

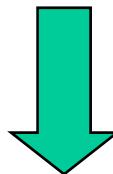
GAPE-DR Project: a combination of peta-scale computing and high-speed networking

Kei Hiraki

Department of Computer Science
The University of Tokyo

Computing System for real Scientists

- **Fast Computation at low cost**
 - High-speed CPU, huge memory, good graphics
- **Very fast data movements**
 - High-speed network, Global file system, Replication facilities
- **Transparency to local computation**
 - No complex middleware, or no modification to existing software



- **Real Scientists are not computer scientists**
- **Computer scientists aren't work forces for real scientists**

Our goal

- HPC system
 - as an infrastructure of scientific research
 - Simulation, data intensive computation, searching and data mining
- Tools for real scientists
 - High-speed, low-cost, and easy to use
 - ⇒ Today's supercomputer cannot
 - High-speed/Low cost computation
 - High-speed global networking
 - OS/ ALL IP network technology

GRAPE-DR
project

GRAPE-DR

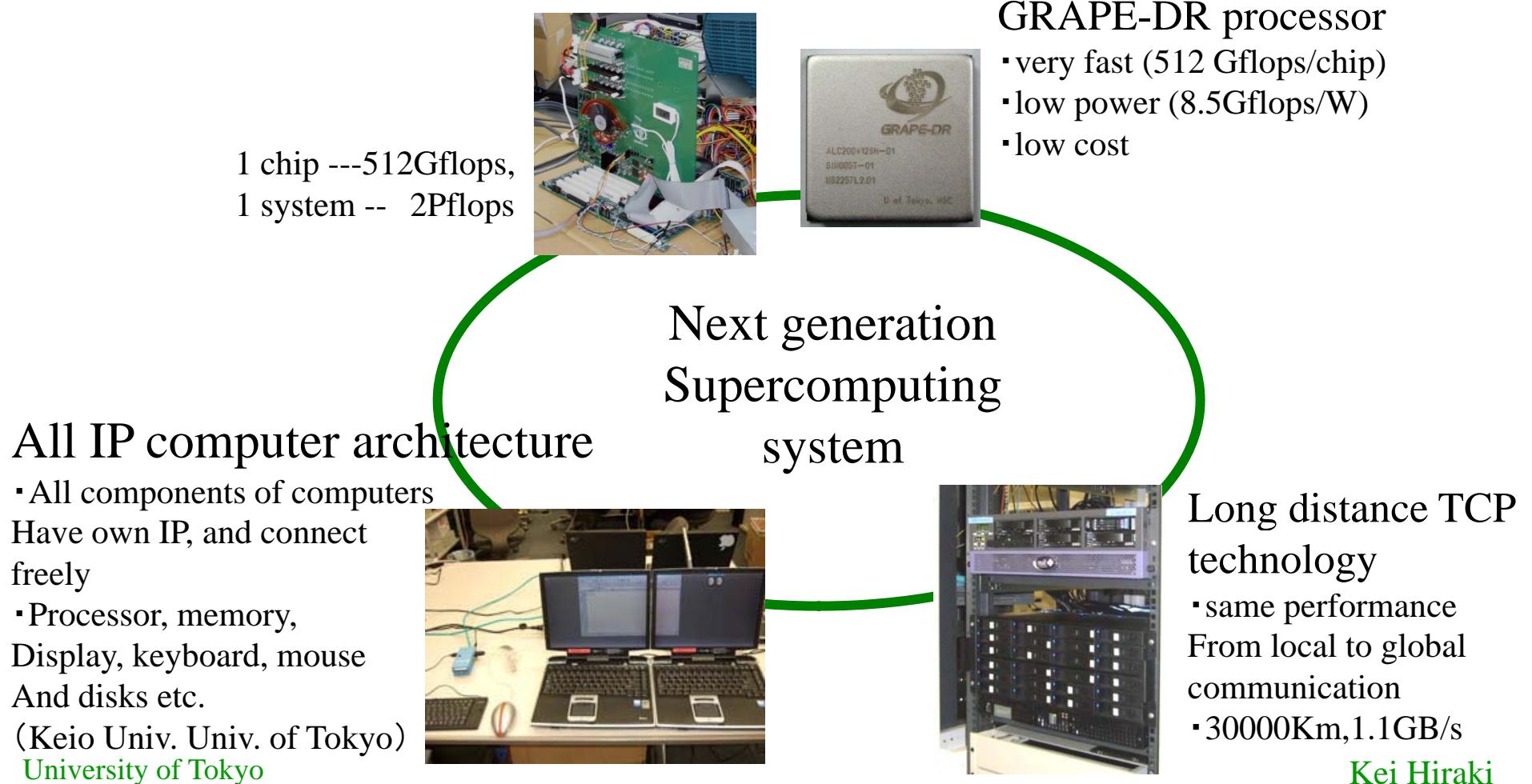
- GRAPE-DR:Very high-speed attached processor
 - 2Pflops(2008) Possibly fastest supercomputer
 - Successor of Grape-6 astronomical simulator
- 2PFLOPS on 512 node cluster system
 - 512 Gflops / chip
 - 2 M processor / system
- Semi-general-purpose processing
 - N-body simulation, gridless fluid dynamics
 - Linear solver, molecular dynamics
 - High-order database searching

Data Reservoir

- Sharing Scientific Data between distant research institutes
 - Physics, astronomy, earth science, simulation data
- Very High-speed single file transfer on Long Fat pipe Network
 - > 10 Gbps, > 20,000 Km, > 400ms RTT
- High utilization of available bandwidth
 - Transferred file data rate > 90% of available bandwidth
 - Including header overheads, initial negotiation overheads
- OS and File system transparency
 - Storage level data sharing (high speed iSCSI protocol on stock TCP)
 - Fast single file transfer

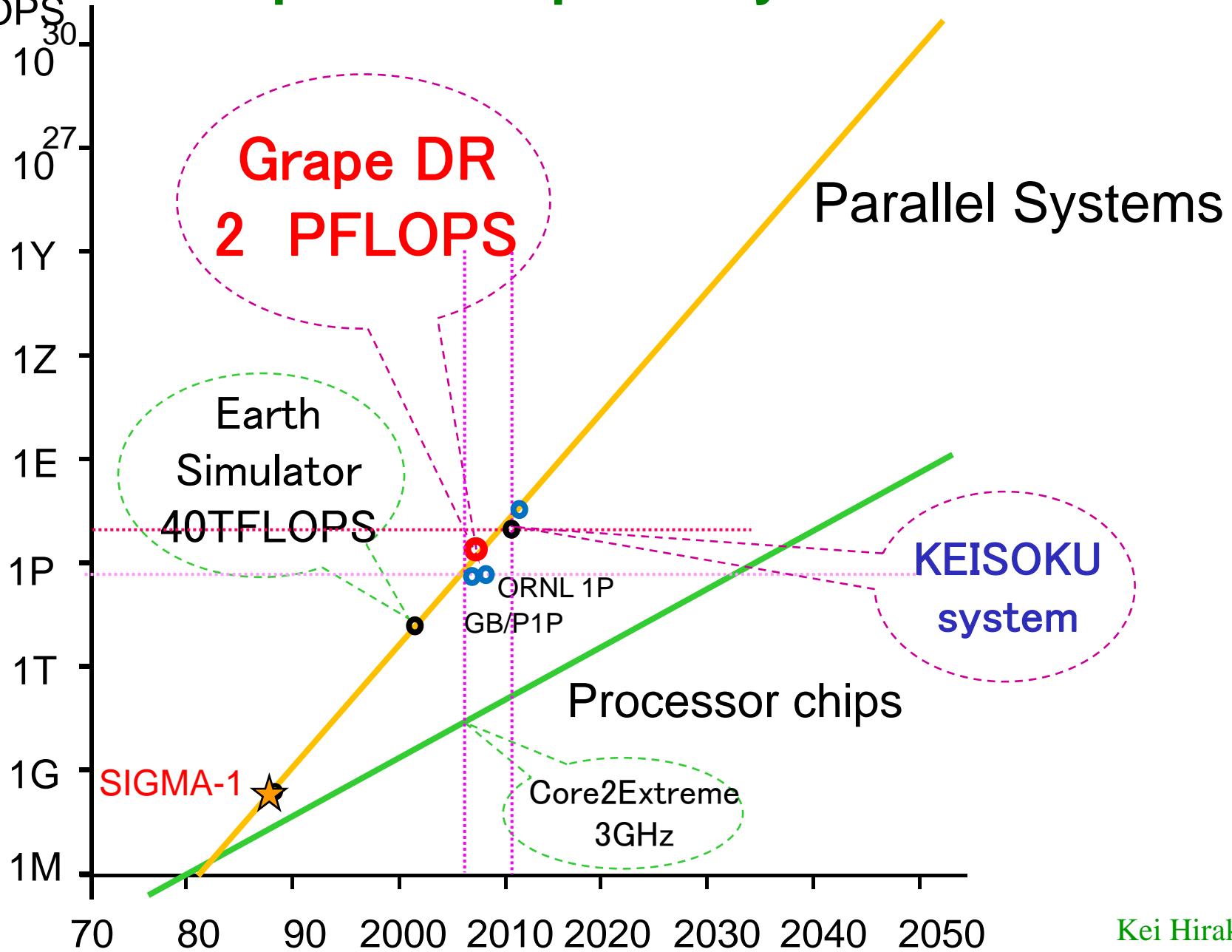
Our Next generation Supercomputing system

- Very high-performance computer
- Supercomputing/networking system for non CS area

