uma breve revisão sobre arquitetura de computadores – parte 1

MCZA020-13 – PROGRAMAÇÃO PARALELA

Emilio Francesquini e.francesquini@ufabc.edu.br 20 de setembro de 2018

Centro de Matemática, Computação e Cognição Universidade Federal do ABC



- Estes slides foram preparados para o curso de **Programaçao Paralela na UFABC**.
- Estes slides são baseados naqueles produzidos por Peter Pacheco como parte do livro An Introduction to Parallel Programming disponíveis em: https://www.cs.usfca.edu/~peter/ipp/
- Este material pode ser usado livremente desde que sejam mantidos, além deste aviso, os créditos aos autores e instituições.

FUNDAMENTOS

HARDWARE E SOFTWARE SEQUENCIAIS



ARQUITETURA DE VON NEUMANN



Main memory

- É uma coleção de posições, cada uma capaz de armazenar **tanto** instruções como dados
- Cada posição consiste de um endereço (usado para acessar aquela posição) e seu conteúdo



UNIDADE CENTRAL DE PROCESSAMENTO CPU

- A **CPU** (*Central Processing Unit*) pode ser dividida em duas partes
- Unidade de controle (Control Unit) é responsável por decidir quais instruções devem ser executadas. (*O chefe*)
- A Unidade Lógica e Aritimética ALU (Arithmetic and Logic Unit) – é responsável por, de fato, executar as instruções. (O Trabalhador)



- **Registrador (***Register***)** Memória com altíssimo desempenho, parte integrada à CPU
- Contador de Programa PC (Program Counter) Armazena o endereço da próxima instrução a ser executada
 - Processadores Intel usam o nome Instruction Pointer (IP)
- Barramento (Bus) hardware que conecta a CPU à memória e aos demais dispositivos.







O GARGALO DA ARQUITETURA DE VON NEUMANN



- Um processo do sistema operacional SO (operating system) é um programa que está sendo executado
- Componentes de um processo
 - O programa em linguagem de máquina
 - Regiões da memória
 - Descritores dos recursos alocados pelo SO ao processo
 - Informações de segurança
 - · Informações sobre o estado do processo

- *Multitasking* cria a ilusão de que um processador simples, com apenas um núcleo (*core*) de processamento, está rodando múltiplos programas simultâneamente
- · Cada processo se reveza no uso do processador (time slice)
- Quando o tempo alocado a um processo acaba ele é colocado em uma fila e espera novamente pela sua vez. Nestes casos dizemos que o processo está bloqueado (*blocked*)
 - A entidade responsável por fazer este trabalho é o escalonador (scheduler) do SO

- Threads elementos que compõem um processo
- Threads permitem os programadores dividirem os seus programas em tarefas virtualmente autônomas e independentes
- A ideia é de que quando um thread bloqueia por estar esperando um recurso, outro thread estara esperando por sua vez
 - A intenção é promove uma utilização mais otimizada da CPU



Modificações no Modelo de Von Neumann

CACHING - O BÁSICO

- Cache Uma coleção de posições de memória que pode ser acessada pelo processador de maneira muito mais rápida que outras posições de memória
- A cache de uma CPU é tipicamente localizada **no mesmo chip** ou em uma memória que pode ser acessada **muito** rapidamente
 - Alguns processadores como o Power8 tem até L4



- O princípio da localidade é uma heurística pela qual os projetistas de hardware esperam que o acesso a uma posição de memória seja seguido por acessos a posições em sua vizinhança
- Há dois tipos principais de localidade:
 - Localidade Espacial Acessos ocorrerão em posições de memória próximas.
 - Localidade Temporal Acessos ocorrerão em um futuro próximo.
- As caches de um processador aproveitam-se desses dois casos

```
1 float z[1000];
2 ...
3 sum = 0.0;
4 for (i = 0; i < 1000; i++)
5 sum += z[i];
```











- Quando uma CPU escreve dados na cache, o valor daquele dado pode ficar inconsistente, ou defasado (*stale*) em relação aos dados na memória principal.
- write-through os dados são escritos na cache e imediatamente enviados para a memória principal
- write-back a cache controla os dados armazenados por ela como sujos (dirty). Quando a linha de cache é substituída por uma nova linha a linha suja é enviada para a memória principal.

