

Всё меняется.

и НРС тоже...

Алексей Перевозчиков **Server Solutions Product Manager**

82189117@ru.ibm.com





2010 Patent Leadership (18 лет лидерства)

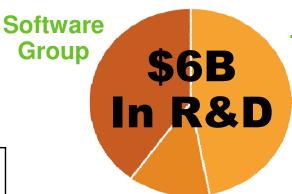
	Totals
IBM	5,896
Samsung	4,551
Microsoft	3,094
Canon	2,552
Panasonic	2,482
Toshiba	2,246
Sony	2,150
Intel	1,653
LG Electronics	1,490
HP	1,480

^{*} Source: IFI Patent Intelligence

IBM Austin the Home of Power Systems development: 950 Patents #1 IBM location for 8th year







Research

Systems & Technology Group



2012 Patent Leadership (20 лет лидерства)

Totals

IBM	6,478
Samsung	5,081
Canon	3,174
Sony	3,032
Panasonic	2,769
Microsoft	2,613
Toshiba	2,447
Hon Hai Precision	2,031
General Electric	1,652
LG Electronics	1,624





Software Group

\$6B In R&D Systems & Technology Group

Source: IFI Patent Intelligence

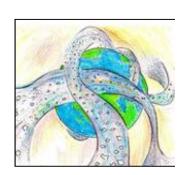
HP #15 1,394 Intel #18 1,290

Oracle is not in the top 50



Research

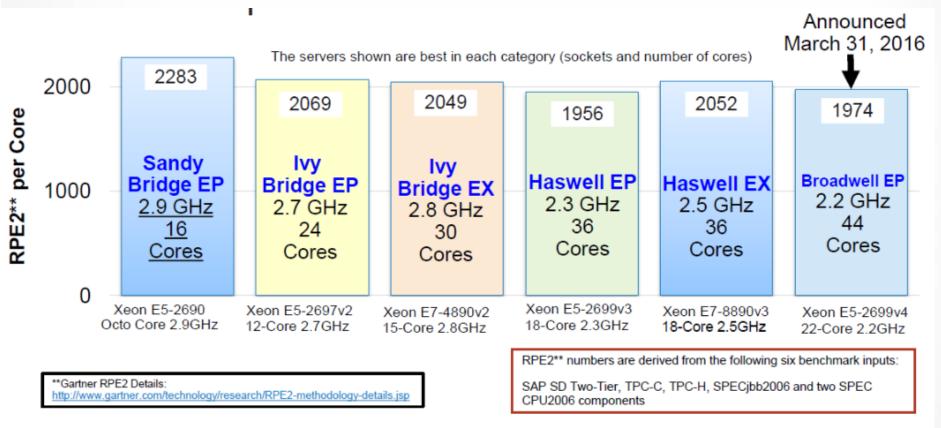
Streaming Analytics



32 Nanometer factory © 2014 IBM Corporation



RPE2



The data on this chart is derived from RPE2 from Gartner, Inc's Competitive Profile tool. © 2016 Gartner, Inc. and/or its affiliates. All rights reserved.



Intel отказалась от своей знаменитой стратегии «тик-так»

Техника

Основная версия

24.03.2016, ЧТ, 11:57, Мск , Текст: Сергей Попсулин

Intel решила сбавить темп смены технологического процесса и увеличить время выпуска процессоров с 14-нм транзисторами до нестандартных трех лет.

Отказ от стратегии

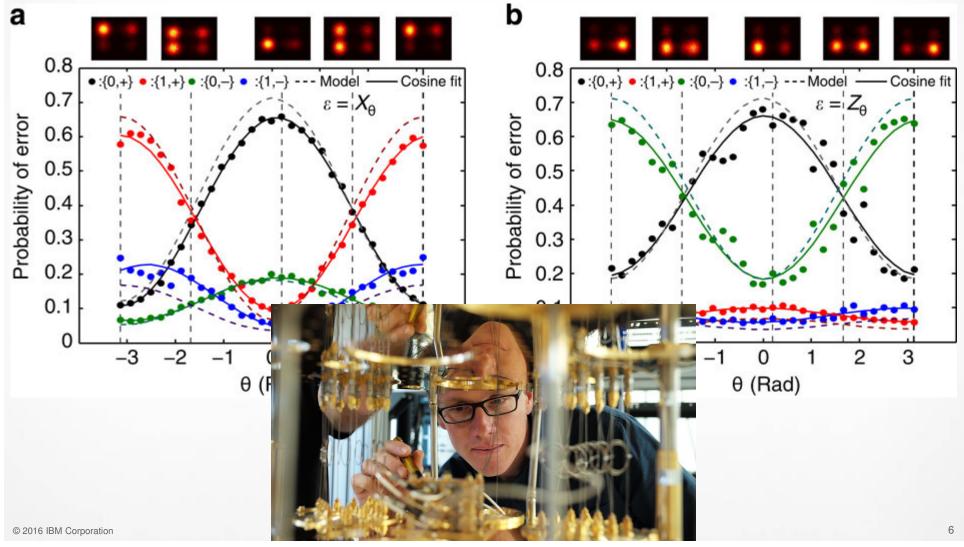
Intel отказалась от стратегии «тик-так» и вместо нее перейдет на стратегию из трех этапов: техпроцесс-архитектура-оптимизация. Так будет на протяжении по крайней мере двух следующих технологических процессов, сообщает AnandTech со ссылкой на годовой отчет, в котором корпорация изложила суть перемен.