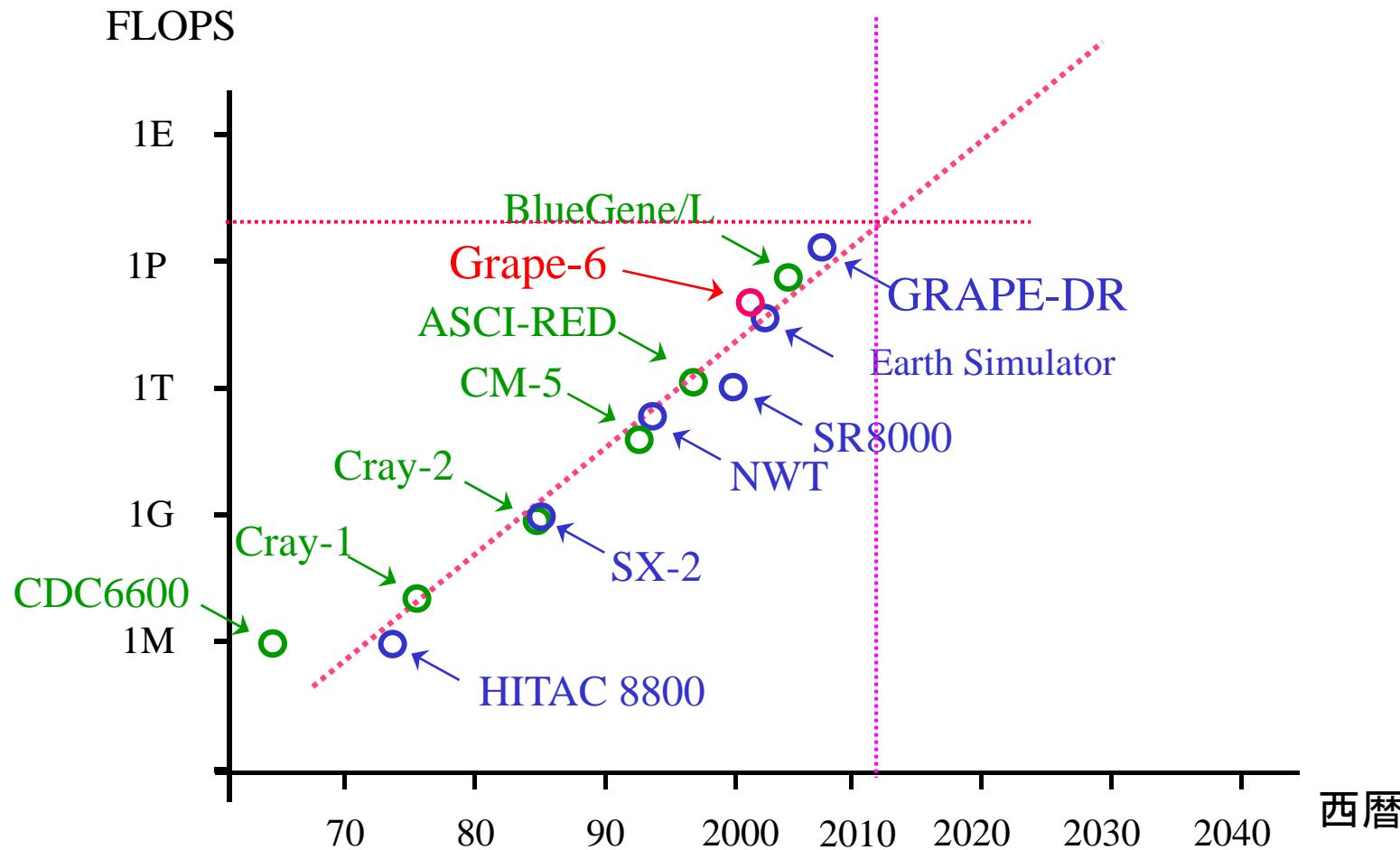


Peak
FLOPS

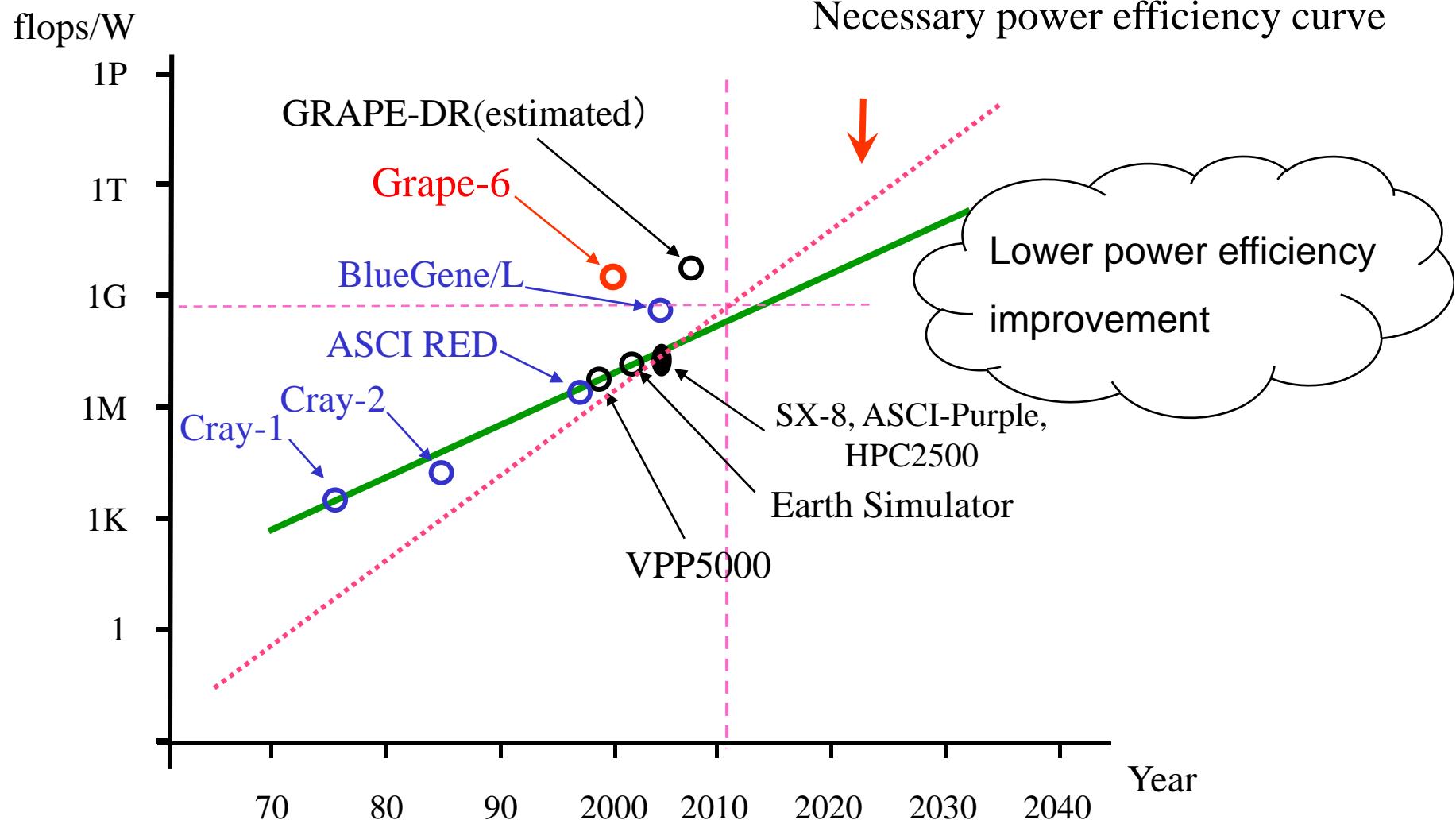
Speed of Top-End systems



History of speedup (some important systems)



History of power efficiency



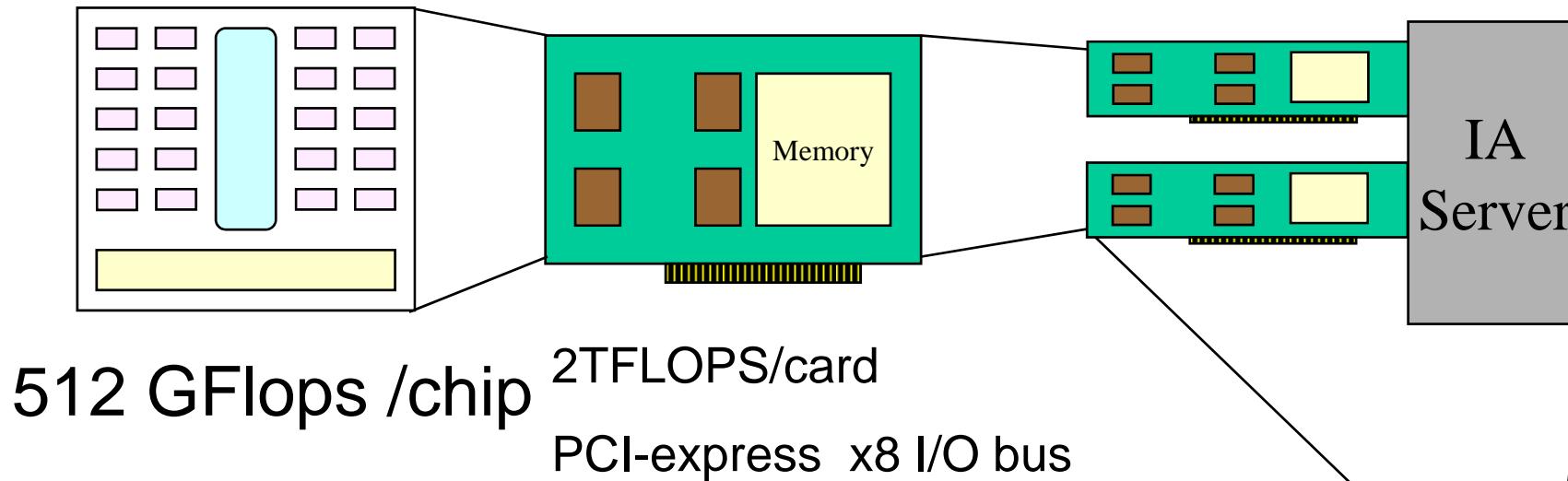
GRAPE-DR project (2004 to 2008)

- Development of practical MPP system
 - Pflops systems for real scientists
- Sub projects
 - Processor system
 - System software
 - Application software
 - Simulation of stars and galaxies, CFD, MD, Linear system
- 2-3 Pflops peak, 20 to 30 racks
- 0.7 – 1Pflops Linpack
- 4000 – 6000 GRAPE-DR chips
- 512 – 768 servers + Interconnect
- 500Kw/Pflops

GRAPE-DR system

- 2 Pflops peak, 40 racks
- 512 servers + Infiniband
- 4000 GRAPE-DR chips
- 0.5MW/Pflops(peak)
- World fastest computer?
 - GRAPE-DR (2008)
 - IBM BlueGene/P(2008?)
 - Cray Baker(ORNL) (2008?)

