# Programming on Parallel Machines

Norman Matloff
University of California, Davis [1]

# Contents

# Chapter 1

# Introduction to Parallel Processing

Parallel machines provide a wonderful opportunity for applications with large computational requirements. Effective use of these machines, though, requires a keen understanding of how they work. This chapter provides an overview.

## 1.1 Overview: Why Use Parallel Systems?

### 1.1.1 Execution Speed

There is an ever-increasing appetite among some types of computer users for faster and faster machines. This was epitomized in a statement by Steve Jobs, founder/CEO of Apple and Pixar. He noted that when he was at Apple in the 1980s, he was always worried that some other company would come out with a faster machine than his. But now at Pixar, whose graphics work requires extremely fast computers, he is always hoping someone produces faster machines, so that he can use them!

A major source of speedup is the parallelizing of operations. Parallel operations can be either within-processor, such as with pipelining or having several ALUs within a processor, or between-processor, in which many processor work on different parts of a problem in parallel. Our focus here is on between-processor operations.

For example, the Registrar's Office at UC Davis uses shared-memory multiprocessors for processing its on-line registration work. Online registration involves an enormous amount of database computation. In order to handle this computation reasonably quickly, the program partitions the work to be done, assigning different portions of the database to different processors. The database field has contributed greatly to the commercial success of large shared-memory machines.

As the Pixar example shows, highly computation-intensive applications like computer graphics also have a

need for these fast parallel computers. No one wants to wait hours just to generate a single image, and the use of parallel processing machines can speed things up considerably. For example, consider **ray tracing** operations. Here our code follows the path of a ray of light in a scene, accounting for reflection and absorbtion of the light by various objects. Suppose the image is to consist of 1,000 rows of pixels, with 1,000 pixels per row. In order to attack this problem in a parallel processing manner with, say, 25 processors, we could divide the image into 25 squares of size 200x200, and have each processor do the computations for its square.

Note, though, that it may be much more challenging than this implies. First of all, the computation will need some communication between the processors, which hinders performance if it is not done carefully. Second, if one really wants good speedup, one may need to take into account the fact that some squares require more computation work than others. More on this below.

In this setting you need the program to run as fast as possible. Thus, in order to write good parallel processing software, you must have a good knowledge of the underlying hardware. You must find clever tricks for **load balancing,** i.e. keeping all the processors busy as much as possible. In the graphics ray-tracing application, for instance, suppose a ray is coming from the "northeast" section of the image, and is reflected by a solid object. Then the ray won't reach some of the "southwest" portions of the image, which then means that the processors assigned to those portions will not have any work to do which is associated with this ray. What we need to do is then try to give these processors some other work to do; the more they are idle, the slower our system will be.

### 1.1.2   Memory

Yes, execution speed is the reason that comes to most people's minds when the subject of parallel processing comes up. But in many applications, an equally important consideration is memory capacity. Parallel processing application often tend to use huge amounts of memory, and in many cases the amount of memory needed is more than can fit on one machine. If we have many machines working together, especially in the message-passing settings described below, we can accommodate the large memory needs.

## 1.2   Parallel Processing Hardware

This is not a hardware course, but since the goal of using parallel hardware is speed, the efficiency of our code is a major issue. That in turn means that we need a good understanding of the underlying hardware that we are programming. In this section, we give an overview of parallel hardware.

### 1.2.1 Shared-Memory Systems

#### 1.2.1.1 Basic Architecture

Here many CPUs share the same physical memory. This kind of architecture is sometimes called MIMD, standing for Multiple Instruction (different CPUs are working independently, and thus typically are executing different instructions at any given instant), Multiple Data (different CPUs are generally accessing different memory locations at any given time).

Until recently, shared-memory systems cost hundreds of thousands of dollars and were affordable only by large companies, such as in the insurance and banking industries. The high-end machines are indeed still quite expensive, but now **dual-core** machines, in which two CPUs share a common memory, are commonplace in the home.

#### 1.2.1.2 Example: SMP Systems

A Symmetric Multiprocessor (SMP) system has the following structure:



Here and below:

- The Ps are processors, e.g. off-the-shelf chips such as Pentiums.

- The Ms are **memory modules**. These are physically separate objects, e.g. separate boards of memory chips. It is typical that there will be the same number of memory modules as processors. In the shared-memory case, the memory modules collectively form the entire shared address space, but with the addresses being assigned to the memory modules in one of two ways:

  - (a)
    High-order interleaving. Here consecutive addresses are in the <u>same</u> M (except at boundaries). For example, suppose for simplicity that our memory consists of addresses 0 through 1023, and that there are four Ms. Then M0 would contain addresses 0-255, M1 would have 256-511, M2 would have 512-767, and M3 would have 768-1023.

    We need 10 bits for addresses (since $1024 = 2^{10}$). The two most-significant bits would be used to select the module number (since $4 = 2^2$); hence the term *high-order* in the name of this design. The remaining eight bits are used to select the word within a module.

- (b)

    Low-order interleaving. Here consecutive addresses are in consecutive memory modules (except when we get to the right end). In the example above, if we used low-order interleaving, then address 0 would be in M0, 1 would be in M1, 2 would be in M2, 3 would be in M3, 4 would be back in M0, 5 in M1, and so on.

    Here the two least-significant bits are used to determine the module number.

- To make sure only one processor uses the bus at a time, standard bus arbitration signals and/or arbitration devices are used.

- There may also be **coherent caches**, which we will discuss later.

### 1.2.2   Message-Passing Systems

#### 1.2.2.1   Basic Architecture

Here we have a number of independent CPUs, each with its own independent memory. The various processors communicate with each other via networks of some kind.

#### 1.2.2.2   Example: Networks of Workstations (NOWs)

Large shared-memory multiprocessor systems are still very expensive. A major alternative today is networks of workstations (NOWs). Here one purchases a set of commodity PCs and networks them for use as parallel processing systems. The PCs are of course individual machines, capable of the usual uniprocessor (or now multiprocessor) applications, but by networking them together and using parallel-processing software environments, we can form very powerful parallel systems.

The networking does result in a significant loss of performance. This will be discussed in a later unit. But even without these techniques, the price/performance ratio in NOW is much superior in many applications to that of shared-memory hardware.

One factor which can be key to the success of a NOW is the use of a fast network, fast both in terms of hardware and network protocol. Ordinary Ethernet and TCP/IP are fine for the applications envisioned by the original designers of the Internet, e.g. e-mail and file transfer, but is slow in the NOW context. A good network for a NOW is, for instance, Infiniband.

NOWs have become so popular that there are now "recipes" on how to build them for the specific purpose of parallel processing. The term **Beowulf** come to mean a cluster of PCs, usually with a fast network connecting them, used for parallel processing. Software packages such as ROCKS (`http://www.rocksclusters.org/wordpress/`) have been developed to make it easy to set up and administer such systems.

### 1.2.3 SIMD

In contrast to MIMD systems, processors in SIMD—Single Instruction, Multiple Data—systems execute in lockstep. At any given time, all processors are executing the same machine instruction on different data.

