HPC-GPU: Large-Scale GPU Accelerated Windows HPC Clusters and its Application to Advanced Bioinformatics and Structural Proteomics (and Climate/Environment)

> Satoshi Matsuoka, Professor/Dr.Sci. & Yutaka Akiyama, Professor with Toshio Endo, Fumikazu Konishi, Akira Nukada, Naoya Maruyama...

> > Tokyo Inst. Technology



### The TSUBAME 1.0 "Supercomputing Grid Cluster" April 2006 at Tokyo Tech 80 TFlops, 1400 Users, 200 Industry Users







## Biggest Problem is Power...

Machine	CPU Cores	Watts	Peak GFLOPS	Peak MFLOPS/ Watt	Watts/ CPU Core	Ratio c.f. TSUBAME
TSUBAME(Opteron)	10480	800,000	50,400	63.00	76.34	
TSUBAME2006 (w/360CSs)	11,200	810,000	79,430	98.06	72.32	
TSUBAME2007 (w/648CSs)	11,776	820,000	102,200	124.63	69.63	1.00
Earth Simulator	5120	6,000,000	40,000	6.67	1171.88	0.05
ASCI Purple (LLNL)	12240	6,000,000	77,824	12.97	490.20	0.10
AIST Supercluster (Opteron)	3188	522,240	14400	27.57	163.81	0.22
LLNL BG/L (rack)	2048	25,000	5734.4	229.38	12.21	1.84
Next Gen BG/P (rack)	4096	30,000	16384	546.13	7.32	4.38
TSUBAME 2.0 (2010Q3/4)	160,000	810,000	1,024,000	1264.20	5.06	10.14

TSUBAME 2.0 x24 improvement in 4.5 years...?  $\rightarrow$  ~ x1000 over 10 years <sup>5</sup>

## NVIDIA CUDA Architecture



240 SPs / Chip (Tesla G200 GPU) Great Relative Power/Performance Possibilities for Scientific Codes GPUs (Tesla, FireStream, Larrabee, ClearSpeed) as Extreme Many-Core and Solution to This Problem



65~55nm(2008) => 15 nm (2016) x20 transitors (30 bil) 5000 Cores 20TF FMA SFP 10TF FMA DFP

nVidia Tesla T10: 65nm, 600m2, 1.4 bil Tr "Massive FMA FPUs" 1.08 TF SFP "Powerful Scalar" 240 Cores 1.5Ghz 1.5Ghz 1.08 TF lops SFP 90 GF lops DFP 102 GBytes/s Tesla Accelerator



w/GPUs





10Gbpc

\$70~80 mil

MSRP

**Fastest SC** 

in Japan

- Whole Cluster Specs (~50 nodes)
  - 50-100 Teraflops
  - 20K-30KW Power
  - Massive FFT Engine
  - \$100,000-\$200,000 MSRP

### **Protein-Protein Interaction**



Protein docking is a central issue both in computational biology and current drug design technique in pharmaceutical industry.

It is becoming a new trend to design an "inhibitor" drug compound which controls specified protein-protein interaction as a target.

#### **Protein Docking is important for:**

- 1) Enzymatic reaction
- 2) Signal Transduction
- 3) Formulating protein scaffold, etc.

![](_page_10_Picture_8.jpeg)

#### All-to-all 3-D Protein Docking Challenge (by Y. Akiyama, in collaboration w/AIST CBRC and other pharma companies)

![](_page_11_Figure_1.jpeg)

1,000 x1,000 all-to-all docking fitness evaluation will take only

1-2 months (15 deg. pitch) with a 32-node HPC-GPGPU cluster (128 GPGPUs).

#### cf.

~ 500 years with single CPU (sus. 1+GF)

> 2 years with 1-rack BlueGene/L

![](_page_11_Picture_7.jpeg)

![](_page_11_Picture_8.jpeg)

## **Rigid Docking vs. Flexible Docking**

![](_page_12_Figure_1.jpeg)

Porcine Pancreatic Trypsin (PDB:1AVX).

**Rigid Docking** 

- Protein is regarded as a rigid body
- Shape complementarity and only simple physicochemical potentials.

<u>No biologist</u> think proteins are rigid. However,

- Flexible docking is prohibitively expensive due to energy local minima.
- Flexible docking usually needs good initial docking structure.
- Some proteins are almost rigid.

