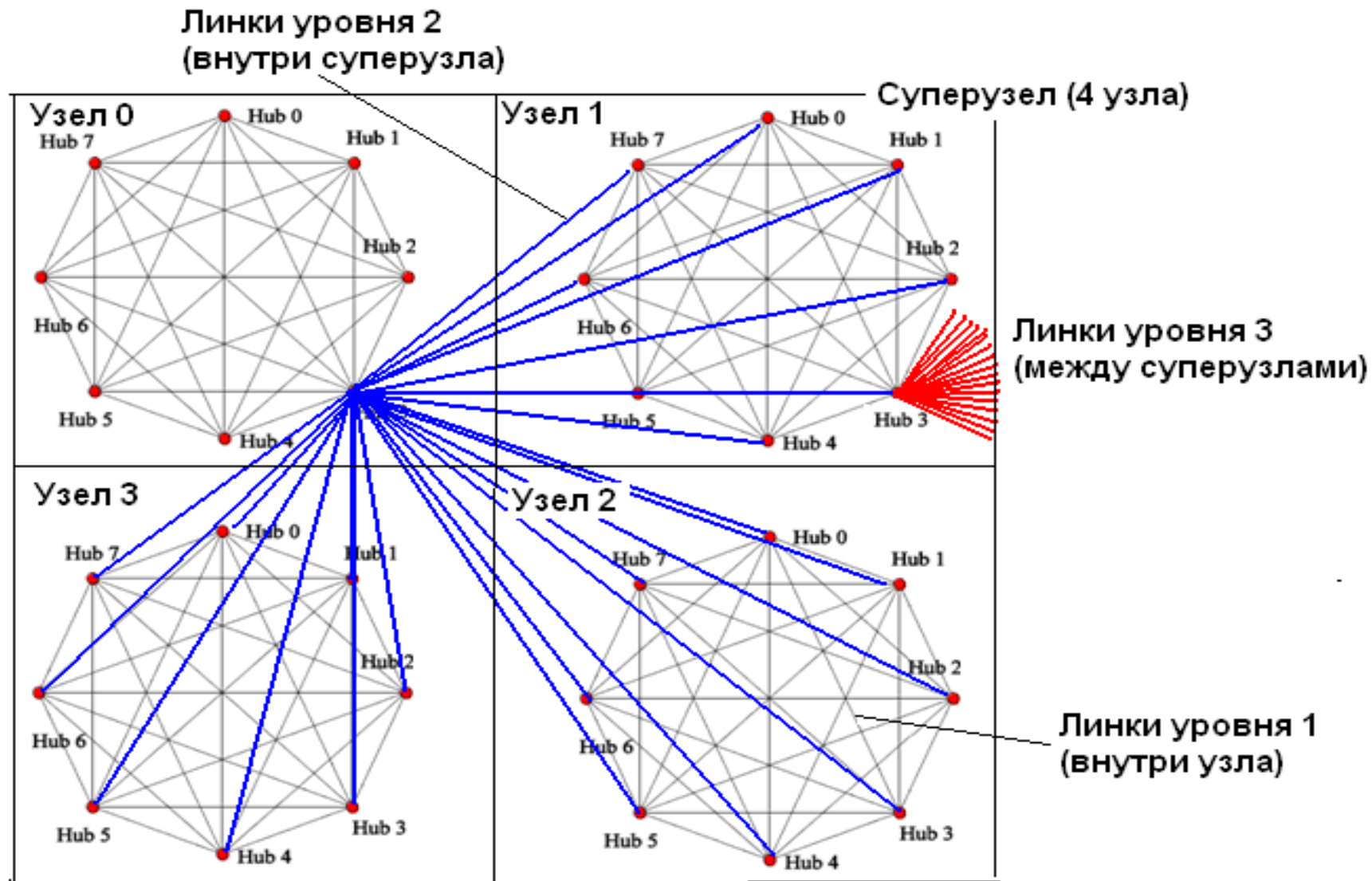


Внешний вид макроузла

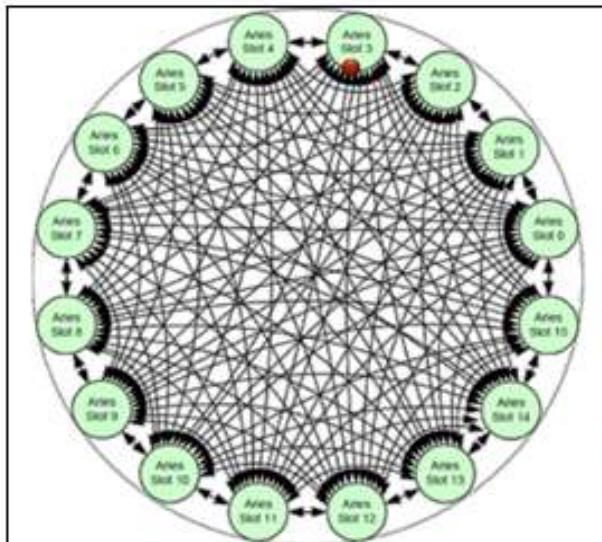


4-х процессорный узел (4 Power7 + HUB)

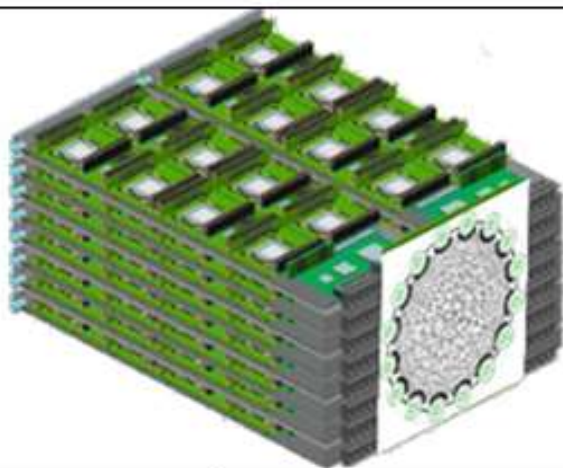
Многоуровневая сеть PERCS суперкомпьютера Power 775



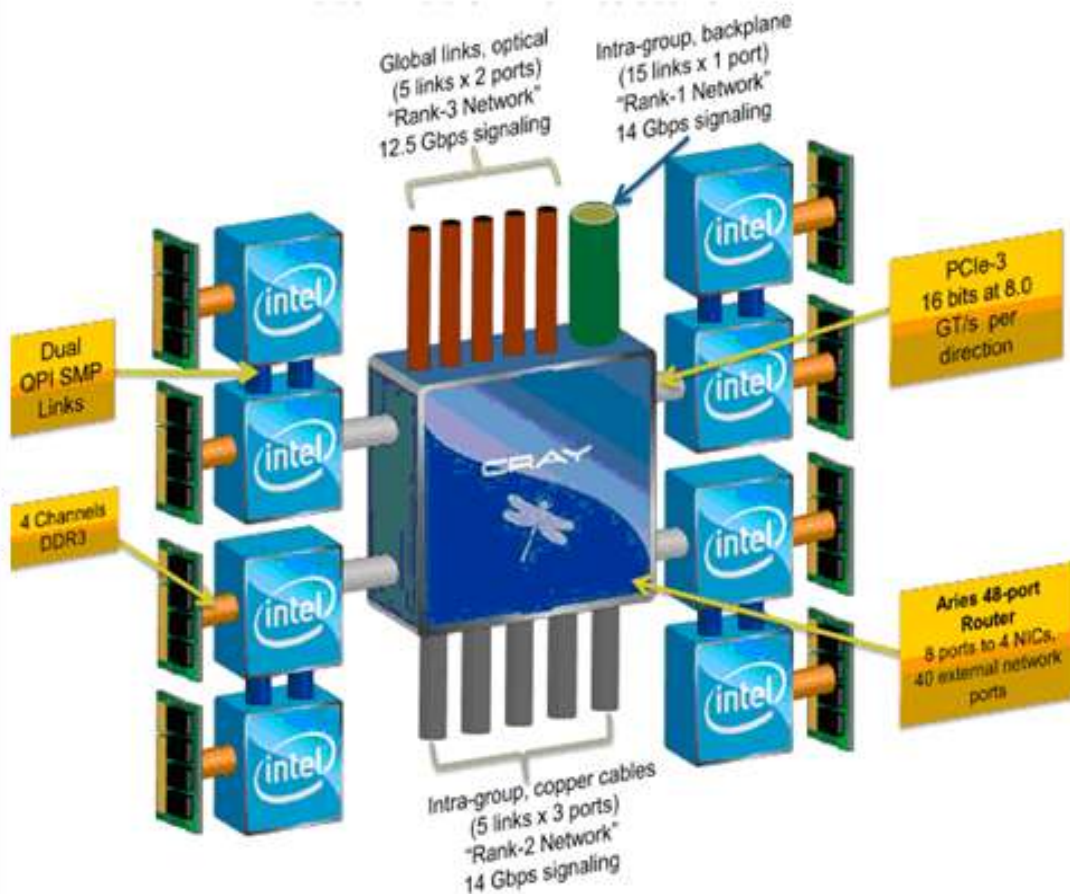
Узел и макроузел Cray XC30



Внутренняя сеть макроузла

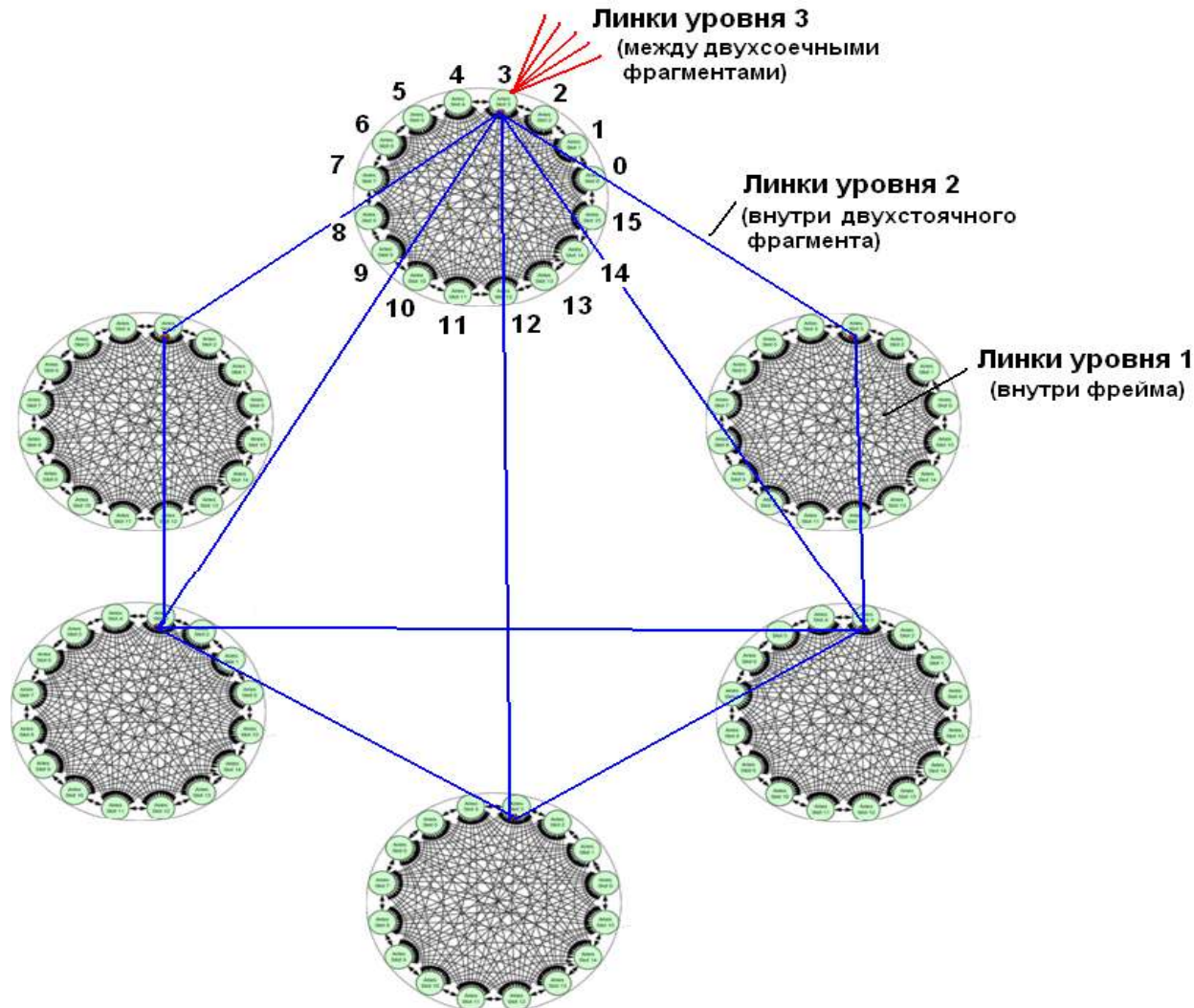


Внешний вид макроузла



8-ми процессорный узел (8 Intel + YARC)

Многоуровневая сеть суперкомпьютера Cray XC30



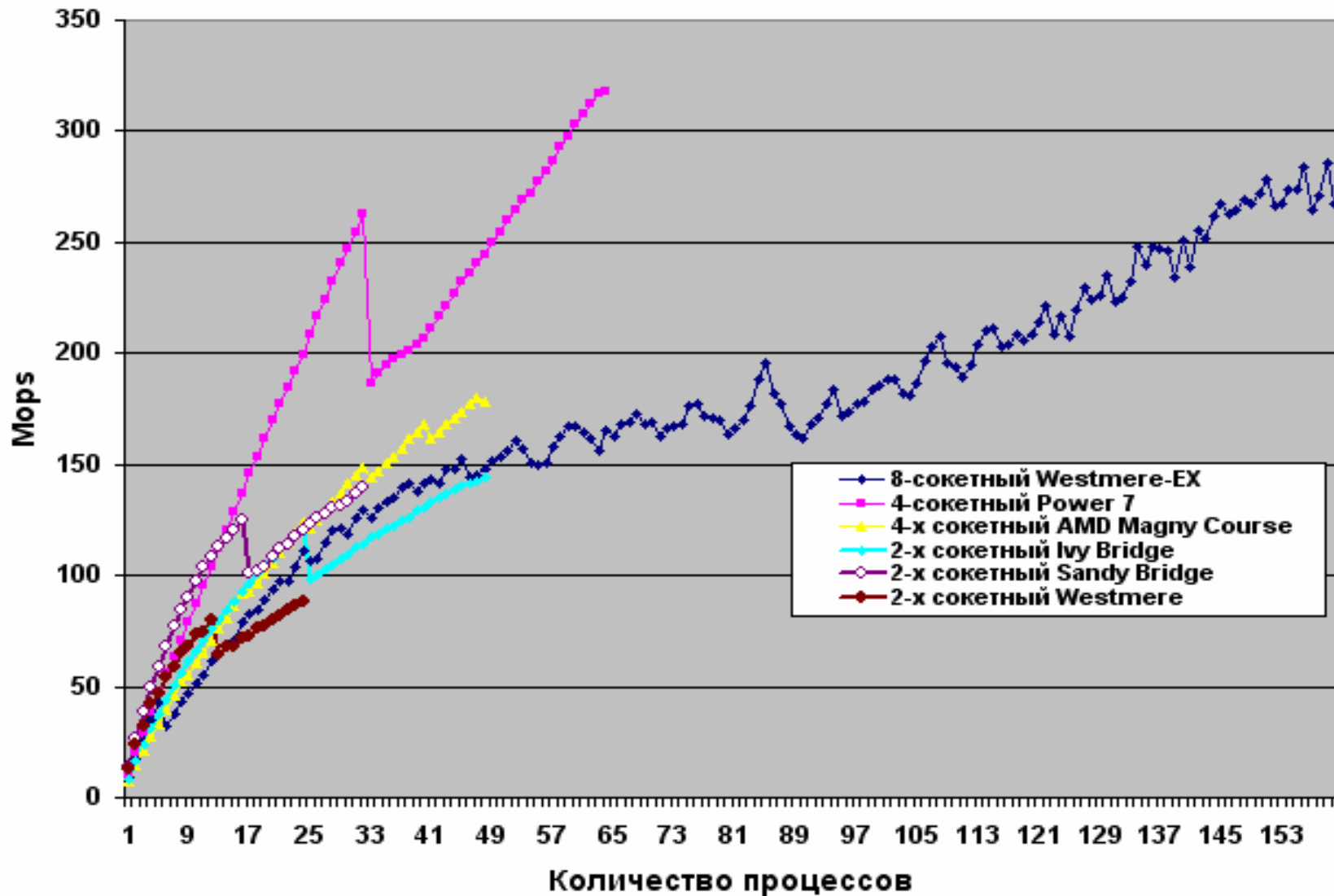
Реальная история и перспектива суперкомпьютеров NERSC (LBNL)