Изменение стратегии приведет к тому, что переход на новую технологическую норму теперь будет происходить не так часто, как прежде. Вместе с тем Intel обещает продолжить выпускать новые продукты ежегодно.





А тем временем... **IBM** сообщила об открытии способа контроля одновременно обоих типов квантовых ошибок и ...









IBM открывает эпоху публичных квантовых вычислений

Автор: Андрей Колесов

05.05.2016

IBM Research сообщило о запуске публичного бесплатного облачного сервиса IBM Quantum Experience, с помощью которого можно на практике познакомиться с возможностями

© 2016 IBM Corporation



Энергоэффективность

GOOGLE'S ENERGY CONSUMPTION

16 February, 2013 · by Arnfinn Oines · in Energy, News. ·

Googles 13 data centres continuously draws 260 million watts. Below is an infographic made by CO2 Sense that highlights what it takes to power Google. It also includes statistics regarding the company's investment in renewable energy.

IT now 10 percent of world's electricity consumption, report finds

New analysis finds IT power suck has eclipsed aviation

By Jack Clark, 16th August 2013

Follow 3,798 followers



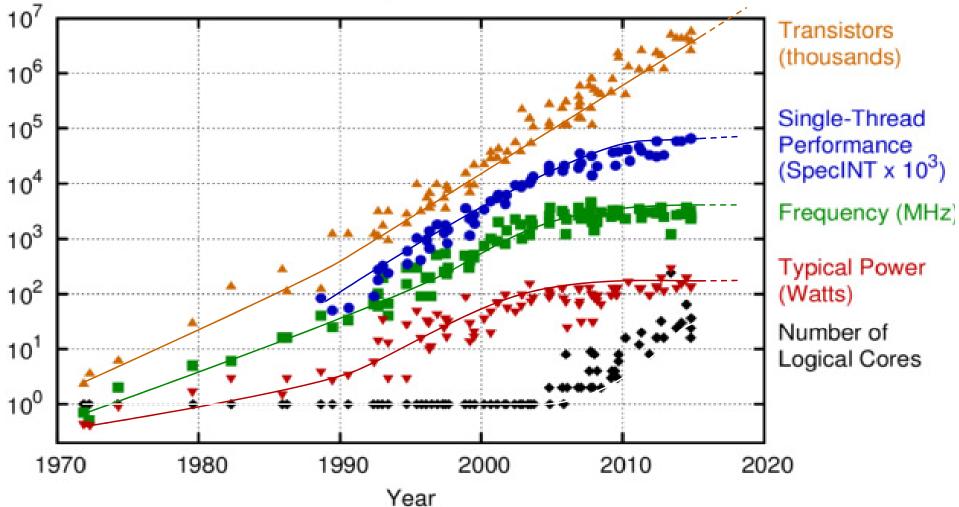
POWER8





40 лет микропроцессорам: тенденции





Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp



Планы развития процессора POWER

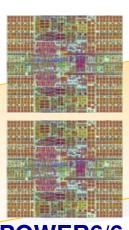






POWER5/5+ 130/90 nm

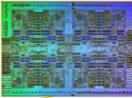
- ✓ Dual Core
- ✓ Enhanced Scaling
- √SMT
- ✓ Distributed Switch +
- √Core Parallelism +
- √FP Performance +
- ✓ Memory Bandwidth +
- √ Virtualization



POWER6/6+ 65/65 nm

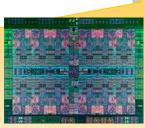
- ✓ Dual Core
- √ High Frequencies
- √ Virtualization +
- ✓ Memory Subsystem +
- ✓ Altivec
- ✓Instruction Retry
- **✓ Dynamic Energy Mgmt**
- ✓SMT +
- ✓ Protection Keys





POWER7/7+ 45/32 nm

- **✓ Eight Cores**
- ✓ On-Chip eDRAM
- **✓ Power-Optimized Cores**
- ✓ Memory Subsystem ++
- ✓SMT++
- ✓ Reliability +
- ✓VSM & VŠX
- ✓Protection Keys+



22nm



POWER8

✓ Extreme Analytics **Optimization**

POWER9

- ✓ Extreme Big Data **Optimization**
- √On-chip accelerators
- ✓ More Cores
- ✓SMT+++
- ✓ Reliability ++ **✓ FPGA Support**
- √ Transactional Memory
- ✓ PCle Acceleration