## 1.3 Programmer World Views

To explain the two paradigms, we will use the term **nodes**, where roughly speaking one node corresponds to one processor, and use the following example:

> Suppose we wish to multiply an nx1 vector X by an nxn matrix A, putting the product in an nx1 vector Y, and we have p processors to share the work.

### 1.3.1 Shared-Memory

#### 1.3.1.1 Programmer View

In the shared-memory paradigm, the arrays for A, X and Y would be held in common by all nodes. If for instance node 2 were to execute

```
Y[3] = 12;
```

and then node 15 were to subsequently execute

```
print("%d\n",Y[3]);
```

then the outputted value from the latter would be 12.

#### 1.3.1.2 Example

Today, programming on shared-memory multiprocessors is typically done via **threading**. (Or, as we will see in other units, by higher-level code that runs threads underneath.) A **thread** is similar to a **process** in an operating system (OS), but with much less overhead. Threaded applications have become quite popular in even uniprocessor systems, and Unix,[1] Windows, Python, Java and Perl all support threaded programming.

---

[1]Here and below, the term *Unix* includes Linux.

In the typical implementation, a thread is a special case of an OS process. One important difference is that the various threads of a program share memory. (One can arrange for processes to share memory too in some OSs, but they don't do so by default.)

On a uniprocessor system, the threads of a program take turns executing, so that there is only an illusion of parallelism. But on a multiprocessor system, one can genuinely have threads running in parallel.

One of the most popular threads systems is Pthreads, whose name is short for POSIX threads. POSIX is a Unix standard, and the Pthreads system was designed to standardize threads programming on Unix. It has since been ported to other platforms.

Following is an example of Pthreads programming, in which we determine the number of prime numbers in a certain range. Read the comments at the top of the file for details; the threads operations will be explained presently.

```
1   // PrimesThreads.c
2
3   // threads-based program to find the number of primes between 2 and n;
4   // uses the Sieve of Eratosthenes, deleting all multiples of 2, all
5   // multiples of 3, all multiples of 5, etc.
6
7   // for illustration purposes only; NOT claimed to be efficient
8
9   // Unix compilation:  gcc -g -o primesthreads PrimesThreads.c -lpthread -lm
10
11  // usage:  primesthreads n num_threads
12
13  #include <stdio.h>
14  #include <math.h>
15  #include <pthread.h>  // required for threads usage
16
17  #define MAX_N 100000000
18  #define MAX_THREADS 25
19
20  // shared variables
21  int nthreads,  // number of threads (not counting main())
22      n,  // range to check for primeness
23      prime[MAX_N+1],  // in the end, prime[i] = 1 if i prime, else 0
24      nextbase;  // next sieve multiplier to be used
25  // lock for the shared variable nextbase
26  pthread_mutex_t nextbaselock = PTHREAD_MUTEX_INITIALIZER;
27  // ID structs for the threads
28  pthread_t id[MAX_THREADS];
29
30  // "crosses out" all odd multiples of k
31  void crossout(int k)
32  {  int i;
33     for (i = 3; i*k <= n; i += 2)  {
34        prime[i*k] = 0;
35     }
36  }
37
38  // each thread runs this routine
39  void *worker(int tn)  // tn is the thread number (0,1,...)
```

```
40   {   int lim,base,
41           work = 0;   // amount of work done by this thread
42       // no need to check multipliers bigger than sqrt(n)
43       lim = sqrt(n);
44       do   {
45           // get next sieve multiplier, avoiding duplication across threads
46           // lock the lock
47           pthread_mutex_lock(&nextbaselock);
48           base = nextbase;
49           nextbase += 2;
50           // unlock
51           pthread_mutex_unlock(&nextbaselock);
52           if (base <= lim)   {
53               // don't bother crossing out if base known composite
54               if (prime[base])   {
55                   crossout(base);
56                   work++;   // log work done by this thread
57               }
58           }
59           else return work;
60       } while (1);
61   }
62
63   main(int argc, char **argv)
64   {   int nprimes,   // number of primes found
65           i,work;
66       n = atoi(argv[1]);
67       nthreads = atoi(argv[2]);
68       // mark all even numbers nonprime, and the rest "prime until
69       // shown otherwise"
70       for (i = 3; i <= n; i++)   {
71           if (i%2 == 0) prime[i] = 0;
72           else prime[i] = 1;
73       }
74       nextbase = 3;
75       // get threads started
76       for (i = 0; i < nthreads; i++)   {
77           // this call says to create a thread, record its ID in the array
78           // id, and get the thread started executing the function worker(),
79           // passing the argument i to that function
80           pthread_create(&id[i],NULL,worker,i);
81       }
82
83       // barrier, to wait for all done
84       for (i = 0; i < nthreads; i++)   {
85           // this call said to wait until thread number id[i] finishes
86           // execution, and to assign the return value of that thread to our
87           // local variable work here
88           pthread_join(id[i],&work);
89           printf("%d values of base done\n",work);
90       }
91
92       // report results
93       nprimes = 1;
94       for (i = 3; i <= n; i++)
95           if (prime[i])   {
96               nprimes++;
97           }
```

```
98      printf("the number of primes found was %d\n",nprimes);
99
100  }
```

To make our discussion concrete, suppose we are running this program with two threads. Suppose also the both threads are running simultaneously most of the time. This will occur if they aren't competing for turns with other big threads, say if there are no other big threads, or more generally if the number of other big threads is less than or equal to the number of processors minus two.

Note the global variables:

```
int nthreads,  // number of threads (not counting main())
    n,  // range to check for primeness
    prime[MAX_N+1],  // in the end, prime[i] = 1 if i prime, else 0
    nextbase;  // next sieve multiplier to be used
pthread_mutex_t nextbaselock = PTHREAD_MUTEX_INITIALIZER;
pthread_t id[MAX_THREADS];
```

This will require some adjustment for those who've been taught that global variables are "evil." All communication between threads is via global variables, so if they are evil, they are a necessary evil. Personally I think the stern admonitions against global variables is overblown anyway. See `http://heather.cs.ucdavis.edu/~matloff/globals.html`.

As mentioned earlier, the globals are shared by all processors.[2] If one processor, for instance, assigns the value 0 to **prime[35]** in the function **crossout()**, then that variable will have the value 0 when accessed by any of the other processors as well. On the other hand, local variables have different values at each processor; for instance, the variable **i** in that function has a different value at each processor.

Note that in the statement

```
pthread_mutex_t nextbaselock = PTHREAD_MUTEX_INITIALIZER;
```

the right-hand side is not a constant. It is a macro call, and is thus something which is executed.

In the code

```
pthread_mutex_lock(&nextbaselock);
base = nextbase
nextbase += 2
pthread_mutex_unlock(&nextbaselock);
```

we see a **critical section** operation which is typical in shared-memory programming. In this context here, it means that we cannot allow more than one thread to execute

---

[2]Technically, we should say "shared by all threads" here, as a given thread does not always execute on the same processor, but at any instant in time each executing thread is at some processor, so the statement is all right.

```
base = nextbase;
nextbase += 2;
```

at the same time. The calls to **pthread_mutex_lock()** and **pthread_mutex_unlock()** ensure this. If thread A is currently executing inside the critical section and thread B tries to lock the lock by calling **pthread_mutex_lock()**, the call will block until thread B executes **pthread_mutex_unlock()**.