Single Rigid docking has only limited validity. But all-to-all Rigid Docking screening can be a good basis for quick survey of potential pairs, and for flexible docking study.

### **Calculation flow and Complexity**

![](_page_13_Figure_1.jpeg)

Calculation for a single protein-protein pair: ~= 200 Tera ops. 3-D complex convolution  $O(N^3 \log N)$ , typically N = 256x Possible rotations R = 54,000 (6 deg. pitch) ~= 200 Exa (10<sup>20</sup>) ops

### **Rigid Docking Example**

- Implemented 6 clustering methods for Post-Docking analysis
- **Developed a confidence level evaluation procedure**

![](_page_14_Picture_3.jpeg)

prediction

**Receptor protein (Trypsin) is shown in center** as a ribbon diagram. Top 2,000 candidate ligand docking sites are shown by small dots. Final prediction site is shown by a red sphere.

![](_page_14_Figure_5.jpeg)

If significant candidate is found, the system outputs predicted docking structure with Trypsin and inhibitor (PDB:1AVX)<sup>a</sup> confidence level, otherwise reports as "No docking". Other clusters

### Bandwidth Intensive 3-D FFT on NVIDIA CUDA GPUs [SC08]

- By Akira Nukada, Tokyo Tech.
- Our 3-D FFT algorithm consists of the following two algorithms
- to maximize the memory bandwidth.
- (1) optimized 1-D FFTs for dimension X,
- (2) *multi-row FFT* for dimension Y & Z.

The multi-row FFT computes multiple 1-D FFTs simultaneously. Adapted from vector algorithms, assuming high memory bandwidth.

![](_page_15_Figure_7.jpeg)

This algorithm accesses multiple streams, but each of them is successive.

Since each thread compute independent set of small FFT, thousands of registers are required Solution: for 256-point FFT, use two- pass 16-point FFT kernels.

![](_page_16_Figure_0.jpeg)

Note: An earlier sample with 1.3GHz is used for Tesla S1070.

![](_page_17_Figure_0.jpeg)

## Performance in Double Precision

![](_page_18_Figure_1.jpeg)

The bottleneck is floating-point operation in double precision.

Both GPUs are running at 1.3GHz, but product version of S1070 will come with higher clock improve the performance.

Performance is competitive with that of a single-node vector supercomputer NEC SX-6 (64 GFLOPS peak).

3-D FFT of size 256<sup>3</sup>

![](_page_19_Figure_0.jpeg)

## Performance including Data

![](_page_20_Figure_1.jpeg)

only support PCI-e 1.1

## Eliminating the Bandwidth Bottleneck---Entire docking in GPU

![](_page_21_Figure_1.jpeg)

## Raccoon: Acclerated WinHPC GPU Prototype Cluster

- 32 compute nodes
- 128 8800GTS GPGPUs
- one head node.
- Gigabit Ethernet network
- Three 40U rack cabinets.
- Linux or Windows HPC 2008
- Visual Studio 2005 SP1
- nVidia CUDA 2.x
- New nodes for GTX280, Infiniband, ...

![](_page_22_Picture_10.jpeg)

### Performance Estimation for 3D PPD Single Node

	Power (W)	Peak (GFLOPS)	3D-FFT (GFLOPS)	Docking (GFLOPS)	Nodes per 40 U rack
Blue Gene/L	20	5.6	-	1.8	1024
TSUBAME	1000 (est.)	76.8 (DP)	18.8 (DP)	26.7 (DP)	10
8800 GTS *4	570	1664	256	207	10~13

#### System Total ! Only CPUs for TSUBAME. DP=double precision.

	# of nodes	Power (kW)	Peak (TFLOPS)	Docking (TFLOPS)	MFLOPS/W
Blue Gene/L (Blue Protein @ AIST)	4096 (4racks)	80	22.9	7.0	87.5
TSUBAME	655 (~70 racks)	~700	50.3 (DP)	17.5 (DP)	25
8800 GTS	32 (3racks)	18	53.2	6.5	361 (x4 B <i>G</i> !)