2007

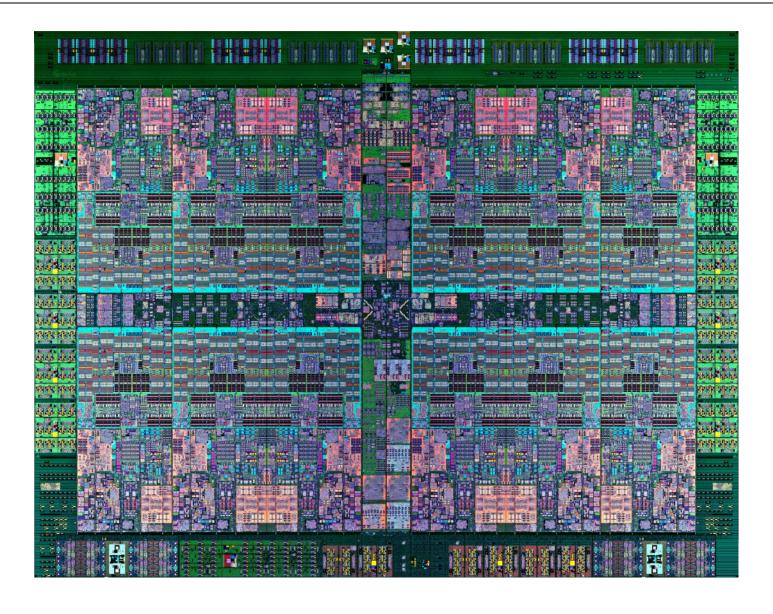
2010

2014



Процессор POWER8





Процессор POWER8



Технология

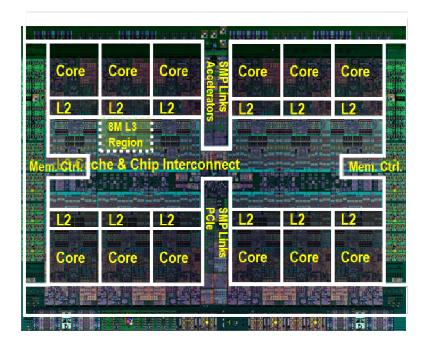
22nm SOI, eDRAM, 650mm2, 4.2B transistors

Ядра

- 12 ядер (SMT8)
- 8 dispatch, 10 issue,
 16 exec pipe
- 2X internal data flows/queues
- Enhanced prefetching
- 64К кэш данных,
 32К кэш инструкций

Акселераторы

- Криптография
- Расширение памяти
- Транзакционная память
- Поддержка VMM
- Перемещение данных / VM



Energy Management

- On-chip Power Management Micro-controller
- Integrated Per-core VRM
- Critical Path Monitors

Увеличенные кэши

- •512 KB SRAM L2 / core
- •96 MB eDRAM shared L3
- •Up to 128 MB eDRAM L4 (off-chip)

Память

•Up to 230 GB/s sustained bandwidth

Шинные интерфейсы

- •Durable open memory attach interface
- •Интегрированный PCle G3
- SMP Interconnect
- •CAPI (Coherent Accelerator Processor Interface)

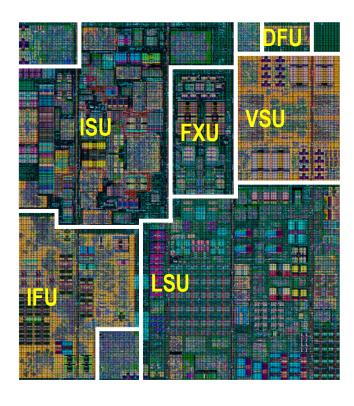


Ядро POWER8 (относительно POWER7)



SMT4 → SMT8

- 8 dispatch
- 10 issue
- 16 execution pipes:
 2 FXU, 2 LSU, 2 LU, 4 FPU,
 2 VMX, 1 Crypto, 1 DFU,
 1 CR, 1 BR
- Larger Issue queues (4 x 16-entry)
- Larger global completion,
 Load/Store reorder
- Improved branch prediction
- Improved unaligned storage access



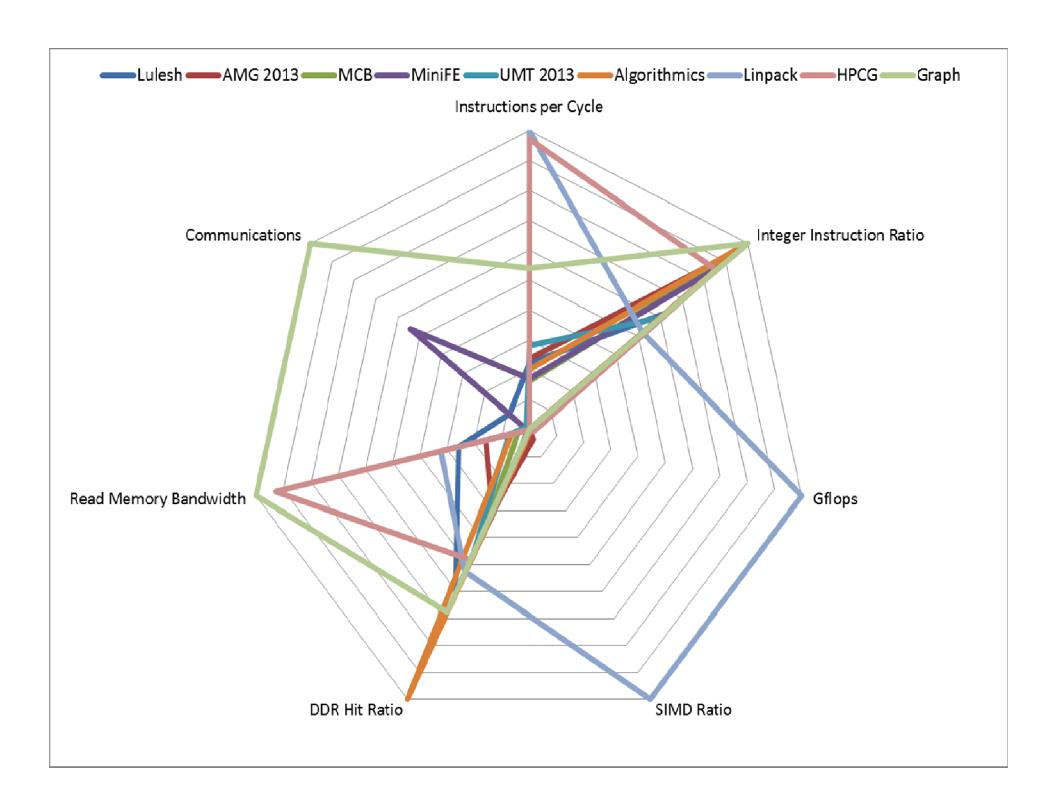
- 2x L1 data cache (64 KB)
- 2x outstanding data cache misses
- 4x translation Cache