System attributes	NERSC-6	NERSC-7	NERSC-8 (proposed)	NERSC-9 (Proposed)
	Hopper	Edison	Cori 2015-2016	? 2019-2020
System peak	1.3 PF	2.6PF	20-40PF	200-300 PF
Power	2.9 MW (Peak) 2.2MW (Typical)	2.3 MW (Peak) 1.6 MW (Typical)	<5 MW (Peak)	< 15 MW (peak)
System memory	0.21 PB	0.35 PB	1-2 PB	~10 PB (128 GB on package, 512-1024 GB DRAM)
Node performance	202GF	460 GF	2-3.5TF	~10 TF
Node memory BW	50 GB/s	90 GB/s	100-500 GB/s	~200 GB/s ? 2-4 TB/s on package
Node concurrency	24 AMD Magnycours cores	24 Intel Ivy Bridge Cores	up to 300 Knight Landing Haswell	Up to 2048
System size (nodes)	6,384 nodes	5,576 nodes	8,000-12,000 nodes	O(10,000)
MPI Node Interconnect BW	~3 GB/s	~9GB/s	~9 GB/s	Up to 50 GB/s

**IBM p775: один QSM
(4-х процессорный узел)**

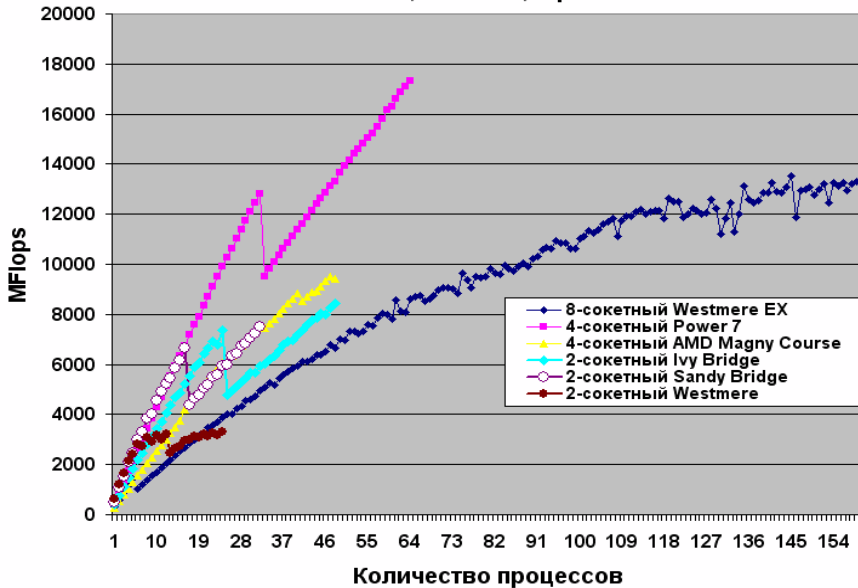
IBM Power 775 - тест UA (класс C)

Тест UA, класс C, OpenMP

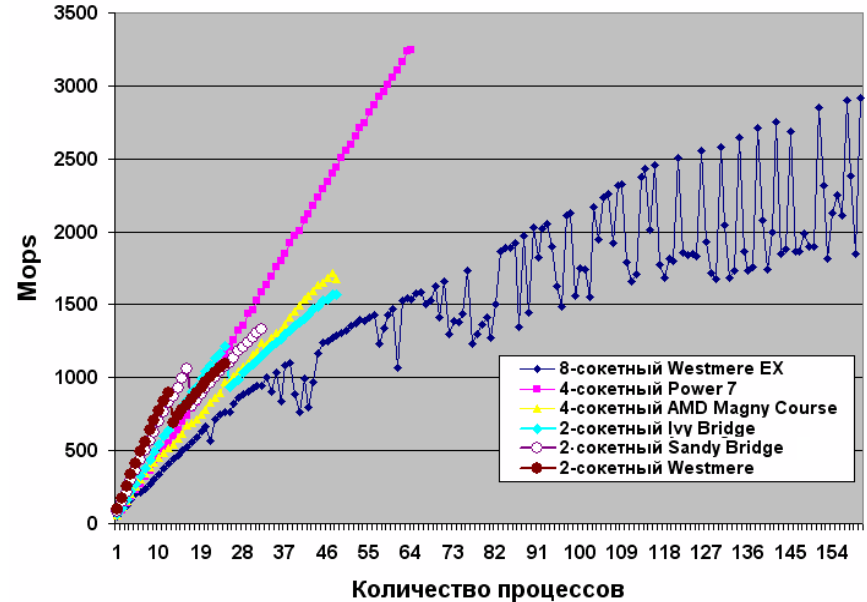


IBM Power 775 - тесты CG, IS, MG, BT

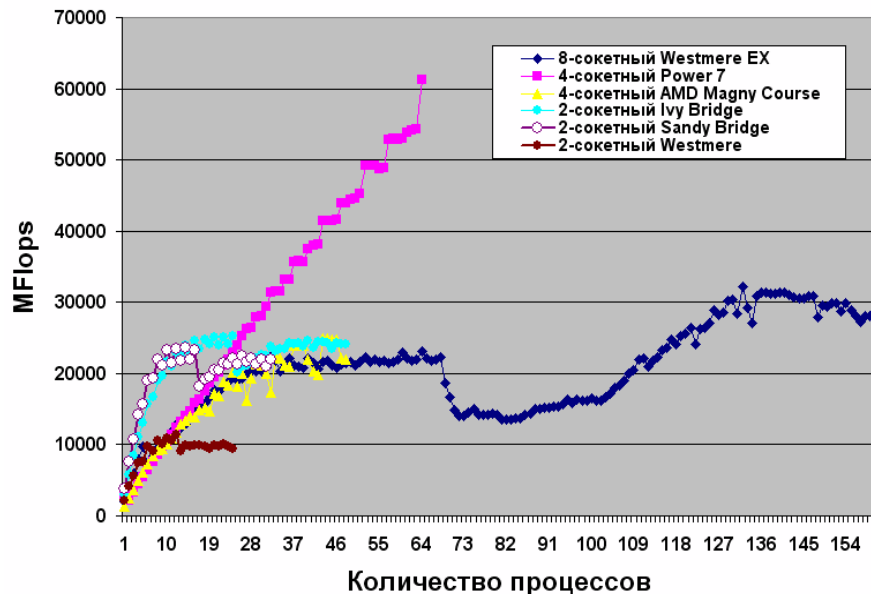
Тест CG, класс C, OpenMP



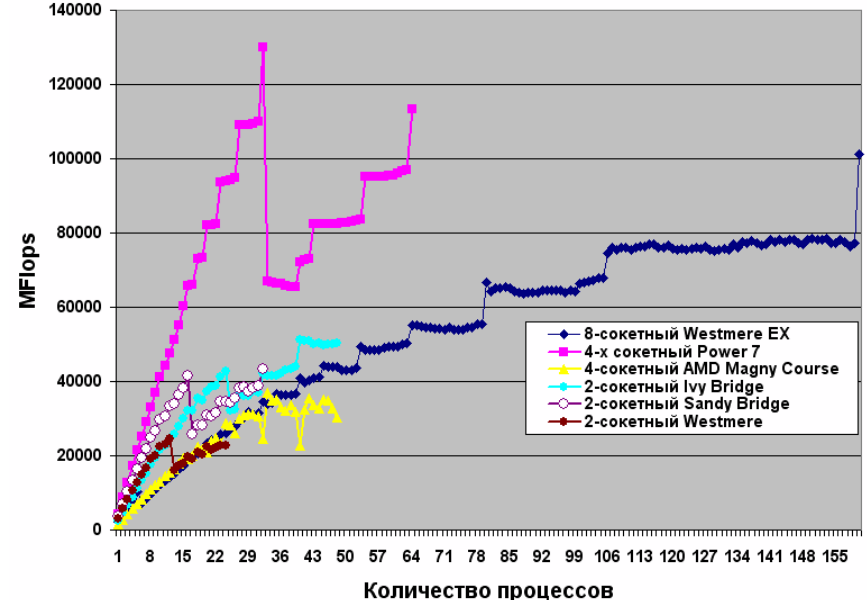
Тест IS, класс C, OpenMP



Тест MG, класс C, OpenMP



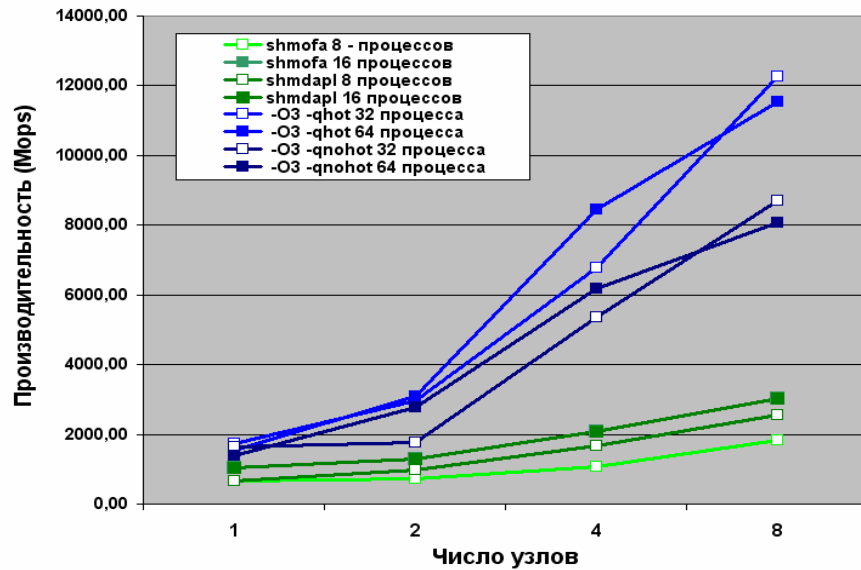
Тест BT, класс C, OpenMP



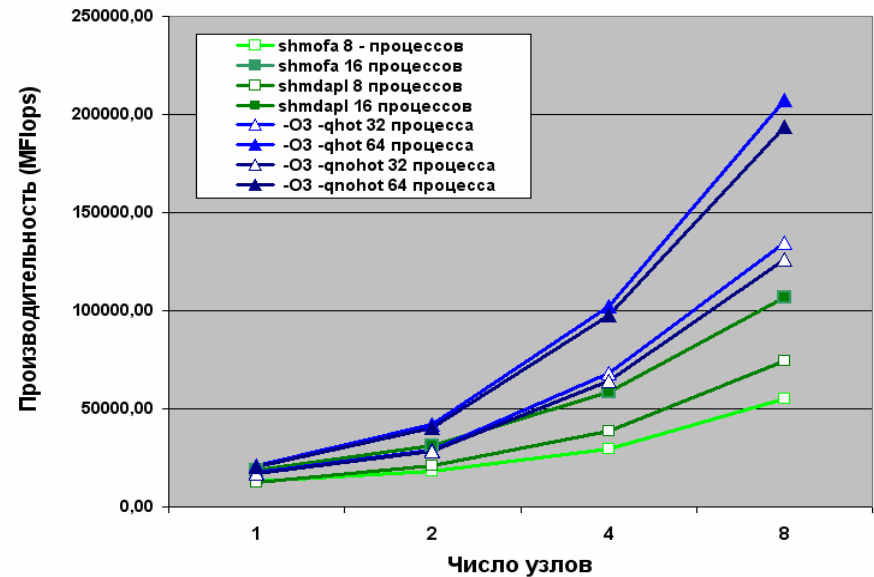
**IBM p775: много QSM одной
серверной платы
(на ней до 8 QSM, это
макроузел первого уровня)**

