


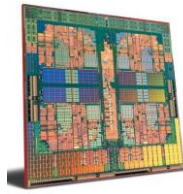

Chapter 6

Parallel Processors from Client to Cloud



Multiprocessor vs. Multicore

- Multiprocessor
 - Several processors connected by an interconnection network
 - Shared memory system
 - Uniform vs. non-uniform memory access
 - Message passing system
 - Processors communicate via send and receive
- Multicore
 - Multiple processor cores on a chip using SoC technology
 - They communicate via on-chip network
 - Homogeneous
 - » IBM Power 5, Intel, Sun T1/T2
 - Heterogeneous
 - » IBM Cell



2016/4/21

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU

2

§6.1 Introduction

Introduction

- Goal: connecting multiple computers to get higher performance
 - Multiprocessors
 - Scalability, availability, power efficiency
- Task-level (process-level) parallelism
 - High throughput for **independent** jobs
- Parallel processing program
 - **Single** program run on multiple processors
- Multicore microprocessors
 - Chips with multiple processors (cores)

632

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU
2016/4/21

3


Hardware and Software

- Hardware
 - Serial: e.g., Pentium 4
 - Parallel: e.g., quad-core Xeon e5345
- Software
 - Sequential: e.g., MatLab matrix multiplication
 - Concurrent: e.g., operating system
- Sequential/concurrent software can run on serial/parallel hardware
 - Challenge: making effective use of parallel hardware

633

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU
2016/4/21

4



What We've Already Covered


- §2.11: Parallelism and Instructions
 - Synchronization (ll 'load linked' and sc 'store conditional')
- §3.6: Parallelism and Computer Arithmetic
 - Associativity
- §4.10: Parallelism and Advanced Instruction-Level Parallelism
- §5.8: Parallelism and Memory Hierarchies
 - Cache Coherence

634

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/21

5



Parallel Programming

- Parallel software is the problem
- Need to get significant performance improvement
 - Otherwise, just use a faster uniprocessor, since it's easier!
- Difficulties
 - Partitioning
 - Coordination
 - Communications overhead


634

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/21

6

§6.2 The Difficulty of Creating Parallel Processing Programs



Amdahl's Law


- Sequential part can limit speedup
- Example: 100 processors, 90× speedup?
 - $T_{\text{new}} = T_{\text{parallelizable}}/100 + T_{\text{sequential}}$
 - $$\text{Speedup} = \frac{1}{(1 - F_{\text{parallelizable}}) + F_{\text{parallelizable}}/100} = 90$$
 - Solving: $F_{\text{parallelizable}} = 0.999$
- Need sequential part to be 0.1% of original time

635

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/21

7



Scaling Example

$a_1 + a_2 + \dots + a_{10}$ $A + B$

- Workload: sum of 10 scalars, and 10 × 10 matrix sum
- Speed up from 10 to 100 processors
 - Single processor: $\text{Time} = (10 + 100) \times t_{\text{add}}$
 - 10 processors
 - $\text{Time} = 10 \times t_{\text{add}} + 100/10 \times t_{\text{add}} = 20 \times t_{\text{add}}$
 - $\text{Speedup} = 110/20 = 5.5$ (55% of potential)
 - 100 processors
 - $\text{Time} = 10 \times t_{\text{add}} + 100/100 \times t_{\text{add}} = 11 \times t_{\text{add}}$
 - $\text{Speedup} = 110/11 = 10$ (10% of potential)
- Assumes load can be balanced across processors

636

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/21

8

636

Scaling Example (cont)

- What if matrix size is 100×100 ? 0.1%
 1. Single processor: $\text{Time} = (10 + 10000) \times t_{\text{add}}$
 2. 10 processors
 - $\text{Time} = 10 \times t_{\text{add}} + 10000/10 \times t_{\text{add}} = 1010 \times t_{\text{add}}$
 - $\text{Speedup} = 10010/1010 = 9.9$ (99% of potential)
 3. 100 processors
 - $\text{Time} = 10 \times t_{\text{add}} + 10000/100 \times t_{\text{add}} = 110 \times t_{\text{add}}$
 - $\text{Speedup} = 10010/110 = 91$ (91% of potential)
- Assuming load balanced

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/21

9

637

Strong vs Weak Scaling

- **Strong** scaling: how the solution time varies with the number of processors for a fixed *total* problem size
 - As in the previous example
- **Weak** scaling: how the solution time varies with the number of processors for a fixed problem size *per processor*.
 - 10 processors, 10×10 matrix sum
 - $\text{Time} = 20 \times t_{\text{add}}$
 - 100 processors, 32×32 matrix sum
 - $\text{Time} = 10 \times t_{\text{add}} + 1000/100 \times t_{\text{add}} = 20 \times t_{\text{add}}$
 - Increases 10x problem size for 10x processors for same performance gain

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/21

10

§6.3 SISD, MIMD, SIMD, SPMD, and Vector

Instruction and Data Streams

- An alternate classification

		Data Streams	
		Single	Multiple
Instruction Streams	Single	SISD: Intel Pentium 4	SIMD: SSE instructions of x86
	Multiple	MISD: No examples today	MIMD: Intel Xeon e5345

- SPMD: Single Program Multiple Data
 - A parallel program on a MIMD computer
 - Conditional code for different processors

648

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU2016/4/2111

Example: DAXPY ($Y = a \times X + Y$)

- Conventional MIPS code

```
l.d    $f0,a($sp)      ;load scalar a
addiu  r4,$s0,#512     ;upper bound of what to load
loop:  l.d    $f2,0($s0) ;load x(i)
      mul.d   $f2,$f2,$f0 ;a x x(i)
      l.d    $f4,0($s1) ;load y(i)
      add.d   $f4,$f4,$f2 ;a x x(i) + y(i)
      s.d    $f4,0($s1) ;store into y(i)
      addiu  $s0,$s0,#8  ;increment index to x
      addiu  $s1,$s1,#8  ;increment index to y
      subu   $t0,r4,$s0  ;compute bound
      bne    $t0,$zero,loop ;check if done
```
- Vector MIPS code

```
l.d    $f0,a($sp)      ;load scalar a
lv     $v1,0($s0)       ;load vector x
mulvs.d $v2,$v1,$f0     ;vector-scalar multiply
lv     $v3,0($s1)       ;load vector y
addv.d $v4,$v2,$v3      ;add y to product
sv     $v4,0($s1)       ;store the result
```