Wider Load/Store

- 32B \rightarrow 64B L2 to L1 data bus
- 2x data cache to execution dataflow

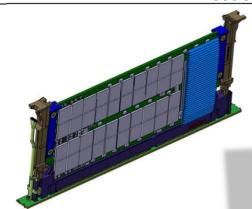
Enhanced Prefetch

- Instruction speculation awareness
- Data prefetch depth awareness
- Adaptive bandwidth awareness
- Topology awareness





Memory Buffer Chip ... with 16MB Cache...



"L4 cache"

Модули памяти наполняются интеллектом

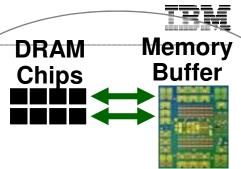
- •Умная система кэширования
- •Оптимизация энергии
- •Надежность

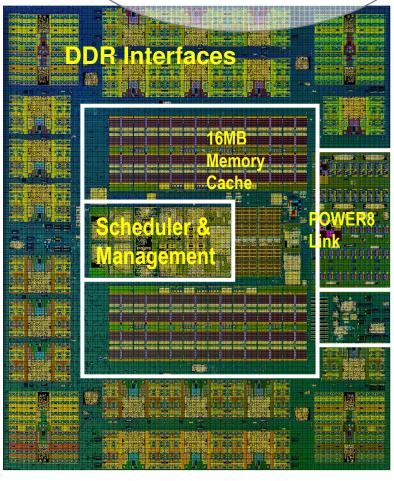
Оптимизированный интерфейс

- •9.6 GB/s high speed interface
- •Интеллектуальная надежность
- •Изоляция сбоев на лету

Уникальная производительность

- •Уменьшенная латентность fastpath
- •Cache → latency/bandwidth, partial updates
- •Логика предсказания
- •22nm SOI for optimal performance / energy
- •15 metal levels (latency, bandwidth)

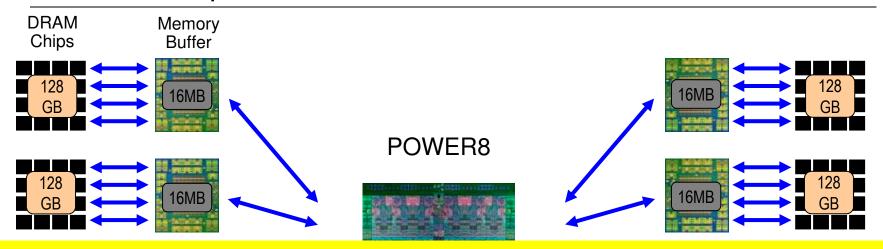






Организация памяти в POWER8





- У Intel нет L4 и они показывают цифры "to the DIMM"
- Наши 230 ГБ/с вполне достижимы в реальных условиях
- Цифры "to-DIMM" теоретические, реально достижимые намного ниже (из-за используемых протоколов DIMM, это справедливо для всех производителей)
 - → 8 скоростных каналов, каждый до 9.6 Гб/с
 до 230 ГБ/с в устойчивом режиме (sustained)
 - → До 32 портов DDR выдающих в пике 410 ГБ/с (на уровне DRAM)
 - → До 1 ТБ памяти на сокет (для старших версий до 2 ТБ на сокет)

CAPI (Coherent Accelerator Processor Interface)

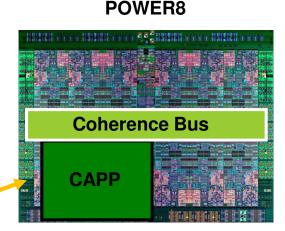


Virtual Addressing

- •Ускоритель работает напрямик с разделяемой памятью
- •Обмен данными с кэшем процессора.
- •Исключает накладные расходы ОС и драйверов.

Hardware Managed Cache Coherence

•Стандартный механизм блокировок.