IBM Power 775 - tests IS, FT, MG, LU

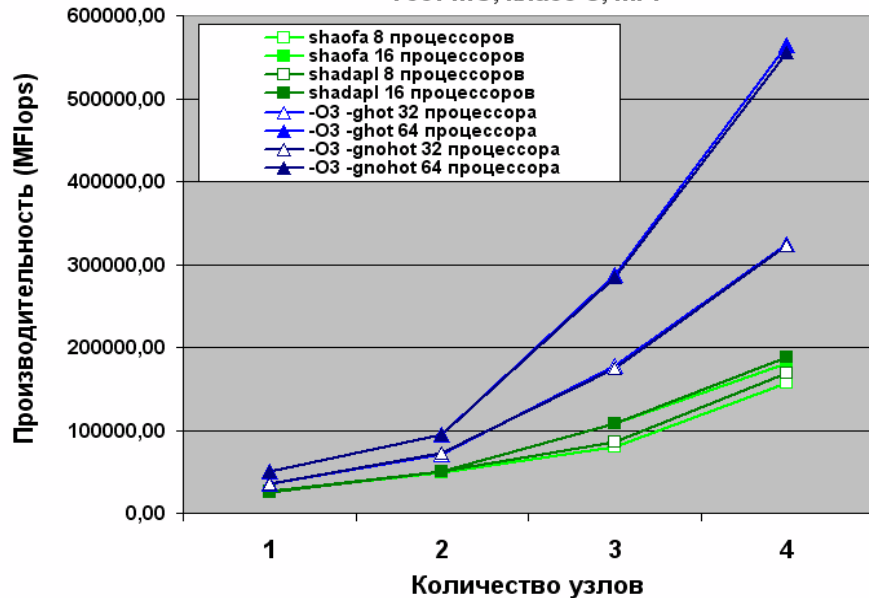
Тест IS, класс C, MPI



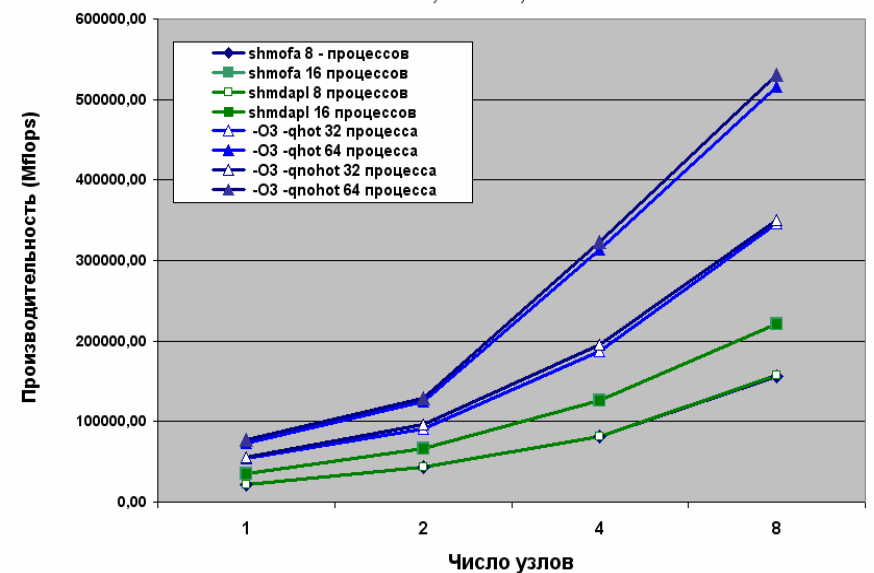
Тест FT, класс C, MPI



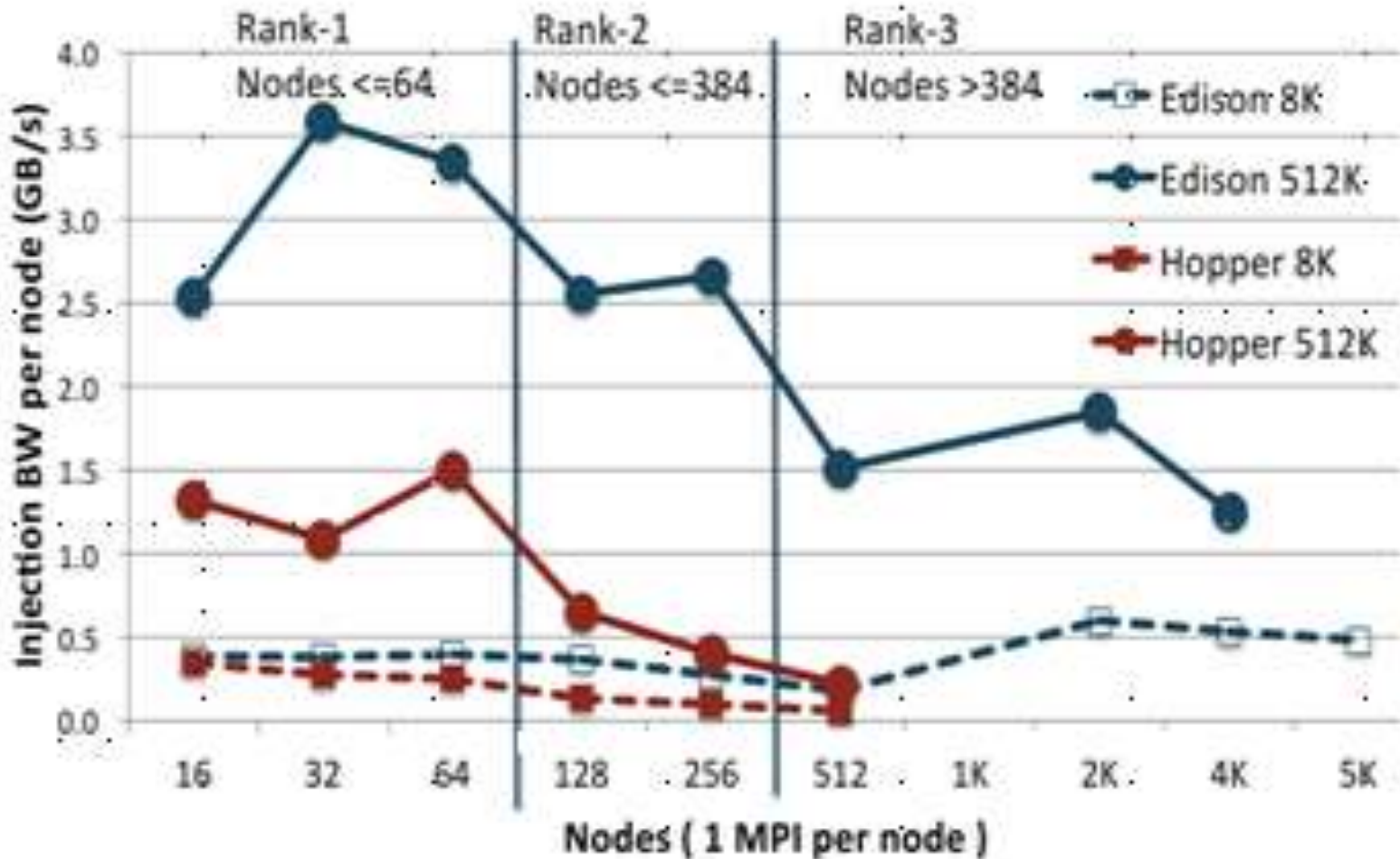
Тест MG, класс C, MPI



Тест LU, класс C, MPI

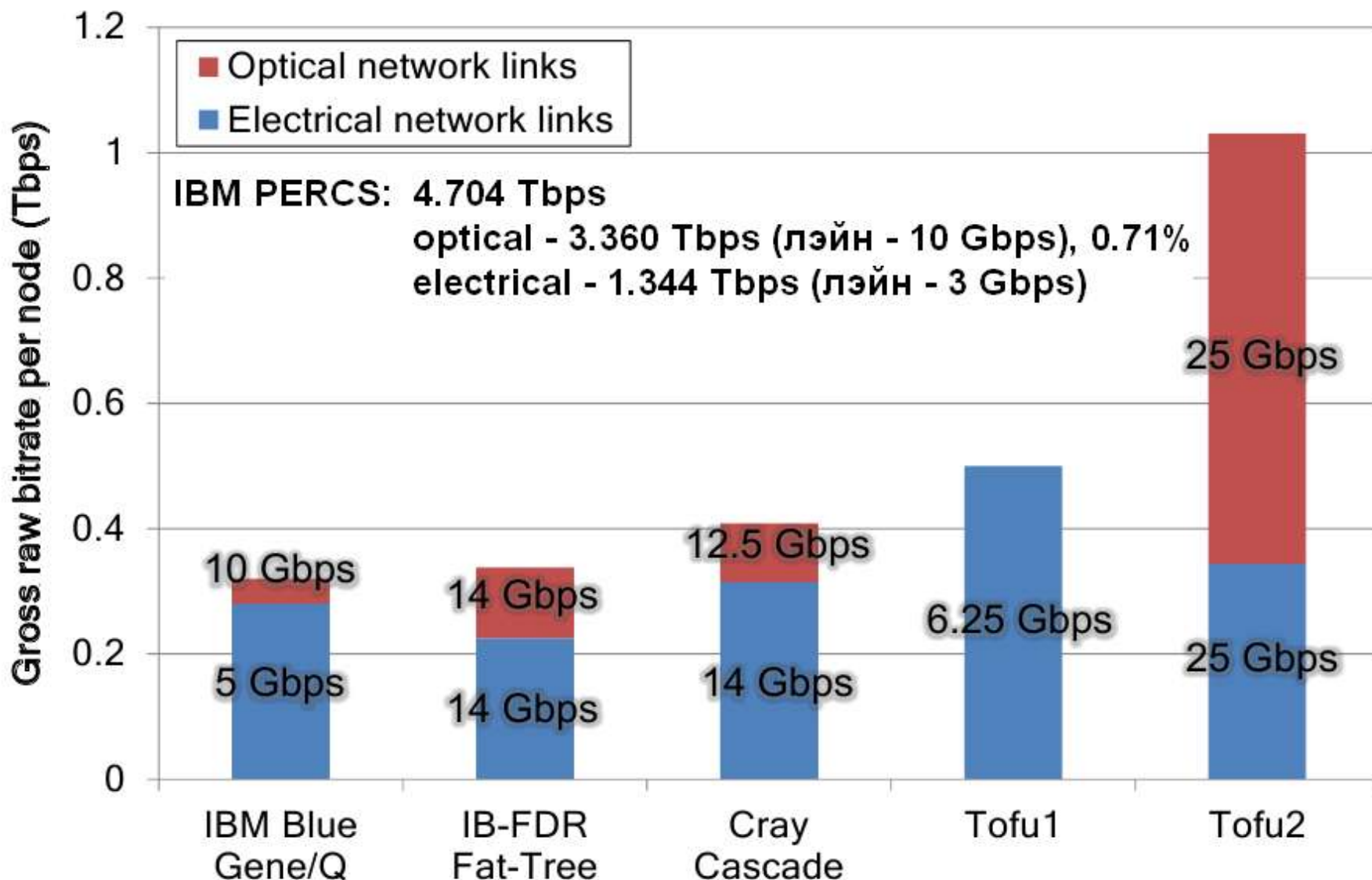


Cray XC30 (Edison, сеть Aries) эффективность Alltoall в сравнении с Cray XE6 (Hopper, сеть Gemini, 3D-top)

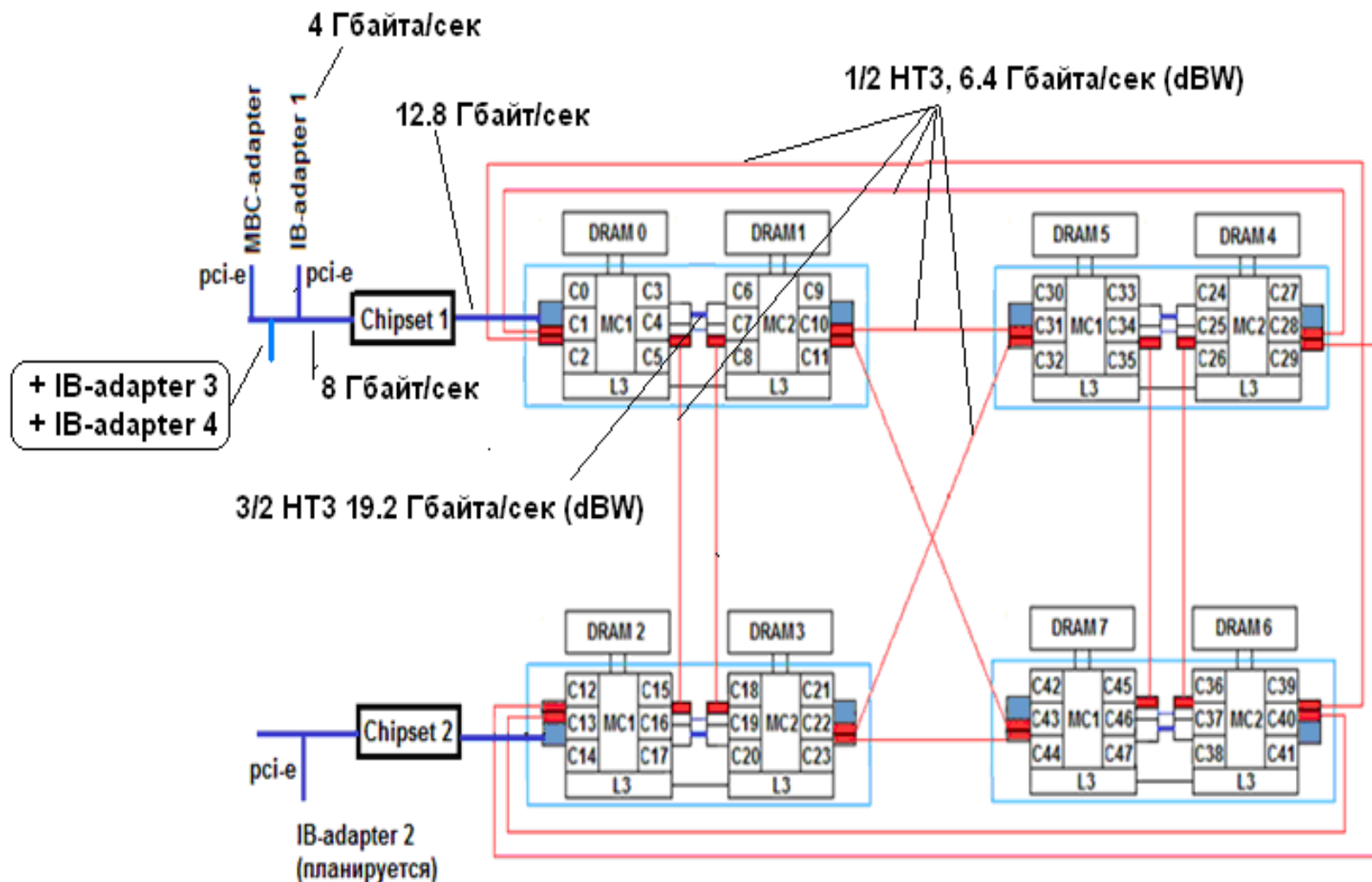


“Петафлопсные возможности” эволюционных СКТ

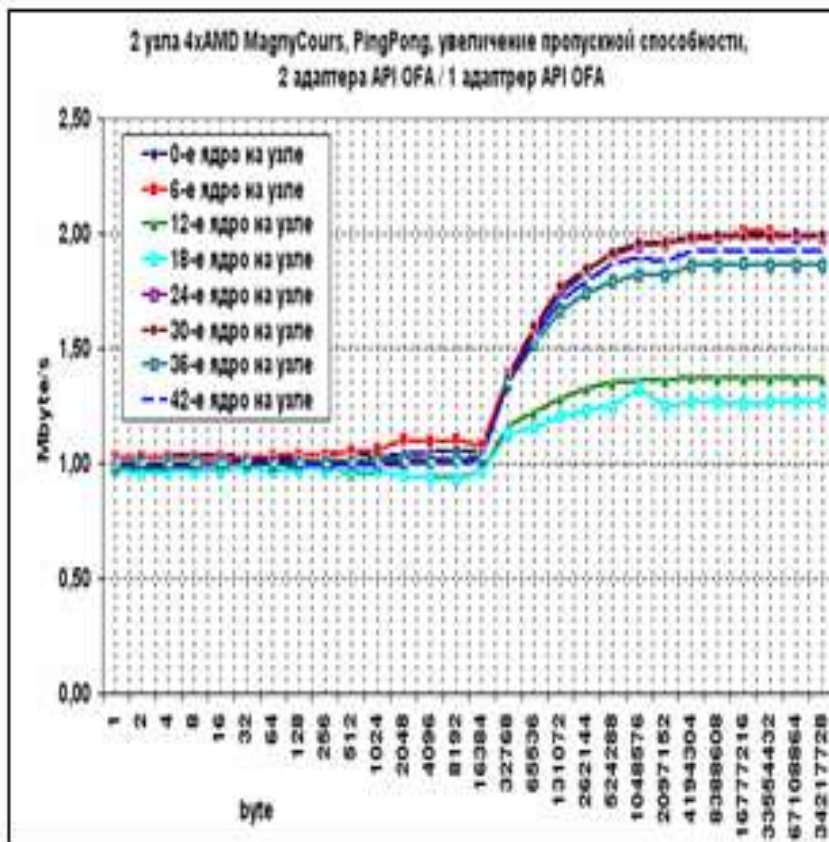
Сравнение пропускных способностей лэйнов, маршрутизаторов



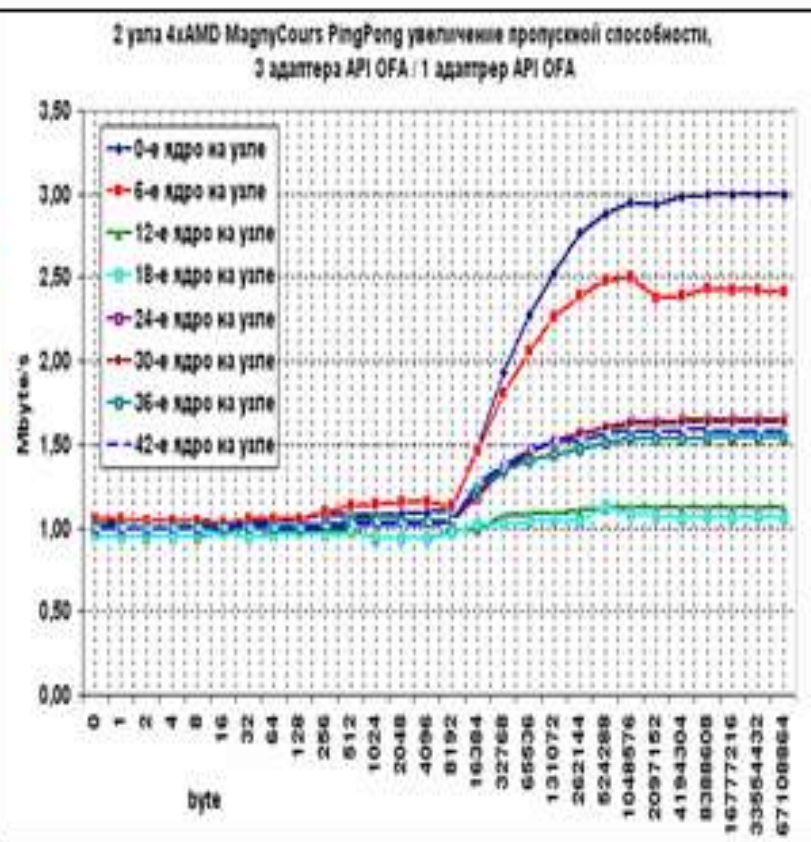
Серверная плата ПТК (СПБГПУ)



Ускорение операций “точка-точка” на 2-х и 3-х адаптерной установке

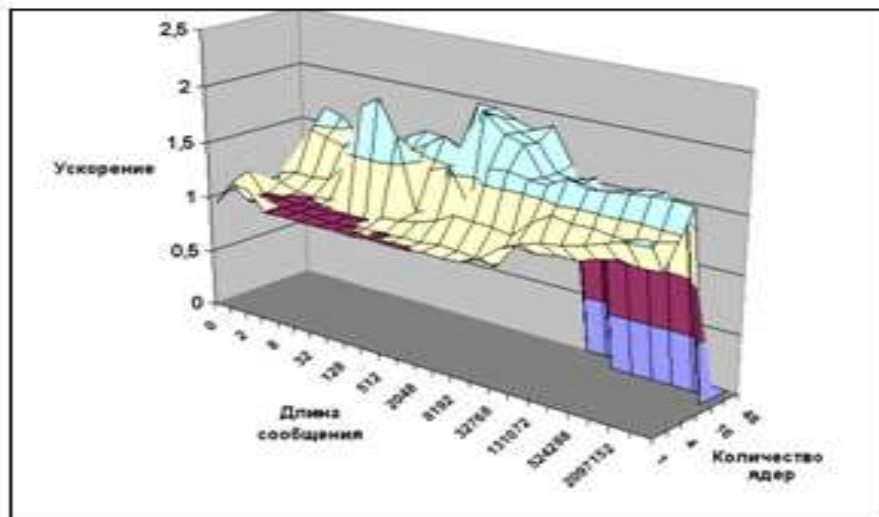


(A) Межузловые взаимодействия, Mellanox, PMPI+ofa, 1 и 2 адаптера на узле.

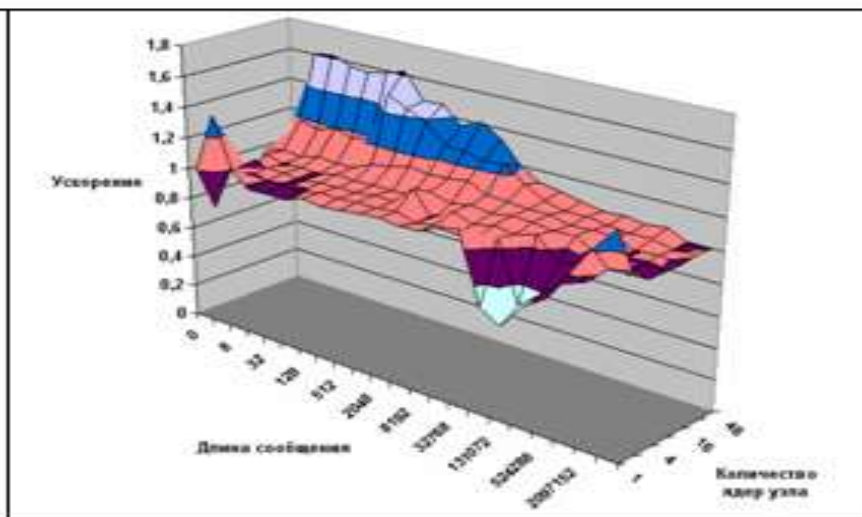


(B) Межузловые взаимодействия, Mellanox, PMPI+ofa, 1 и 3 адаптера на узле.