Here is why this is so important: Say currently **nextbase** has the value 11. What we want to happen is that the next thread to read **nextbase** will "cross out" all multiples of 11. But if we allow two threads to execute the critical section at the same time, the following may occur:

- thread A reads **nextbase**, setting its value of **base** to 11

- thread B reads **nextbase**, setting its value of **base** to 11

- thread A adds 2 to **nextbase**, so that **nextbase** becomes 13

- thread B adds 2 to **nextbase**, so that **nextbase** becomes 15

Two problems would then occur:

- Both threads would do "crossing out" of multiples of 11, duplicating work and thus slowing down execution speed.

- We will never "cross out" multiples of 13.

Thus the lock is crucial to the correct (and speedy) execution of the program.

Note that these problems could occur either on a uniprocessor or multiprocessor system. In the uniprocessor case, thread A's turn might end right after it reads **nextbase**, followed by a turn by B which executes that same instruction. In the multiprocessor case, A and B could literally be running simultaneously, but still with the action by B coming an instant after A.

This problem frequently arises in parallel database systems. For instance, consider an airline reservation system. If a flight has only one seat left, we want to avoid giving it to two different customers who might be talking to two agents at the same time. The lines of code in which the seat is finally assigned (the **commit** phase, in database terminology) is then a critical section.

A critical section is always a potential bottlement in a parallel program, because its code is serial instead of parallel. In our program here, we may get better performance by having each thread work on, say, five values of **nextbase** at a time. Our line

```
nextbase += 2;
```

would become

```
nextbase += 10;
```

That would mean that any given thread would need to go through the critical section only one-fifth as often, thus greatly reducing overhead. On the other hand, near the end of the run, this may result in some threads being idle while other threads still have a lot of work to do.

Note the **barrier**:

```
for (i = 0; i < nthreads; i++)  {
   pthread_join(id[i],&work);
   printf("%d values of base done\n",work);
}
```

A barrier is a point in the code that all threads must reach before continuing. In this case, a barrier is needed in order to prevent premature execution of the later code

```
for (i = 3; i <= n; i++)
   if (prime[i])  {
      nprimes++;
   }
```

which would result in possibly wrong output if we start counting primes before some threads are done.

The **pthread_join()** function actually causes the given thread to exit, so that we then "join" the thread that created it, i.e. **main()**. Thus some may argue that this is not really a true barrier.

Barriers are very common in shared-memory programming, and will be discussed in more detail in a later unit.

### 1.3.2   Message Passing

#### 1.3.2.1   Programmer View

By contrast, in the message-passing paradigm, all nodes would have separate copies of A, X and Y. In this case, in our example above, in order for node 2 to send this new value of Y[3] to node 15, it would have to execute some special function, which would be something like

```
  send(15,12,"Y[3]");
```

and node 15 would have to execute some kind of **receive()** function.

The conventional wisdom is that the shared-memory paradigm is much easier to program in than the message-passing paradigm. The latter, however, may be easier to implement, and in some settings may have greater speed.

### 1.3.3  Example

Here we use the MPI system, with our hardware being a NOW.

MPI is a popular public-domain set of interface functions, callable from C/C++, to do message passing. We are again counting primes, though in this case using a **pipelining** method. It is similar to hardware pipelines, but in this case it is done in software, and each "stage" in the pipe is a different computer.

The program is self-documenting, via the comments.

```
1
2  /* MPI sample program; NOT INTENDED TO BE EFFICIENT as a prime
3     finder, either in algorithm or implementation
4
5     MPI (Message Passing Interface) is a popular package using
6     the "message passing" paradigm for communicating between
7     processors in parallel applications; as the name implies,
8     processors communicate by passing messages using "send" and
9     "receive" functions
10
11    finds and reports the number of primes less than or equal to N
12
13    uses a pipeline approach:  node 0 looks at all the odd numbers
14    (i.e. has already done filtering out of multiples of 2) and
15    filters out those that are multiples of 3, passing the rest
16    to node 1; node 1 filters out the multiples of 5, passing
17    the rest to node 2; in this simple example, we just have node
18    2 filter out all the rest and then report the number of primes
19
20    note that we should NOT have a node run through all numbers
21    before passing them on to the next node, since we would then
22    have no parallelism at all; on the other hand, passing on just
23    one number at a time isn't efficient either, due to the high
24    overhead of sending a message if it is a network (tens of
25    microseconds until the first bit reaches the wire, due to
26    software delay); thus efficiency would be greatly improved if
27    each node saved up a chunk of numbers before passing them to
28    the next node */
29
30  // this include file is mandatory
31  #include <mpi.h>
32
33  #define MAX_N 100000
34  #define PIPE_MSG 0  // type of message containing a number to
35                         be checked
36  #define END_MSG 1  // type of message indicating no more data will
```

```
37                         be coming
38
39   int NNodes,  /* number of nodes in computation*/
40       N,  /* find all primes from 2 to N */
41       Me,  /* my node number */
42       ToCheck;  /* current number to check for passing on to next node;
43                     stylistically this might be nicer as a local in
44                     Node*(), but I have placed it here to dramatize
45                     the fact that the globals are NOT shared among
46                     the nodes */
47
48   double T1,T2;  /* start and finish times */
49
50   Init(Argc,Argv)
51       int Argc;  char **Argv;
52
53   {  int DebugWait;
54
55       N = atoi(Argv[1]);
56       DebugWait = atoi(Argv[2]);
57
58       /* this loop is here to synchronize all nodes for debugging;
59          if DebugWait is specified as 1 on the command line, all nodes
60          wait here until the debugging programmer starts GDB at all
61          nodes and within GDB sets DebugWait to 0 to then proceed */
62       while (DebugWait) ;
63
64       /* mandatory to begin any MPI program */
65       MPI_Init(&Argc,&Argv);
66
67       /* puts the number of nodes in NNodes */
68       MPI_Comm_size(MPI_COMM_WORLD,&NNodes);
69       /* puts the node number of this node in Me */
70       MPI_Comm_rank(MPI_COMM_WORLD,&Me);
71
72       /* OK, get started; first record current time in T1 */
73       if (Me == 2) T1 = MPI_Wtime();
74   }
75
76   Node0()
77
78   {  int I,Dummy,
79          Error;  /* not checked in this example */
80       for (I = 1; I <= N/2; I++)  {
81          ToCheck = 2 * I + 1;
82          if (ToCheck > N) break;
83          /* MPI_Send  --  send a message
84                   parameters:
85                   pointer to place where message is to be drawn from
86                   number of items in message
87                   item type
88                   destination node
89                   message type ("tag") programmer-defined
90                   node group number (in this case all nodes) */
91          if (ToCheck % 3 > 0)
92              Error = MPI_Send(&ToCheck,1,MPI_INT,1,PIPE_MSG,MPI_COMM_WORLD);
93       }
94       Error = MPI_Send(&Dummy,1,MPI_INT,1,END_MSG,MPI_COMM_WORLD);
```

```
95   }
96
97   Node1()
98
99   {  int Error,  /* not checked in this example */
100         Dummy;
101      MPI_Status Status;  /* see below */
102
103      while (1)  {
104         /* MPI_Recv  --  receive a message
105                   parameters:
106                   pointer to place to store message
107                   number of items in message (see notes on
108                      this at the end of this file)
109                   item type
110                   accept message from which node(s)
111                   message type ("tag"), programmer-defined (in this
112              case any type)
113                   node group number (in this case all nodes)
114                   status (see notes on this at the end of this file) */
115         Error = MPI_Recv(&ToCheck,1,MPI_INT,0,MPI_ANY_TAG,
116                  MPI_COMM_WORLD,&Status);
117         if (Status.MPI_TAG == END_MSG) break;
118         if (ToCheck % 5 > 0)
119            Error = MPI_Send(&ToCheck,1,MPI_INT,2,PIPE_MSG,MPI_COMM_WORLD);
120      }
121      /* now send our end-of-data signal, which is conveyed in the
122         message type, not the message (we have a dummy message just
123         as a placeholder */
124      Error = MPI_Send(&Dummy,1,MPI_INT,2,END_MSG,MPI_COMM_WORLD);
125   }
126
127   Node2()
128
129   {  int ToCheck,  /* current number to check from Node 0 */
130         Error,  /* not checked in this example */
131         PrimeCount,I,IsComposite;
132      MPI_Status Status;  /* see below */
133
134      PrimeCount = 3;  /* must account for the primes 2, 3 and 5, which
135                      won't be detected below */
136      while (1)  {
137         Error = MPI_Recv(&ToCheck,1,MPI_INT,1,MPI_ANY_TAG,
138                  MPI_COMM_WORLD,&Status);
139         if (Status.MPI_TAG == END_MSG) break;
140         IsComposite = 0;
141         for (I = 7; I*I <= ToCheck; I += 2)
142            if (ToCheck % I == 0)  {
143       IsComposite = 1;
144       break;
145   }
146         if (!IsComposite) PrimeCount++;
147      }
148      /* check the time again, and subtract to find run time */
149      T2 = MPI_Wtime();
150      printf("elapsed time = %f\n",(float)(T2-T1));
151      /* print results */
152      printf("number of primes = %d\n",PrimeCount);
```

```
153   }
154
155   main(argc,argv)
156      int argc; char **argv;
157
158   {  Init(argc,argv);
159      /* note:  instead of having a switch statement, we could write
160         three different programs, each running on a different node */
161      switch (Me)  {
162         case 0:  Node0();
163                  break;
164         case 1:  Node1();
165                  break;
166         case 2:  Node2();
167      };
168      /* mandatory for all MPI programs */
169      MPI_Finalize();
170   }
171
172   /* explanation of "number of items" and "status" arguments at the end
173      of MPI_Recv():
174
175      when receiving a message you must anticipate the longest possible
176      message, but the actual received message may be much shorter than
177      this; you can call the MPI_Get_count() function on the status
178      argument to find out how many items were actually received
179
180      the status argument will be a pointer to a struct, containing the
181      node number, message type and error status of the received
182      message
183
184      say our last parameter is Status; then Status.MPI_SOURCE
185      will contain the number of the sending node, and
186      Status.MPI_TAG will contain the message type; these are
187      important if used MPI_ANY_SOURCE or MPI_ANY_TAG in our
188      node or tag fields but still have to know who sent the
189      message or what kind it is */
```