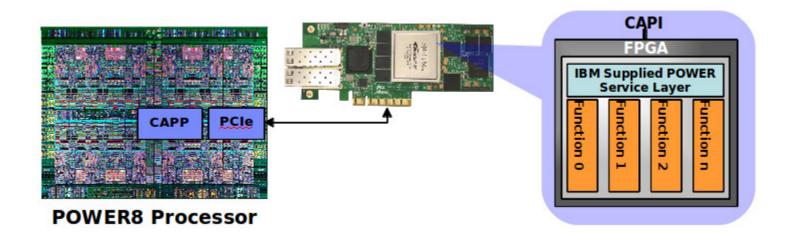
PCle Gen 3
Transport for encapsulated messages

Специализированные контроллеры Программные ускорители





Coherent Accelerator Processor Interface (CAPI) Flow



Типичный процесс работы І/О



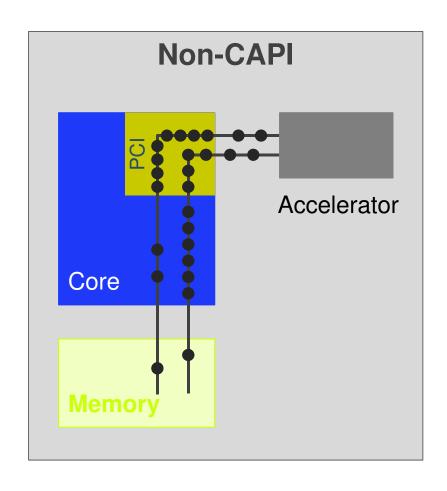
Процесс при использовании когерентной памяти

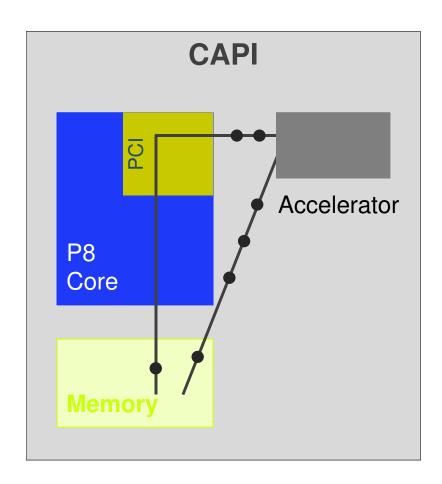






Coherent Accelerator Processor Interface









Несколько слов о стратегии в области НРС





Развитие стратегии аппаратных средств для НРС

- Общий дизайн платформы для высокопроизводительных вычислений и высокопроизводительной аналитики
- Углубление отношений с технологическими партнёрами
- Серверы для данного сегмента в основном 2 сокета
- Усиление поддержки InfiniBand и Ethernet
- Большая часть производительности на операциях с плавающей точкой будет достигаться за счёт GPU
- Стандартные индустриальные стойки и корпуса
 - Варианты воздушного и водяного охлаждения

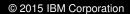


Стратегия развития процессоров архитектуры POWER

- Консолидация усилий и фокус на одном процессоре (чипе) общего назначения для каждого поколения
 - ❖Дизайн для более плотной интеграции с вспомогательным оборудованием
 - ❖Множественный дизайн модулей обеспечивает различные комбинации памяти и шин I/O
- Использование ускорителей подключаемых к процессору для соответствующих платформ и приложений
 - ❖FPGA для коммерческих задач, таких как Java, СУБД, аналитика
 - ❖GPU для научных и вычислительных задач



OpenPOWER Foundation – что, как, зачем.





Основные особенности OpenPOWER

- Это общественная организация, деятельность которой не регулируется кем бы то ни было. Ни коммерческими, ни государственными структурами
- Идея близка к концепции ПО с открытым кодом, но в применении к аппаратуре
- Отличие от мира СПО участники консорциума кооперируются, а не конкурируют.
- Каждый участник делает свою часть или создаёт свои изделия используя наработки остальных участников сообщества.