- Full associative, completamente associativa uma nova linha pode ser colocada em qualquer posição da cache.
- **Direct mapped**, mapeamento direto cada linha da cache tem uma posição única onde ela pode ser colocada.
- n-way set associative, associativa de n vias cada linha da cache pode ser colocada em n diferentes posições na cache.

 Quando mais de uma linha na mamória pode ser mapeada para diversas posições diferentes, acaba aparecendo a necessidade de definir uma política de substituição (*replacement policy/eviction policy* para decidir qual das linhas precisa ser substituida (replaced/evicted)



EXEMPLO

	Cache Location					
Memory Index	Fully Assoc	Direct Mapped	2-way			
0	0, 1, 2, or 3	0	0 or 1			
1	0, 1, 2, or 3	1	2 or 3			
2	0, 1, 2, or 3	2	0 or 1			
3	0, 1, 2, or 3	3	2 or 3			
4	0, 1, 2, or 3	0	0 or 1			
5	0, 1, 2, or 3	1	2 or 3			
6	0, 1, 2, or 3	2	0 or 1			
7	0, 1, 2, or 3	3	2 or 3			
8	0, 1, 2, or 3	0	0 or 1			
9	0, 1, 2, or 3	1	2 or 3			
10	0, 1, 2, or 3	2	0 or 1			
11	0, 1, 2, or 3	3	2 or 3			
12	0, 1, 2, or 3	0	0 or 1			
13	0, 1, 2, or 3	1	2 or 3			
14	0, 1, 2, or 3	2	0 or 1			
15	0.1.2 or 3	3	2 or 3			

Associações de uma memória de 16 linhas a uma cache de 4 linhas.

CACHES E PROGRAMAS

```
double A[MAX][MAX], x[MAX], y[MAX];
/* Initialize A and x, assign y = 0 * /
/* First pair of loops */
for (i = 0; i < MAX; i++)
  for (j = 0; j < MAX; j++)
     v[i] += A[i][i]*x[i]:
/* Assign y = 0 */
. . .
/* Second pair of loops */
for (j = 0; j < MAX; j++)
   for (i = 0; i < MAX; i++)
     y[i] += A[i][i]*x[i];
```

Cache Line	Elements of A						
0	A[0][0]	A[0][1]	A[0][2]	A[0][3]			
1	A[1][0]	A[1][1]	A[1][2]	A[1][3]			
2	A[2][0]	A[2][1]	A[2][2]	A[2][3]			
3	A[3][0]	A[3][1]	A[3][2]	A[3][3]			

Virtual memory (1)

- If we run a very large program or a program that accesses very large data sets, all of the instructions and data may not fit into main memory.
- Virtual memory functions as a cache for secondary storage.



Virtual memory (2)

 It exploits the principle of spatial and temporal locality.

 It only keeps the active parts of running programs in main memory.



Virtual memory (3)

- Swap space those parts that are idle are kept in a block of secondary storage.
- Pages blocks of data and instructions.
 - Usually these are relatively large.
 - Most systems have a fixed page size that currently ranges from 4 to 16 kilobytes.





Virtual memory (4)





Virtual page numbers

- When a program is compiled its pages are assigned *virtual* page numbers.
- When the program is run, a table is created that maps the virtual page numbers to physical addresses.
- A page table is used to translate the virtual address into a physical address.



Virtual Address									
Virtual Page Number			Byte Offset						
31	30		13	12	11	10		1	0
1	0	••••	1	1	0	0	••••	1	1

Table 2.2: Virtual Address Divided into Virtual Page Number and Byte Offset



Translation-lookaside buffer (TLB)

 Using a page table has the potential to significantly increase each program's overall run-time.