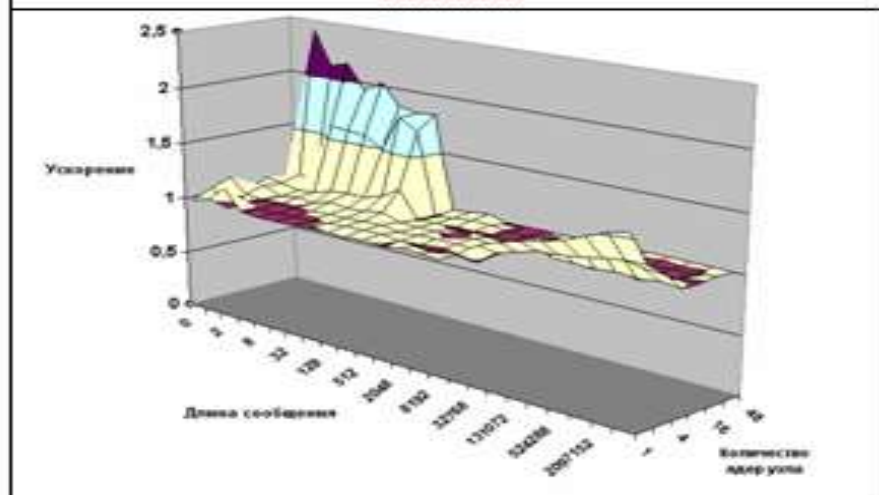
Ускорение коллективных операций на 2-х адаптерной установке



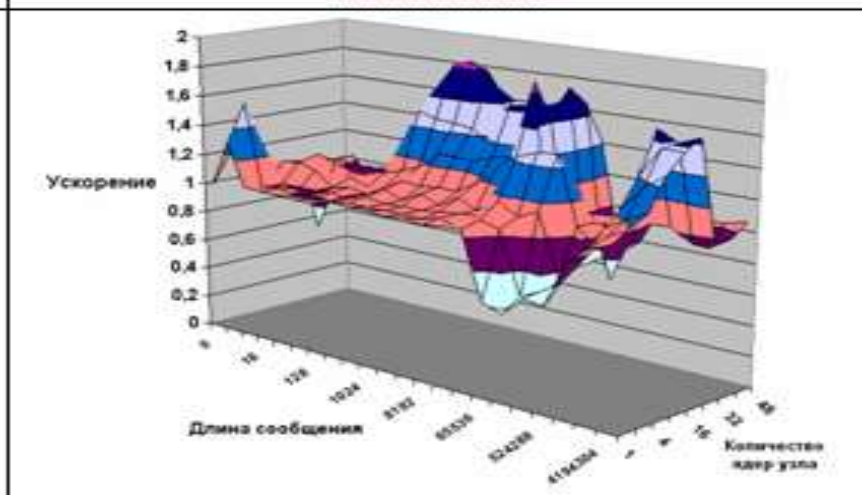
Alltoall



Allreduce

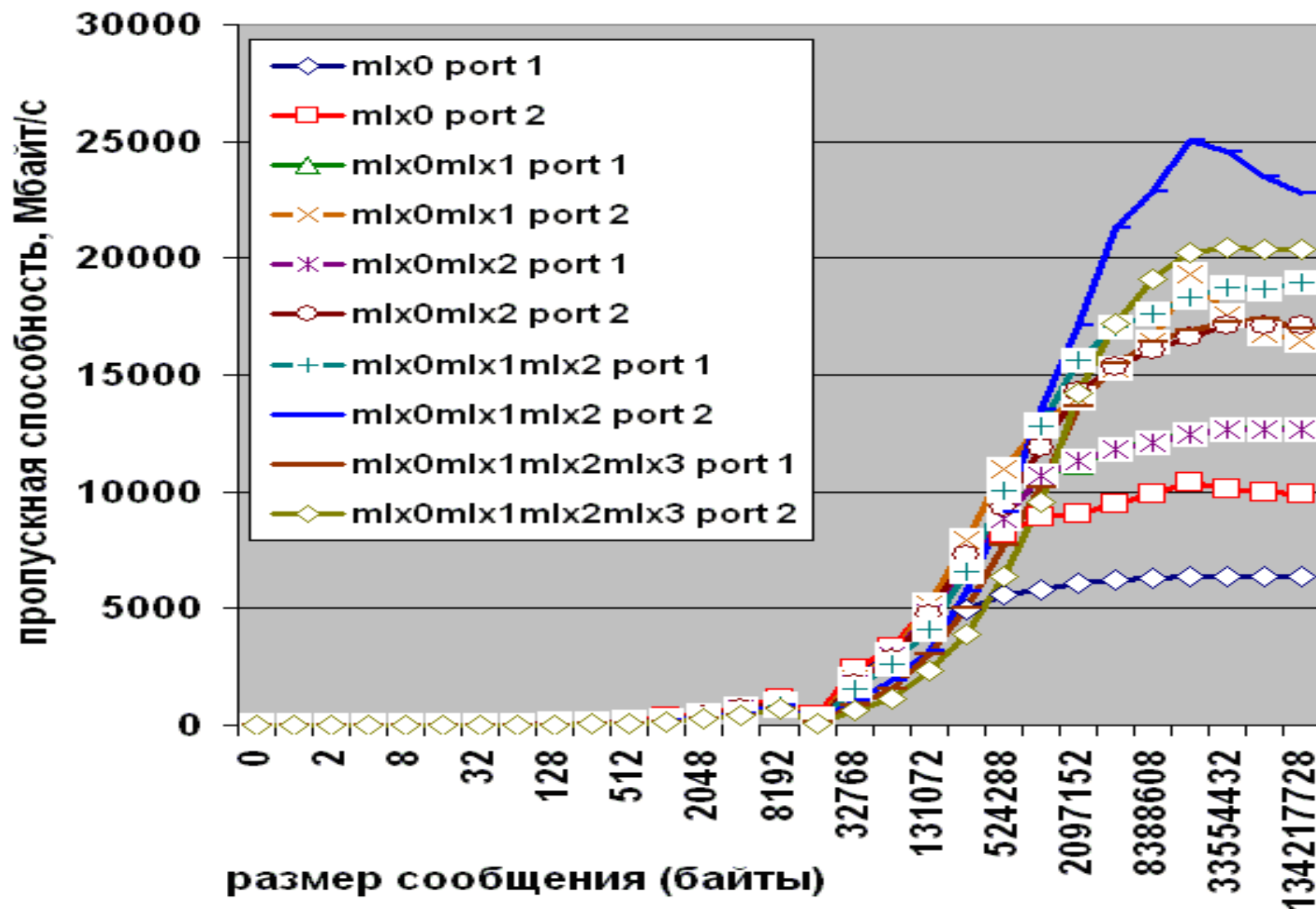


Allgather



ReduceScatter

PingPong на 2-х, 3-х и 4-х адаптерной установке (2-х портовые адаптеры Connect-IB)



Эксафлопс

(эволюционные и инновационные СКТ)

Россия – только эволюционные СКТ

?

10 актуальных проблем разработки экзафлопсных суперкомпьютеров в США

1. Энергоэффективность - создание энергоэффективных схем процессоров, маршрутизаторов коммуникационных сетей, технологий обеспечения питанием и технологий охлаждения.

2. Технологии соединений как внутри вычислительных узлов, так и между вычислительными узлами – увеличение производительности и сокращение задержек передачи данных, достижение энергоэффективности в линиях связи коммуникационных сетей и интерфейсах, наиболее важный показатель – снижение энергопотребления и повышения быстродействия обращений к памяти удаленных узлов (RDMA).

3. Технологии оперативной памяти – интеграция новых улучшенных технологий памяти для повышения емкости при повышении плотности размещения запоминающих элементов и сокращения их стоимости, снижение задержек выполнения операций с памятью, повышение их сложности и локализации, повышение пропускной способности даже для обращений с большой мелкозернистостью, т.е. когда обращения происходят к небольшим участкам памяти, а не к блокам большого объема.

10 актуальных проблем разработки экзафлопсных суперкомпьютеров в США

4. Создание масштабируемого (при увеличении параллелизма) системного программного обеспечения в виде операционных систем нового типа и систем поддержки выполнения программ (run-time систем), обеспечивающего высокий параллелизм уровня 10^9 (основной прирост, до трех порядков, ожидается непосредственно внутри вычислительного узла, а межузловой – на порядок), энергоэффективность за счет глубокого проникновения в управление работой оборудования и отказоустойчивость.

5. Системы программирования для пользователей – создание новых систем программирования, которые обеспечивают: создание эффективных с массовым параллелизмом программ; прозрачную для пользователя работу с иерархической глобально адресуемой памятью с обеспечением как эффективных удаленных обращений к памяти, так и локализацию данных при вычислениях и вычислений при данных; высокую многоуровневую отказоустойчивость.

6. Управление данными – создание хранилищ данных, программного обеспечения, которые бы справились с объемами и интенсивностью поступающих данных, их обработкой и хранением, ожидаемым разнообразием типов данных

10 актуальных проблем разработки экзафлопсных суперкомпьютеров в США

7. Создание экзамасштабных алгоритмов – переформулирование научных проблем и реконструирование или переработка алгоритмов их решения с целью эффективного выполнения на создаваемых суперкомпьютерах экзафлопсного класса.

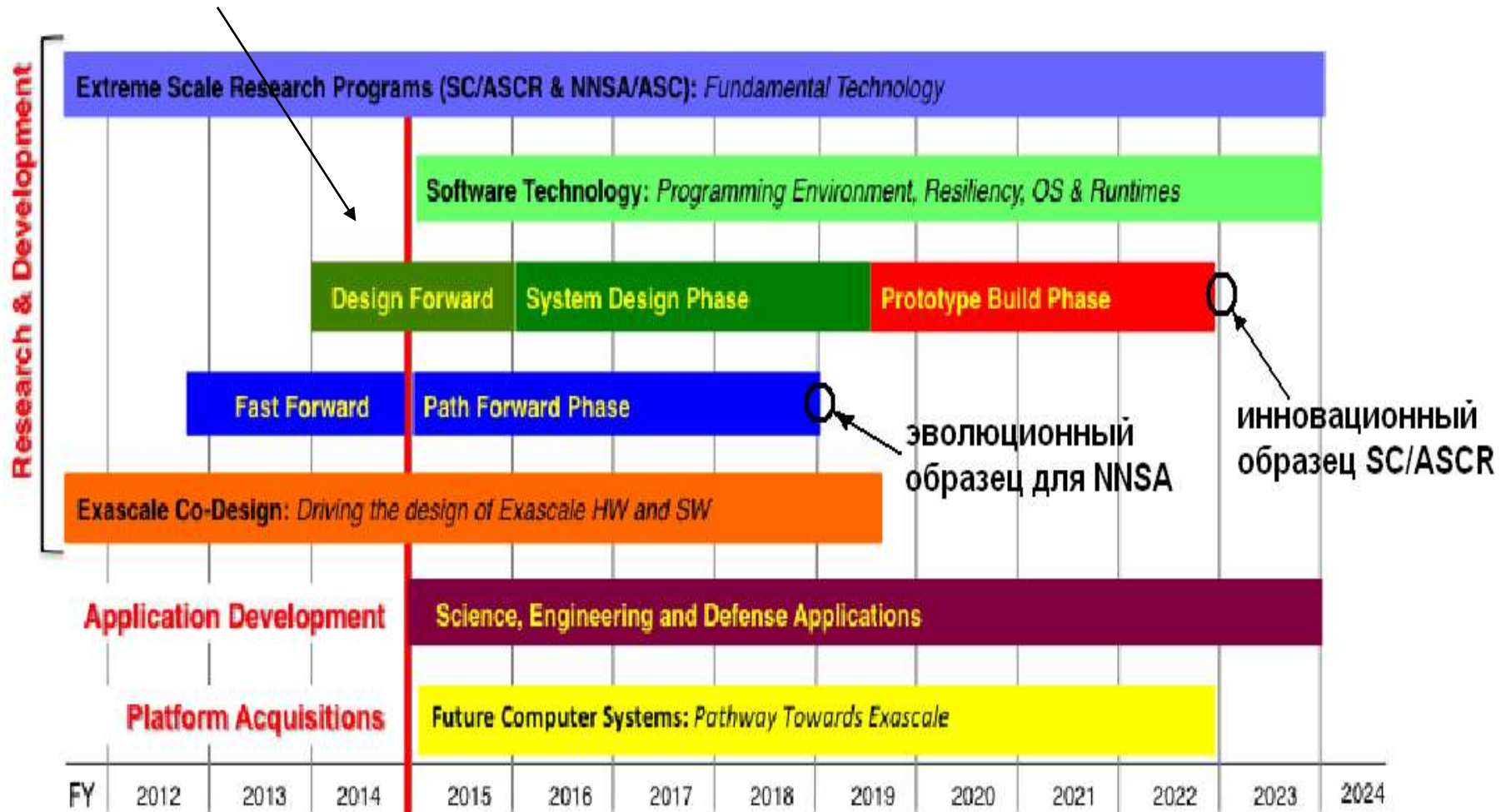
8. Создания алгоритмов автоматизации принятия решений по получаемым на экзафлопсных суперкомпьютерах результатам, что связано с оптимизацией принятия инженерных решений при создании сложных технических изделий, извлечением знаний из результатов научных расчетов и поступающих данных от сенсоров и физических установок.

9. Обеспечение как отказоустойчивости, так и достоверности вычислений в условиях наличия сбоев и отказов оборудования, ошибок программ и информационных не повторяемости результатов вычислений.

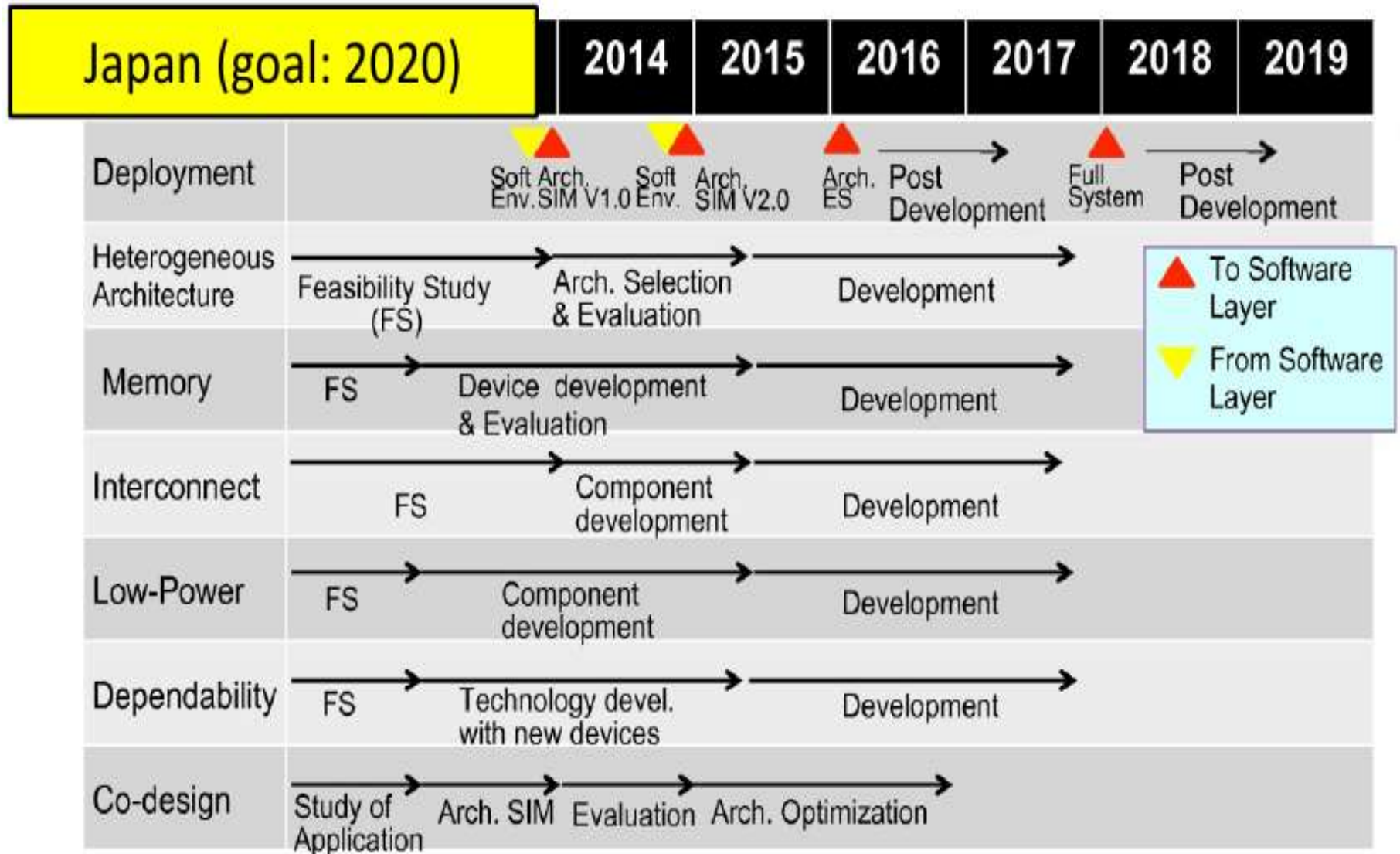
10. Повышение продуктивности разработки прикладных программ.