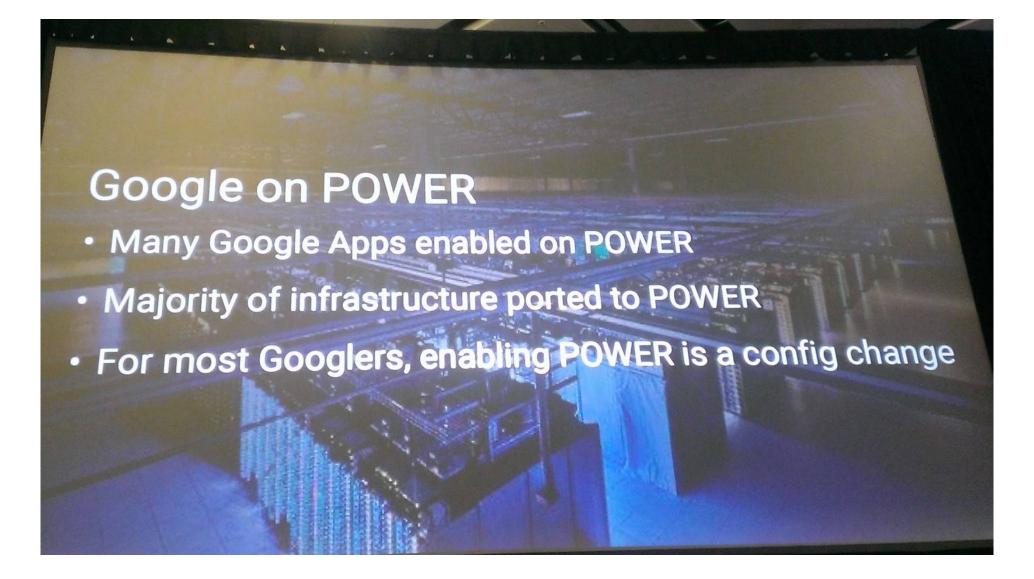


Открытое сообщество разработчиков



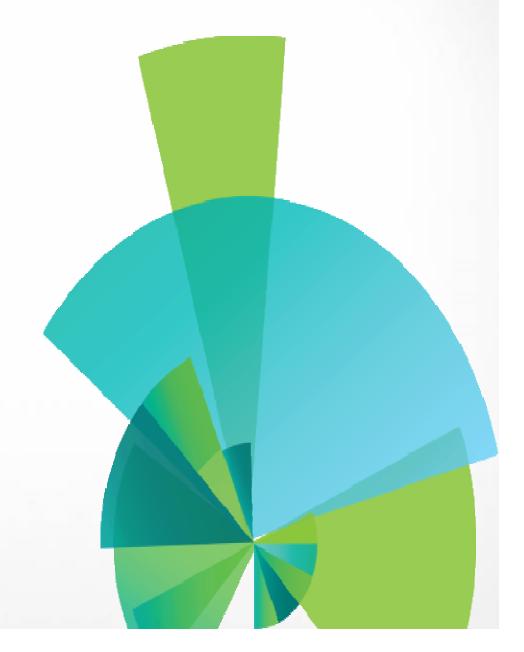




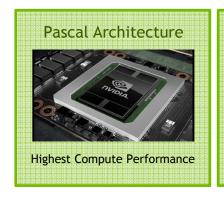


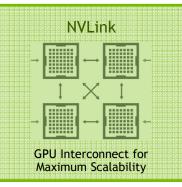


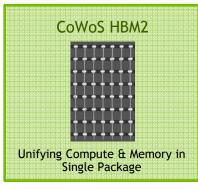
Специализированный сервер для **HPC**

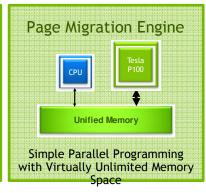


TESLA P100 ACCELERATOR



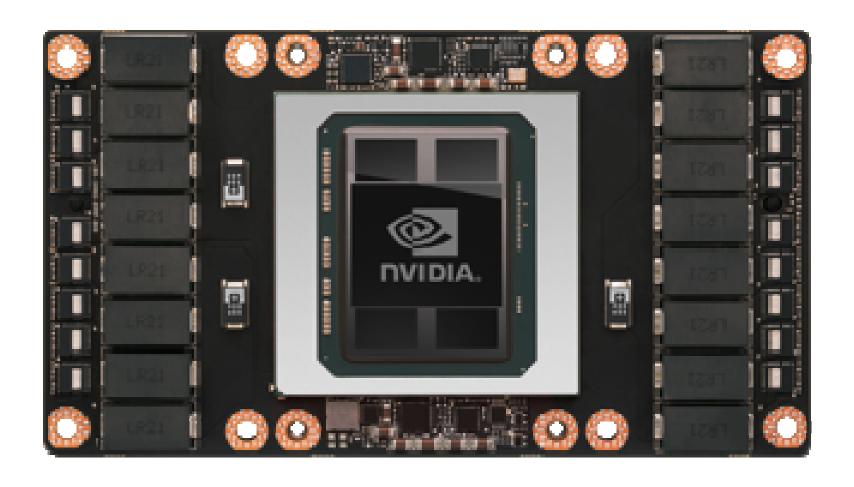






Compute	5.3 TF DP · 10.6 TF SP · 21.2 TF HP
Memory	HBM2: 720 GB/s · 16 GB
Interconnect	NVLink (up to 8 way) + PCIe Gen3
Programmability Page Migration Engine Unified Memory	
Availability	Ships in IBM "Minsky" System: September 2016