A special address translation cache in the processor.



Translation-lookaside buffer (2)

 It caches a small number of entries (typically 16–512) from the page table in very fast memory.

 Page fault – attempting to access a valid physical address for a page in the page table but the page is only stored on disk.



Instruction Level Parallelism (ILP)

 Attempts to improve processor performance by having multiple processor components or **functional units** simultaneously executing instructions.


Instruction Level Parallelism (2)

 Pipelining - functional units are arranged in stages.

 Multiple issue - multiple instructions can be simultaneously initiated.



Pipelining





Pipelining example (1)

Time	Operation	Operand 1	Operand 2	Result
1	Fetch operands	9.87×10^{4}	6.54×10^{3}	
2	Compare exponents	9.87×10^{4}	6.54×10^{3}	
3	Shift one operand	9.87×10^{4}	0.654×10^{4}	
4	Add	9.87×10^{4}	0.654×10^{4}	10.524×10^{4}
5	Normalize result	9.87×10^{4}	0.654×10^{4}	1.0524×10^{5}
6	Round result	9.87×10^{4}	0.654×10^{4}	1.05×10^{5}
7	Store result	9.87×10^{4}	0.654×10^{4}	1.05×10^{5}

Add the floating point numbers 9.87×10^4 and 6.54×10^3



Pipelining example (2)

```
float x[1000], y[1000], z[1000];
. . .
for (i = 0; i < 1000; i++)
z[i] = x[i] + y[i];
```

- Assume each operation takes one nanosecond (10⁻⁹ seconds).
- This for loop takes about 7000 nanoseconds.



Pipelining (3)

- Divide the floating point adder into 7 separate pieces of hardware or functional units.
- First unit fetches two operands, second unit compares exponents, etc.
- Output of one functional unit is input to the next.



Pipelining (4)

Time	Fetch	Compare	Shift	Add	Normalize	Round	Store
0	0						
1	1	0					
2	2	1	0				
3	3	2	1	0			
4	4	3	2	1	0		
5	5	4	3	2	1	0	
6	6	5	4	3	2	1	0
	:	:	:	:	:	:	
999	999	998	997	996	995	994	993
1000		999	998	997	996	995	994
1001			999	998	997	996	995
1002				999	998	997	996
1003					999	998	997
1004						999	998
1005							999

Table 2.3: Pipelined Addition.

Numbers in the table are subscripts of operands/results.



Pipelining (5)

 One floating point addition still takes 7 nanoseconds.

 But 1000 floating point additions now takes 1006 nanoseconds!



Multiple Issue (1)

 Multiple issue processors replicate functional units and try to simultaneously execute different instructions in a program.





Multiple Issue (2)

 static multiple issue - functional units are scheduled at compile time.

 dynamic multiple issue – functional units are scheduled at run-time.

superscalar



Speculation (1)

 In order to make use of multiple issue, the system must find instructions that can be executed simultaneously.



 In speculation, the compiler or the processor makes a guess about an instruction, and then executes the instruction on the basis of the guess.



Speculation (2)



If the system speculates incorrectly,

it must go back and recalculate w = y.



Hardware multithreading (1)

- There aren't always good opportunities for simultaneous execution of different threads.
- Hardware multithreading provides a means for systems to continue doing useful work when the task being currently executed has stalled.
 - Ex., the current task has to wait for data to be loaded from memory.



Hardware multithreading (2)

- Fine-grained the processor switches between threads after each instruction, skipping threads that are stalled.
 - <u>Pros</u>: potential to avoid wasted machine time due to stalls.
 - <u>Cons</u>: a thread that's ready to execute a long sequence of instructions may have to wait to execute every instruction.



Hardware multithreading (3)

 Coarse-grained - only switches threads that are stalled waiting for a timeconsuming operation to complete.

- <u>Pros</u>: switching threads doesn't need to be nearly instantaneous.
- <u>Cons</u>: the processor can be idled on shorter stalls, and thread switching will also cause delays.