Дорожная карта создания экзафлопсных суперкомпьютеров DoE США

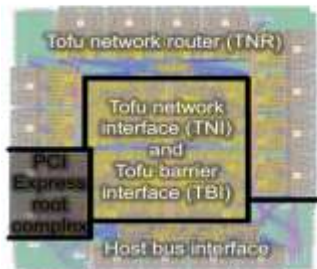
- ORNL Titan, Cray XK7, 27 PF
- LLNL Sequoia, IBM BG/Q 20 PF
- ANL Mira, IBM BG/Q, 10 PF
- LBNL Edison, Cray XC30, 2 PF
- LANL Cielo, Cray XE6, 1.1 PF



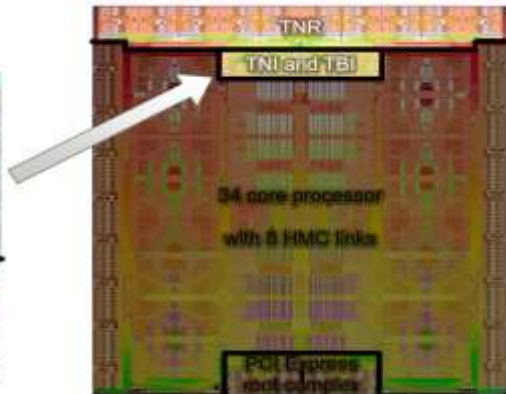
Дорожная карта создания экзафлопсных суперкомпьютеров Японии



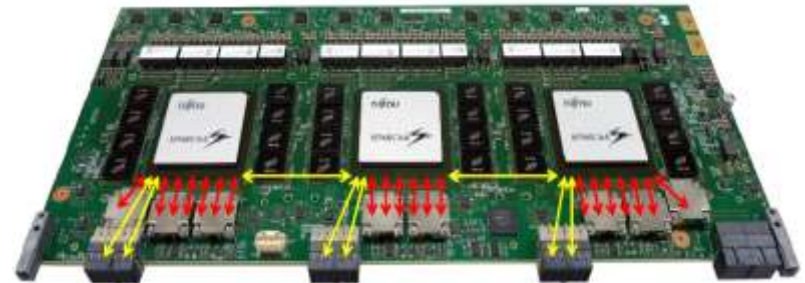
Новый вариант К-компьютера (Fujitsu, 2015 год)



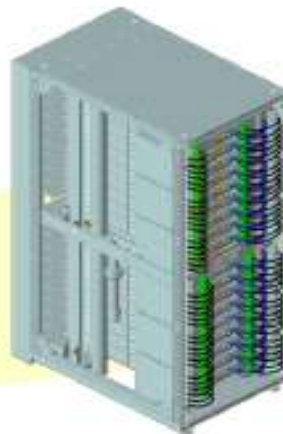
Tofu1 < 300mm²
InterConnect Controller (65nm)



Tofu2 < 100mm² – SPARC64™ Xlfx (20nm)



2U chassis
12 nodes

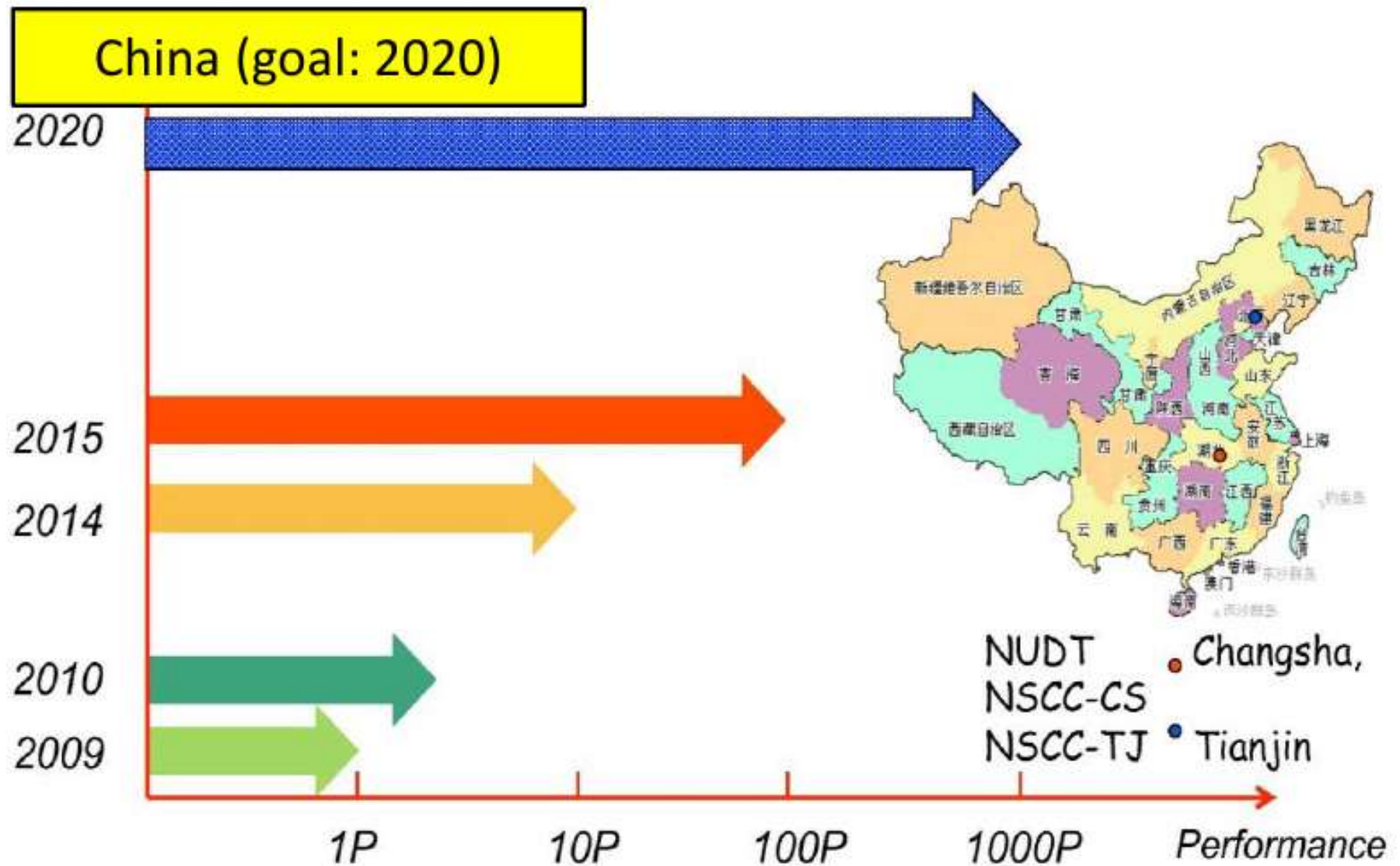


19-inch rack
216 nodes



Post-FX10 System
Petaflops per 5 racks

Дорожная карта создания экзафлопсных суперкомпьютеров Китая



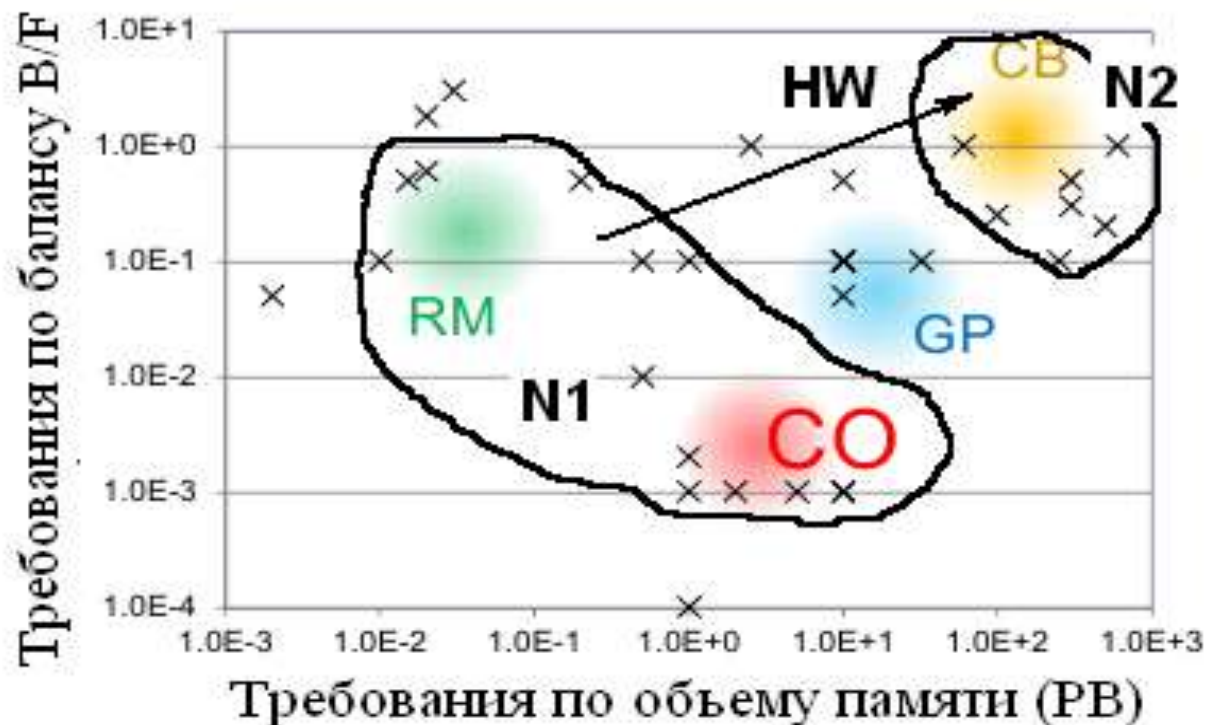
Участники работ по экзафлопсной тематике в Европе

European Union (goal: 2020)



**Российские
инновационные СКТ
(примеры из проекта
МГВС РАН)**

Предложения по стратегии и тактике импортозамещения в контексте развития инновационных СКТ

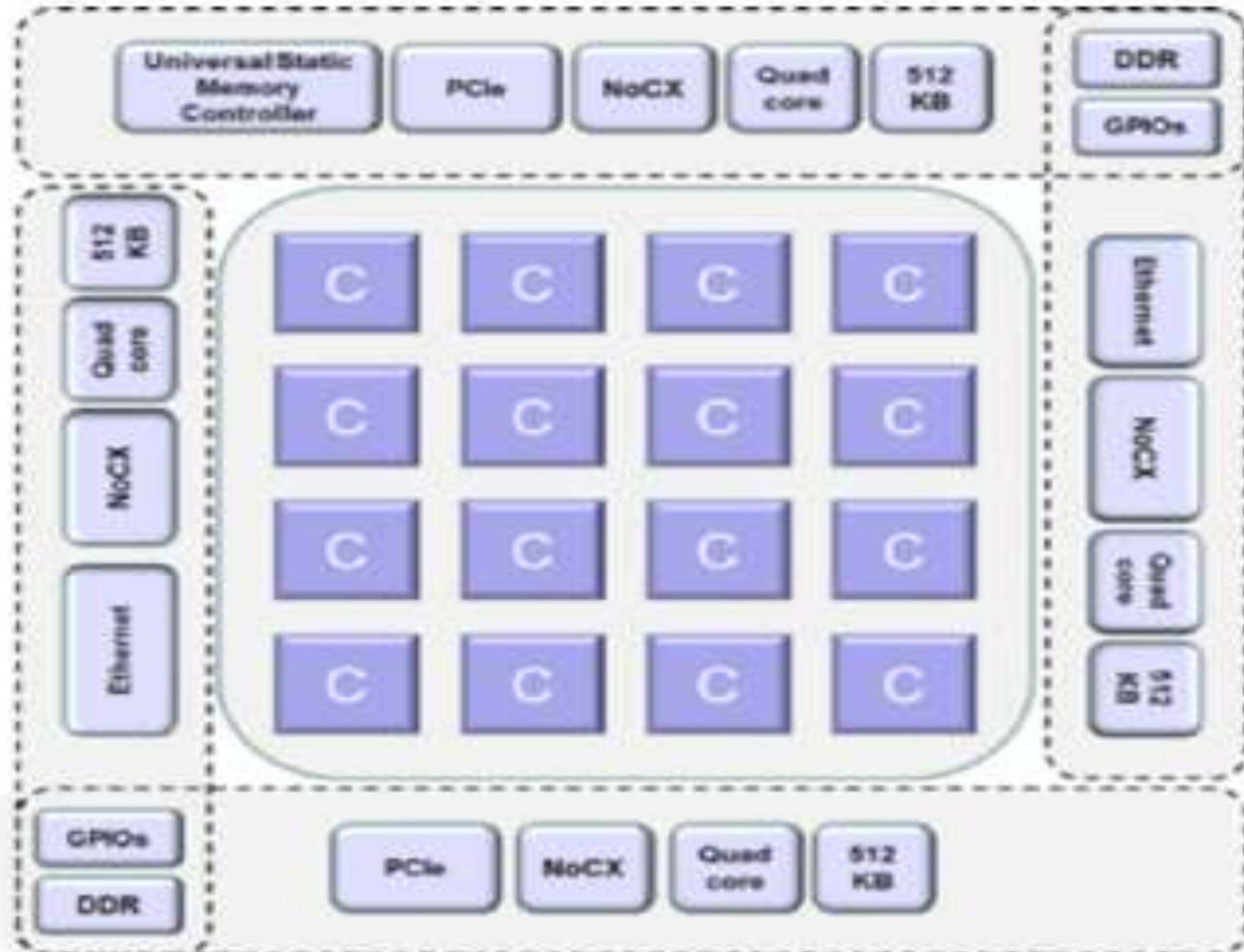


N1 - упрощенные проблемно-ориентированные КМОП-микропроцессоры, далее - их реализация на пост-Муровской ЭКБ (TSV, нанофотоника, RSFQ, QCA (???))

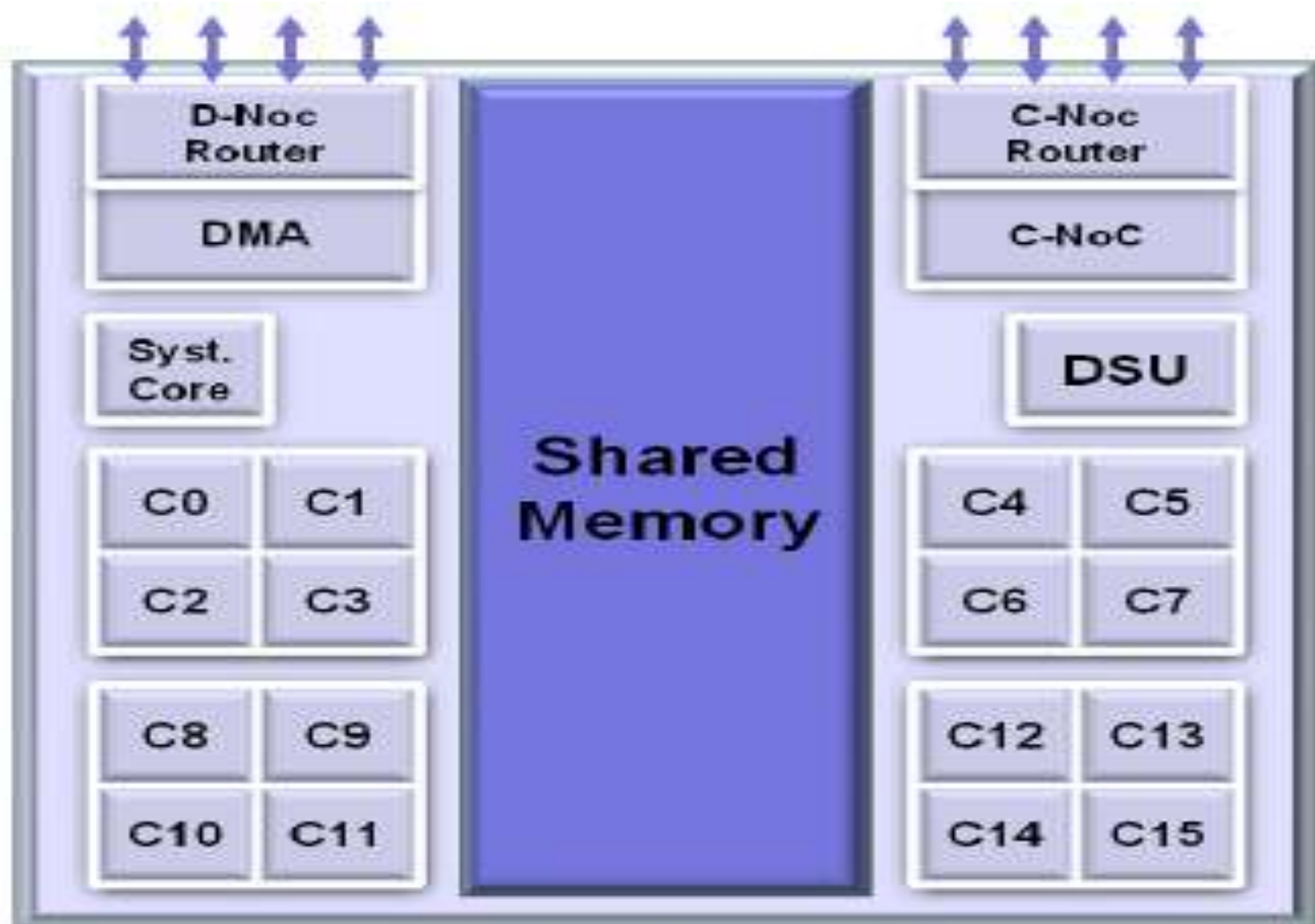
N2 - эмуляция HPGAS и массовой мультитредовости на GP-классе, микропроцессоры для СВ-класса можно создать посредством модернизации проблемно-ориентированных микропроцессоров для суперкомпьютеров RM- и CO-класса