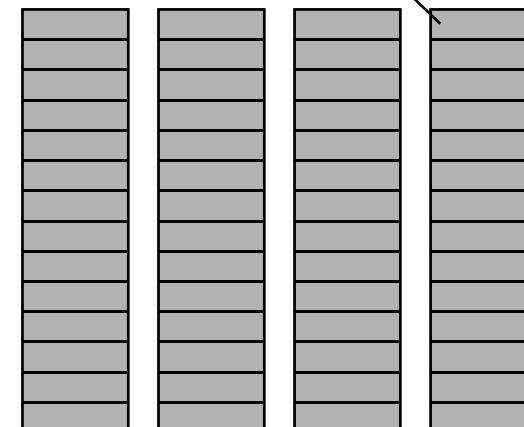
GRAPE-DR chip



2PFLOPS/system

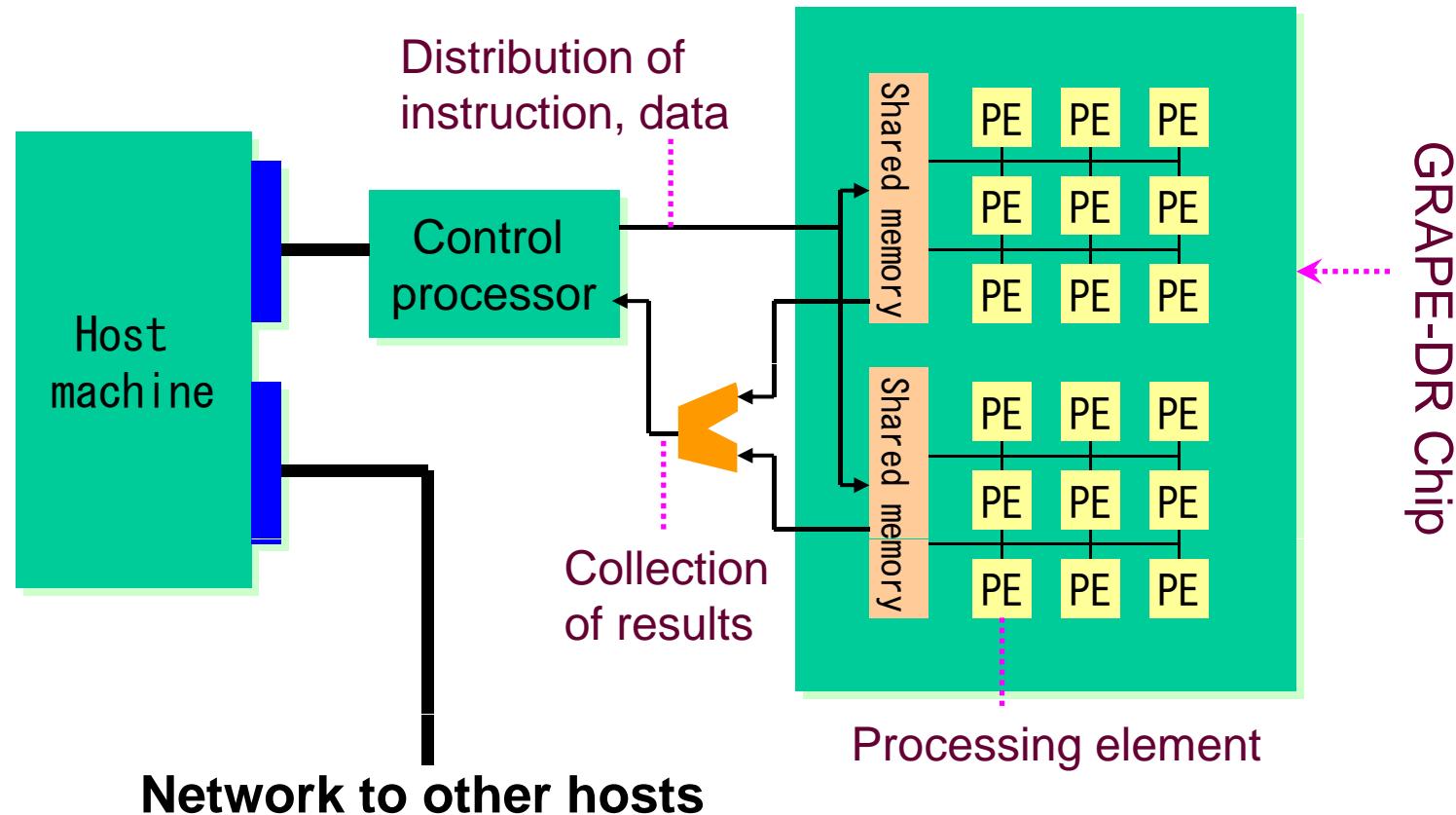
2M PE/system

20Gbps Infiniband



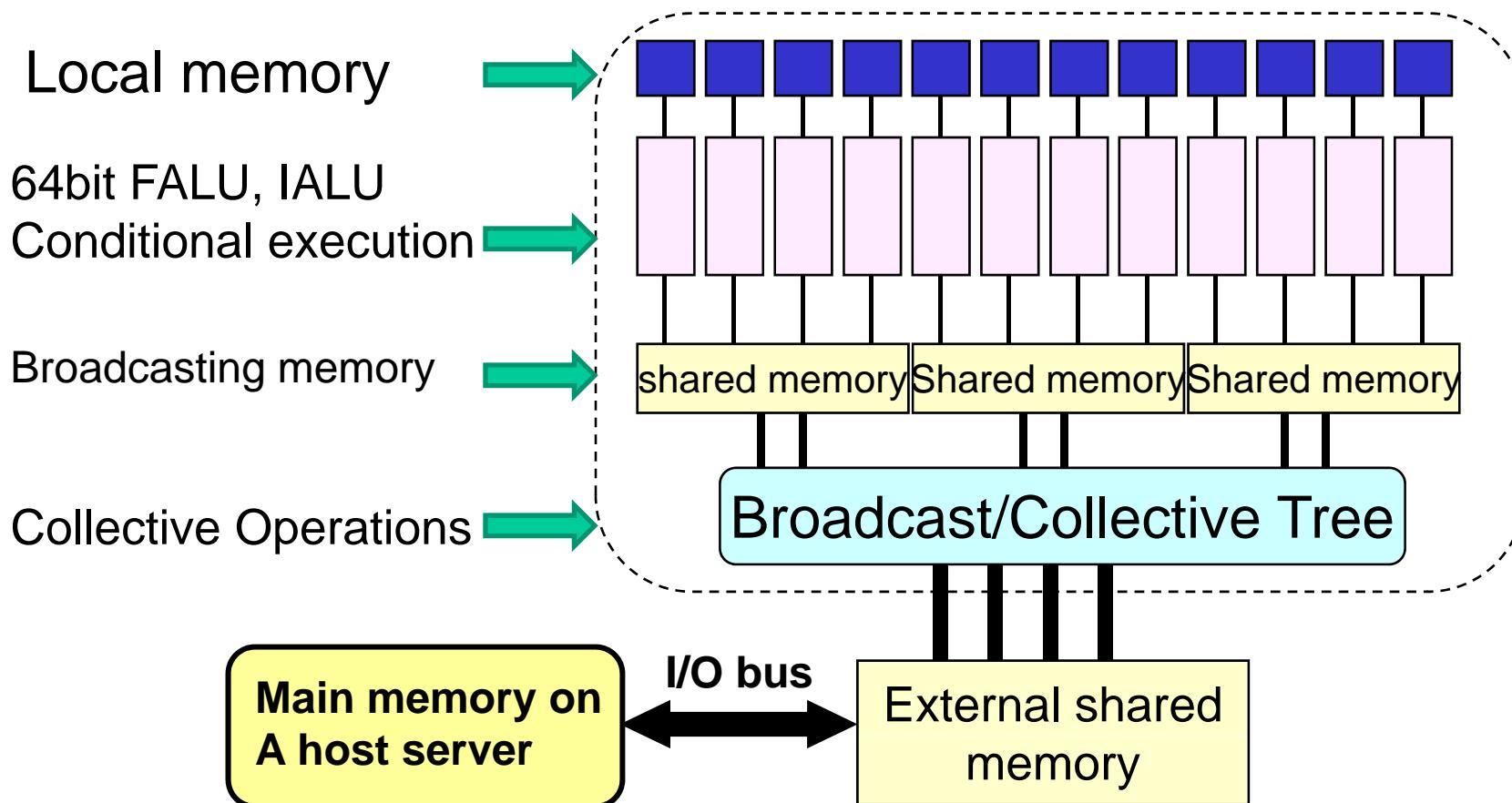
Node Architecture

- Accelerator attached to a host machine via I/O bus
- Instruction execution between shared memory and local data
- All the PE in a chip execute the same instruction (SIMD)