Hardware multithreading (3)

- Simultaneous multithreading (SMT) a variation on fine-grained multithreading.
- Allows multiple threads to make use of the multiple functional units.



A programmer can write code to exploit.

PARALLEL HARDWARE



Flynn's Taxonomy

clae	sic von Neumann SISD Single instruction stream Single data stream	(SIMD) Single instruction stream Multiple data stream
	MISD Multiple instruction stream Single data stream	(MIMD) Multiple instruction stream Multiple data stream
	- Covered	





- Parallelism achieved by dividing data among the processors.
- Applies the same instruction to multiple data items.

Called data parallelism.



SIMD example







- What if we don't have as many ALUs as data items?
- Divide the work and process iteratively.
- Ex. m = 4 ALUs and n = 15 data items.

Round3	ALU ₁	ALU ₂	ALU ₃	ALU ₄
1	X[0]	X[1]	X[2]	X[3]
2	X[4]	X[5]	X[6]	X[7]
3	X[8]	X[9]	X[10]	X[11]
4	X[12]	X[13]	X[14]	



SIMD drawbacks

- All ALUs are required to execute the same instruction, or remain idle.
- In classic design, they must also operate synchronously.
- The ALUs have no instruction storage.
- Efficient for large data parallel problems, but not other types of more complex parallel problems.



Vector processors (1)

 Operate on arrays or vectors of data while conventional CPU's operate on individual data elements or scalars.

- Vector registers.
 - Capable of storing a vector of operands and operating simultaneously on their contents.



Vector processors (2)

Vectorized and pipelined functional units.
The same operation is applied to each element in the vector (or pairs of elements).

Vector instructions.

Operate on vectors rather than scalars.



Vector processors (3)

Interleaved memory.

- Multiple "banks" of memory, which can be accessed more or less independently.
- Distribute elements of a vector across multiple banks, so reduce or eliminate delay in loading/storing successive elements.
- Strided memory access and hardware scatter/gather.
 - The program accesses elements of a vector located at fixed intervals.



Vector processors - Pros

- Fast.
- Easy to use.



- Vectorizing compilers are good at identifying code to exploit.
- Compilers also can provide information about code that cannot be vectorized.
 - Helps the programmer re-evaluate code.
- High memory bandwidth.
- Uses every item in a cache line.



Vector processors - Cons

 They don't handle irregular data structures as well as other parallel architectures.



 A very finite limit to their ability to handle ever larger problems. (scalability)



Graphics Processing Units (GPU)

 Real time graphics application programming interfaces or API's use points, lines, and triangles to internally represent the surface of an object.





GPUs

 A graphics processing pipeline converts the internal representation into an array of pixels that can be sent to a computer screen.

 Several stages of this pipeline (called shader functions) are programmable.



• Typically just a few lines of C code.



 Shader functions are also implicitly parallel, since they can be applied to multiple elements in the graphics stream.

- GPU's can often optimize performance by using SIMD parallelism.
- The current generation of GPU's use SIMD parallelism.
 - Although they are not pure SIMD systems.





 Supports multiple simultaneous instruction streams operating on multiple data streams.

 Typically consist of a collection of fully independent processing units or cores, each of which has its own control unit and its own ALU.



Shared Memory System (1)

- A collection of autonomous processors is connected to a memory system via an interconnection network.
- Each processor can access each memory location.
- The processors usually communicate implicitly by accessing shared data structures.



Shared Memory System (2)

- Most widely available shared memory systems use one or more multicore processors.
 - (multiple CPU's or cores on a single chip)





Shared Memory System



Figure 2.3



UMA multicore system





NUMA multicore system



A memory location a core is directly connected to can be accessed faster than a memory location that must be accessed through another chip.