651

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU2016/4/2112

650

Vector Processors

- Highly pipelined function units
- Stream data from/to vector registers to units
 - Data collected from memory into registers
 - Results stored from registers to memory
- Example: Vector extension to MIPS
 - 32×64 -element registers (64-bit elements)
 - Vector instructions
 - `lv, sv`: load/store vector
 - `addv.d`: add vectors of double
 - `addvs.d`: add scalar to each element of vector of double
- Significantly reduces instruction-fetch bandwidth

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU2016/4/2113

652

Vector vs. Scalar

- Vector architectures and compilers
 - Simplify data-parallel programming
 - Explicit statement of absence of loop-carried dependences
 - Reduced checking in hardware
 - Regular access patterns benefit from interleaved and burst memory
 - Avoid control hazards by avoiding loops
- More general than ad-hoc media extensions (such as MMX, SSE)
 - Better match with compiler technology

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU2016/4/2114

SIMD

- Operate elementwise on **vectors of data**
 - E.g., MMX and SSE instructions in x86
 - Multiple data elements in 128-bit wide registers
- All processors execute the same instruction at the same time
 - Each with different data address, etc.
- Simplifies synchronization
- Reduced instruction control hardware
- Works best for highly data-parallel applications

649

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU
2016/4/21

15

Vector vs. Multimedia Extensions

- Vector instructions have a variable vector width, multimedia extensions have a fixed width
- Vector instructions support strided access, multimedia extensions do not
- Vector units can be combination of pipelined and arrayed functional units:

Diagram illustrating elementwise addition of two vectors A and B to produce result C. The vectors are represented as columns of elements A[0] to A[9] and B[0] to B[9]. The result C is shown as a single element C[0] for the first addition, and a group of four elements C[0] to C[3] for the second addition, labeled 'Element group'.

Diagram illustrating a four-lane SIMD architecture. Each lane (Lane 0 to Lane 4) contains a pipeline of FP mul pipe, Vector registers, and FP add pipe. The Vector registers are labeled with their element indices: Lane 0 (0, 4, 8, ...), Lane 1 (1, 5, 9, ...), Lane 2 (2, 6, 10, ...), and Lane 3 (3, 7, 11, ...). A shared Vector load store unit is shown at the bottom.

2016/4/21

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU

16

§6.4 Hardware Multithreading

Multithreading

- Performing multiple threads of execution in parallel (share functional units)
 - Replicate registers, PC, etc. (state)
 - Fast switching between threads
- 1. Fine-grain multithreading
 - Switch threads after each cycle
 - Interleave instruction execution
 - If one thread stalls, others are executed
- 2. Coarse-grain multithreading
 - Only switch on long stall (e.g., L2-cache miss)
 - Simplifies hardware, but doesn't hide short stalls (e.g., data hazards)

645

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU2016/4/21

17

Simultaneous Multithreading (SMT)

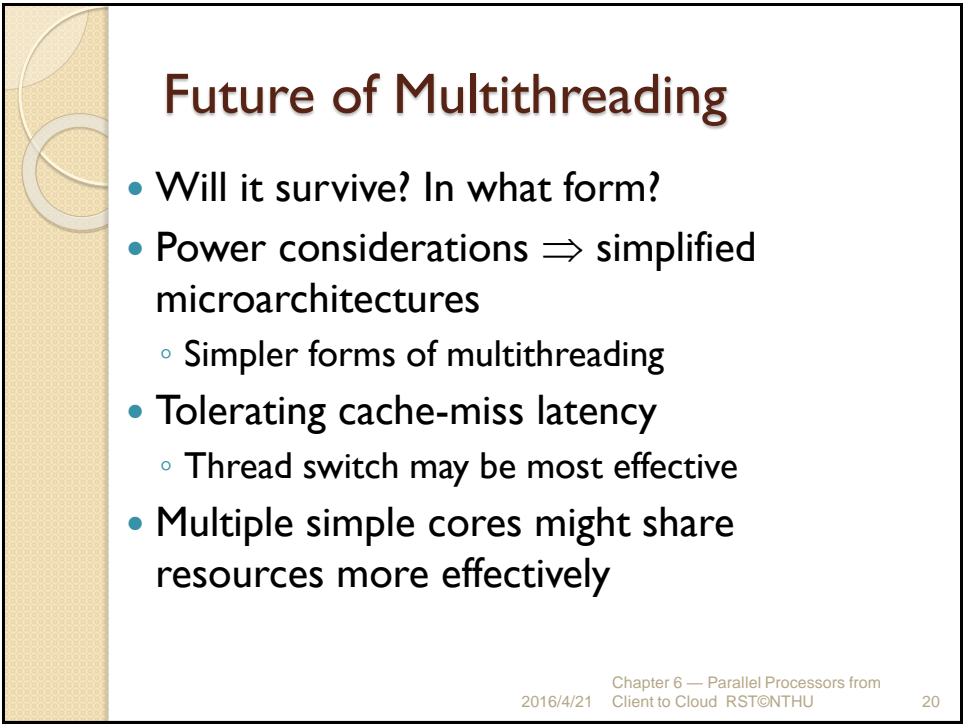
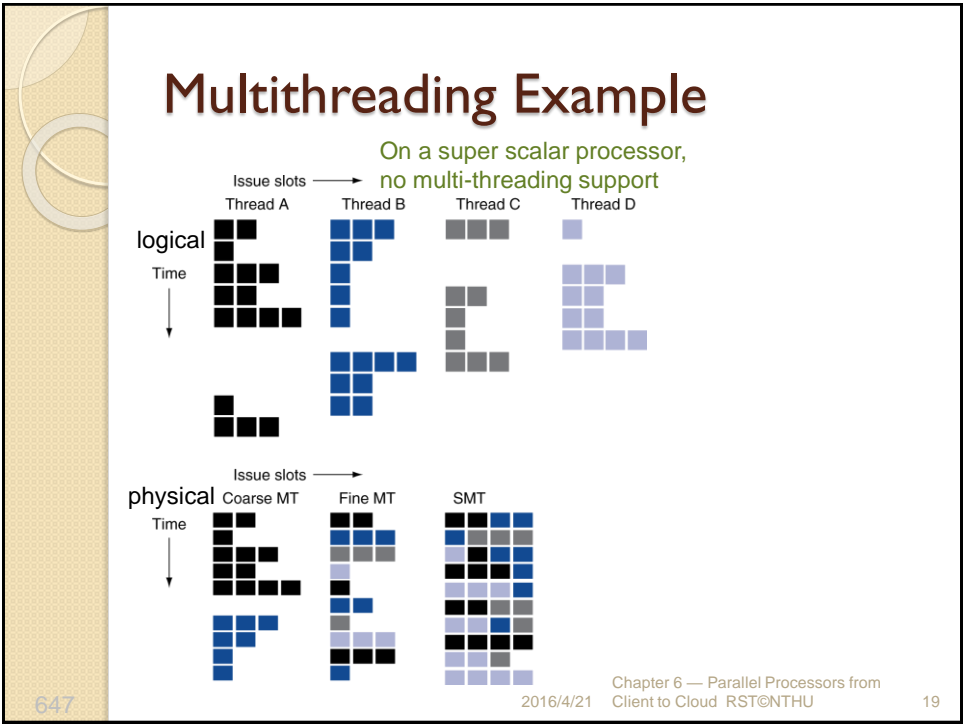
- In multiple-issue dynamically scheduled processor
 - Schedule instructions from multiple threads
 - Instructions from independent threads execute when function units are available
 - Within threads, dependencies handled by scheduling and register renaming
- Example: Intel Pentium-4 HT
 - Two threads: duplicated registers, shared function units and caches