**Разработка упрощенных
проблемно-ориентированных
микропроцессорных СБИС
(примеры зарубежных изделий)**

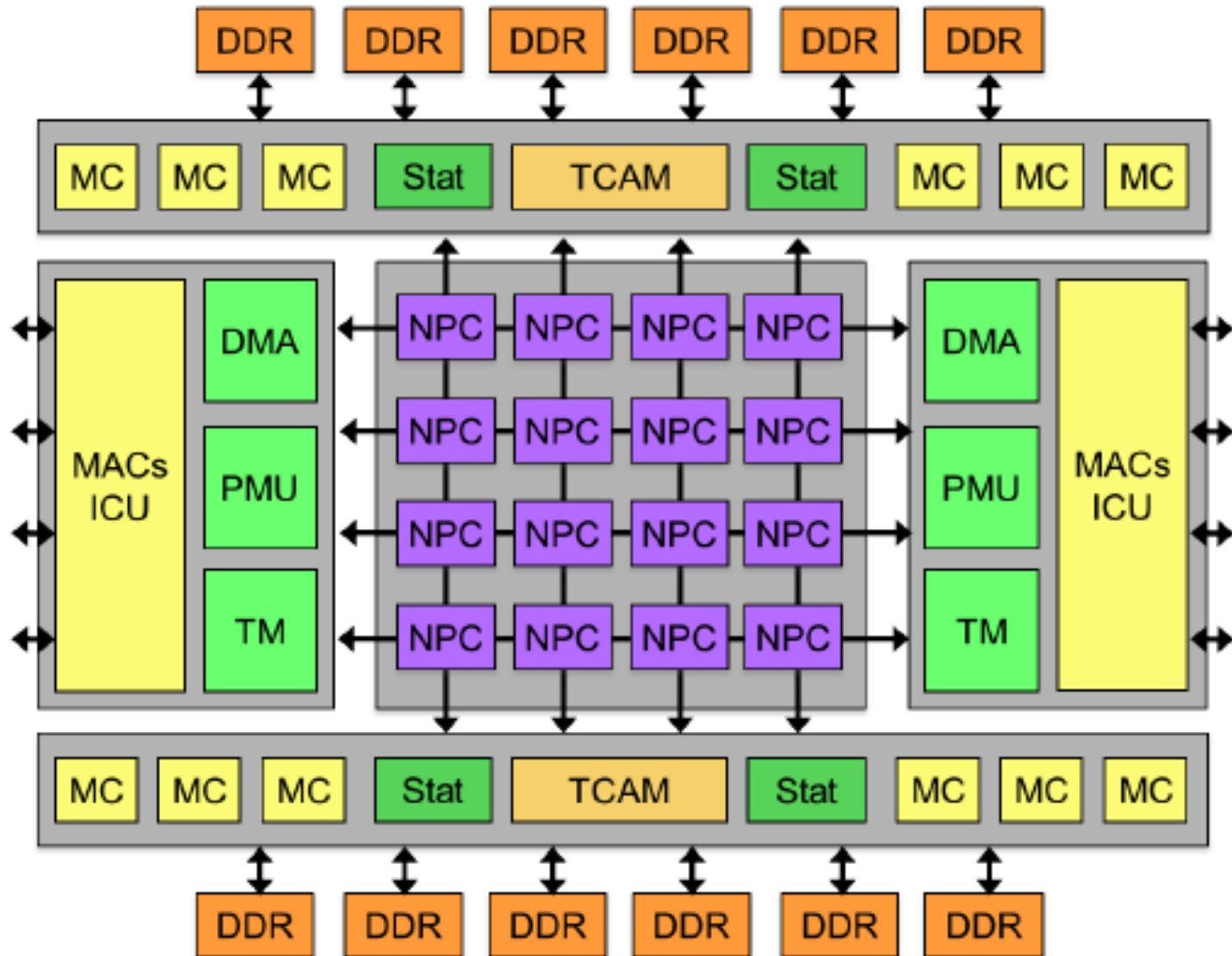
Микропроцессор МРРА-256, фирма Kalray (Франция)



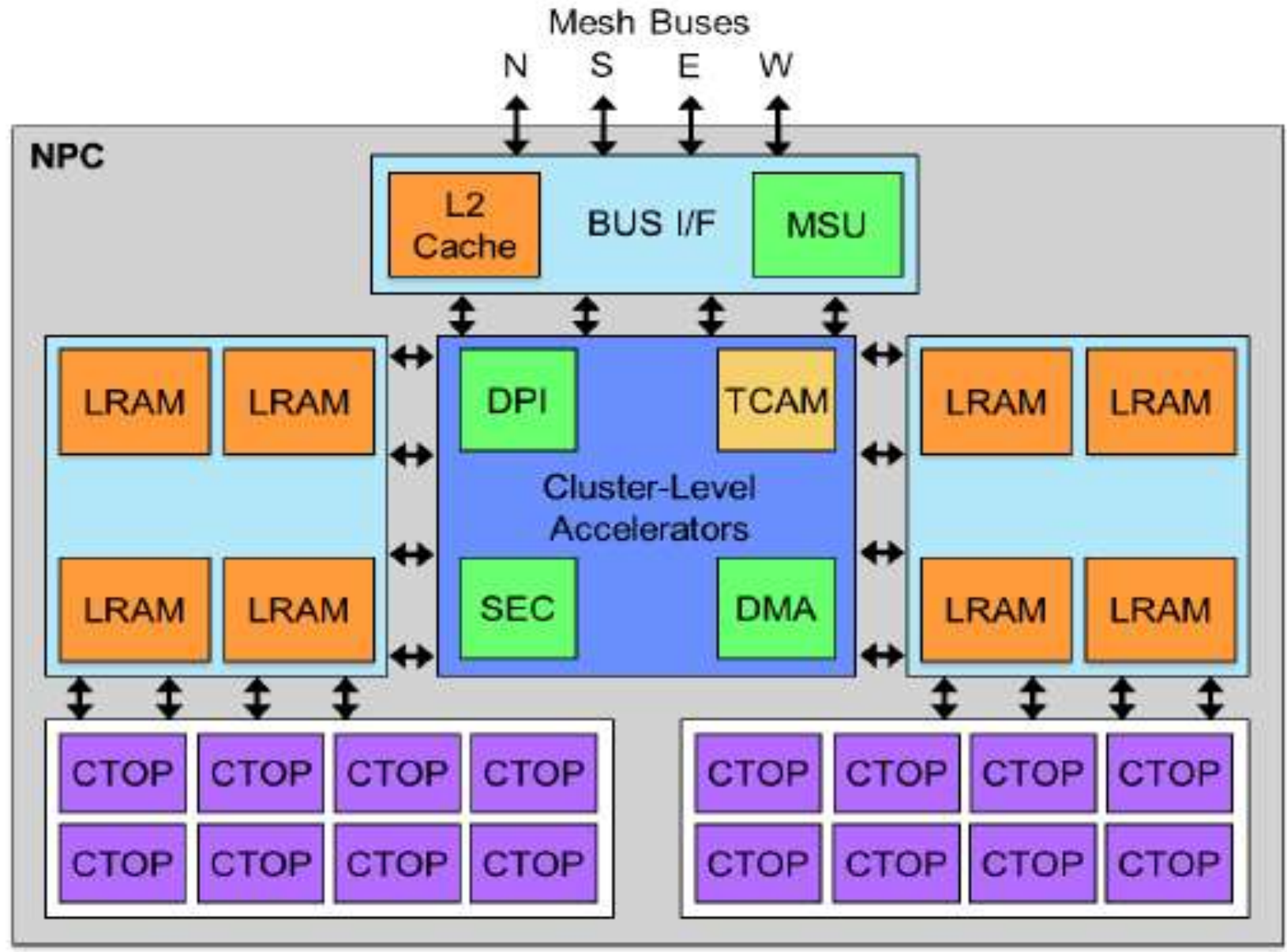
Кластер микропроцессора MPRA-256 (16 VLIW-ядер + ядро управления)



Сетевой микропроцессор NPS-400 фирмы EZChip (Израиль)



Кластер NPC микропроцессора NPS-400

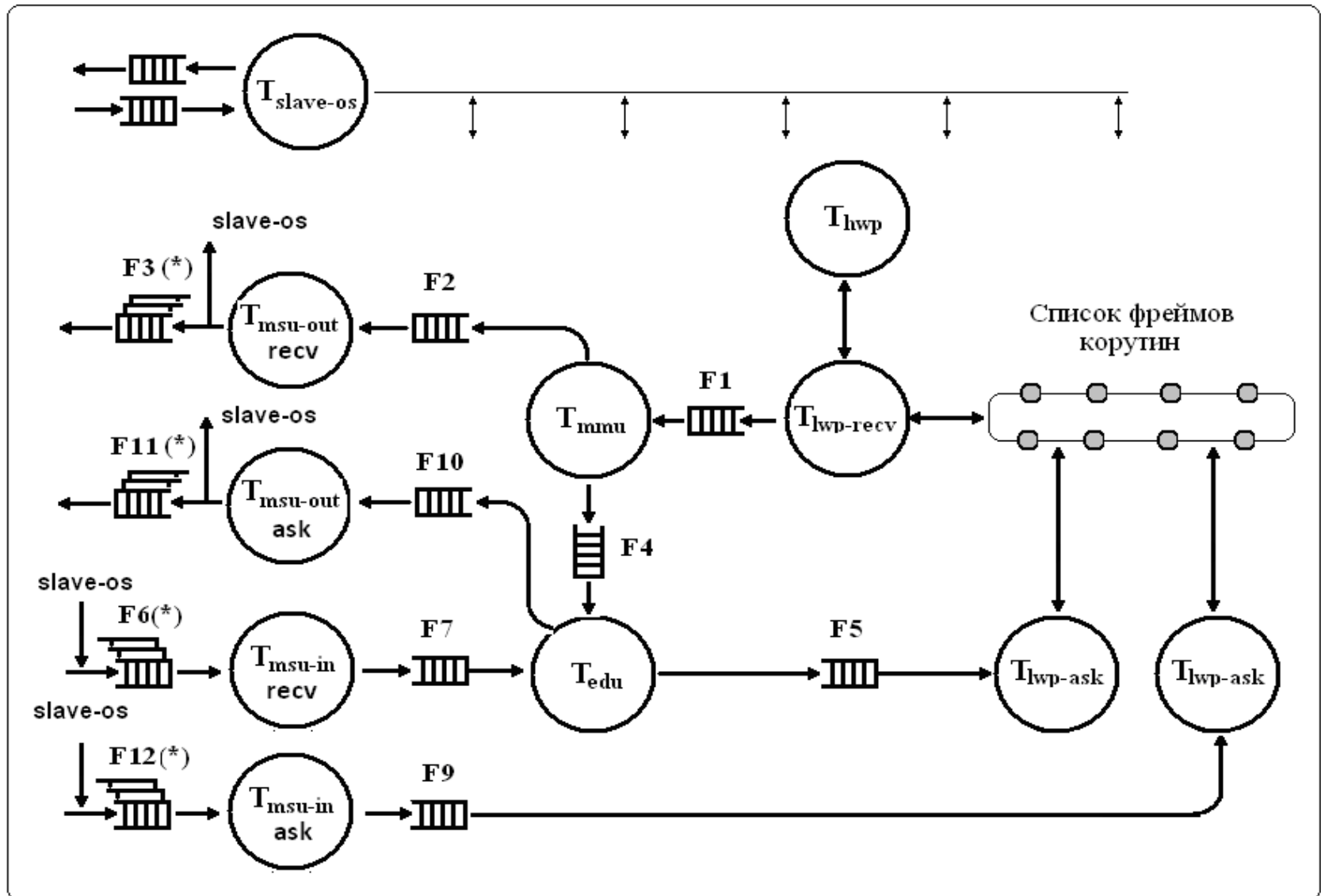


**Эмуляция HPGAS и массовой
мультитредовости на кластерных
суперкомпьютерах.**

**Модели параллельных программ
для экзафлопсных машин.**

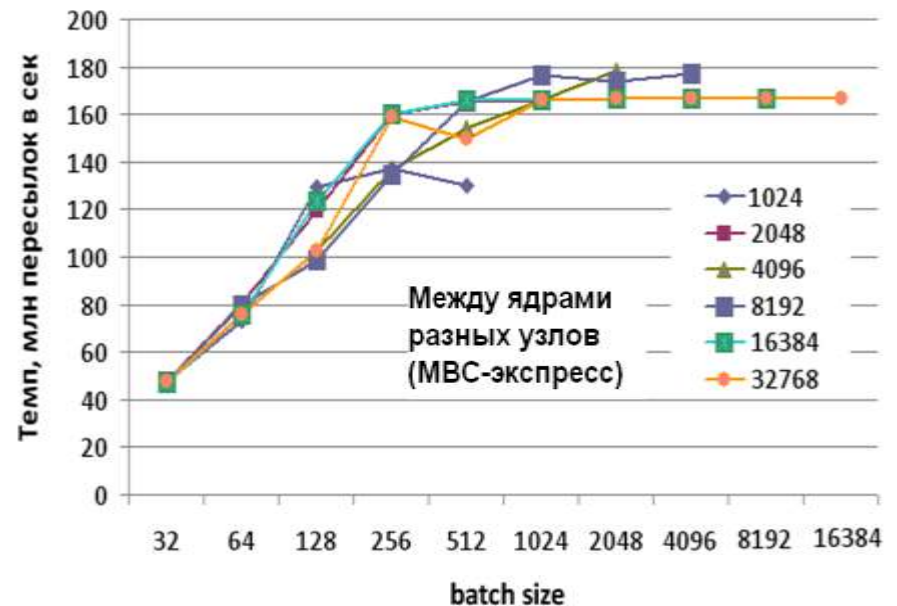
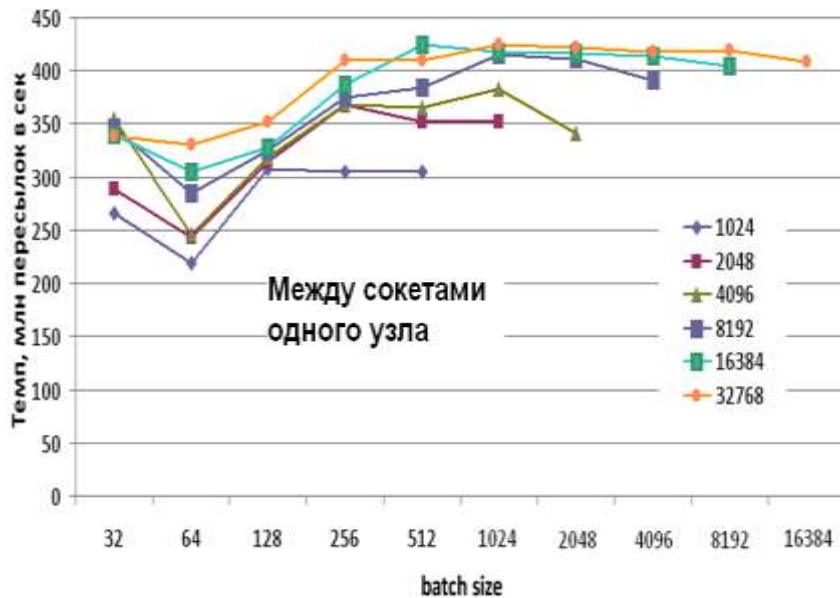
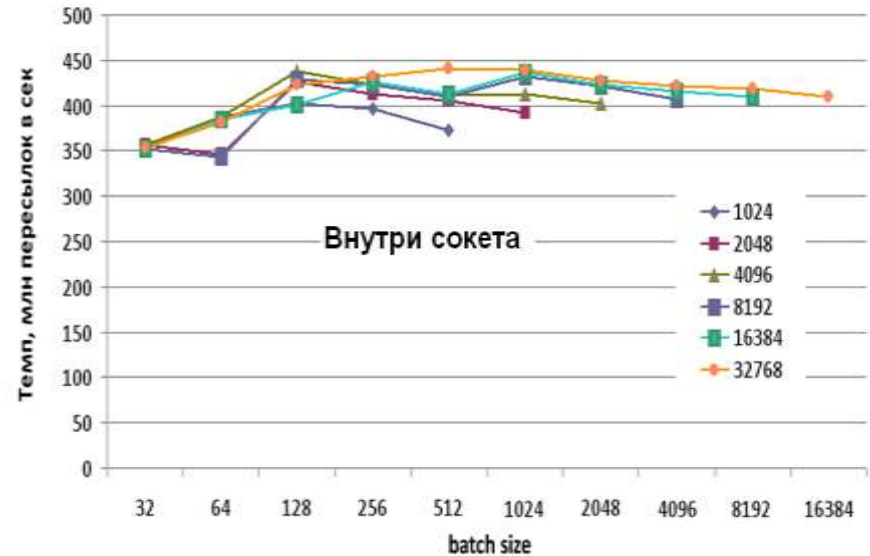
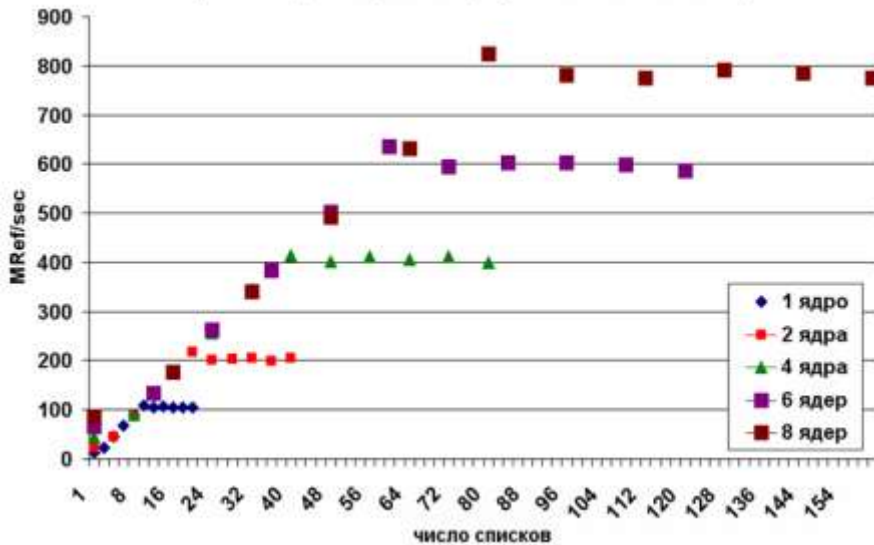
Функциональный суперкомпьютер

Уточненная схема логического узла проекта HPGAS/MT



Результаты экспериментов по HPGAS/MT

Простой обход списков, 2 x E5 2660, iss (165150720 элементов на поток)



Крупнозернистое распараллеливание программ – одновременное выполнение функций.