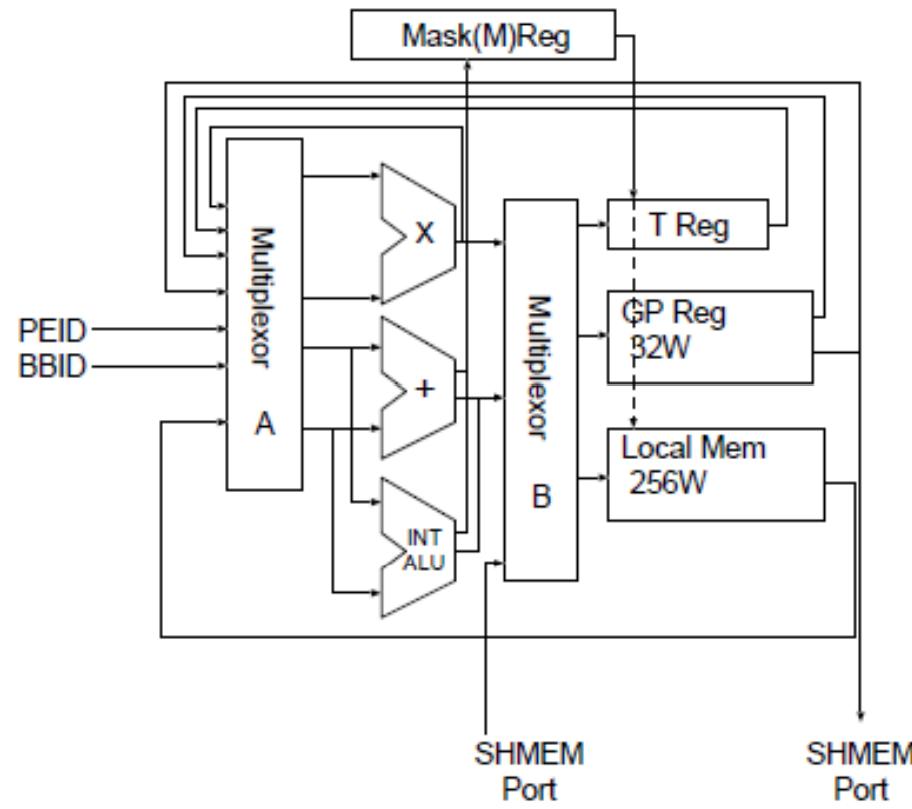
GRAPE-DR Processor chip

- SIMD architecture
- 512PEs
- No interconnections between PEs



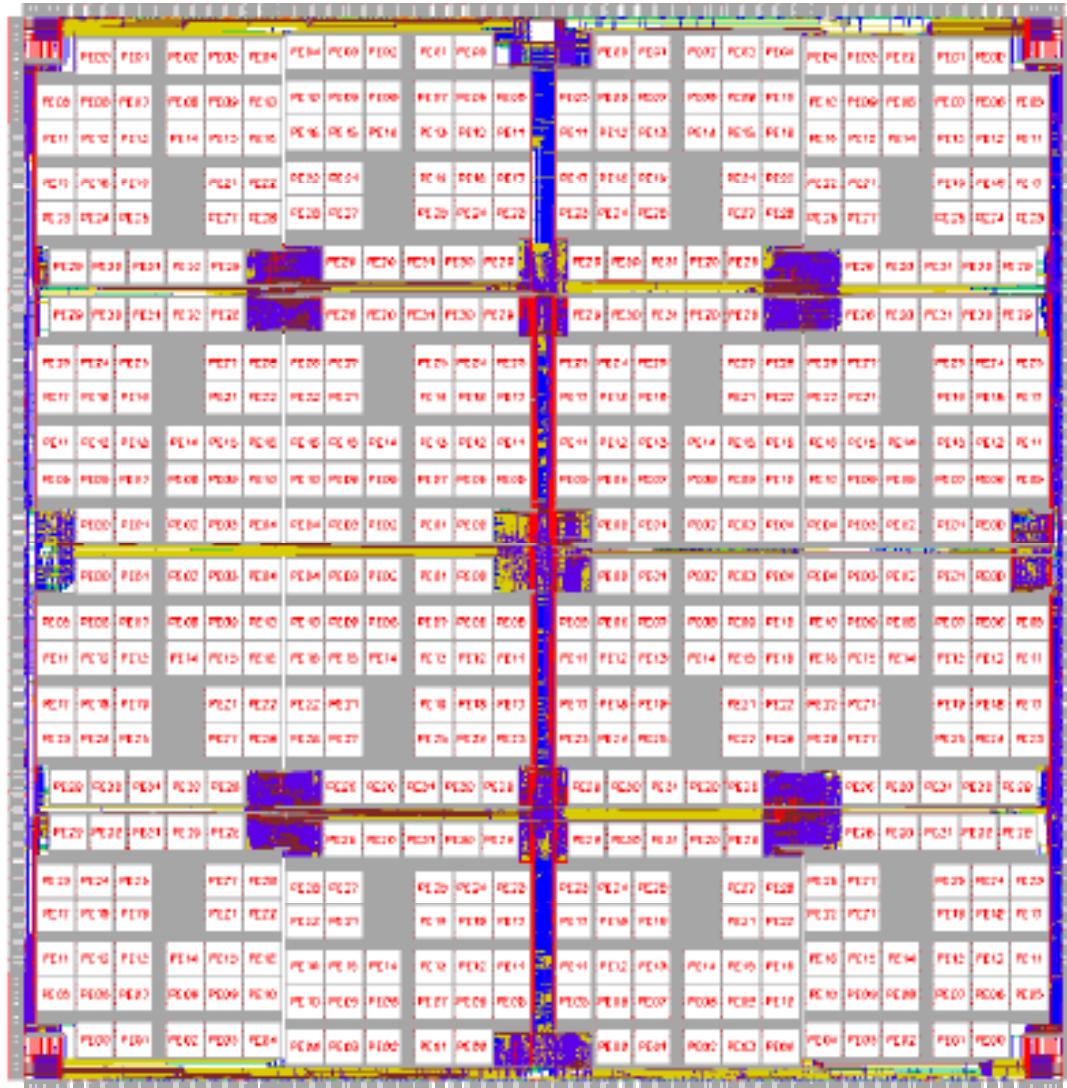
Processing Element

- 512 PE in a chip



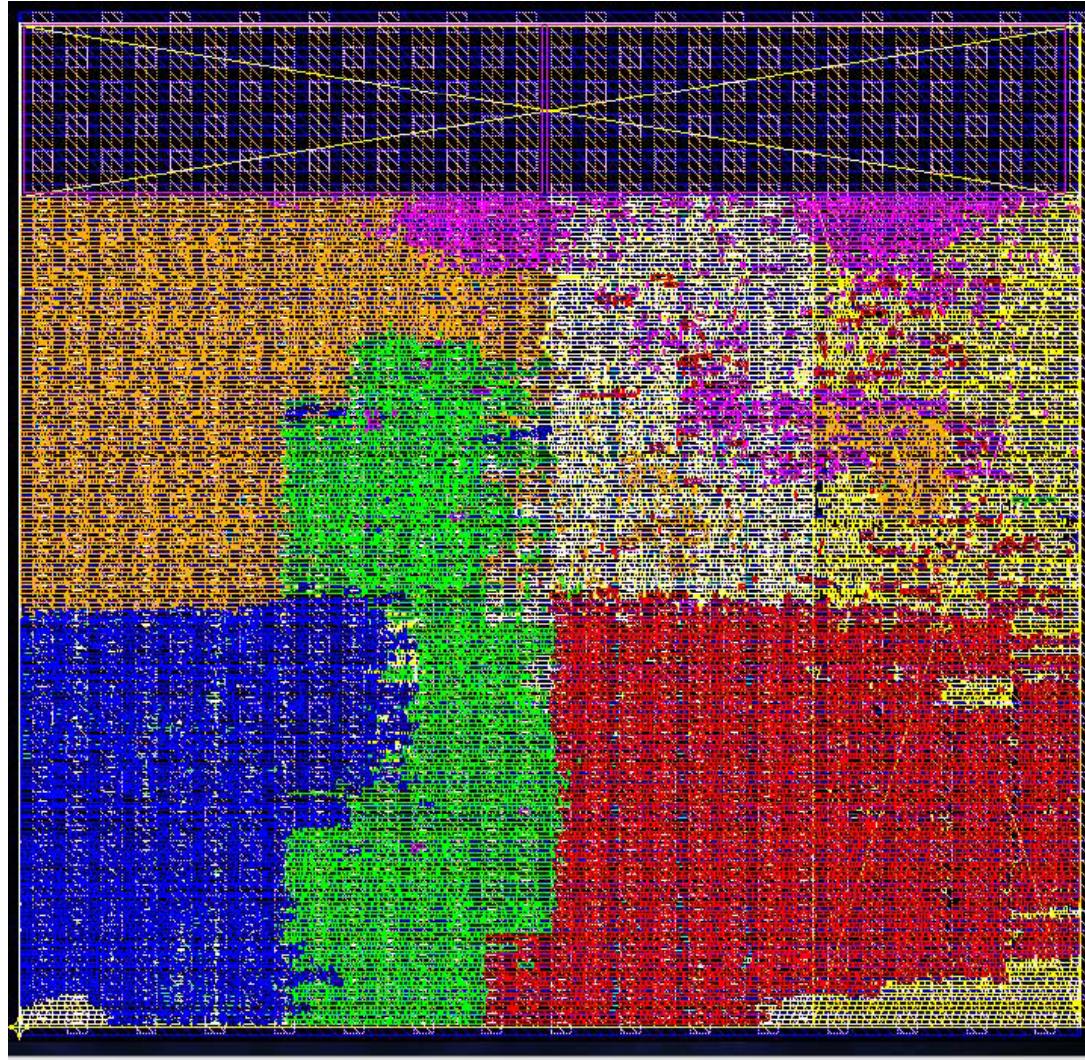
- Float Mult
- Float add/sub
- Integer ALU
- 32-word registers
- 256-word memory
- communication port

- 90nm CMOS
- 18mm x18mm
- 400M Tr
- BGA 725
- 512Gflops(32bit)
- 256Gflops(64bit)
- 60W(max)



Chip layout in a processing element

Module Name	Color
grf0	red
fmul0	orange
fadd0	green
alu0	blue
lm0	magenta
others	yellow



GRAPE-DR Processor chip

Engineering Sample (2006)

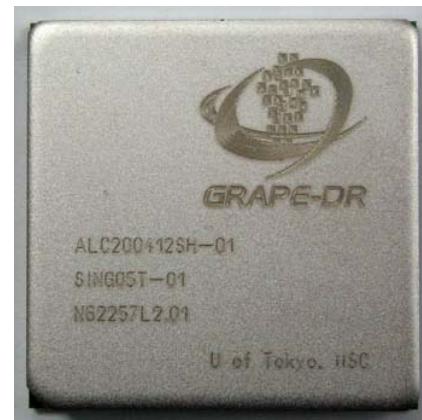
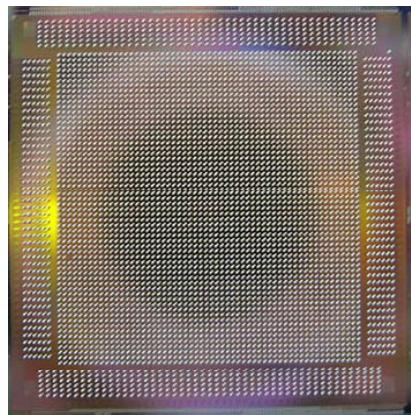
512 Processing Elements/Chip (**Largest number in a chip**)

Working on a prototype board (at designed speed)

500MHz, 512Gflops (**Fastest chip for production**)

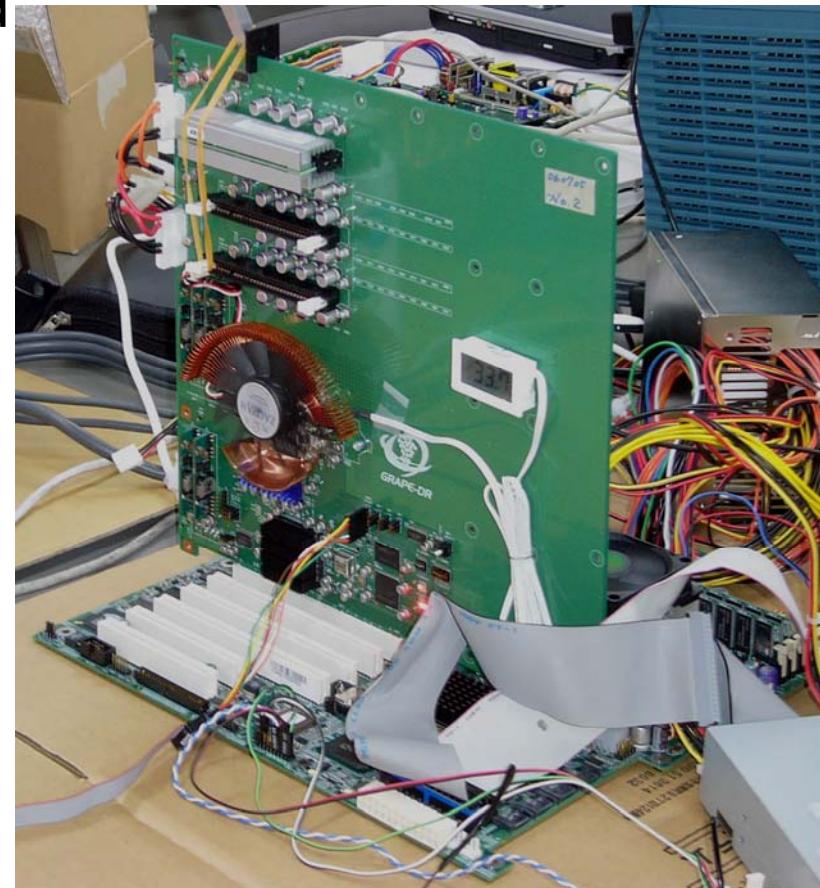
Power Consumption max. 60W Idle 30W

(**Lowest power per computation**)