Figure 2.6



Distributed Memory System

- Clusters (most popular)
 - A collection of commodity systems.
 - Connected by a commodity interconnection network.
- Nodes of a cluster are individual computations units joined by a communication network.

a.k.a. hybrid systems


Distributed Memory System



Figure 2.4



Interconnection networks

- Affects performance of both distributed and shared memory systems.
- Two categories:
 - Shared memory interconnects
 - Distributed memory interconnects



Shared memory interconnects

Bus interconnect

- A collection of parallel communication wires together with some hardware that controls access to the bus.
- Communication wires are shared by the devices that are connected to it.
- As the number of devices connected to the bus increases, contention for use of the bus increases, and performance decreases.



Shared memory interconnects

Switched interconnect

- Uses switches to control the routing of data among the connected devices.
- Crossbar
 - Allows simultaneous communication among different devices.
 - Faster than buses.
 - But the cost of the switches and links is relatively high.



Figure 2.7

(a)

A crossbar switch connecting 4 processors (P_i) and 4 memory modules (M_i)

(b) Configuration of internal switches in a crossbar

(c) Simultaneous memory accesses by the processors









Copyright © 2010, Elsevier Inc. All rights Reserved

Distributed memory interconnects

Two groups

- Direct interconnect
 - Each switch is directly connected to a processor memory pair, and the switches are connected to each other.
- Indirect interconnect
 - Switches may not be directly connected to a processor.



Direct interconnect





Bisection width

- A measure of "number of simultaneous communications" or "connectivity".
- How many simultaneous communications can take place "across the divide" between the halves?





Two bisections of a ring



Figure 2.9



A bisection of a toroidal mesh



Figure 2.10



Definitions

Bandwidth

- The rate at which a link can transmit data.
- Usually given in megabits or megabytes per second.
- Bisection bandwidth
 - A measure of network quality.
 - Instead of counting the number of links joining the halves, it sums the bandwidth of the links.



Fully connected network

 Each switch is directly connected to every other switch.



Figure 2.11



Hypercube

- Highly connected direct interconnect.
- Built inductively:
 - A one-dimensional hypercube is a fullyconnected system with two processors.
 - A two-dimensional hypercube is built from two one-dimensional hypercubes by joining "corresponding" switches.
 - Similarly a three-dimensional hypercube is built from two two-dimensional hypercubes.



Hypercubes





Indirect interconnects

Simple examples of indirect networks:

- Crossbar
- Omega network

 Often shown with unidirectional links and a collection of processors, each of which has an outgoing and an incoming link, and a switching network.



A generic indirect network



Figure 2.13



Crossbar interconnect for distributed memory





An omega network



Figure 2.15



A switch in an omega network



Figure 2.16



Copyright © 2010, Elsevier Inc. All rights Reserved

More definitions

- Any time data is transmitted, we're interested in how long it will take for the data to reach its destination.
- Latency
 - The time that elapses between the source's beginning to transmit the data and the destination's starting to receive the first byte.

Bandwidth

 The rate at which the destination receives data after it has started to receive the first byte.







Cache coherence

 Programmers have no control over caches and when they get updated.



Figure 2.17

A shared memory system with two cores and two caches



Cache coherence

y0 privately owned by Core 0 y1 and z1 privately owned by Core 1

 $\mathbf{x} = 2$; /* shared variable */

Time	Core 0	Core 1
0	y0 = x;	y1 = 3*x;
1	x = 7;	Statement(s) not involving x
2	Statement(s) not involving x	z1 = 4*x;

y0 eventually ends up = 2 y1 eventually ends up = 6 z1 = ???



Snooping Cache Coherence

- The cores share a bus .
- Any signal transmitted on the bus can be "seen" by all cores connected to the bus.
- When core 0 updates the copy of x stored in its cache it also broadcasts this information across the bus.
- If core 1 is "snooping" the bus, it will see that x has been updated and it can mark its copy of x as invalid.



Directory Based Cache Coherence

- Uses a data structure called a directory that stores the status of each cache line.
- When a variable is updated, the directory is consulted, and the cache controllers of the cores that have that variable's cache line in their caches are invalidated.