The set of machines can be heterogeneous, but MPI "translates" for you automatically. If say one node has a big-endian CPU and another has a little-endian CPU, MPI will do the proper conversion.

# Chapter 2

# Shared Memory Parallelism

Shared-memory programming is considered by many in the parallel processing community as being the clearest of the various parallel paradigms available.

## 2.1 What Is Shared?

The term **shared memory** means that the processors all share a common address space. Say this is occurring at the hardware level, and we are using Intel Pentium CPUs. Suppose processor P3 issues the instruction

```
movl 200, %eabx
```

which reads memory location 200 and places the result in the EAX register in the CPU. If processor P4 does the same, they both will be referring to the same physical memory cell. In non-shared-memory machines, each processor has its own private memory, and each one will then have its own location 200, completely independent of the locations 200 at the other processors' memories.

Say a program contains a global variable **X** and a local variable **Y** on share-memory hardware (and we use shared-memory software). If for example the compiler assigns location 200 to the variable **X**, i.e. `&X = 200`, then the point is that all of the processors will have that variable in common, because any processor which issues a memory operation on location 200 will access the same physical memory cell.

On the other hand, each processor will have its own separate run-time stack. All of the stacks are in shared memory, but they will be accessed separately, since each CPU has a different value in its SP (Stack Pointer) register. Thus each processor will have its own independent copy of the local variable **Y**.

To make the meaning of "shared memory" more concrete, suppose we have a bus-based system, with all the processors and memory attached to the bus. Let us compare the above variables **X** and **Y** here. Suppose

15

again that the compiler assigns **X** to memory location 200.  Then in the machine language code for the program, every reference to **X** will be there as 200. Every time an instruction that writes to **X** is executed by a CPU, that CPU will put 200 into its Memory Address Register (MAR), from which the 200 flows out on the address lines in the bus, and goes to memory.  This will happen in the same way no matter which CPU it is.  Thus the same physical memory location will end up being accessed, no matter which CPU generated the reference.

By contrast, say the compiler assigns **Y** to something like ESP+8, the third item on the stack.  (It is 8 bytes past the word pointed to by the stack pointer, ESP.) Each CPU will have its own current value for ESP, so the stacks of the various CPUs will be separate. Note that the stacks <u>will</u> be in the physical shared memory, and thus P3, say, could theoretically access P8's stack, say if there were an erroneous pointer value.  But even that would not occur if we are using virtual memory and thus have protections against this.

## 2.2   Structures for Sharing

### 2.2.1   SMP Systems

A Symmetric Multiprocessor (SMP) system has the following structure:



Here and below:

- The Ps are processors, e.g. off-the-shelf chips such as Pentiums.

- The Ms are **memory modules**. These are physically separate objects, e.g. separate boards of memory chips.  It is typical that there will be the same number of Ms as Ps, but it does not have to be this way.  In the shared-memory case, the Ms collectively form the entire shared address space, but with the addresses being assigned to the Ms in one of two ways:

  - (a)
    High-order interleaving. Here consecutive addresses are in the <u>same</u> M (except at boundaries). For example, suppose for simplicity that our memory consists of addresses 0 through 1023, and that there are four Ms. Then M0 would contain addresses 0-255, M1 would have 256-511, M2 would have 512-767, and M3 would have 768-1023.

  - (b)

> Low-order interleaving. Here consecutive addresses are in consecutive M's (except when we get to the right end). In the example above, if we used low-order interleaving, then address 0 would be in M0, 1 would be in M1, 2 would be in M2, 3 would be in M3, 4 would be back in M0, 5 in M1, and so on.

- To make sure only one P uses the bus at a time, standard bus arbitration signals and/or arbitration devices are used.

- There may also be **coherent caches**, which we will discuss later.

### 2.2.2 NUMA Systems

In a **Nonuniform Memory Access** (NUMA) architecture, each CPU has a memory module physically next to it, and these processor/memory (P/M) pairs are connected by some kind of network.