GRAPE-DR Prototype boards

- For evaluation(1 GRAPE-DR chip, PCI-X bus)
- Next step is 4 chip board



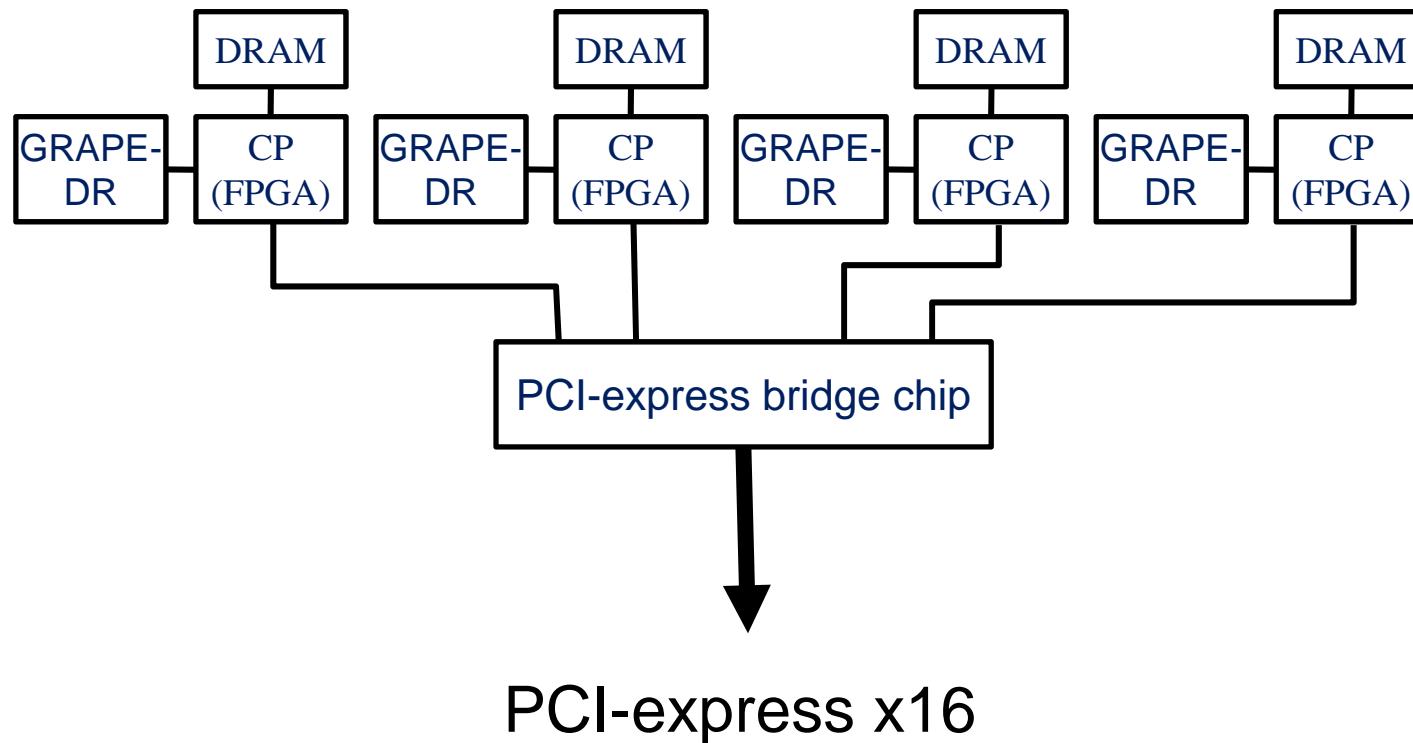
GRAPE-DR Prototype boards(2)

- 2nd prototype
- Single-chip board
- PCI-Express x8 interface
- On-board DRAM



Block diagram of a board

- 4 GRAPE-DR chip ⇒ 2 Tflops(single), 1Tflops(double)
- 4 FPGA for control processors
- On-board DDR2-DRAM 4GB
- Interconnection by a PCI-express bridge chip



Compiler for GRAPE-DR

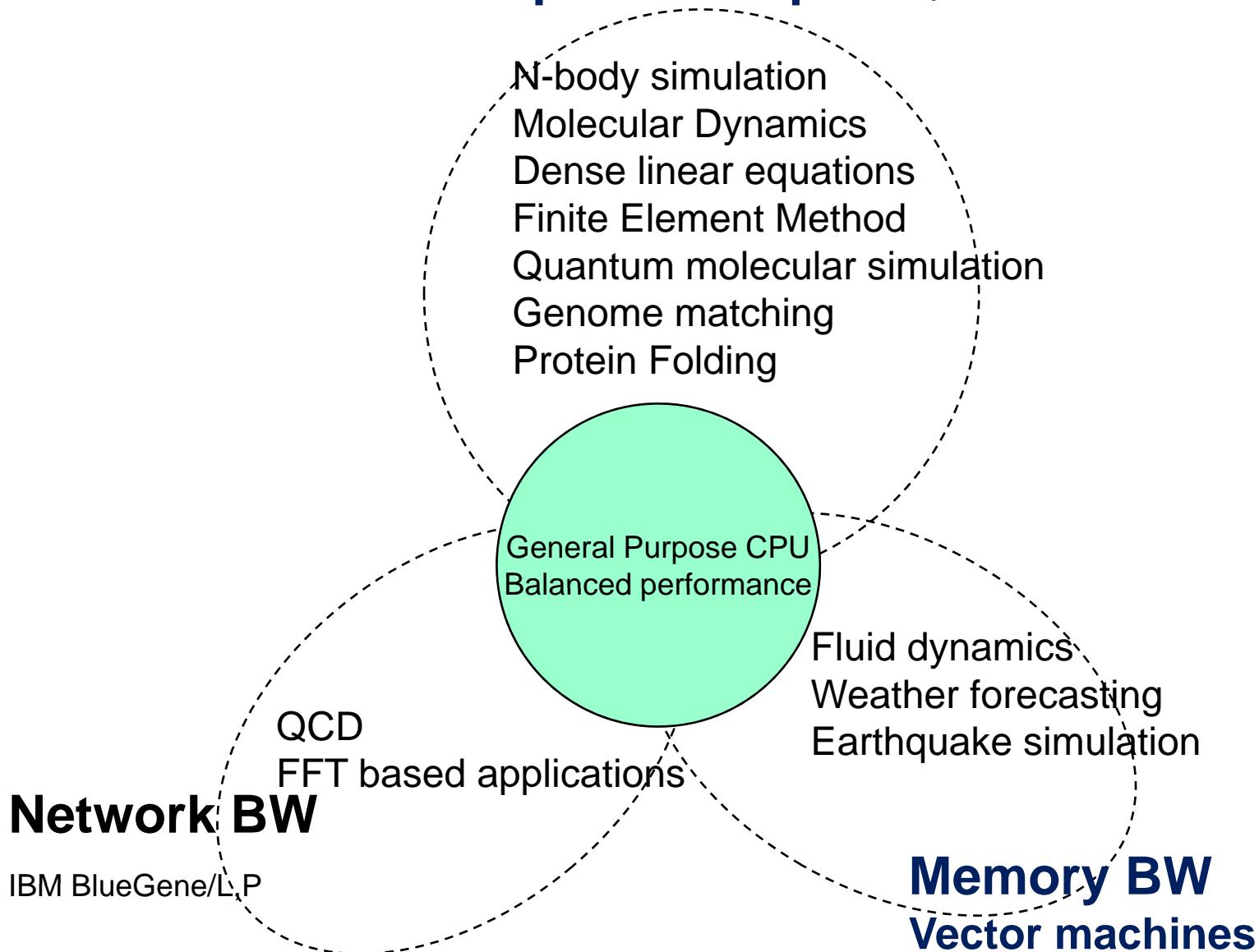
- GRAPE-DR optimizing compiler
 - Parallelization with global analysis
 - Generation of multi-threaded code for complex operation
- Sakura-C compiler
 - Basic optimization
(PDG,flow analysys、pointer analysis etc.)
 - Souce program in C ⇒ intermediate language
⇒ GRAPE-DR machine code
- Currently, 2x to 3x slower than hand optimized codes

GRAPE-DR Application software

- Astronomical Simulation
- SPH(Smoothed Particle Hydrodynamics)
- Molecular dynamics
- Quantum molecular simulation
- Genome sequence matching
- Linear equations on dense matrices
 - Linpack

Application fields for GRAPE-DR

Computation Speed (GRAPE-DR, GPGPU etc.)



Eflops system

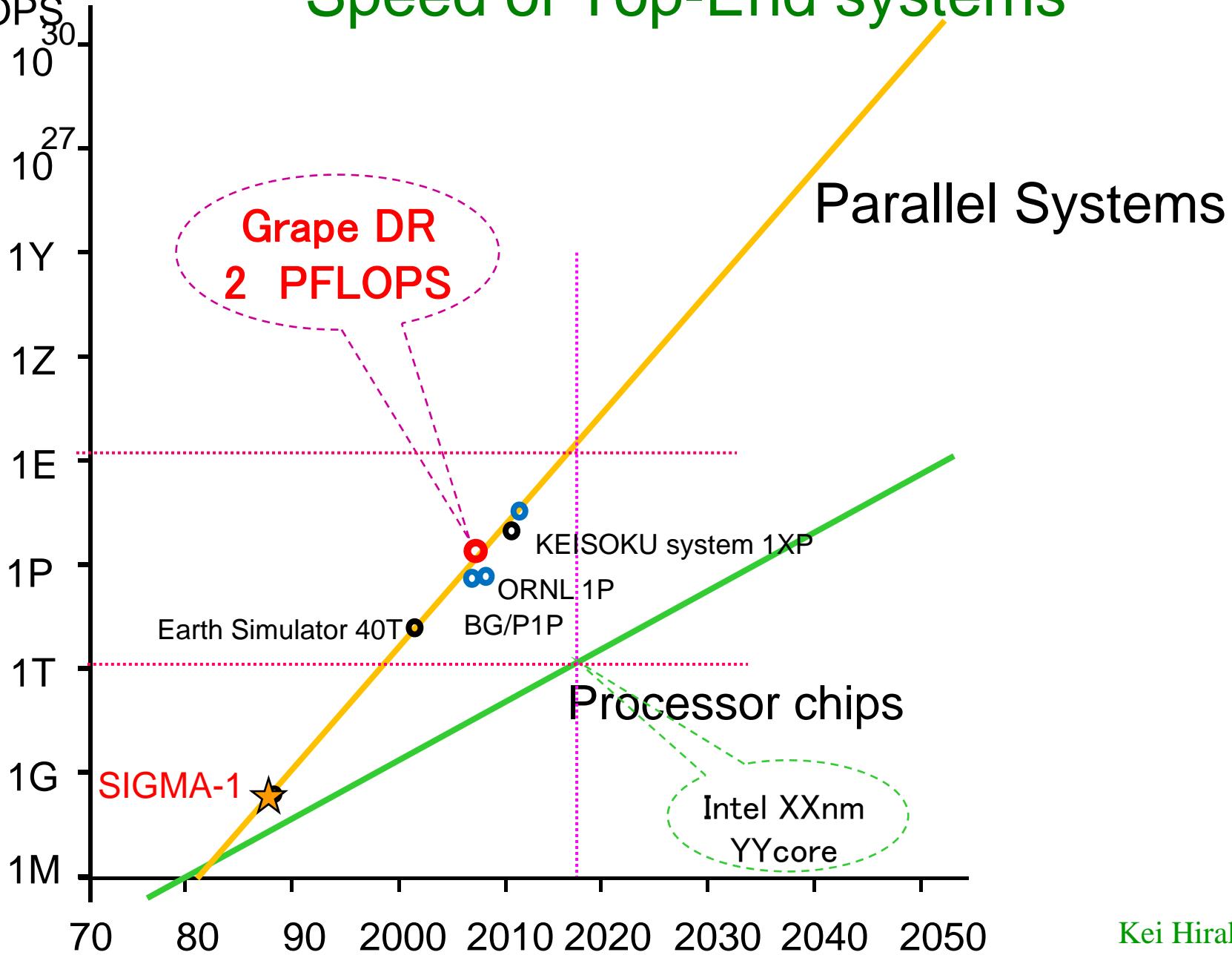
- 2011 10Pflops systems will be appeared
 - Next Generation Supercomputer by Japanese Government
 - IBM BlueGene/Q system
 - Cray Cascade system (ORNL)
 - IBM HPCS