Программа

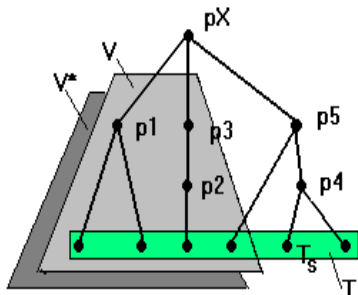
$$\begin{aligned} \S 1 kfe_1 + e_2 &> kfe_1 _ , kfe_2 _ , + \\ \S 2 kfe_1 * e_2 &> kfe_1 _ , kfe_2 _ , * \\ \S 3 kfe_1 &> e_1 \end{aligned} \quad (1)$$

Последовательное выполнение

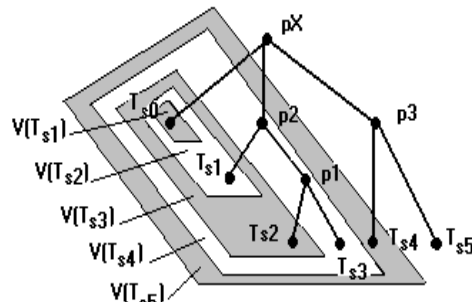
```

шаг 1.  kfa * B + C * D \_
шаг 2.  kfa * B \_ , kfc * D \_ , +
шаг 3.  kfa \_ , kfb \_ , * , kfc * D \_ , +
шаг 4.  A , kfb \_ , * , kfc * D \_ , +
шаг 5.  A , B , * , kfc * D \_ , +
шаг 6.  A , B , * , kfc \_ , kfd \_ , * , +
шаг 7.  A , B , * , C , kfd \_ , * , +
шаг 8.  A , B , * , C , D , * , +
    
```

(2)



Дерево вызовов функций
(конкретизаций).



Информационные зависимости
на дереве конкретизаций.

Параллельное выполнение

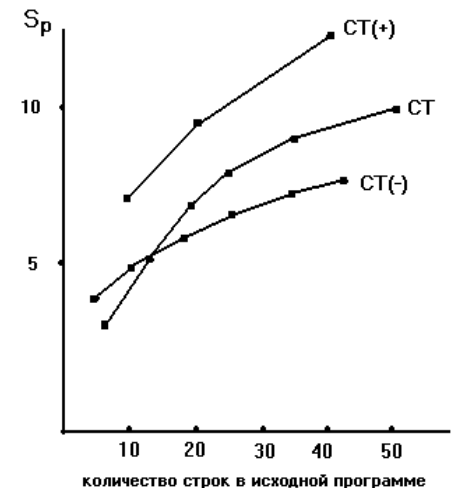
```

шаг 1.  kfa * B + C * D \_
шаг 2.  kfa * B \_ , kfc * D \_ , +
шаг 3.  kfa \_ , kfb \_ , * , kfc \_ , kfd \_ , * , +
шаг 4.  A , B , * , C , D , * , +
    
```

(3)

```

Тест CT(-)  real A,B,C
             B = C+A
             ...
             B = C+A
             end
Тест CT     real A,B,C
             A = B + C/(C+2.3-B)
             B = C+A
             ...
             A = B + C/(C+2.3-B)
             B = C+A
             end
Тест CT(+)  real A,B,C
             A = B + C/(C+2.3-B)
             ...
             A = B + C/(C+2.3-B)
             end
    
```

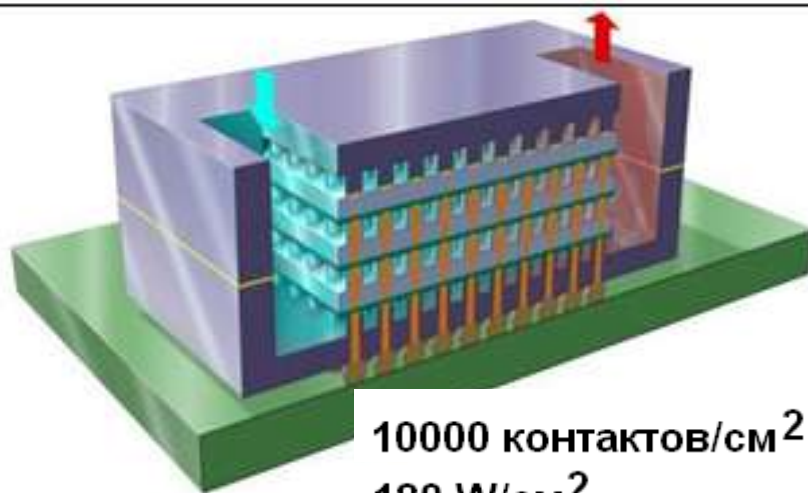


Пример совмещения проектирования элементов одной левой части

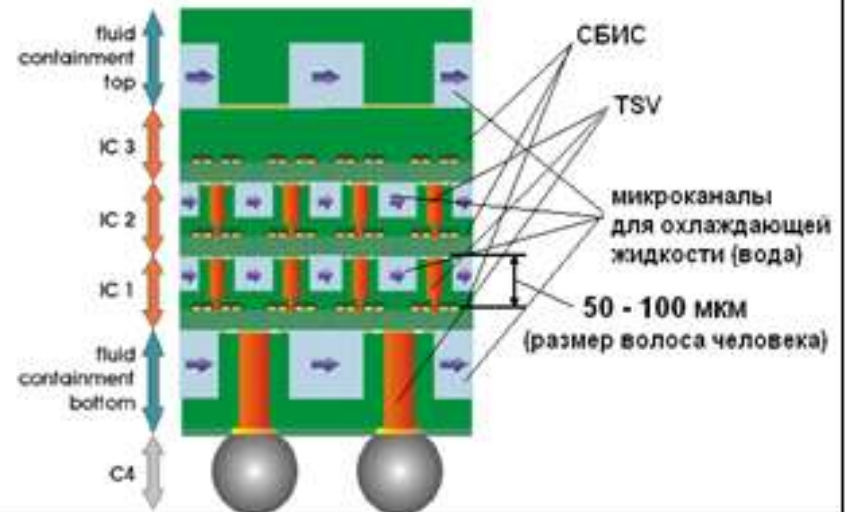
Такты	k	f	e_1^L	e_1^R	$+$	e_2^L	e_2^R	$($	w_3^L	w_3^R	A	$($	$*$	$*$	$)$	e_4^L	e_4^R	B	C	$)$	$>$
1	s_1	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_3
2	s_1	s_5	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_6	s_3
3	s_1	s_7	s_5	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_6	s_8	s_3
4	s_1	s_1	s_7	s_5	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_6	s_8	s_1	s_3
5	s_4	s_1	s_1	s_7	s_5	s_0	s_0	s_1	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_8	s_3	s_4	s_3
6	s_3	s_4	s_3	s_1	s_7	s_5	s_6	s_1	s_5	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_0	s_3	s_3	s_3	s_3
7	s_3	s_3	s_3	s_1	s_2	s_7	s_6	s_1	s_7	s_5	s_0	s_0	s_0	s_0	s_0	s_0	s_6	s_3	s_3	s_3	s_3
8	s_3	s_3	s_3	s_9	s_0	s_2	s_6	s_1	s_1	s_7	s_5	s_0	s_0	s_0	s_0	s_0	s_6	s_3	s_3	s_3	s_3
9	s_3	s_3	s_1	s_9	s_0	s_0	s_6	s_4	s_1	s_2	s_7	s_5	s_0	s_0	s_0	s_0	s_6	s_3	s_3	s_3	s_3
10	s_3	s_3	s_1	s_1	s_0	s_0	s_6	s_3	s_1	s_0	s_2	s_7	s_5	s_0	s_0	s_0	s_6	s_3	s_3	s_3	s_3
11	s_3	s_3	s_3	s_1	s_5	s_0	s_6	s_3	s_1	s_1	s_0	s_2	s_7	s_5	s_0	s_0	s_6	s_3	s_3	s_3	s_3
12	s_3	s_3	s_3	s_1	s_7	s_5	s_6	s_3	s_3	s_1	s_5	s_0	s_2	s_7	s_5	s_0	s_6	s_3	s_3	s_3	s_3
13	s_3	s_3	s_3	s_1	s_2	s_7	s_6	s_3	s_2	s_1	s_7	s_5	s_0	s_2	s_7	s_5	s_6	s_3	s_3	s_3	s_3
14	s_3	s_3	s_3	s_9	s_0	s_2	s_6	s_3	s_3	s_1	s_1	s_7	s_5	s_0	s_2	s_7	s_6	s_3	s_3	s_3	s_3
15	s_3	s_3	s_1	s_9	s_0	s_0	s_6	s_3	s_3	s_4	s_1	s_1	s_7	s_5	s_0	s_2	s_6	s_3	s_3	s_3	s_3
16	s_3	s_3	s_1	s_1	s_0	s_0	s_6	s_3	s_3	s_3	s_4	s_1	s_1	s_7	s_5	s_0	s_6	s_3	s_3	s_3	s_3
17	s_3	s_3	s_3	s_1	s_6	s_0	s_6	s_3	s_3	s_3	s_3	s_3	s_4	s_1	s_1	s_7	s_5	s_6	s_3	s_3	s_3
18	s_3	s_3	s_3	s_1	s_7	s_5	s_6	s_3	s_3	s_3	s_3	s_3	s_4	s_1	s_1	s_7	s_6	s_3	s_3	s_3	s_3
19	s_3	s_3	s_3	s_1	s_1	s_7	s_6	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_1	s_6	s_3	s_3	s_3	s_3
20	s_3	s_3	s_3	s_4	s_1	s_1	s_6	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_4	s_6	s_3	s_3	s_3	s_3
21	s_3	s_3	s_3	s_3	s_4	s_3	s_6	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_6	s_3	s_3	s_3	s_3
22	s_3	s_3	s_3	s_3	s_3	s_3	s_6	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_8	s_3	s_3	s_3
23	s_3	s_3	s_3	s_3	s_3	s_3	s_8	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3
24	s_3	s_3	s_3	s_3	s_3	s_3	s_1	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_8	s_3	s_3	s_3	s_3
25	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_6	s_3	s_3	s_3	s_3
26	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_8	s_3	s_3	s_3	s_3
27	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_1	s_3	s_3	s_3	s_3
28	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_4	s_3	s_3	s_3	s_3
29	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3
30	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3
31	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3	s_3

Варианты пост-Муровских технологий ЭКБ

3D сборка – IBM TSV

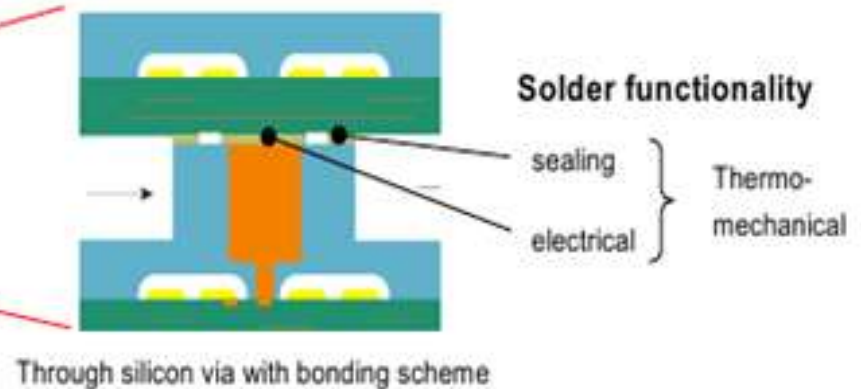
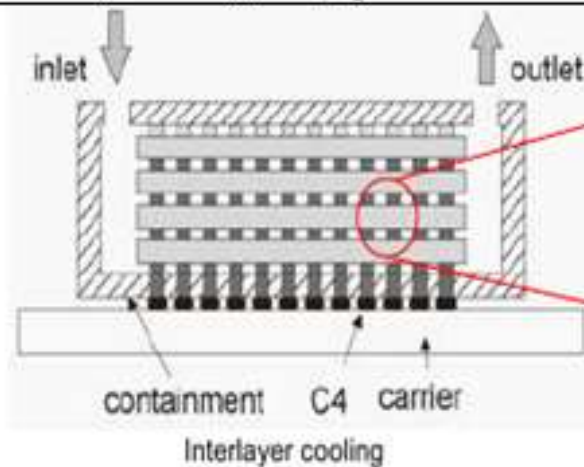


10000 контактов/см²
180 W/см²



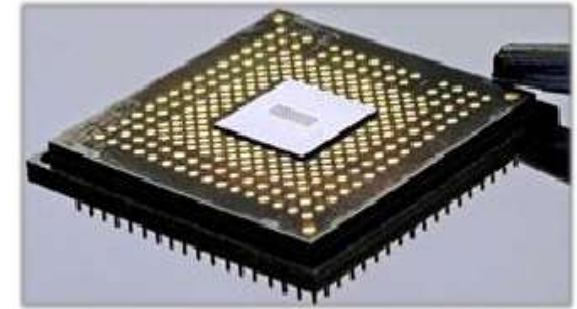
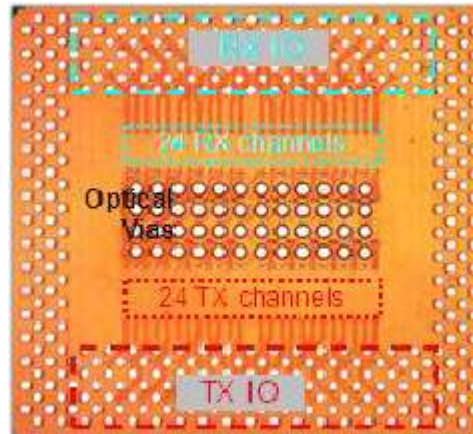
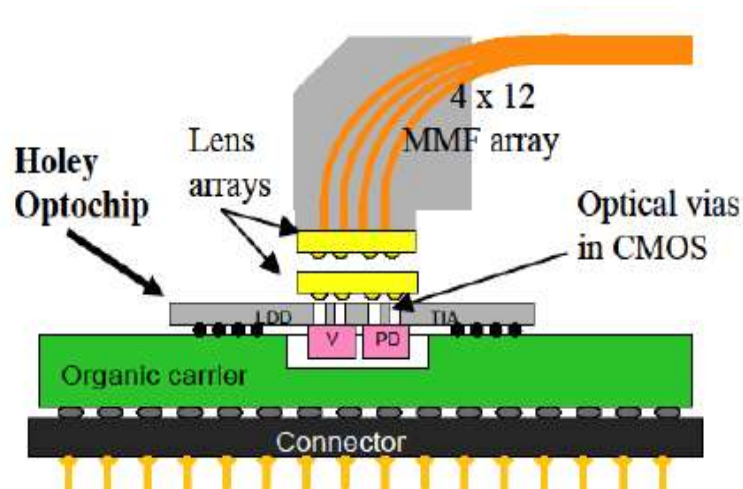
Общий вид модуля с охлаждением

Вид модуля сбоку.

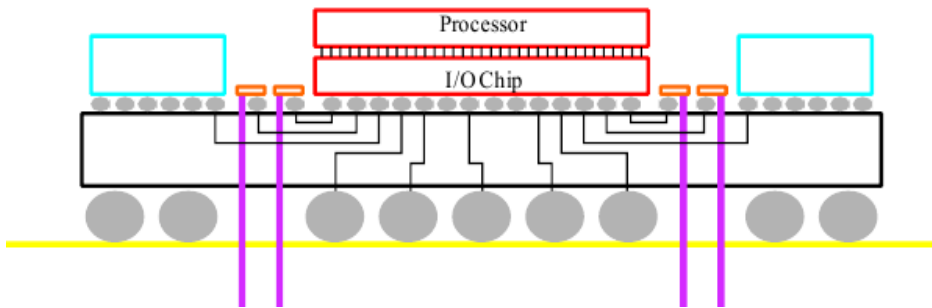


Уточненный вид сбоку с выделением соединений и изоляции.