The diagram illustrates Simultaneous Multithreading (SMT). It shows two yellow boxes labeled 'R' (representing registers) at the top. Below them is a pink arrow pointing downwards, labeled 'P' (representing the processor or functional unit). This indicates that multiple threads share the same processor unit while having their own registers.

646

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU2016/4/21

18



§6.5 Shared Memory Multiprocessors

Shared Memory Multiprocessor

- SMP: shared memory multiprocessor
 - Hardware provides single physical address space for all processors
 - Synchronize shared variables using locks
 - Memory access time
 - UMA (uniform) vs. NUMA (nonuniform)

```
graph TD
    P1[Processor] <--> C1[Cache]
    P2[Processor] <--> C2[Cache]
    Pn[Processor] <--> Cn[Cache]
    C1 <--> IN[Interconnection Network]
    C2 <--> IN
    Cn <--> IN
    IN <--> M[Memory]
    IN <--> IO[I/O]
```

638

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU
2016/4/21

21

Example: Sum Reduction

- Sum 100,000 numbers on 100 processor UMA
 - Each processor has ID: $0 \leq P_n \leq 99$
 - Partition 1000 numbers per processor ($=100,000/100$)
 - Initial summation on each processor

```
sum[Pn] = 0;
for ( i = 1000*Pn; i < 1000*(Pn+1); i = i + 1)
    sum[Pn] = sum[Pn] + A[i];
```
- Now need to add these partial sums
 - Reduction: divide and conquer
 - Half the processors add pairs, then quarter, ...
 - Need to synchronize between reduction steps

639

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU
2016/4/21

22

Example: Sum Reduction

```
half = 100;
repeat
  synch();
  if ( half%2 != 0 && Pn == 0)
    sum[0] = sum[0] + sum[half-1];
    /* Conditional sum needed when half is odd;
       Processor_0 gets missing element */
  half = half/2; /* dividing line on who sums */
  if (Pn < half) sum[Pn] = sum[Pn] + sum[Pn+half];
until (half == 1);
for(Pn=0; Pn<half; Pn++)
```

640

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU 23

History of GPUs

- Early video cards
 - Frame buffer memory with address generation for video output
- 3D graphics processing
 - Originally high-end computers (e.g., SGI)
 - Moore's Law \Rightarrow lower cost, higher density
 - 3D graphics cards for PCs and game consoles
- Graphics Processing Units
 - Processors oriented to 3D graphics tasks
 - Vertex/pixel processing, shading, texture mapping, rasterization

654

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU 24

§6.5 Introduction to Graphics Processing Units

Texture Mapping

From Computer Desktop Encyclopedia
Reproduced with permission.
© 2001 Intergraph Computer Systems