Architecture

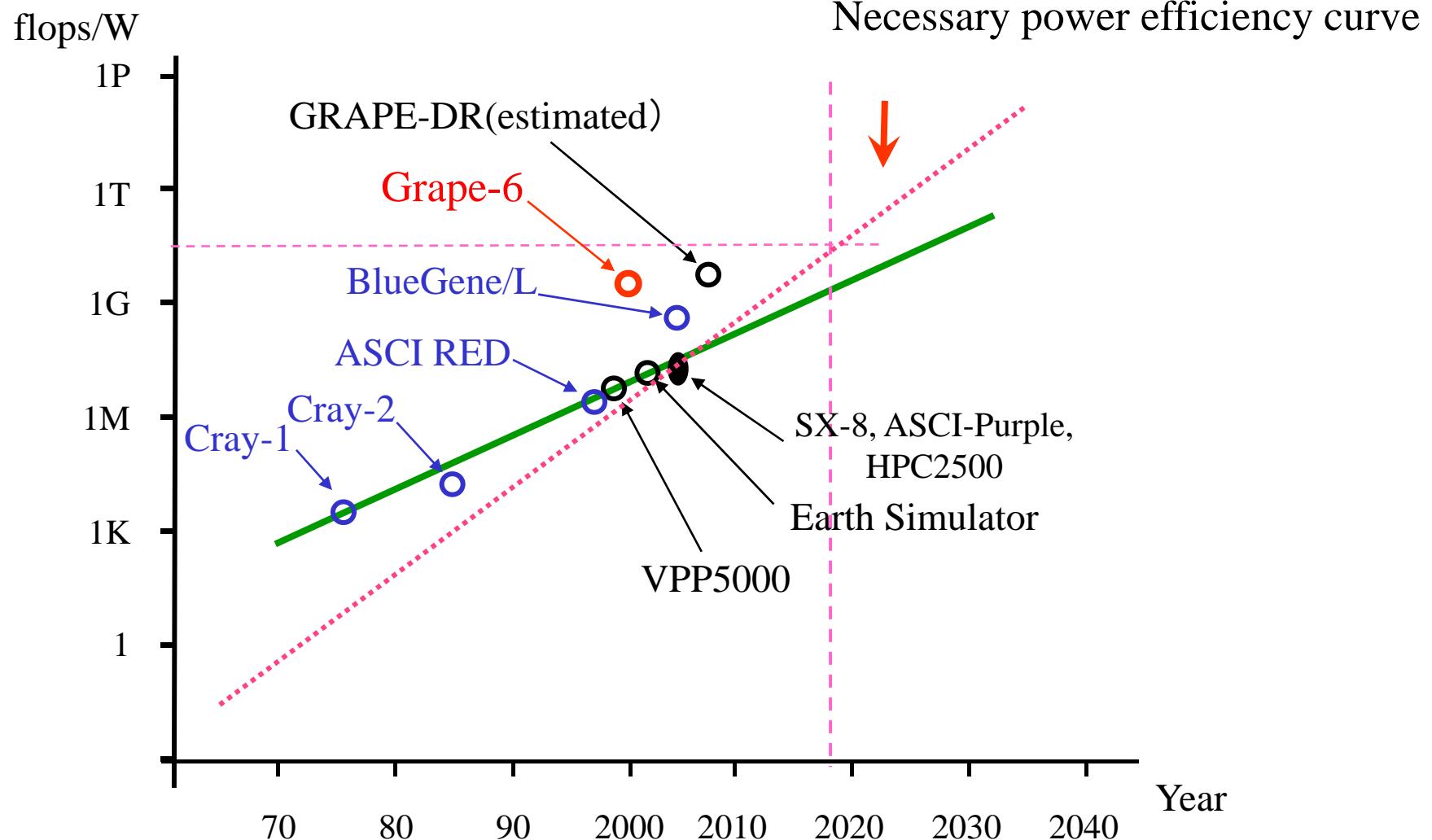
- Cluster of Commodity MPU (Sparc/Intel/AMD/IBM Power)
- Vector processor
- Densely integrated low-power processors
- 1 ~ 1.5 ~2.5 MW/Pflops power efficiency

Peak
FLOPS

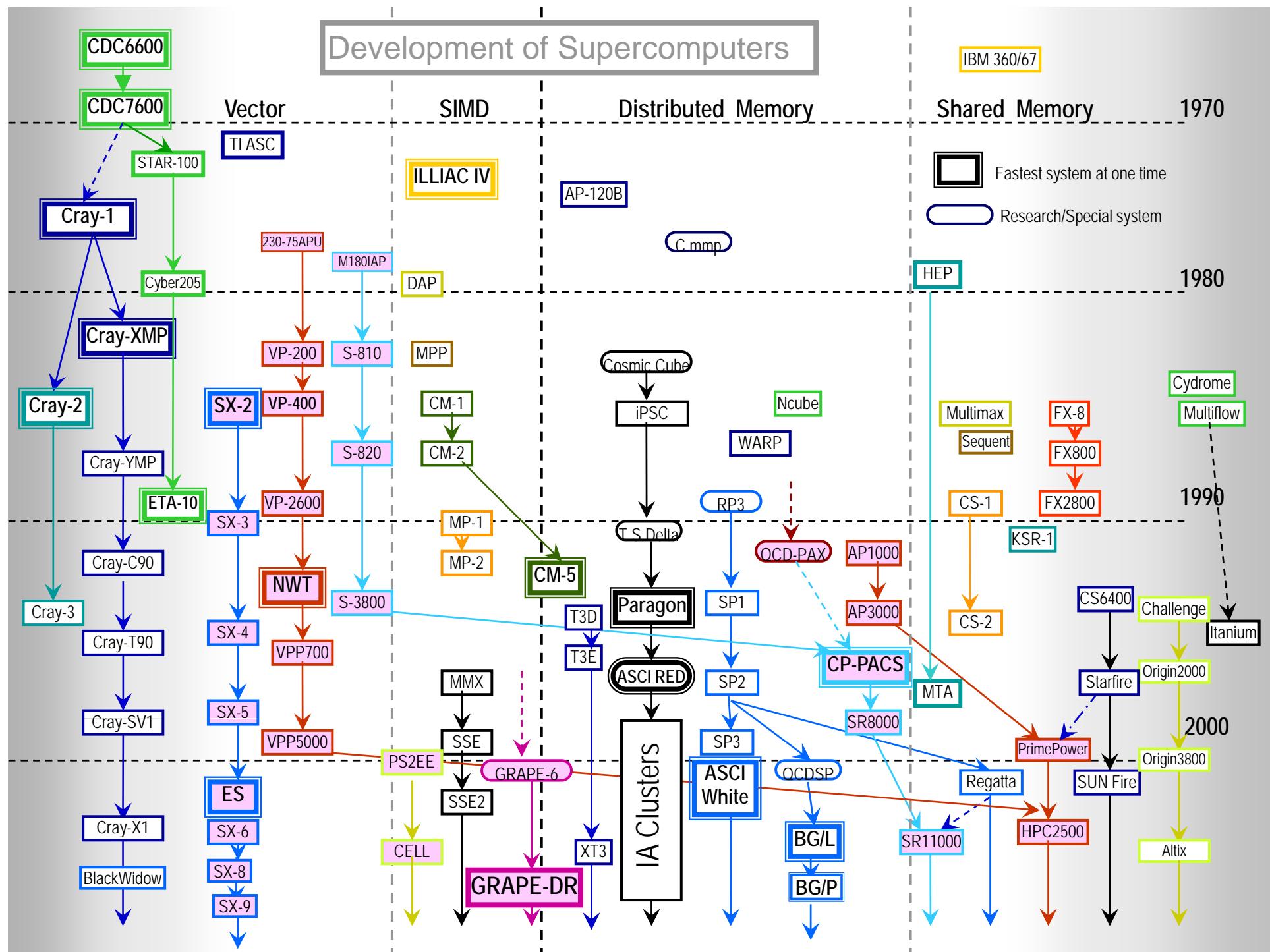
Speed of Top-End systems



History of power efficiency



Development of Supercomputers



Eflops system

- Requirements for building Eflops systems
 - Power efficiency
 - MAX Power \sim 30 MW
 - Efficiency must be smaller than 30 Gflops/W
 - Number of processor chips
 - 300,000 chips ? (3 Tflops/chip)
 - Memory system
 - Smaller memory / processor chip \Rightarrow Power consumption
 - Some kind of shared memory mechanism is a must
 - Limitation of cost
 - 1 billion \$/ Eflops \Rightarrow 1000 \$/Tflops

Effective approaches

- General purpose MPUs (clusters)
- Vectors
- GPGPU
- SIMD (GRAPE-DR)
- FPGA based simulating engine

Unsuitable architecture

- Vectors
 - Too much power for numerical computation
 - Too much power for wide memory interface
 - Less effective to wide range of numerical algorithms
 - Too expensive for computing speed
- General purpose MPUs (clusters)
 - Too low FPU speed per die size
 - Too much power for numerical computation
 - Too many processor chips for stable operation

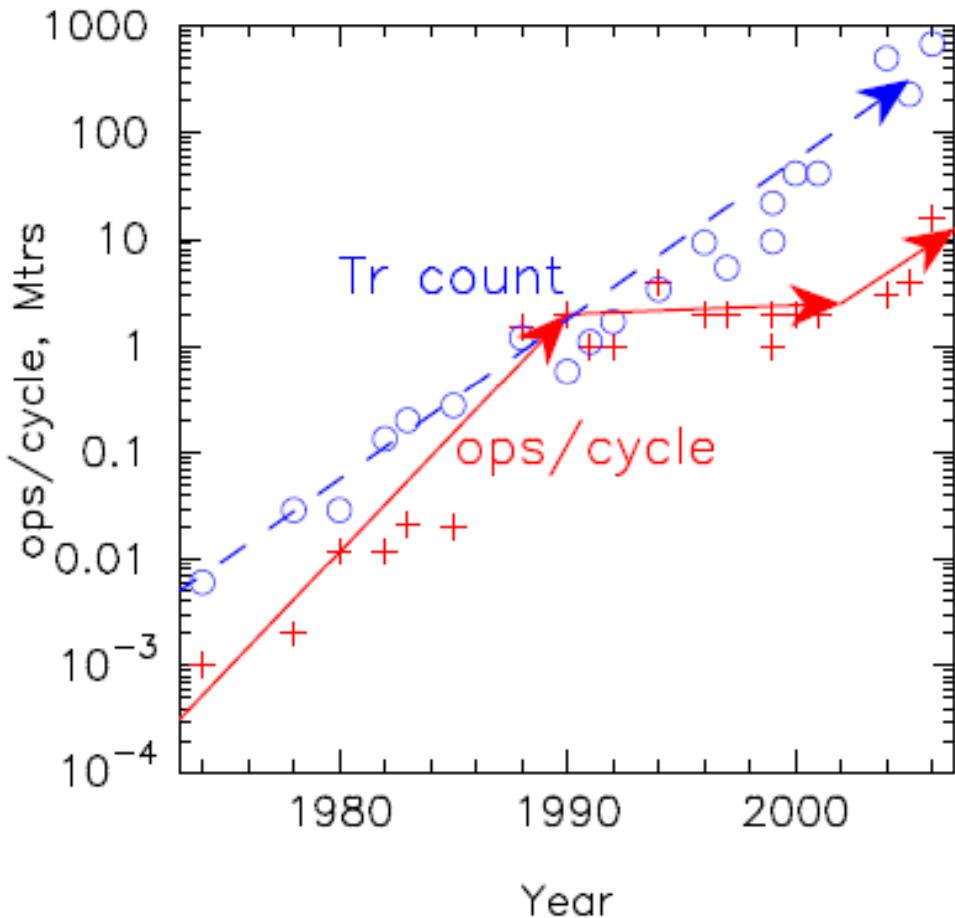
Unsuitable architecture

- Vectors
 - Current version 65nm Technology
 - 0.06 Gflops/W peak
 - 2,000,000\$/Tflops
- General purpose MPUs (clusters)
 - Current version 65nm Technology
 - 0.3 Gflops/W peak (BlueGene/P, TACC etc.)
 - 200,000\$/Tflops