PARALLEL SOFTWARE



Copyright © 2010, Elsevier Inc. All rights Reserved

The burden is on software

- Hardware and compilers can keep up the pace needed.
- From now on...
 - In shared memory programs:
 - Start a single process and fork threads.
 - Threads carry out tasks.
 - In distributed memory programs:
 - Start multiple processes.
 - Processes carry out tasks.



SPMD – single program multiple data

 A SPMD programs consists of a single executable that can behave as if it were multiple different programs through the use of conditional branches.





Writing Parallel Programs

- 1. Divide the work among the processes/threads
 - (a) so each process/thread gets roughly the same amount of work
 - (b) and communication is minimized.

- 2. Arrange for the processes/threads to synchronize.
- 3. Arrange for communication among processes/threads.



Shared Memory

Dynamic threads

- Master thread waits for work, forks new threads, and when threads are done, they terminate
- Efficient use of resources, but thread creation and termination is time consuming.
- Static threads
 - Pool of threads created and are allocated work, but do not terminate until cleanup.
 - Better performance, but potential waste of system resources.



Nondeterminism





Nondeterminism

Time	Core 0	Core 1
0	Finish assignment to my_val	In call to Compute_val
1	Load $x = 0$ into register	Finish assignment to my_val
2	Load my_val = 7 into register	Load $x = 0$ into register
3	Add my_val = $7 \text{ to } x$	Load my_val = 19 into register
4	Store $x = 7$	Add my_val to x
5	Start other work	Store $x = 19$



Nondeterminism

- Race condition
- Critical section
- Mutually exclusive
- Mutual exclusion lock (mutex, or simply lock)





message-passing

```
char message [100];
my_rank = Get_rank ( );
i f (my_rank == 1) 
  sprintf (message, "Greetings from process 1");
  Send (message, MSG_CHAR, 100, 0);
e l s e i f (my_rank == 0) 
  Receive (message, MSG_CHAR, 100, 1);
  printf ("Process 0 > \text{Received: } \% \text{s/n"}, message);
```



}

Partitioned Global Address Space Languages

shared in t $n = \ldots$: shared double x [n] , y [n] ; private i n t i , my_first_element , my_last_element ; my_first_element = . . . ; my_last_element = . . . ; /* Initialize x and y */ . . . f or (i = my_first_element ; i <= my_last_element ; i++) x[i] += y[i];


Input and Output

 In distributed memory programs, only process 0 will access *stdin*. In shared memory programs, only the master thread or thread 0 will access *stdin*.

 In both distributed memory and shared memory programs all the processes/threads can access stdout and stderr.



Input and Output

 However, because of the indeterminacy of the order of output to *stdout*, in most cases only a single process/thread will be used for all output to *stdout* other than debugging output.

 Debug output should always include the rank or id of the process/thread that's generating the output.



Input and Output

 Only a single process/thread will attempt to access any single file other than *stdin*, *stdout*, or *stderr*. So, for example, each process/thread can open its own, private file for reading or writing, but no two processes/threads will open the same file.





PERFORMANCE



Copyright © 2010, Elsevier Inc. All rights Reserved

Speedup

- Number of cores = p
- Serial run-time = T_{serial}
- Parallel run-time = T_{parallel}







Speedup of a parallel program





Efficiency of a parallel program





Speedups and efficiencies of a parallel program





Speedups and efficiencies of parallel program on different problem sizes

	p	1	2	4	8	16
Half	S	1.0	1.9	3.1	4.8	6.2
	E	1.0	0.95	0.78	0.60	0.39
Original	S	1.0	1.9	3.6	6.5	10.8
	E	1.0	0.95	0.90	0.81	0.68
Double	S	1.0	1.9	3.9	7.5	14.2
	E	1.0	0.95	0.98	0.94	0.89



Speedup





Efficiency





Effect of overhead

$T_{parallel} = T_{serial} / p + T_{overhead}$



Copyright © 2010, Elsevier Inc. All rights Reserved

Amdahl's Law

 Unless virtually all of a serial program is parallelized, the possible speedup is going to be very limited — regardless of the number of cores available.