Here is a simple version:



Each P/M/R set here is called a **processing element** (PE). Note that each PE has its own local bus, and is also connected to the global bus via R, the router.

Suppose for example that P3 needs to access location 200, and suppose that high-order interleaving is used. If location 200 is in M3, then P3's request is satisfied by the local bus.[1] On the other hand, suppose location 200 is in M8. Then the R3 will notice this, and put the request on the global bus, where it will be seen by R8, which will then copy the request to the local bus at PE8, where the request will be satisfied. (E.g. if it was a read request, then the response will go back from M8 to R8 to the global bus to R3 to P3.)

It should be obvious now where NUMA gets its name. P8 will have much faster access to M8 than P3 will to M8, if none of the buses is currently in use—and if say the global bus is currently in use, P3 will have to wait a long time to get what it wants from M8.

---

[1]This sounds similar to the concept of a cache. However, it is very different. A cache contains a local copy of some data stored elsewhere. Here it is the data itself, not a copy, which is being stored locally.

Today almost all high-end MIMD systems are NUMAs. One of the attractive features of NUMA is that by good programming we can exploit the nonuniformity. In matrix problems, for example, we can write our program so that, for example, P8 usually works on those rows of the matrix which are stored in M8, P3 usually works on those rows of the matrix which are stored in M3, etc. In order to do this, we need to make use of the C language's & address operator, and have some knowledge of the memory hardware structure, i.e. the interleaving.

### 2.2.3   NUMA Interconnect Topologies

The problem with a bus connection, of course, is that there is only one pathway for communication, and thus only one processor can access memory at the same time. If one has more than, say, two dozen processors are on the bus, the bus becomes saturated, even if traffic-reducing methods such as adding caches are used. Thus multipathway topologies are used for all but the smallest systems. In this section we look at two alternatives to a bus topology.

#### 2.2.3.1   Crossbar Interconnects

Consider a shared-memory system with n processors and n memory modules. Then a crossbar connection would provide $n^2$ pathways. E.g. for n = 8:

M 7

M 6

M 5

M 4

M 3

M 2

M 1

M 0

P0    P 1    P 2    P3    P 4    P 5    P 6    P7

Generally serial communication is used from node to node, with a packet containing information on both source and destination address. E.g. if P2 wants to read from M5, the source and destination will be 3-bit strings in the packet, coded as 010 and 101, respectively. The packet will also contain bits which specify which word within the module we wish to access, and bits which specify whether we wish to do a read or a write. In the latter case, additional bits are used to specify the value to be written.

Each diamond-shaped node has two inputs (bottom and right) and two outputs (left and top), with buffers at the two inputs. If a buffer fills, there are two design options: (a) Have the node from which the input comes block at that output. (b) Have the node from which the input comes discard the packet, and retry later, possibly outputting some other packet for now. If the packets at the heads of the two buffers both need to go out the same output, the one (say) from the bottom input will be given priority.

There could also be a return network of the same type, with this one being memory $\rightarrow$ processor, to return

the result of the read requests.[2]

Another version of this is also possible. It is not shown here, but the difference would be that at the bottom edge we would have the PEi and at the left edge the memory modules Mi would be replaced by lines which wrap back around to PEi, similar to the Omega network shown below.

Crossbar switches are too expensive for large-scale systems, but are useful in some small systems. The 16-CPU Sun Microsystems Enterprise 10000 system includes a 16x16 crossbar.

### 2.2.3.2  Omega (or Delta) Interconnects

These are multistage networks similar to crossbars, but with fewer paths. Here is an example of a NUMA 8x8 system:



Recall that each PE is a processor/memory pair. PE3, for instance, consists of P3 and M3.

Note the fact that at the third stage of the network (top of picture), the outputs are routed back to the PEs, each of which consists of a processor and a memory module.[3]

At each network node (the nodes are the three rows of rectangles), the output routing is done by destination bit. Let's number the stages here 0, 1 and 2, starting from the bottom stage, number the nodes within a stage 0, 1, 2 and 3 from left to right, number the PEs from 0 to 7, left to right, and number the bit positions in a destination address 0, 1 and 2, starting from the most significant bit. Then at stage i, bit i of the destination address is used to determine routing, with a 0 meaning routing out the left output, and 1 meaning the right one.

Say P2 wishes to read from M5. It sends a read-request packet, including 5 = 101 as its destination address, to the switch in stage 0, node 1. Since the first bit of 101 is 1, that means that this switch will route the packet out its right-hand output, sending it to the switch in stage 1, node 3. The latter switch will look at the next bit in 101, a 0, and thus route the packet out its left output, to the switch in stage 2, node 2. Finally, that switch will look at the last bit, a 1, and output out its right-hand output, sending it to PE5, as desired. M5 will process the read request, and send a packet back to PE2, along the same

Again, if two packets at a node want to go out the same output, one must get priority (let's say it is the one

---

[2]For safety's sake, i.e. fault tolerance, even writes are typically acknowledged in multiprocessor systems.

[3]The picture may be cut off somewhat at the top and left edges. The upper-right output of the rectangle in the top row, leftmost position should connect to the dashed line which leads down to the second PE from the left. Similarly, the upper-left output of that same rectangle is a dashed lined, possibly invisible in your picture, leading down to the leftmost PE.

from the left input).

Here is how the more general case of N = $2^n$ PEs works. Again number the rows of switches, and switches within a row, as above. So, $S_{ij}$ will denote the switch in the i-th row from the bottom and j-th column from the left (starting our numbering with 0 in both cases). Row i will have a total of N input ports $I_{ik}$ and N output ports $O_{ik}$, where k = 0 corresponds to the leftmost of the N in each case. Then if row i is not the last row ($i < n - 1$), $O_{ik}$ will be connected to $I_{jm}$, where j = i+1 and

$$m = (2k + \lfloor (2k)/N \rfloor) \, mod \, N \tag{2.1}$$

If row i is the last row, then $O_{ik}$ will be connected to, PE k.

### 2.2.4   Comparative Analysis

In the world of parallel architectures, a key criterion for a proposed feature is **scalability**, meaning how well the feature performs as we go to larger and larger systems. Let n be the system size, either the number of processors and memory modules, or the number of PEs. Then we are interested in how fast the latency, bandwidth and cost grow with n:

| criterion | bus | Omega | crossbar |
|---|---|---|---|
| latency | O(1) | $O(\log_2 n)$ | O(n) |
| bandwidth | O(1) | O(n) | O(n) |
| cost | O(1) | $O(n \log_2 n)$ | $O(n^2)$ |

Let us see where these expressions come from, beginning with a bus: No matter how large n is, the time to get from, say, a processor to a memory module will be the same, thus O(1). Similarly, no matter how large n is, only one communication can occur at a time, thus again O(1).[4]

Again, we are interested only in "O( )" measures, because we are only interested in growth rates as the system size n grows. For instance, if the system size doubles, the cost of a crossbar will quadruple; the $O(n^2)$ cost measure tells us this, with any multiplicative constant being irrelevant.