2D texture

3D object

2D texture "draped" over 3D object

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU

2016/4/21

25

Rasterization

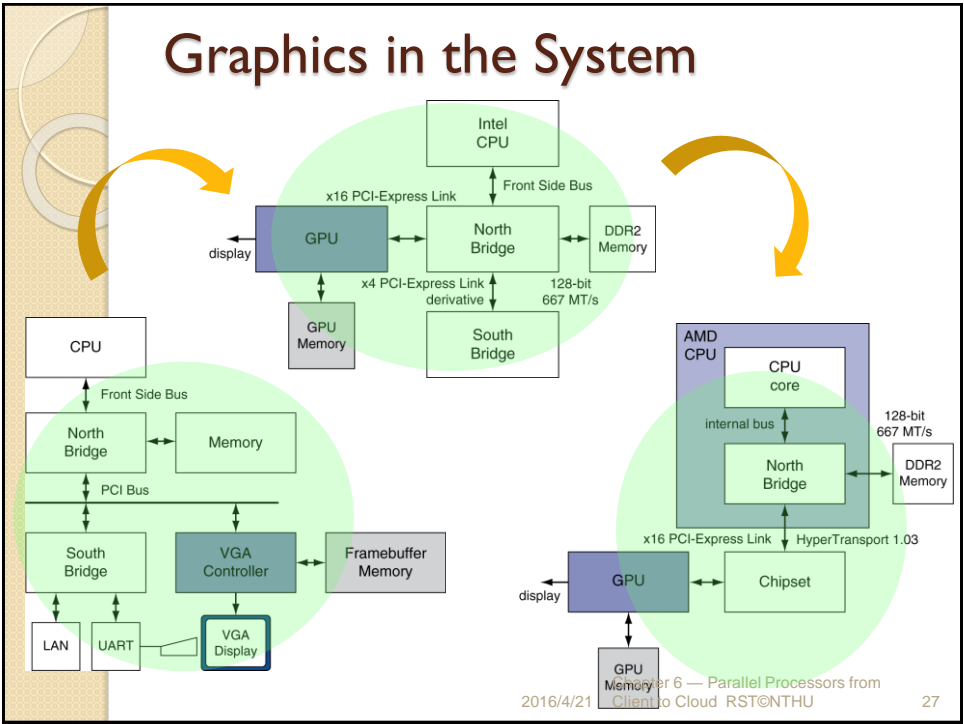
Rasterized

Ray traced

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU

2016/4/21

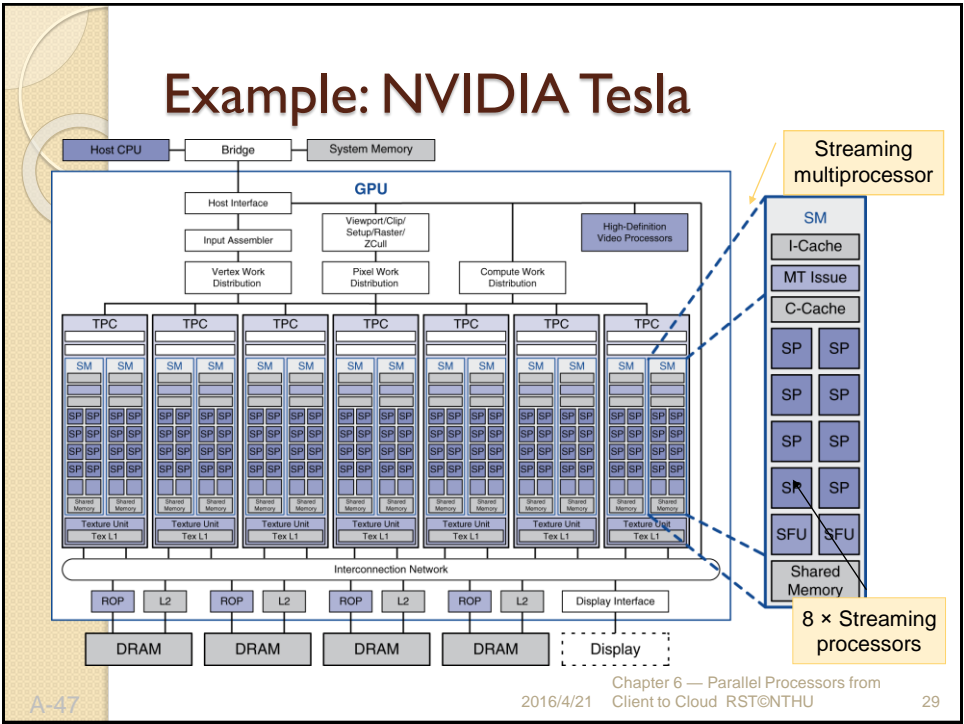
26



GPU Architectures

- Processing is highly data-parallel
 - GPUs are highly **multithreaded**
 - Use thread switching to hide memory latency
 - Less reliance on multi-level caches
 - Graphics memory is wide and high-bandwidth
- Trend toward general purpose GPUs
 - Heterogeneous CPU/GPU systems
 - CPU for sequential code, GPU for parallel code
- Programming languages/APIs
 - DirectX, OpenGL
 - C for Graphics (Cg), High Level Shader Language (HLSL)
 - Compute Unified Device Architecture (CUDA)

2016/4/21 Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU 28



Example: NVIDIA Tesla

- Streaming Processors (SPs)
 - Single-precision FP and integer units
 - Each SP is fine-grained multithreaded
- Warp: group of 32 threads
 - Executed in parallel, SIMD style
 - 8 SPs × 4 clock cycles
 - Hardware contexts for 24 warps
 - Registers, PCs, ...

The diagram shows a comparison between UltraSPARC T2 and Tesla Multiprocessor. On the left, a vertical stack of 8 threads (Thread0 through Thread7) is shown, labeled 'UltraSPARC T2' and 'Hardware Supported Threads'. On the right, a grid of threads is shown, labeled 'Tesla Multiprocessor'. The grid is divided into 24 warps (Warp0 through Warp23), each containing 32 threads. The threads are arranged in a 4x8 grid, with 8 threads per warp and 4 threads per clock cycle.

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU 30

Classifying GPUs

- Don't fit nicely into SIMD/MIMD model
 - Conditional execution in a thread allows an illusion of MIMD
 - But with performance degradation
 - Need to write general purpose code with care

	Static: Discovered at Compile Time	Dynamic: Discovered at Runtime
Instruction-Level Parallelism	VLIW	Superscalar
Data-Level Parallelism	SIMD or Vector	Tesla Multiprocessor

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU

2016/4/2131

GPU Memory Structures

The diagram illustrates the memory hierarchy of a GPU. At the top, a 'CUDA Thread' is shown with a wavy line representing its execution path, connected to a box labeled 'Per-CUDA Thread Private Memory'. Below this, a 'Thread block' is depicted as a group of threads, connected to a box labeled 'Per-Block Local Memory'. The next level shows two 'Grid' boxes, 'Grid 0' and 'Grid 1', each containing multiple thread blocks. These grids are connected to a large box labeled 'GPU Memory'. A dashed line between the grids is labeled 'Inter-Grid Synchronization'. A vertical arrow on the right side of the GPU Memory box is labeled 'Sequence'.