Goals are
30Gflops/W, 1000\$/Tflops, 3Tflops/chip

General purpose CPU

Evolution of Microprocessors



- Transistors: 1000 times in last 15 years
- FPUs: 8 times more in the same period
- a factor of 100 “lost”

Architectural Candidates

- SIMD (GRAPE-DR)

- 30% FALU
- 20% IALU
- 20% Register file
- 20% Memory
- 10% Interconnect

About 50% is used for computation

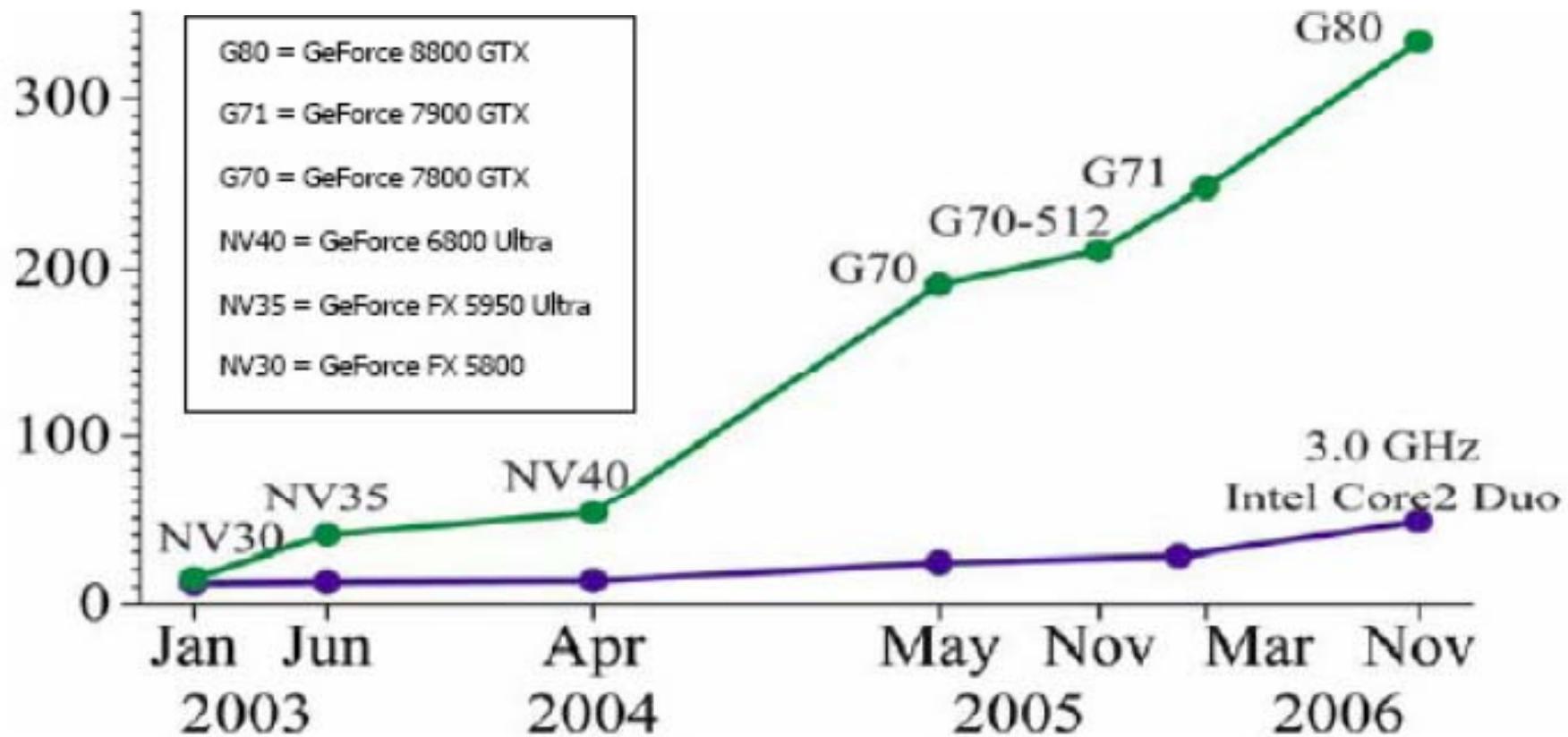
Very efficient

- GPGPU

- Less expensive due to the number of production
- Less efficient than properly designed SIMD processor
 - Larger die size
 - Larger power consumption

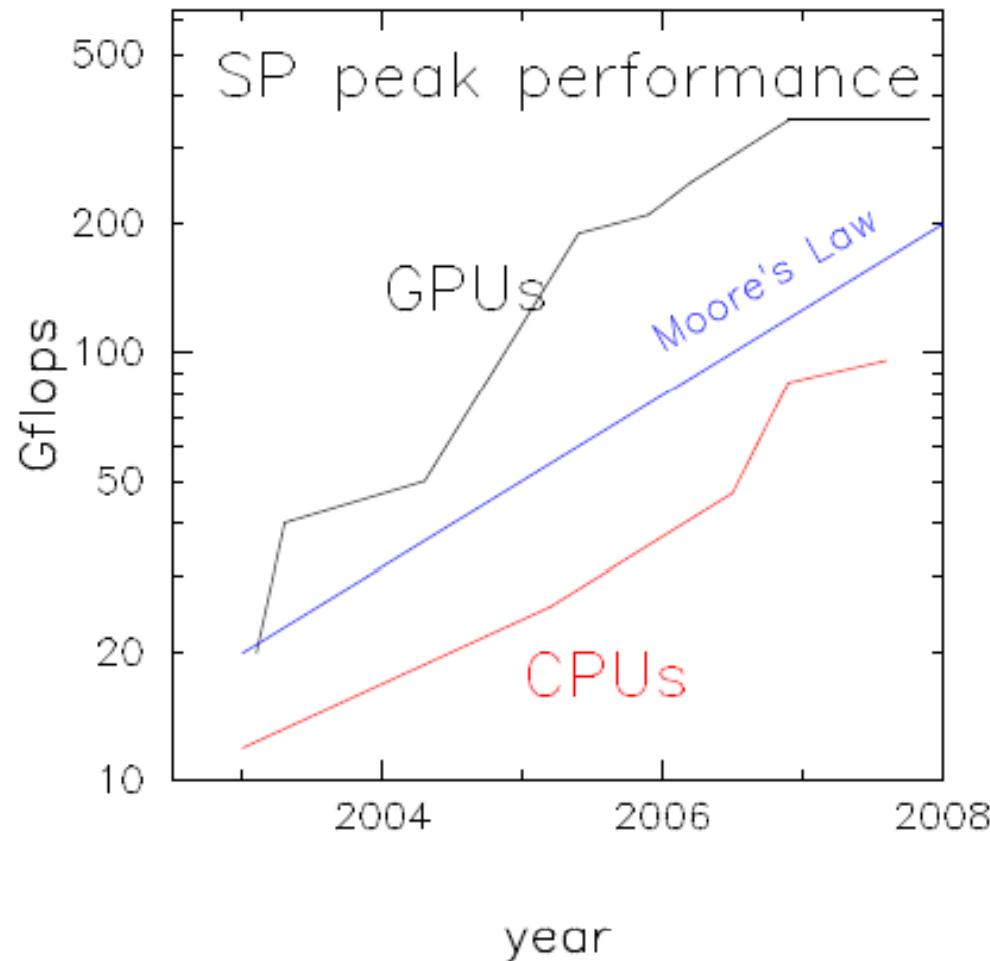
GPGPU(1)

GFLOPS



“GPUs beat Moore’s Law!”

GPGPU(2)



- Faster-than-Moore period ended in 2005
- Microprocessors are catching up
- DP performance?
- Design limit with memory bandwidth

Comparison with GPGPU

Pros:

- Significantly better silicon usage
(512PEs with 90nm)
- Designed for scientific applications
reduction, small communication overhead, etc

Cons:

- Higher cost per silicon area...
(small production quantity)
- Longer product cycle... 5 years vs 1 year

Good implementations of N -body code on GPGPU
are coming (Hamada, Nitadori, Portegies Zwart,
Harris, ...)

Comparison with GPGPU(2)

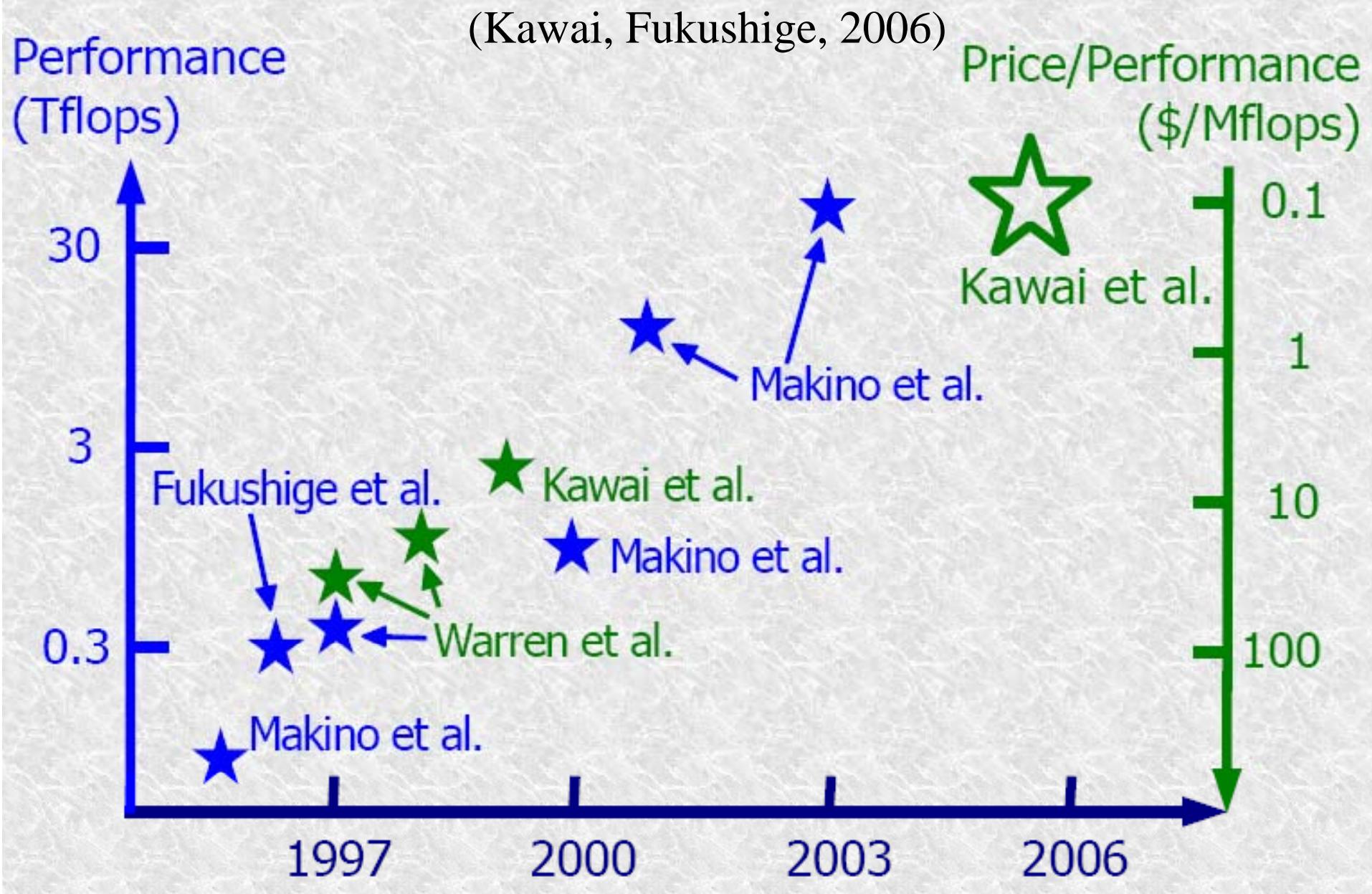
	GRAPE-DR	nV G92	AMD FS9170
Design rule	90	65	55
Clock(GHz)	0.5	1.5	0.8
# FPUs	512	112	320
SP peak(GF)	512	336	512
DP peak(GF)	256	—	?
Power(W)	65	70?	150?