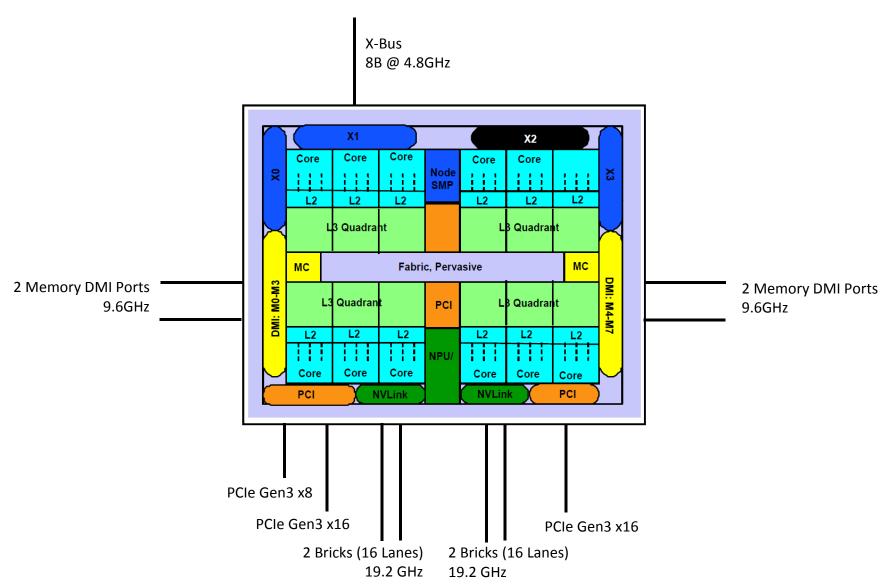
P100







POWER8 with NVLink Module Interfaces



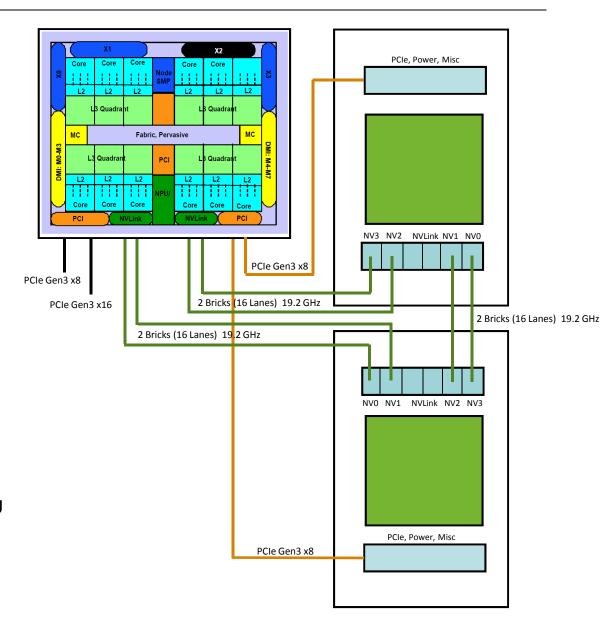
© 2014 IBM Corporation





GPU Interconnect

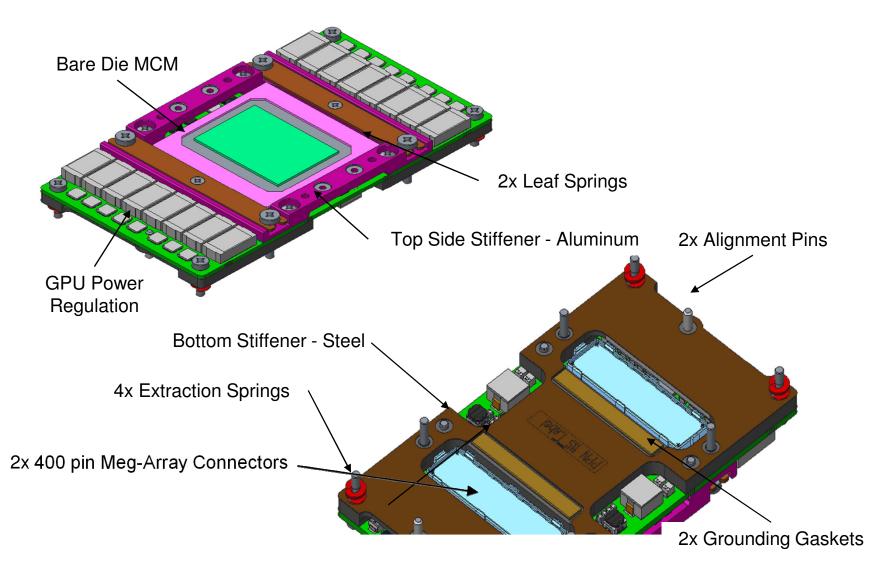
- NVLink Interface
 - High bandwidth interface far exceeding any existing or planned future PCIe interface
 - 16 Lanes CPU to GPU
 - 16 Lanes GPU to GPU
- PCle Interface
 - Used for initiation, control and in band reporting of GPU status
 - PCIe x16 interface on the GPU, Garrison uses this interface in x8 mode