Can compute 1000x1000 in 1 month (15 deg.) or 1 year (6 deg.)<sup>24</sup>

# DNA Giga Sequencing Using GPUs Personalized DNA-based Diagnostics and

base pairs

- Medicine => Giga Sequencers
- Proposal: GPUs for giga-size short fragment DNA sequencing, for 50,000,000 fragments
  - Rigourous Smith-Watherman gapped alighnmenter-phosphate
  - Months on a conventional cluster, just 17 hours on our test cluster

Type of GPU	<i>#</i> 9f	Pruning	# of Probes	total CPU (hours)	yclopædia Britann
8800GTS	128	off	10,000,000	22.6	
Tesla10	128	off	10,000,000	12.2	
8800GTS	128	on	50,000,000	16.7	
Tesla 10	128	on	50,000,000	8.9	5 

## CFD, esp. Climate and Disaster Prevention on GPUs

(Material from Prof. Takayuki Aoki, Tokyo Tech.) Safety Nuclear (Cooling)

![](_page_25_Picture_2.jpeg)

#### Weather/Environmental

![](_page_25_Picture_4.jpeg)

![](_page_25_Picture_5.jpeg)

#### Civil Engineering

![](_page_25_Picture_7.jpeg)

## **Rayleigh-Taylor Instability**

Heavy fluid lays on light fluid and unstable. Euler equation:

$$\frac{\partial Q}{\partial t} + \frac{\partial E}{\partial x} + \frac{\partial F}{\partial y} = 0$$

$$\boldsymbol{Q} = \begin{bmatrix} \boldsymbol{\rho} \\ \boldsymbol{\rho} \boldsymbol{u} \\ \boldsymbol{\rho} \boldsymbol{v} \\ \boldsymbol{e} \end{bmatrix} \boldsymbol{E} = \begin{bmatrix} \boldsymbol{\rho} \boldsymbol{u} \\ \boldsymbol{\rho} \boldsymbol{u}^{2} + \boldsymbol{p} \\ \boldsymbol{\rho} \boldsymbol{u} \boldsymbol{v} \\ \boldsymbol{\rho} \boldsymbol{u} \boldsymbol{v} \\ \boldsymbol{e} \boldsymbol{u} + \boldsymbol{p} \boldsymbol{u} \end{bmatrix} \boldsymbol{F} = \begin{bmatrix} \boldsymbol{\rho} \boldsymbol{v} \\ \boldsymbol{\rho} \boldsymbol{u} \boldsymbol{v} \\ \boldsymbol{\rho} \boldsymbol{v}^{2} + \boldsymbol{p} \\ \boldsymbol{e} \boldsymbol{v} + \boldsymbol{p} \boldsymbol{v} \end{bmatrix}$$

## 42 GFLOPS using GTX280

![](_page_26_Figure_5.jpeg)

## Two-Stream Instability in Plasma Physics

**Vlasov-Poisson Equation:** 

![](_page_27_Picture_2.jpeg)

$$\frac{\partial f}{\partial t} + v \frac{\partial f}{\partial x} - \frac{eE}{m_e} \frac{\partial f}{\partial v} = 0 \quad \frac{\partial^2 \phi}{\partial x^2} = \frac{e(n_e - n_i)}{\varepsilon_0}$$
$$\left(E = -\frac{\partial \phi}{\partial x}, \quad n_e = \int f dv\right)$$

f : electron distribution function

*n* : electron number density

![](_page_28_Figure_0.jpeg)

## **Real-time Tsunami Simulation**

Collaboration with ADPC (Asian Disaster Preparedness Center) and Japan Meteorological Agency

### **Early Warning System:**

![](_page_29_Picture_3.jpeg)

![](_page_29_Figure_4.jpeg)

![](_page_29_Figure_5.jpeg)

**Shallow-Water Eq.** 

**Conservative Form:** 

Assuming hydrostatic balance in the vertical direction,

3D → 2D equation

$$\frac{\partial h}{\partial t} + \frac{\partial hu}{\partial x} + \frac{\partial hv}{\partial y} = 0$$

$$\frac{\partial hu}{\partial t} + \frac{\partial}{\partial x} \left( hu^2 + \frac{1}{2}gh^2 \right) + \frac{\partial huv}{\partial y} = -gh\frac{\partial z}{\partial x}$$
$$\frac{\partial hv}{\partial t} + \frac{\partial huv}{\partial x} + \frac{\partial}{\partial y} \left( hv^2 + \frac{1}{2}gh^2 \right) = -gh\frac{\partial z}{\partial y}$$