FPGA based simulator

- Commercial FPGAs cannot do FP operations faster than microprocessors.
 - Has been so in last decade
 - Will remain so for foreseeable future
- Custom reconfigurable processors cannot compete with commercial FPGAs in price-performance ratio.
 - Longer development cycle
 - Far smaller quantity

Good for applications with short-wordlength

Astrophysical Simulations in GB History

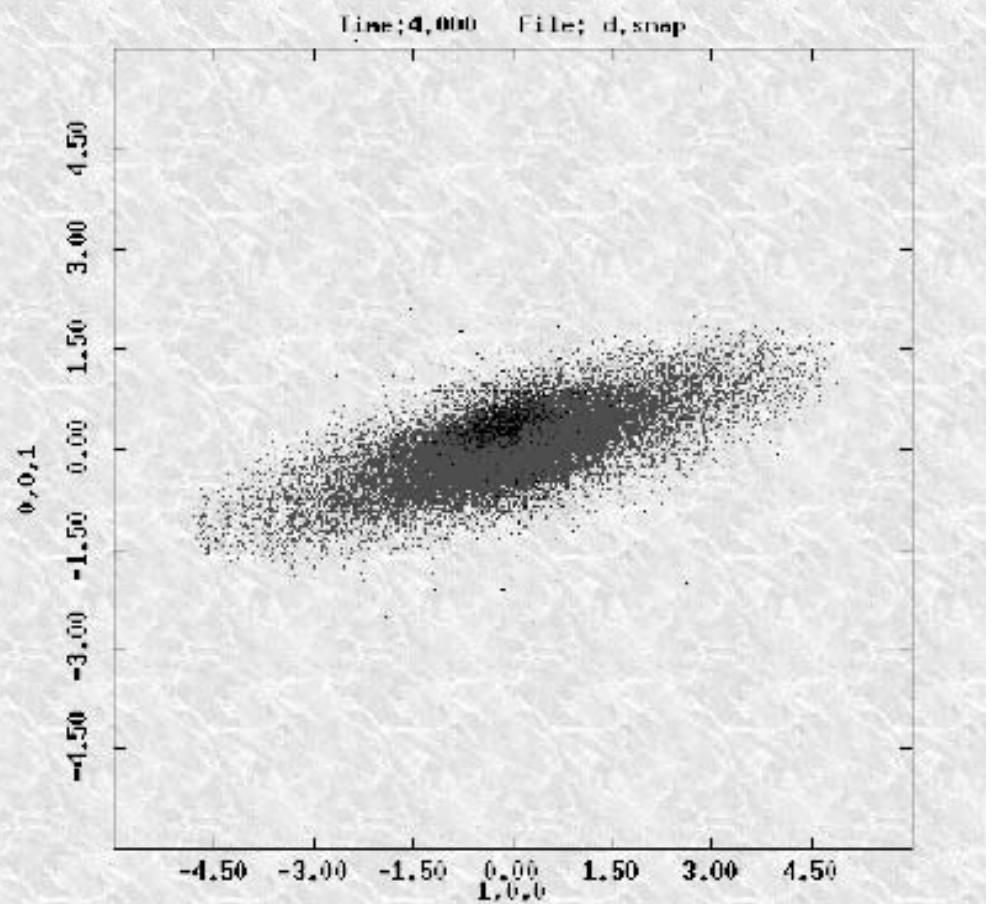


Astronomical Simulation

(Kawai, Fukushige, 2006)

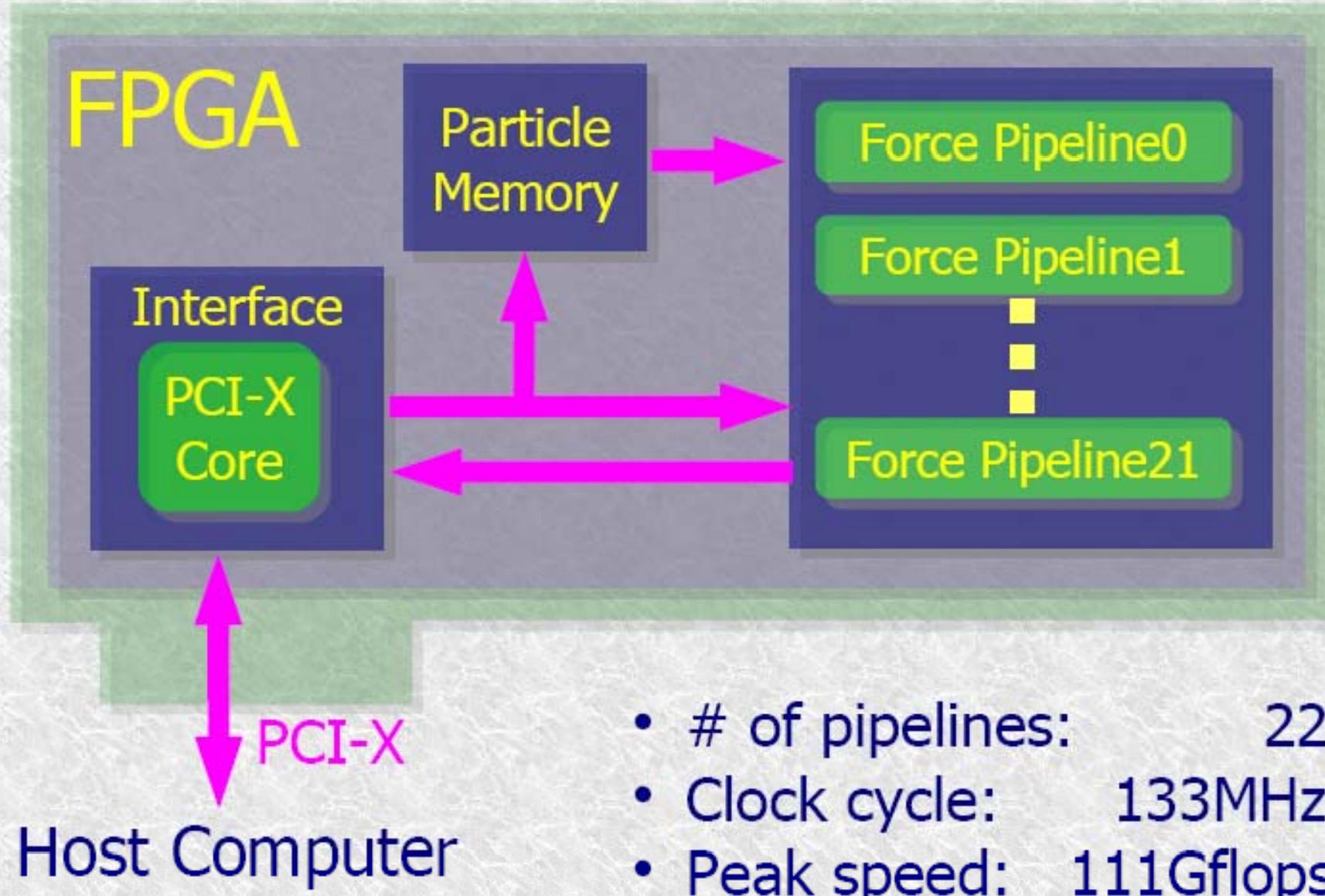


An astronomical object



A representation by particles

Add-in Card Configuration



(Kawai, Fukushige, 2006)

Grape-7 board (Fukushige et al.)

(Kawai, Fukushige, 2006)



Price/Performance

Price (JPY) (Kawai, Fukushige, 2006)

Add-in Card: 187,619

Host PC: 91,238

Total: 278,857 → \$2,363

\$1=118JPY

Corrected Performance: 22.59 Gflops

Price/Performance: \$105/Gflops

(improved from \$158/Gflops marked in July)



Previous record: \$246/Gflops (Kim et al. 2001)

SIMD v.s. FPGA

- High-efficiency of FPGA based simulation
 - Small precision calculation
 - Bit manipulation
 - Use of general purpose FPGA is essential to reduce cost.
 - Current version 65nm Technology
 - 8 Gflops/W peak
 - 105,000\$/Tflops
 - 30 Gflops/W peak will be feasible at 2017
 - 1000 \$/ Tflops will be a tough target

SIMD v.s. FPGA

- GRAPE-DRx based simulation
 - Double precision calculation
 - Programming is much easier than FPGA based approach
 - HDL based programming v.s. C programming
 - Use of latest semiconductor technology is the key
 - Current version 90nm Technology
 - 4 Gflops/W peak
 - 10,000\$/Tflops
 - 30 Gflops/W peak will be feasible at 2017
 - 1000 \$/ Tflops will be feasible at 2017

Systems in 2017

- Eflops system
 - Necessary condition
 - Cost ~1000\$/Tflops
 - Power consumption (30Gflops/W -> 30MW)
 - GRAPE-DR technology
 - 45nm CMOS -> 4Tflops/chip, 40Gflops/W
 - 22nm CMOS -> 16Tflops/chip, 160Gflops/W
 - GPGPU Low double precision performance
 - Power consumption, Die size
 - ClearSpeed
 - Low performance (25Gflops/chip)

To FPGA community

- Key issues to achieve next generation HPC
 - Low cost. Special purpose FPGA will be useless to survive in HPC
 - Low power. Low power is the most important feature of FPGA based simulation
 - High-speed. 500 Gflops/chip? (short precision)
- Programming environment
 - More important to HPC users (not skilled in hardware design)
 - Programming using conventional programming languages
 - Good compiler, run-time and debugger
 - New algorithm that utilize short-precision arithmetic is necessary