© 2014 IBM Corporation



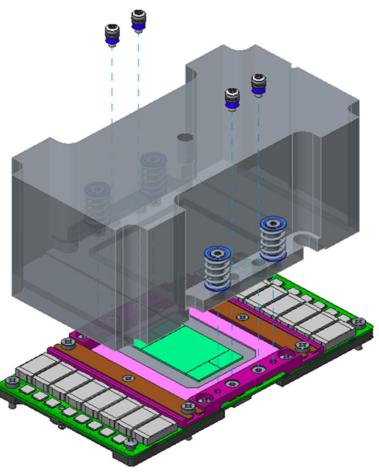
NVIDIA Pascal GPU





NVIDIA Pascal GPU с радиатором

IBM FRU Creation from NVidia PPN

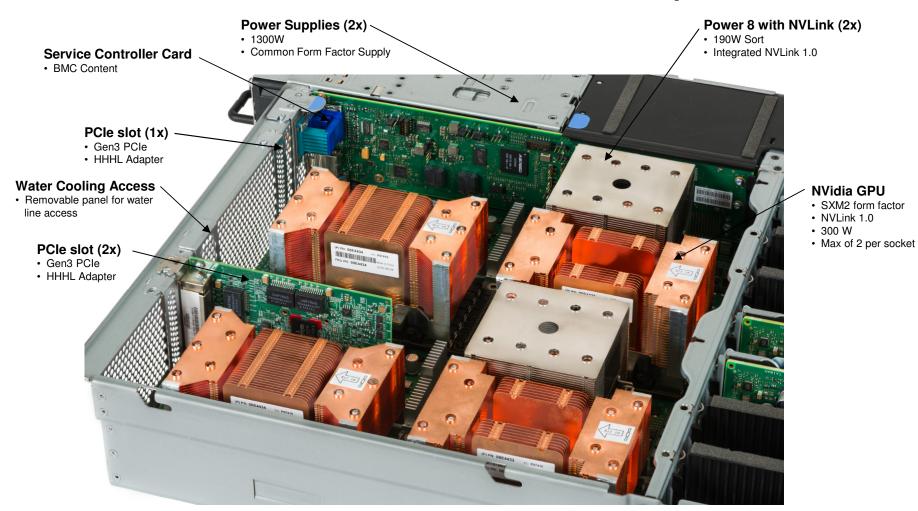


"Assemble TIM, heatsink & NIFs to NVidia PPN"





2 Socket P8 with NVLink + 4 GPU P100 = 21.5TFlops in 2U



© 2014 IBM Corporation



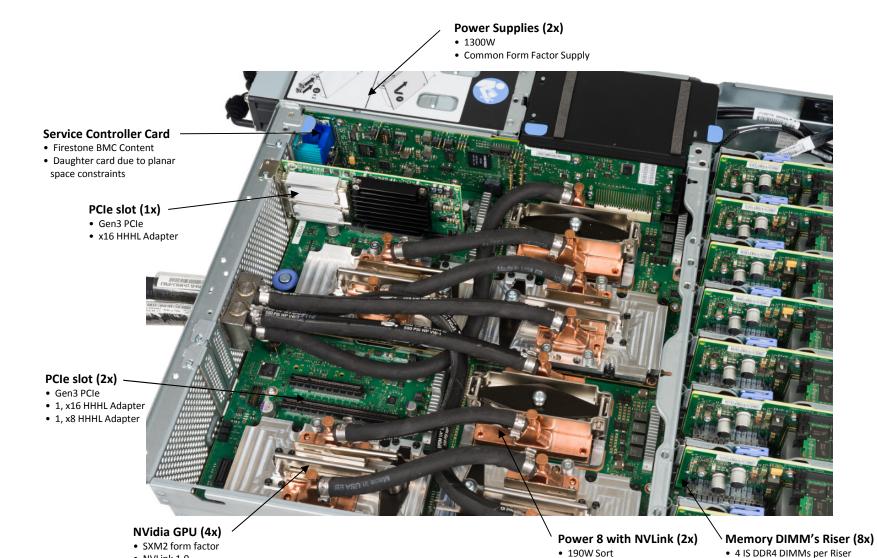
• NVLink 1.0

• 2 per socket

• 300 W

37





• Integrated NVLink 1.0

© 2014 IBM Corporation

Single Centaur per Riser

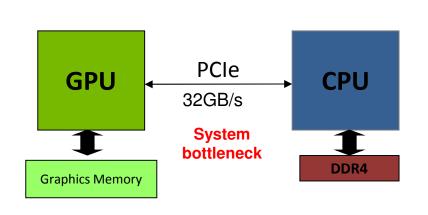
• 32 IS DIMM's total



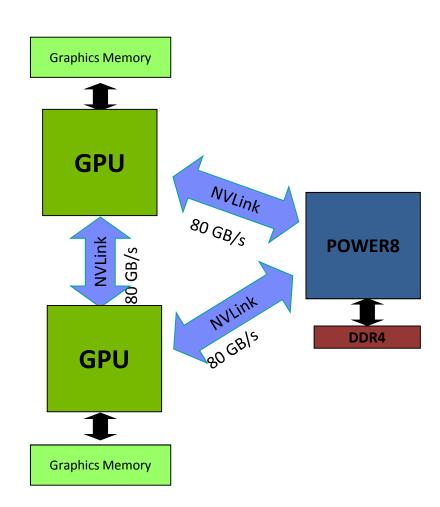


NVLink: формально в 2.5 раза быстрее связь CPU-GPU

Реально: 3-4ГБ/с vs 17-18ГБ/с на PCle3 на NVLink



GPUs Limited by PCIe Bandwidth From CPU-System Memory



NVLink Enables Fast Unified Memory Access between CPU & GPU Memories



Design: Flat and Fat

Дизайн "flat and fat"

- Данные свободно протекают в системе
- Полоса CPU: GPU почти такая же как Системная Память: CPU
- Широкие каналы между GPU подключенными к тому же сокету

Устраняет ограничения PCI-е для многих типов задач

- Пики на старте / сброс итогов
- Обеспечение непрерывного потока данных Host-Device
- Постоянные пересылки между 2 GPU
- Скрытые пересылки по шине в направлении Host-Device