## Numerical Methods of Tsunami Simulation

- 2-dimensional Problem : Directional-Splitting Fractional Method
- Point Value Comp. : Characteristic-based Method using Multi-moment Interpolation
- Integral Value Comp. : Conservative Semi-Lagrangian CIP + IDO

Run-up to dry area: thin water layer and artificial viscosities

![](_page_31_Picture_0.jpeg)

## **GPU** Performance

Speed Comparison

#### x-direction : y-direction = 10 : 7

**Current Speed-up** 

### **GPU** : **CPU** = **62** : **1**

GPU – GeForce GTX280 (sp = 240, clock 1.3Ghz)

CPU – Xeon 2.4GHz 6MB Cache Memory

## **Multi-GPU Estimation**

### 3000km x 3000km (500m mesh)

![](_page_33_Figure_2.jpeg)

Colioris Force Tidal Potential Wind effect

![](_page_33_Figure_4.jpeg)

covering Indian Ocean

![](_page_34_Picture_0.jpeg)

### 680 Unit Tesla Installation... While TSUBAME in Production Service (!)

![](_page_35_Picture_1.jpeg)

Sm

## TSUBAME 1.2. The most Heterogeneous Supercomputer in the world

 Three node configurations with four different processors → >30,000 cores, ~170TFlops system

![](_page_36_Picture_2.jpeg)

#### SunFire X4600+ 2 TESLAs + ClearSpeed

- Opteron 2.4GHz 16 cores
- TESLA S1070 (30cores) 2boards
- ClearSpeed X620 (2cores) 1board
- $\rightarrow$  78 cores, 330 Gflops peak

![](_page_36_Picture_8.jpeg)

![](_page_36_Picture_9.jpeg)

#### SunFire X4600+ClearSpeed

- Opteron 2.4GHz 16 cores
- ClearSpeed X620 (2cores) 1board
- → 18 cores, 157 Gflops peak

#### SunBlade X6250 (TSUBASA cluster)

- Xeon 2.83GHz 8 cores
- $\rightarrow$  8 cores, 90.7 Gflops peak

x 318nodes

x 330nodes

x 90nodes

## TSUBAME in Top500 Ranking w/our Hetero-Linpack

	Jun 06	Nov 06	Jun 07	Nov 07	Jun 08	Nov 08
Speed(TF )	38.18	47.38	48.88 [HPDC 2008]	56.43	67.70	77.48
Rank	7	9	14	16	24	29
		CS	Opteron x 360 ≻		S x 648 -	

Six consecutive improvements (world's first)

 The 2<sup>nd</sup> fastest heterogeneous supercomputer in the world (No.1 is RoadRunner)

 Achieved through extremely intricate heteroalgorithmic advances and tuning

Xeon<sup>I</sup>

Tesla

## Portfolio of our GPU Computing Base Technologies for HPC & eScience

- TSUBAME 1.2 (680 Teslas) & 2.0
- Kernels(FFT, Dense/Sparse Matrix)
- Parallel Algorithms (Large FFT, LINPACK, CG)
- Task & Resource Mgmt (Heterogeneity, Scheduling, BQ Scheduling, etc.)
- Fault Tolerance (ECC, redundant computation GPU checkpointing)
- Languages(OpenMP on GPU, Accelerator, MP)
- GPU Low Power computing (power modeling, measurement, optimization)

GPU Leadership from research to deployment(!)

Software-Based ECC for GPUs (N. Maruyama) Possible Collab. w/MSR Vivian Sewelsen and HPC Cluster

![](_page_39_Figure_1.jpeg)

## Fast Conjugate Gradient Solver on Multi-GPU Systems (Ali Cevahir)

Iterative CG on multi-CUDA <u>GPUs w/mixed precision</u> Level 1 BLAS => OK for GPUs Sparse Matrix\*Vector => incurs random access NG for GPUs

Previous Methods CRS Method (Cache-based CPUs) Small # of memory accesses JDS Method (for vector CPUs) Continueous memory accesses

<u>GPU-friendly Implementation</u> 1) JDS-like Storage Format •Efficient memory access Adjust alignment for coalescing 2) CRS-like computing order Reduce # of memory accesses Utilize various memory on CUDA GPUs 1) Matrix element and index (accessed only once) => global memory 2) JDS method offset table (accessed several times) => constant memory (cached) 3) Vector element (accessed many times) => texture memory (cached)