Copyright © 2010, Elsevier Inc. All rights Reserved

Example

- We can parallelize 90% of a serial program.
- Parallelization is "perfect" regardless of the number of cores p we use.
- T_{serial} = 20 seconds
- Runtime of parallelizable part is

 $0.9 \times T_{serial} / p = 18 / p$



Example (cont.)

Runtime of "unparallelizable" part is

 $0.1 \times T_{serial} = 2$

Overall parallel run-time is

$$T_{parallel} = 0.9 \text{ x } T_{serial} / p + 0.1 \text{ x } T_{serial} = 18 / p + 2$$



Example (cont.)

Speed up





Scalability

- In general, a problem is *scalable* if it can handle ever increasing problem sizes.
- If we increase the number of processes/threads and keep the efficiency fixed without increasing problem size, the problem is *strongly scalable*.
- If we keep the efficiency fixed by increasing the problem size at the same rate as we increase the number of processes/threads, the problem is *weakly scalable*.



- What is time?
- Start to finish?
- A program segment of interest?
- CPU time?
- Wall clock time?









```
private double start, finish;
....
start = Get_current_time();
/* Code that we want to time */
....
finish = Get_current_time();
printf("The elapsed time = %e seconds\n", finish-start);
```



```
shared double global_elapsed;
private double my_start, my_finish, my_elapsed;
/* Synchronize all processes/threads */
Barrier();
my_start = Get_current_time();
/* Code that we want to time */
. . .
my finish = Get current time():
my elapsed = my finish - my start;
/* Find the max across all processes/threads */
global_elapsed = Global_max(my_elapsed);
if (my rank == 0)
   printf("The elapsed time = %e seconds\n", global_elapsed);
```





PARALLEL PROGRAM DESIGN



Copyright © 2010, Elsevier Inc. All rights Reserved

1. Partitioning: divide the computation to be performed and the data operated on by the computation into small tasks.

The focus here should be on identifying tasks that can be executed in parallel.



2. Communication: determine what communication needs to be carried out among the tasks identified in the previous step.





Copyright © 2010, Elsevier Inc. All rights Reserved

3. Agglomeration or aggregation: combine tasks and communications identified in the first step into larger tasks.

For example, if task A must be executed before task B can be executed, it may make sense to aggregate them into a single composite task.



4. Mapping: assign the composite tasks identified in the previous step to processes/threads.

This should be done so that communication is minimized, and each process/thread gets roughly the same amount of work.



Example - histogram

1.3,2.9,0.4,0.3,1.3,4.4,1.7,0.4,3.2,0.3,4.9,2
.4,3.1,4.4,3.9,0.4,4.2,4.5,4.9,0.9





Serial program - input

- 1. The number of measurements: data_count
- 2. An array of data_count floats: data
- 3. The minimum value for the bin containing the smallest values: min_meas
- 4. The maximum value for the bin containing the largest values: max_meas
- 5. The number of bins: bin_count



Serial program - output

1. bin_maxes : an array of bin_count floats

2. bin_counts : an array of bin_count ints







Alternative definition of tasks and communication





Adding the local arrays





Concluding Remarks (1)

- Serial systems
 - The standard model of computer hardware has been the von Neumann architecture.
- Parallel hardware
 - Flynn's taxonomy.
- Parallel software
 - We focus on software for homogeneous MIMD systems, consisting of a single program that obtains parallelism by branching.
 - SPMD programs.



Concluding Remarks (2)

Input and Output

- We'll write programs in which one process or thread can access stdin, and all processes can access stdout and stderr.
- However, because of nondeterminism, except for debug output we'll usually have a single process or thread accessing stdout.



Concluding Remarks (3)

- Performance
 - Speedup
 - Efficiency
 - Amdahl's law
 - Scalability
- Parallel Program Design
 - Foster's methodology