For Omega networks, it is clear that $log_2 n$ network rows are needed, hence the latency value given. Also, each row will have n/2 switches, so the number of network nodes will be O(n $log_2 n$). This figure then gives the cost (in terms of switches, the main expense here). It also gives the bandwidth, since the maximum number of simultaneous transmissions will occur when all switches are sending at once.

Similar considerations hold for the crossbar case.

---

[4] Note that the '1' in "O(1)" does not refer to the fact that only one communication can occur at a time. If we had, for example, a two-bus system, the bandwidth would still be O(1), since multiplicative constants do not matter. What O(1) means, again, is that as n grows, the bandwidth stays at a multiple of 1, i.e. stays constant.

The crossbar's big advantage is that it is guaranteed that n packets can be sent simultaneously, providing they are to distinct destinations.

That is <u>not</u> true for Omega-networks. If for example, PE0 wants to send to PE3, and at the same time PE4 wishes to sent to PE2, the two packets will clash at the leftmost node of stage 1, where the packet from PE0 will get priority.

On the other hand, a crossbar is very expensive, and thus is dismissed out of hand in most modern systems. Note, though, that an equally troublesom aspect of crossbars is their high latency value; this is a big drawback when the system is not heavily loaded.

The bottom line is that Omega-networks amount to a compromise between buses and crossbars, and for this reason have become popular.

## 2.3   Test-and-Set Type Instructions

Consider a bus-based system. In addition to whatever memory read and memory write instructions the processor included, there would also be a TAS instruction.[5] This instruction would control a TAS pin on the processor chip, and the pin in turn would be connected to a TAS line on the bus.

Applied to a location L in memory and a register R, say, TAS does the following:

```
copy L to R
if R is 0 then write 1 to L
```

And most importantly, these operations are done in an **atomic** manner; no bus transactions by other processors may occur between the two steps.

The TAS operation is applied to variables used as **locks**. Let's say that 1 means locked and 0 unlocked. Then the guarding of a critical section C by a lock variable L would be done by having the following code in the program being run:

```
TRY:  TAS R,L
      JNZ TRY
   C: ...   ; start of critical section
      ...
      ...   ; end of critical section
      MOV L,0  ; unlock
```

where of course JNZ is a jump-if-nonzero instruction, and we are assuming that the copying from the Memory Data Register to R results in the processor N and Z flags (condition codes) being affected.

---

[5]This discussion is for a mythical machine, but any real system works in this manner.

On Pentium machines, the LOCK prefix can be used to get atomicity for certain instructions.[6] For example,

```
lock add $2, x
```

would add the constant 2 to the memory location labeled **x** in an atomic manner.

The LOCK prefix locks the bus for the entire duration of the instruction. Note that the ADD instruction here involves two memory transactions—one to read the old value of **x**, and the second the write the new, incremented value back to **x**. So, we are locking for a rather long time, but the benefits can be huge.

A good example of this kind of thing would be our program **PrimesThreads.c** in Chapter 1, where our critical section consists of adding 2 to **nextbase**. There we surrounded the add-2 code by Pthreads lock and unlock operations. These involve system calls, which are very time consuming, involving hundreds of machine instructions. Compare that to the one-instruction solution above! The very heavy overhead of pthreads would be thus avoided.

In crossbar or $\Omega$-network systems, some 2-bit field in the packet must be devoted to transaction type, say 00 for Read, 01 for Write and 10 for TAS. In a sytem with 16 CPUs and 16 memory modules, say, the packet might consist of 4 bits for the CPU number, 4 bits for the memory module number, 2 bits for the transaction type, and 32 bits for the data (for a write, this is the data to be written, while for a read, it would be the requested value, on the trip back from the memory to the CPU).

But note that the atomicity here is best done at the memory, i.e. some hardware should be added at the memory so that TAS can be done; otherwise, an entire processor-to-memory path (e.g. the bus in a bus-based system) would have to be locked up for a fairly long time, obstructing even the packets which go to other memory modules.

The Intel LOCK prefix just locks the bus, so you can see that Intel doesn't want to get into the large multi-processor business in a big way. Their commercial projections are probably right.

There are many variations of test-and-set, so don't expect that all processors will have an instruction with this name, but they all will have some kind of synchronization instruction like it.

Note carefully that in many settings it may not be crucial to get the most up-to-date value of a variable. For example, a program may have a data structure showing work to be done. Some processors occasionally add work to the queue, and others take work from the queue. Suppose the queue is currently empty, and a processor adds a task to the queue, just as another processor is checking the queue for work. As will be seen later, it is possible that even though the first processor has written to the queue, the new value won't be visible to other processors for some time. But the point is that if the second processor does not see work in the queue (even though the first processor has put it there), the program will still work correctly, albeit with some performance loss.

---

[6]The instructions ADD, ADC, AND, BTC, BTR, BTS, CMPXCHG, DEC, INC, NEG, NOT, OR, SBB, SUB, XOR, XADD. Also, XCHG asserts the LOCK# bus signal even if the LOCK prefix is specified. Locking only applies to these instructions in forms in which there is an operand in memory.

## 2.4    Cache Issues

### 2.4.1    Cache Coherency

Consider, for example, a bus-based system.  Relying purely on TAS for interprocessor synchronization would be unthinkable: As each processor contending for a lock variable spins in the loop shown above, it is adding tremendously to bus traffic.

An answer is to have caches at each processor.[7]  These will to store copies of the values of lock variables. (Of course, non-lock variables are stored too.  However, the discussion here will focus on effects on lock variables.)  The point is this:  Why keep looking at a lock variable L again and again, using up the bus bandwidth?  L may not change value for a while, so why not keep a copy in the cache, avoiding use of the bus?

The answer of course is that eventually L <u>will</u> change value, and this causes some delicate problems. Say for example that processor P5 wishes to enter a critical section guarded by L, and that processor P2 is already in there.  During the time P2 is in the critical section, P5 will spin around, always getting the same value for L (1) from C5, P5's cache. When P2 leaves the critical section, P2 will set L to 0—and now C5's copy of L will be incorrect.  This is the **cache coherency problem**, inconsistency between caches.

A number of solutions have been devised for this problem.  For bus-based systems, **snoopy** protocols of various kinds are used, with the word "snoopy" referring to the fact that all the caches monitor ("snoop on") the bus, watching for transactions made by <u>other</u> caches.

The most common protocols are the **invalidate** and **update** types.  This relation between these two is somewhat analogous to the relation between **write-back** and **write-through** protocols for caches in uniprocessor systems:

- Under an invalidate protocol, when a processor writes to a variable in a cache, it first (i.e. before actually doing the write) tells each other cache to mark as invalid its cache line (if any) which contains a copy of the variable.[8]  Those caches will be updated only later, the next time their processors need to access this cache line.

- For an update protocol, the processor which writes to the variable tells all other caches to immediately update their cache lines containing copies of that variable with the new value.

Let's look at an outline of how one implementation (many variations exist) of an invalidate protocol would operate:

---

[7]The reader may wish to review the basics of caches.  See for example `http://heather.cs.ucdavis.edu/~matloff/50/PLN/CompOrganization.pdf`.