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU

2016/4/2132

Putting GPUs into Perspective

Feature	Multicore with SIMD	GPU
SIMD processors	4 to 8	8 to 16
SIMD lanes/processor	2 to 4	8 to 16
Multithreading hardware support for SIMD threads	2 to 4	16 to 32
Typical ratio of single precision to double-precision performance	2:1	2:1
Largest cache size	8 MB	0.75 MB
Size of memory address	64-bit	64-bit
Size of main memory	8 GB to 256 GB	4 GB to 6 GB
Memory protection at level of page	Yes	Yes
Demand paging	Yes	No
Integrated scalar processor/SIMD processor	Yes	No
Cache coherent	Yes	No

Chapter 6 — Parallel Processors from Client to Cloud

2016/4/21 RST@NTHU

33

Guide to GPU Terms

type	More descriptive name	Closest old term outside of GPUs	Official CUDA/ NVIDIA GPU term	Book definition
Program abstractions	Vectorizable Loop	Vectorizable Loop	Grid	A vectorizable loop, executed on the GPU, made up of one or more Thread Blocks (bodies of vectorized loop) that can execute in parallel.
	Body of Vectorized Loop	Body of a (Strip-Mined) Vectorized Loop	Thread Block	A vectorized loop executed on a multithreaded SIMD Processor, made up of one or more threads of SIMD instructions. They can communicate via Local Memory.
	Sequence of SIMD Lane Operations	One iteration of a Scalar Loop	CUDA Thread	A vertical cut of a thread of SIMD instructions corresponding to one element executed by one SIMD Lane. Result is stored depending on mask and predicate register.
Machine object	A Thread of SIMD Instructions	Thread of Vector Instructions	Warp	A traditional thread, but it contains just SIMD instructions that are executed on a multithreaded SIMD Processor. Results stored depending on a per-element mask.
	SIMD Instruction	Vector Instruction	PTX Instruction	A single SIMD instruction executed across SIMD Lanes.
Processing hardware	Multithreaded SIMD Processor	(Multithreaded) Vector Processor	Streaming Multiprocessor	A multithreaded SIMD Processor executes threads of SIMD instructions, independent of other SIMD Processors.
	Thread Block Scheduler	Scalar Processor	Giga Thread Engine	Assigns multiple Thread Blocks (bodies of vectorized loop) to multithreaded SIMD Processors.
	SIMD Thread Scheduler	Thread scheduler in a Multithreaded CPU	Warp Scheduler	Hardware unit that schedules and issues threads of SIMD instructions when they are ready to execute; includes a scoreboard to track SIMD Thread execution.
	SIMD Lane	Vector lane	Thread Processor	A SIMD Lane executes the operations in a thread of SIMD instructions on a single element. Results stored depending on mask.
Memory/hardware	GPU Memory	Main Memory	Global Memory	DRAM memory accessible by all multithreaded SIMD Processors in a GPU.
	Local Memory	Local Memory	Shared Memory	Fast local SRAM for one multithreaded SIMD Processor, unavailable to other SIMD Processors.
	SIMD Lane Registers	Vector Lane Registers	Thread Processor Registers	Registers in a single SIMD Lane allocated across a full thread block (body of vectorized loop).

Chapter 6 — Parallel Processors from Client to Cloud

2016/4/21 RST@NTHU

34

§6.7 Clusters and Other Message-Passing Multiprocessors

Message Passing

- Each processor has private physical address space
- Hardware sends/receives messages between processors

The diagram illustrates a Message Passing architecture. It shows three identical processing units arranged horizontally, separated by ellipses. Each unit consists of a 'Processor' box at the top, connected by a double-headed arrow to a 'Cache' box below it, which is in turn connected by a double-headed arrow to a 'Memory' box at the bottom. The 'Memory' boxes are highlighted in pink. All 'Memory' boxes are connected to a single horizontal 'Interconnection Network' box at the bottom via double-headed arrows.

641

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU
2016/4/21

35

Loosely Coupled Clusters

- Network of independent computers
 - Each has private memory and OS
 - Connected using I/O system
 - E.g., Ethernet/switch, Internet
- Suitable for applications with independent tasks
 - Web servers, databases, simulations, ...
- High availability, scalable, affordable
- Problems
 - Administration cost (prefer virtual machines)
 - Low interconnect bandwidth
 - c.f. processor/memory bandwidth on an SMP

641

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU
2016/4/21

36

642

Sum Reduction (Again)

- Sum 100,000 on 100 processors
- First distribute 1000 numbers to each
 - The do partial sums
sum = 0;
for (i = 0; i<1000; i = i + 1)
sum = sum + AN[i];
- Reduction
 - Half the processors send, other half receive and add
 - The quarter send, quarter receive and add, ...