![](_page_40_Figure_6.jpeg)

14.5GFlops 4 GPUs (nVidia 8800 GTS) vs. 0.54GFlops 4 Core CPU (Phenom 2.5Ghz) (Sparse Matrix collection from UFlorida)

![](_page_41_Figure_0.jpeg)

## Dynamic Scheduling of Matrix Kernels on Heterogeneous GPUs (Y. Watanabe)

Rapid GPU Performance Progress and Updates → New GPU add-ons will differ in performance Dynamic Load distribution to heterogeneous GPUs

Thesis: Dynamic task allocation works effectively for load distribution

<u>Characteristics</u> Good: No a-priori info necessary Bad: GPU control API Overhead Good: natural distribution of Host-GPU transfers, less collision

![](_page_42_Figure_4.jpeg)

![](_page_42_Figure_5.jpeg)

### 94% efficiency c.f. summation of single GPU performances

![](_page_42_Figure_7.jpeg)

#### Better performance c.f. static allocation

## GPU Computing: Power modeling (R. Suda@U-Tokyo)

- Problem: To our knowledge there are no tools to predict the power consumption of large-scale parallel program in the algorithm design stage.
- Approach: BSP model provides parallel framework; BM Model incorporates physical constraints. Extending them with powerrelevant characters, and , power model is built, and performance is examined.

![](_page_43_Figure_3.jpeg)

Intel Core 2 Extreme Quad-core Processor

Level	<i>p</i> <sub>t</sub>	<i>e</i> <sup><i>p</i></sup> (W/I)	g <sub>i</sub> (GT/s)	e <sub>i</sub> <sup>m</sup> (W/word)	т <sub>і</sub> кв
1	4	4.1E-8	300	1.5E-5	64
2	Threads	2.8E-8	300	2.1E-5	256
3	Threads	2.8E-8	75	3.0E-5	12K

Results: The power measurement results of *m*atrix computation agrees with the estimated result. Further refinement to the model is needed.

Power Consumption Estimation vs. Measurement

![](_page_43_Figure_8.jpeg)

### Incompressible Fluid Application Power Measurements thru Sleep Insertions (T. Aoki)

45

![](_page_44_Figure_1.jpeg)

}

### Power Efficiency in 3-D FFD (A. Nukada)

GPU	Computation	Idle	Power	GFLOPS	GFLOPS/W
RIVA128	On CPU	126 W	140 W	10.3	0.074
8800 GT	On GPU	180 W	215 W	62.2	0.289
8800 GTS	On GPU	196 W	238 W	67.2	0.282
8800 GTX	On GPU	224 W	290 W	84.4	0.291

CUDA GPUs have four times higher power efficiency than CPU in high-performance FFT.

#### me

RIVA128 is an old, low-power GPU, to measure pure power consumption of host system (CPU, chipset, memory). The interface is legacy PCI.

![](_page_45_Picture_5.jpeg)

## Onto TSUBAME 2.0 & 2.1

- TSUBAME 2.0 will be
  - Deployed 1H 2010
  - A Petascale machine w/GPU acceleration
  - Will be VERY GREEN (same energy envelope)
  - Will be a Cloudy supercomputer, virtualization/multi-OS dynamic deployment
  - Will be a Supercomputer for Everyone
- TSUBAME will run Windows HPC (!)
  - Users will have a choice along w/Linux
  - The first Petascale Windows HPC Cluster?

#### Road to Exascale and Personal Petascale 2012 10PF (Japanese NLP SC etc.) 2019 1ExaFlops (Leadership machines) TSUBAME 4.0 > 100PF Desktop ~= 1PF 1ExF FSUBANE Japanese NLp 10pf (2012) 100PF SUBAME RoadRunner **10PF** 1.5TF (2008) .3.0 IS BAME2.0 BlueGener 1PF 360TF(2005 <sup>Yersonal</sup> 0 Peta 100TF-SUPAME 1.0 85TF 2016 (10006) TSUBAME Earth Simulator 40TF (2002) b<sub>ecomes</sub> **10TF** deskside Titech GPGPU-Nor **KEK 59TF** Campus 400TE 1.3TF BG/L+SR11100 1TF Grid 2014 2016 2018 2020 2006 2008 2012 2002 2004 2010