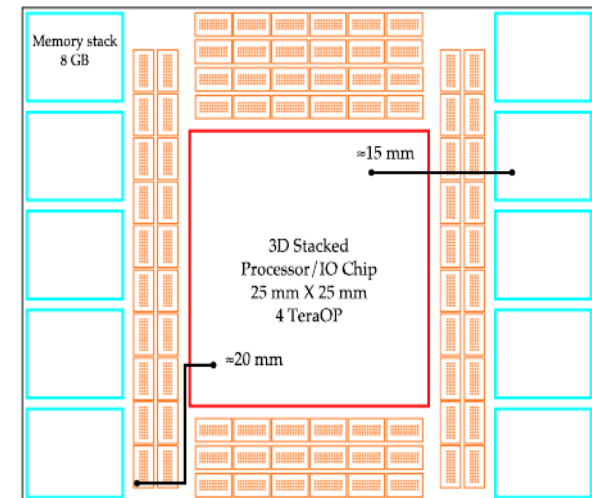
Коммуникационный опточип IBM Holley и вариант перспективной компоновки вычислительного модуля



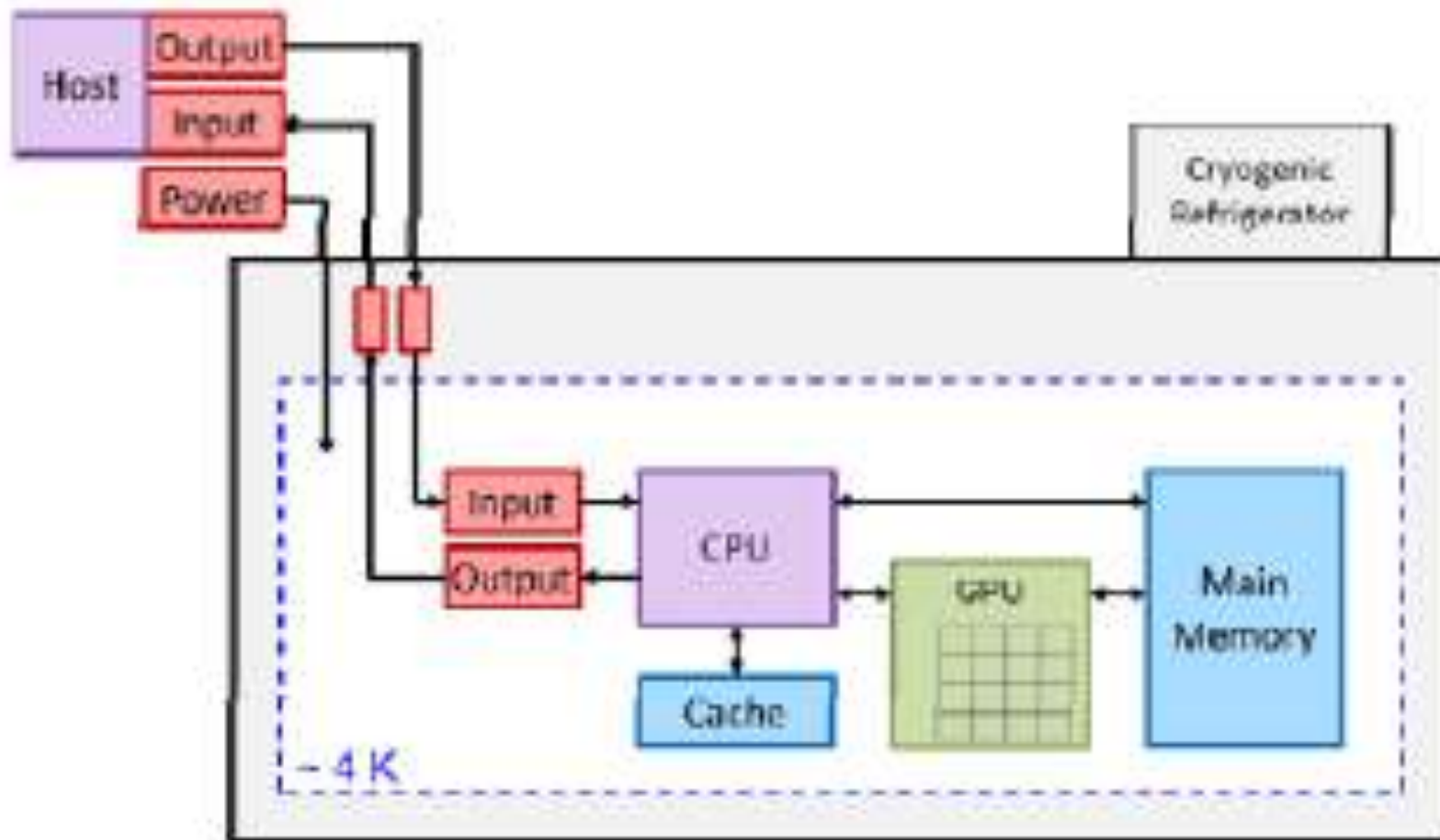
90 нм, 48 линков по 12.5 Gb/s (24(in)+24(out)), ~5x5мм, 8 pJ/bit.
Перспектива (45 нм) – 40 Gb/s, 1 pJ/bit (0.8 – E, 0.2 – O)



Узел - 5x6 см, 82 (OE), 1968 VCSL + 1968 PD,
6 узлов в группе, 4 группы в модуле



Вычислительный узел криогенного суперкомпьютера (RSFQ), создаваемого по программе IARPA C3 (5 лет)



Характеристики модуля и компонентов

CPU

#	Parameter	Goal
1	Throughput (bit-op/s)	10^{13}
1	Efficiency @ 4 K (bit-op/J)	10^{15}
2	Main memory, total (B)	2^{27}
2	Cache memory, total (B)	2^{15}
2	Input/Output (bit/s)	10^8

Общие

GPU

#	Parameter	Goal
1	Word size (bit)	64
2	Efficiency @ 4 K (bit-op/J)	5×10^{15}
2	Processor class	ARM™ or Intel Atom™
2	Instruction set	ARM™ or simple x86
2	ALUs	1 integer
2	Throughput (bit-op/s)	10^{12}

#	Parameter	Goal
2	Efficiency @ 4 K (bit-op/J)	10^{16}
2	Instruction size (bit)	64
2	Data size (bit)	64
2	PU count	8
2	PU register size (bit)	64×128
2	PU ALUs	1 integer
2	PU throughput (bit-op/s)	2×10^{12}
3	PU area (mm ²)	10

Cach

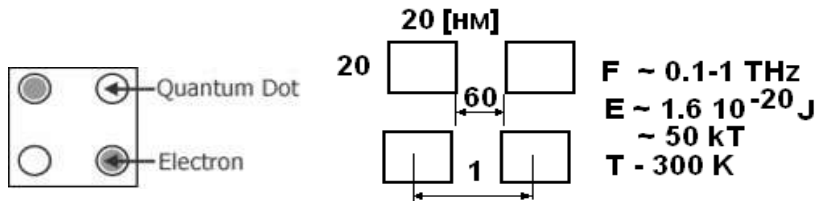
#	Parameter	Goal
1	Access energy @ 4 K, ave. (J/bit)	$5e-18$
1	Power, static (W/bit)	$< 5e-19$
2	Write time (ps)	200
2	Read time (ps)	100
2	Read rate, burst mode (Gbit/s)	300
2	Capacity per memory chip (bit)	$\geq 2^{16}$
2	Density (bit/cm ²)	$1e+7$
3	Read/Write error rate	$1e-12$
3	Read/Write disturb rate	$1e-12$

Memory

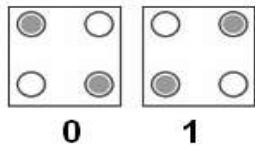
#	Parameter	Goal
1	Access energy @ 4 K, ave. (J/bit)	$5e-17$
1	Power, static (W/bit)	$< 5e-19$
2	Write time (ps)	2,000
2	Read time (ps)	500
2	Read rate, burst mode (Gbit/s)	50
2	Capacity per memory chip (bit)	$\geq 2^{26}$
2	Density (bit/cm ²)	$1e+8$
3	Read/Write error rate	$1e-12$
3	Read/Write disturb rate	$1e-12$

Элементы QCA-логики

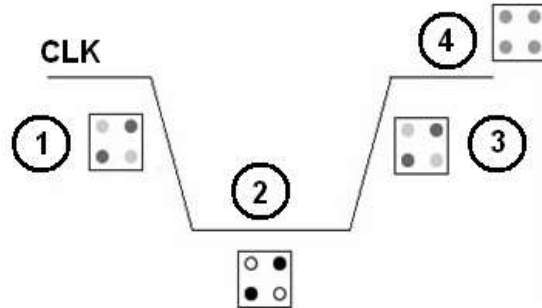
4-х точечная QCA-ячейка



QCA-ячейки в состоянии 0 и 1



Управление QCA-ячейкой



Сигнал CLK влияет на уровень тунnelирования:

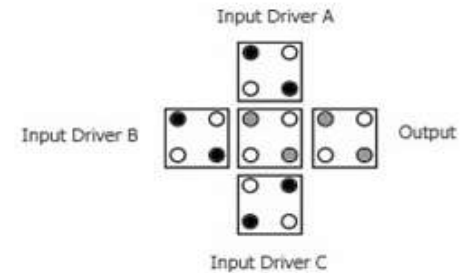
- 1 - переходное состояние (меняется уровень)
- 2- стабильное состояние (уровень высокий)
хранение информации
- 3 - переходное состояние (меняется уровень)
- 4- рабочее состояние (уровень низкий)
обработка информации, логика клеточного автомата

QCA majority gate

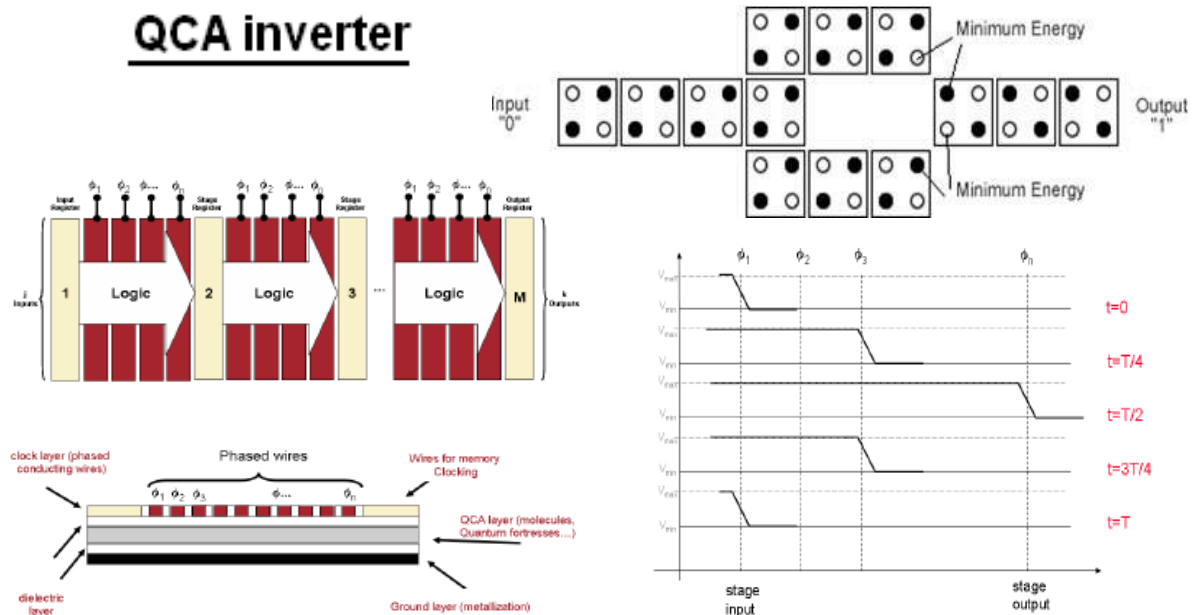
$$M(a,b,c) = ab + bc + ca$$

$$a.b = M(a,b,0)$$

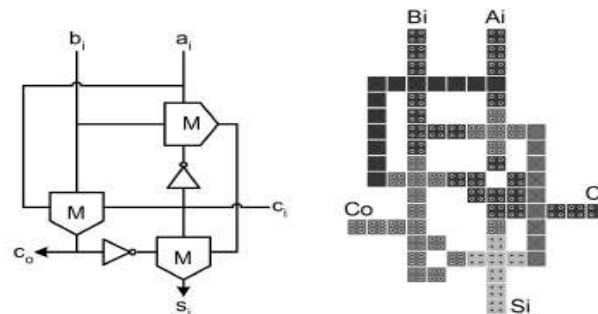
$$a+b = M(a,b,1)$$



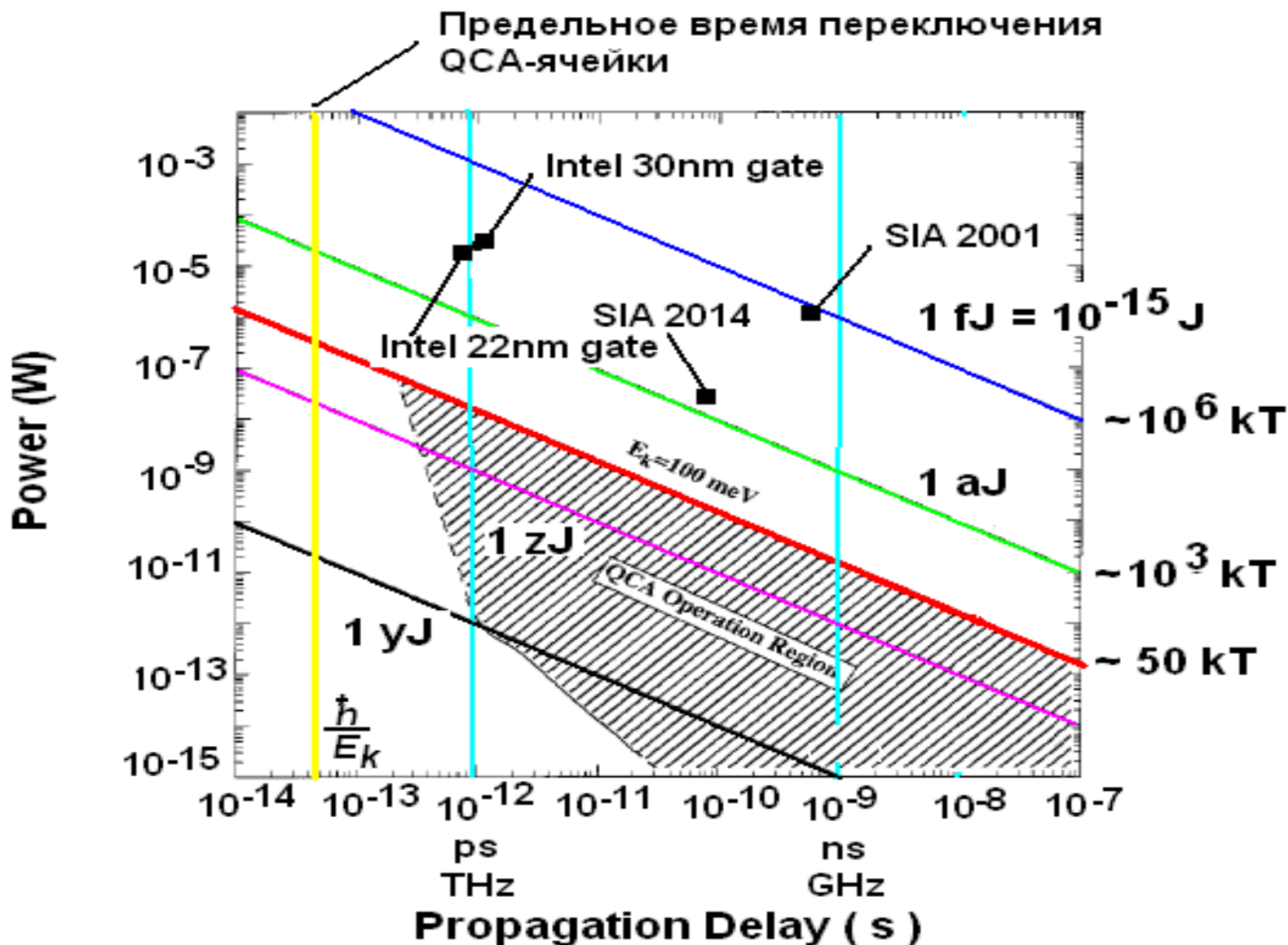
QCA inverter



Разряд сумматора



Прогноз 2006 года по перспективам создания зеттафлопсной машины на КМОП и QCA - 4



Выводы **(первоочередные задачи в области** **инновационных СКТ)**

- **Разработка проблемно-ориентированных (co-design) микропроцессоров на базе перспективных архитектурных принципов**
- **Проведение исследований и разработок по новым моделям вычислений и организации памяти, а также эмуляции массово-мультитредовых суперкомпьютеров с глобально адресуемой памятью**
- **Организация работы экспертного сообщества для планирования и оценки работ, формирование нескольких десятков исследовательских групп для работы в области инновационных суперкомпьютерных технологий**
- **Активизация работ по сверхпроводниковой электронике и нанофотонике, для целенаправленной работы в области ЭКБ перспективных суперЭВМ (криогенная электроника, квантовая электроника)**

Вопросы ?

Эйсымонт Леонид Константинович

(ФГУП"НИИ"Квант", ФГБНУ НИИ РИНКЦЭ Минобрнауки РФ

verger-lk@yandex.ru . eismont@rdi-kvant.ru)

Горбунов Виктор Станиславович

(ФГУП"НИИ"Квант", gorbunov@rdi-kvant.ru)