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/2137

643

Sum Reduction (Again)

- Given send() and receive() operations

```
limit = 100; half = 100; /* 100 processors */
repeat
  half = (half+1)/2; /* send vs. receive dividing line */
  if (Pn >= half && Pn < limit)
    send(Pn - half, sum);
  if (Pn < (limit/2))
    sum = sum + receive();
  limit = half; /* upper limit of senders */
until (half == 1); /* exit with final sum */
```
- Send/receive also provide synchronization
- Assumes send/receive take similar time to addition

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/2138

Grid Computing

- Separate computers interconnected by long-haul networks
 - E.g., Internet connections
 - Work units farmed out, results sent back
- Can make use of idle time on PCs
 - E.g., SETI@home, World Community Grid

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/2139

Interconnection Networks

- Network topologies
 - Arrangements of processors, switches, and links

Bus

Ring

2D Mesh

N-cube (N = 3)

Fully connected

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/2140

§6.8 Introduction to Multiprocessor Network Topologies

Multistage Networks

The diagrams illustrate different multistage network topologies. Diagram (a) shows a crossbar network with 8 processors (P₀ to P₇) on the left and 8 output lines on the right, with a grid of switches. Diagram (b) shows an Omega network with 8 processors (P₀ to P₇) on the left and 8 output lines on the right, with a butterfly-like internal structure. Diagram (c) shows a single Omega network switch box with four inputs (A, B, C, D) and four outputs, with a specific internal routing path.

a. Crossbar b. Omega network

c. Omega network switch box

663

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU
2016/4/21

41

Network Characteristics

- Performance
 - Latency per message (unloaded network)
 - Throughput
 - Link bandwidth
 - Total network bandwidth
 - Bisection bandwidth
 - Congestion delays (depending on traffic)
- Cost
- Power
- Routability in silicon

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU
2016/4/21

42

Parallel Benchmarks

- Linpack: matrix linear algebra
- SPECrate: parallel run of SPEC CPU programs
 - Job-level parallelism
- SPLASH: Stanford Parallel Applications for Shared Memory
 - Mix of kernels and applications, strong scaling
- NAS (NASA Advanced Supercomputing) suite
 - computational fluid dynamics kernels
- PARSEC (Princeton Application Repository for Shared Memory Computers) suite
 - Multithreaded applications using Pthreads and OpenMP

664

Chapter 6 — Parallel Processors from
2016/4/21 Client to Cloud RST@NTHU

43

Code or Applications?

- Traditional benchmarks
 - Fixed code and data sets
- Parallel programming is evolving
 - Should algorithms, programming languages, and tools be part of the system?
 - Compare systems, provided they implement a given application
 - E.g., Linpack, Berkeley Design Patterns
- Would foster innovation in approaches to parallelism

Chapter 6 — Parallel Processors from
2016/4/21 Client to Cloud RST@NTHU

44

§6.10 Roofline: A Simple Performance Model

Modeling Performance

- Assume performance metric of interest is achievable GFLOPs/sec
 - Measured using computational kernels from Berkeley Design Patterns
- Arithmetic intensity of a kernel
 - FLOPs per byte of memory accessed
- For a given computer, determine
 - Peak GFLOPS (from data sheet)
 - Peak memory bytes/sec (using Stream benchmark)

667

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU2016/4/21

45

Arithmetic Intensity

- Arithmetic Intensity: the ratio of floating-point operations per byte of memory accessed.

Complexity labels above the arrow: $O(1)$, $O(\log(N))$, $O(N)$

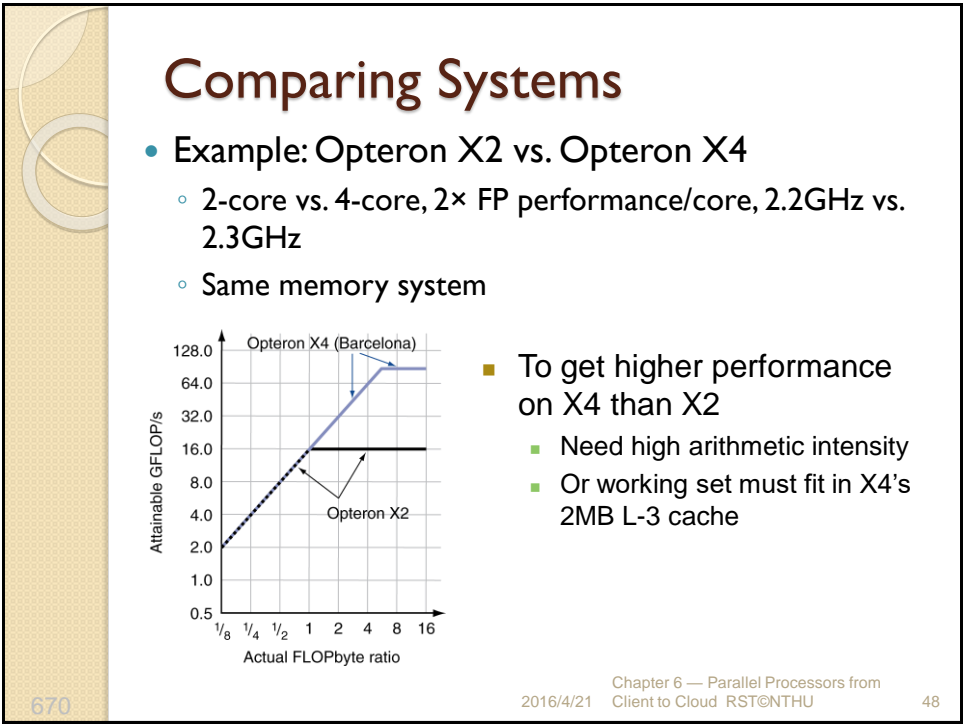
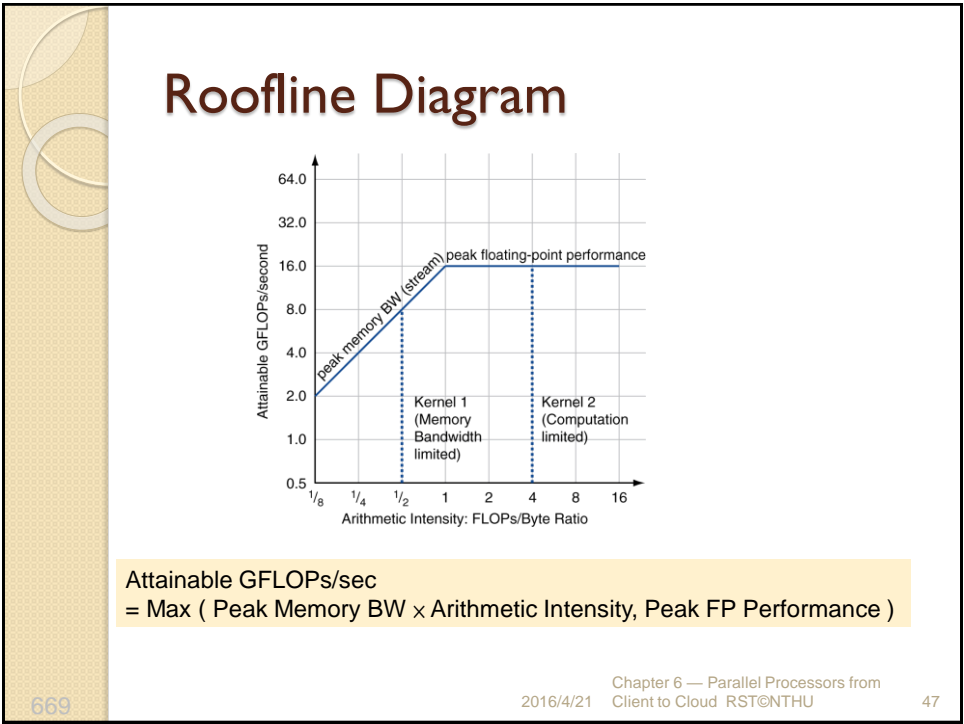
Computational models below the arrow:

- Sparse Matrix (SpMV)
- Structured Grids (Stencils, PDEs)
- Structured Grids (Lattice Method)
- Spectral Methods (FFTs)
- Dense Matrix (BLASS)
- N-Body (Particle Methods)

668

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU2016/4/21

46



Optimizing Performance

- Optimize FP performance
 - Balance adds & multiplies
 - Improve superscalar ILP and use of SIMD instructions
- Optimize memory usage
 - Software prefetch
 - Avoid load stalls
 - Memory affinity
 - Avoid non-local data accesses

AMD Opteron

Attainable GFLOPs/second

Arithmetic Intensity: FLOPs/Byte Ratio

1. Fl. Pt. imbalance

2. Without ILP or SIMD

3. peak memory BW (green)

4. w/out Memory Affinity

peak floating-point performance

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU

2016/4/21

49

Optimizing Performance

- Choice of optimization depends on arithmetic intensity of code

Attainable GFLOPs/second

Arithmetic Intensity: FLOPs/Byte Ratio

1. Fl. Pt. imbalance

2. Without ILP or SIMD

3. peak memory BW (green)

4. w/out Memory Affinity

peak floating-point performance

Kernel 1

Kernel 2

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU

2016/4/21

50

- Arithmetic intensity is not always fixed
 - May scale with problem size
 - Caching reduces memory accesses
 - Increases arithmetic intensity

i7-960 vs. NVIDIA Tesla 280/480

	Core i7-960	GTX 280	GTX 480	Ratio 280/17	Ratio 480/17
Number of processing elements (cores or SMs)	4	30	15	7.5	3.8
Clock frequency (GHz)	3.2	1.3	1.4	0.41	0.44
Die size	263	576	520	2.2	2.0
Technology	Intel 45 nm	TCMS 65 nm	TCMS 40 nm	1.6	1.0
Power (chip, not module)	130	130	167	1.0	1.3
Transistors	700 M	1400 M	3100 M	2.0	4.4
Memory bandwidth (GBytes/sec)	32	141	177	4.4	5.5
Single precision SIMD width	4	8	32	2.0	8.0
Dobule precision SIMD with	2	1	16	0.5	8.0
Peak Single precision scalar FLOPS (GFLOP/sec)	26	117	63	4.6	2.5
Peak Single precision s SIMD FLOPS (GFLOP/Sec)	102	311 to 933	515 to 1344	3.0-9.1	6.6-13.1
(SP 1 add or multiply)	N.A.	(311)	(515)	(3.0)	(6.6)
(SP 1 instruction fused)	N.A.	(622)	(1344)	(6.1)	(13.1)
(face SP dual issue fused)	N.A.	(933)	N.A	(9.1)	-
Peal double precision SIMD FLOPS (GFLOP/sec)	51	78	515	1.5	10.1

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

51

Rooflines

Intel Core i7-960

NVIDIA GTX280

Intel Core i7-960

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

52

Benchmarks

Kernel	Units	Core i7-960	GTX 280	GTX 280/ i7-960
SGEMM	GFLOP/sec	94	364	3.9
MC	Billion paths/sec	0.8	1.4	1.8
Conv	Million pixels/sec	1250	3500	2.8
FFT	GFLOP/sec	71.4	213	3.0
SAXPY	GBytes/sec	16.8	88.8	5.3
LBM	Million lookups/sec	85	426	5.0
Solv	Frames/sec	103	52	0.5
SpMV	GFLOP/sec	4.9	9.1	1.9
GJK	Frames/sec	67	1020	15.2
Sort	Million elements/sec	250	198	0.8
RC	Frames/sec	5	8.1	1.6
Search	Million queries/sec	50	90	1.8
Hist	Million pixels/sec	1517	2583	1.7
Bilat	Million pixels/sec	83	475	5.7

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU

2016/4/21

53

Performance Summary

- GPU (480) has 4.4 X the memory bandwidth
 - Benefits memory bound kernels
- GPU has 13.1 X the single precision throughput, 2.5 X the double precision throughput
 - Benefits FP compute bound kernels
- CPU cache prevents some kernels from becoming memory bound when they otherwise would on GPU
- GPUs offer scatter-gather, which assists with kernels with strided data
- Lack of synchronization and memory consistency support on GPU limits performance for some kernels