[8]We will follow commonly-used terminology here, distinguishing between a *cache line* and a *memory block*. Memory is divided in blocks, some of which have copies in the cache.  The cells in the cache are called *cache lines*.  So, at any given time, a given cache line is either empty or contains a copy (valid or not) of some memory block.

In the scenario outlined above, when P2 leaves the critical section, it will write the new value 0 to L. Under the invalidate protocol, P2 will post an invalidation message on the bus. All the other caches will notice, as they have been monitoring the bus. They then mark their cached copies of the line containing L as invalid.

Now, the next time P5 executes the TAS instruction—which will be very soon, since it is in the loop shown above—P5 will find that the copy of L in C5 is invalid. It will respond to this cache miss by going to the bus, and requesting P2 to supply the "real" (and valid) copy of the line containing L.

But there's more. Suppose that all this time P6 had also been executing the loop shown above, along with P5. Then P5 and P6 may have to contend with each other. Say P6 manages to grab possession of the bus first.[9] P6 then executes the TAS again, which finds L = 0 and changes L back to 1. P6 then relinquishes the bus, and enters the critical section. Note that in changing L to 1, P6 also sends an invalidate signal to all the other caches. So, when P5 tries its execution of the TAS again, it will have to ask P6 to send a valid copy of the block. P6 does so, but L will be 1, so P5 must resume executing the loop. P5 will then continue to use its valid local copy of L each time it does the TAS, until P6 leaves the critical section, writes 0 to L, and causes another cache miss at P5, etc.

At first the update approach seems obviously superior, and actually, if our shared, cacheable[10] variables were only lock variables, this might be true.

But consider a shared, cacheable vector. Suppose the vector fits into one block, and that we write to each vector element sequentially. Under an update policy, we would have to send a new message on the bus/network for each component, while under an invalidate policy, only one message (for the first component) would be needed. If during this time the other processors do not need to access this vector, all those update messages, and the bus/network bandwidth they use, would be wasted.

Or suppose for example we have code like

```
Sum += X[I];
```

in the middle of a **for** loop. Under an update protocol, we would have to write the value of Sum back many times, even though the other processors may only be interested in the final value when the loop ends. (This would be true, for instance, if the code above were part of a critical section.)

Thus the invalidate protocol works well for some kinds of code, while update works better for others. The CPU designers must try to anticipate which protocol will work well across a broad mix of applications.[11]

Now, how is cache coherency handled in non-bus shared-memory systems, say crossbars? Here the problem is more complex. Think back to the bus case for a minute: The very feature which was the biggest negative feature of bus systems—the fact that there was only one path between components made bandwidth very

---

[9]Again, remember that ordinary bus arbitration methods would be used.

[10] Many modern processors, including Pentium and MIPS, allow the programmer to mark some blocks as being noncacheable.

[11]Some protocols change between the two modes dynamically.

limited—is a very <u>positive</u> feature in terms of cache coherency, because it makes <u>broadcast</u> very easy: Since everyone is attached to that single pathway, sending a message to all of them costs no more than sending it to just one—we get the others for free. That's no longer the case for multipath systems. In such systems, extra copies of the message must be created for each path, adding to overall traffic.

A solution is to send messages only to "interested parties." In **directory-based** protocols, a list is kept of all caches which currently have valid copies of all blocks. In one common implementation, for example, while P2 is in the critical section above, it would be the **owner** of the block containing L. (Whoever is the latest node to write to L would be considered its current owner.) It would maintain a directory of all caches having valid copies of that block, say C5 and C6 in our story here. As soon as P2 wrote to L, it would then send either invalidate or update packets (depending on which type was being used) to C5 and C6 (and <u>not</u> to other caches which didn't have valid copies).

There would also be a directory at the memory, listing the current owners of all blocks. Say for example P0 now wishes to "join the club," i.e. tries to access L, but does not have a copy of that block in its cache C0. C0 will thus not be listed in the directory for this block. So, now when it tries to access L and it will get a cache miss. P0 must now consult the **home** of L, say P14. The home might be determined by L's location in main memory according to high-order interleaving; it is the place where the main-memory version of L resides. A table at P14 will inform P0 that P2 is the current owner of that block. P0 will then send a message to P2 to add C0 to the list of caches having valid copies of that block. Similarly, a cache might "resign" from the club, due to that cache line being replaced, e.g. in a LRU setting, when some other cache miss occurs.

### 2.4.2   Example: the MESI Cache Coherency Protocol

Many types of cache coherency protocols have been proposed and used, some of them quite complex. A relatively simple one for snoopy bus systems which is widely used is MESI, which for example is the protocol used in the Pentium series.

MESI is an invalidate protocol for bus-based systems. Its name stands for the four states a given cache line can be in for a given CPU:

- Modified

- Exclusive

- Shared

- Invalid

Note that *each memory block* has such a state at *each cache*. For instance, block 88 may be in state S at P5's and P12's caches but in state I at P1's cache.

Here is a summary of the meanings of the states:

| state | meaning |
|:-:|:-:|
| M | written to more than once; no other copy valid |
| E | valid; no other cache copy valid; memory copy valid |
| S | valid; at least one other cache copy valid |
| I | invalid (block either not in the cache or present but incorrect) |

Following is a summary of MESI state changes.[12] When reading it, keep in mind again that there is a separate state for each cache/memory block combination.

In addition to the terms **read hit**, **read miss**, **write hit**, **write miss**, which you are already familiar with, there are also **read snoop** and **write snoop**. These refer to the case in which our CPU observes on the bus a block request by another CPU that has attempted a read or write action but encountered a miss in its own cache; if our cache has a valid copy of that block, we must provide it to the requesting CPU (and in some cases to memory).

So, here are various events and their corresponding state changes:

**If our CPU does a read:**

| present state | event | new state |
|:-:|:-:|:-:|
| M | read hit | M |
| E | read hit | E |
| S | read hit | S |
| I | read miss; no valid cache copy at any other CPU | E |
| I | read miss; at least one valid cache copy in some other CPU | S |

**If our CPU does a memory write:**

| present state | event | new state |
|:-:|:-:|:-:|
| M | write hit; do not put invalidate signal on bus; do not update memory | M |
| E | same as M above | M |
| S | write hit; put invalidate signal on bus; update memory | E |
| I | write miss; update memory but do nothing else | I |

**If our CPU does a read or write snoop:**

---

[12]See *Pentium Processor System Architecture*, by D. Anderson and T. Shanley, Addison-Wesley, 1995. We have simplified the presentation here, by eliminating certain programmable options.

| present state | event | newstate |
|---|---|---|
| M | read snoop; write line back to memory, picked up by other CPU | S |
| M | write snoop; write line back to memory, signal other CPU now OK to do its write | I |
| E | read snoop; put shared signal on bus; no memory action | S |
| E | write snoop; no memory action | I |
| S | read snoop | S |
| S | write snoop | I |
| I | any snoop | I |

Note that a write miss does NOT result in the associated block being brought in from memory.

Example: Suppose a given memory block has state M at processor A but has state I at processor B, and B attempts to write to the block. B will see that its copy of the block is invalid, so it notifies the other CPUs via the bus that it intends to do this write. CPU A sees this announcement, tells B to wait, writes its own copy of the block back to memory, and then tells B to go ahead with its write. The latter action means that A's copy of the block is not correct anymore, so the block now has state I at A. B's action does not cause loading of that block from memory to its cache, so the block still has state I at B.