Chapter 6 — Parallel Processors from Client to Cloud RST©NTHU

2016/4/21

54

§6.12 Going Faster: Multiple Processors and Matrix Multiply

Multi-threading DGEMM

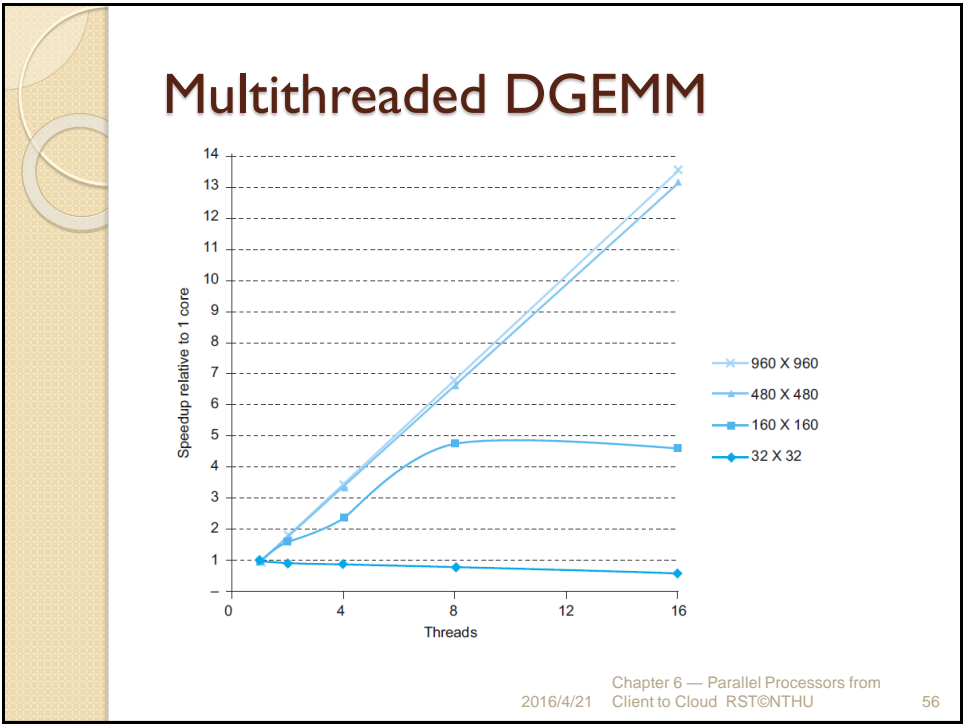
- Use OpenMP:

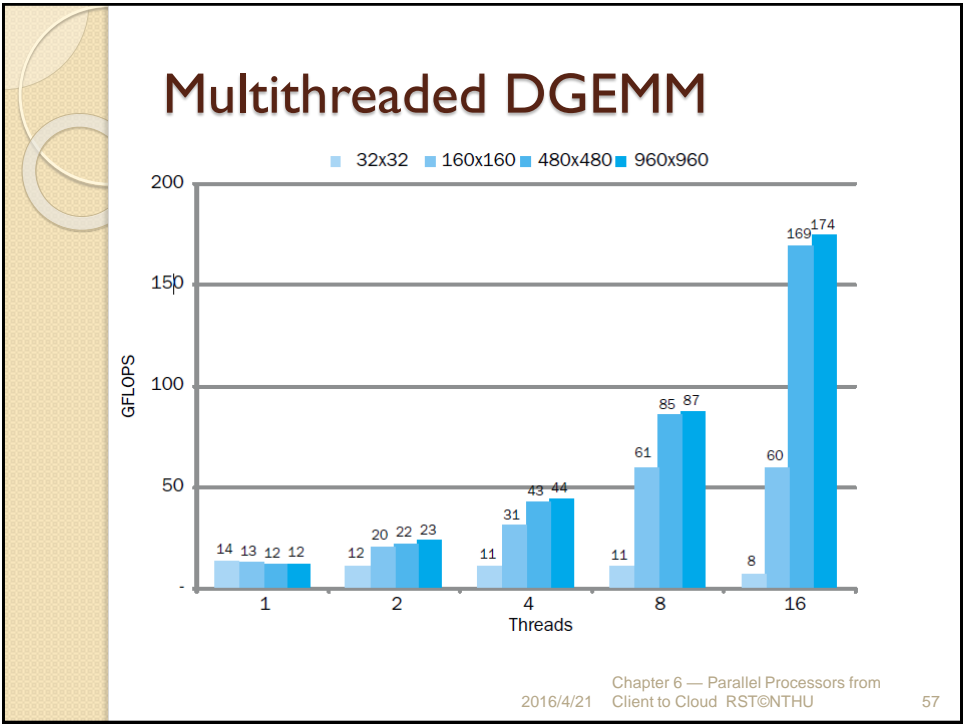
```
void dgemm (int n, double* A, double* B, double* C)
{
    #pragma omp parallel for
    for ( int sj = 0; sj < n; sj += BLOCKSIZE )
        for ( int si = 0; si < n; si += BLOCKSIZE )
            for ( int sk = 0; sk < n; sk += BLOCKSIZE )
                do_block(n, si, sj, sk, A, B, C);
}
```

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/21

55






Fallacies

- ✗ Amdahl's Law doesn't apply to parallel computers
 - Since we can achieve linear speedup
 - But only on applications with weak scaling
- ✗ Peak performance tracks observed performance
 - Marketers like this approach!
 - But compare Xeon with others in example
 - Need to be aware of bottlenecks

684

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU
2016/4/21 58

§6.12 Fallacies and Pitfalls



Pitfalls


- Not developing the software to take account of a multiprocessor architecture
 - Example: using a single lock for a shared composite resource
 - Serializes accesses, even if they could be done in parallel
 - Use finer-granularity locking

685

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/21

59



Concluding Remarks

- Goal: higher performance by using multiple processors
- Difficulties
 - Developing parallel software
 - Devising appropriate architectures
- SaaS importance is growing and clusters are a good match
- Performance per dollar and performance per Joule drive both mobile and WSC

§6.14 Concluding Remarks

Chapter 6 — Parallel Processors from Client to Cloud RST@NTHU

2016/4/21

60