### 2.4.3   The Problem of "False Sharing"

Consider the C declaration

```
int W,Z;
```

Since **W** and **Z** are declared adjacently, most compilers will assign them contiguous memory addresses. Thus, unless one of them is at a memory block boundary, when they are cached they will be stored in the same cache line. Suppose the program writes to **Z**, and our system uses an invalidate protocol. Then **W** will be considered invalid at the other processors, even though its values at those processors' caches are correct. This is the **false sharing** problem, alluding to the fact that the two variables are sharing a cache line even though they are not related.

This can have very adverse impacts on performance. If for instance our variable **W** is now written to, then **Z** will suffer unfairly, as its copy in the cache will be considered invalid even though it is perfectly valid. This can lead to a "ping-pong" effect, in which alternate writing to two variables leads to a cyclic pattern of coherency transactions.

## 2.5   Memory-Access Consistency Policies

Though the word *consistency* in the title of this section may seem to simply be a synonym for *coherency* from the last section, and though there actually is some relation, the issues here are quite different. In this

case, it is a timing issue: After one processor changes the value of a shared variable, when will that value be visible to the other processors?

There are various reasons why this is an issue. For example, many processors, especially in multiprocessor systems, have **write buffers**, which save up writes for some time before actually sending them to memory. (For the time being, let's suppose there are no caches.) The goal is to reduce memory access costs. Sending data to memory in groups is generally faster than sending one at a time, as the overhead of, for instance, acquiring the bus is amortized over many accesses. Reads following a write may proceed, without waiting for the write to get to memory, except for reads to the same address. So in a multiprocessor system in which the processors use write buffers, there will often be some delay before a write actually shows up in memory.

A related issue is that operations may occur, or appear to occur, out of order. As noted above, a read which follows a write in the program may execute before the write is sent to memory. Also, in a multiprocessor system with multiple paths between processors and memory modules, two writes might take different paths, one longer than the other, and arrive "out of order." In order to simplify the presentation here, we will focus on the case in which the problem is due to write buffers, though.

The designer of a multiprocessor system must adopt some **consistency model** regarding situations like this. The above discussion shows that the programmer must be made aware of the model, or risk getting incorrect results. Note also that different consistency models will give different levels of performance. The "weaker" consistency models make for faster machines but require the programmer to do more work.

The strongest consistency model is Sequential Consistency. It essentially requires that memory operations done by one processor are observed by the other processors to occur in the same order as executed on the first processor. Enforcement of this requirement makes a system slow, and it has been replaced on most systems by weaker models.

One such model is **release consistency**. Here the processors' instruction sets include instructions ACQUIRE and RELEASE. Execution of an ACQUIRE instruction at one processor involves telling all other processors to flush their write buffers. However, the ACQUIRE won't execute until pending RELEASEs are done. Execution of a RELEASE basically means that you are saying, "I'm done writing for the moment, and wish to allow other processors to see what I've written." An ACQUIRE waits for all pending RELEASEs to complete before it executes.[13]

A related model is **scope consistency**. Say a variable, say **Sum**, is written to within a critical section guarded by LOCK and UNLOCK instructions. Then under scope consistency any changes made by one processor to **Sum** within this critical section would then be visible to another processor when the latter next enters this critical section. The point is that memory update is postponed until it is actually needed. Also, a barrier operation (again, executed at the hardware level) forces all pending memory writes to complete.

All modern processors include instructions which implement consistency operations. For example, Sun Microsystems' SPARC has a MEMBAR instruction. If used with a STORE operand, then all pending writes

---

[13]There are many variants of all of this, especially in the software distibuted shared memory realm, to be discussed later.

at this processor will be sent to memory. If used with the LOAD operand, all writes will be made visible to this processor.

Now, how does cache coherency fit into all this? There are many different setups, but for example let's consider a design in which there is a write buffer between each processor and its cache. As the processor does more and more writes, the processor saves them up in the write buffer. Eventually, some programmer-induced event, e.g. a MEMBAR instruction,[14] will cause the buffer to be flushed. Then the writes will be sent to "memory"—actually meaning that they go to the cache, and then possibly to memory.

The point is that (in this type of setup) before that flush of the write buffer occurs, the cache coherency system is quite unaware of these writes. Thus the cache coherency operations, e.g. the various actions in the MESI protocol, won't occur until the flush happens.

To make this notion concrete, again consider the example with **Sum** above, and assume release or scope consistency. The CPU currently executing that code (say CPU 5) writes to **Sum**, which is a memory operation—it affects the cache and thus eventually the main memory—but that operation will be invisible to the cache coherency protocol for now, as it will only be reflected in this processor's write buffer. But when the unlock is finally done (or a barrier is reached), the write buffer is flushed and the writes are sent to this CPU's cache. That then triggers the cache coherency operation (depending on the state). The point is that the cache coherency operation would occur only now, not before.

What about reads? Suppose another processor, say CPU 8, does a read of **Sum**, and that page is marked invalid at that processor. A cache coherency operation will then occur. Again, it will depend on the type of coherency policy and the current state, but in typical systems this would result in **Sum**'s cache block being shipped to CPU 8 from whichever processor the cache coherency system thinks has a valid copy of the block. That processor may or may not be CPU 5, but even if it is, that block won't show the recent change made by CPU 5 to **Sum**.

The analysis above assumed that there is a write buffer between each processor and its cache. There would be a similar analysis if there were a write buffer between each cache and memory.

Note once again the performance issues. Instructions such as ACQUIRE or MEMBAR will use a substantial amount of interprocessor communication bandwidth. A consistency model must be chosen carefully by the system designer, and the programmer must keep the communication costs in mind in developing the software.

The recent Pentium models use Sequential Consistency, with any write done by a processor being immediately sent to its cache as well.

---

[14]We call this "programmer-induced," since the programmer will include some special operation in her C/C++ code which will be translated to MEMBAR.

## 2.6 Fetch-and-Add and Packet-Combining Operations

Another form of interprocessor synchronization is a **fetch-and-add** (FA) instruction. The idea of FA is as follows. For the sake of simplicity, consider code like

```
LOCK(K);
Y = X++;
UNLOCK(K);
```

Suppose our architecture's instruction set included an F&A instruction. It would add 1 to the specified location in memory, and return the old value (to **Y**) that had been in that location before being incremented. And all this would be an atomic operation.

We would then replace the code above by a library call, say,

```
FETCH_AND_ADD(X,1);
```

The C code above would compile to, say,

```
F&A X,R,1
```

where **R** is the register into which the old (pre-incrementing) value of **X** would be returned.

There would be hardware adders placed at each memory module. That means that the whole operation could be done in one round trip to memory. Without F&A, we would need two round trips to memory just for the

```
X++;
```

(we would load **X** into a register in the CPU, increment the register, and then write it back to **X** in memory), and then the LOCK() and UNLOCK() would need trips to memory too. This could be a huge time savings, especially for long-latency interconnects.

In addition to read and write operations being specifiable in a network packet, an F&A operation could be specified as well (a 2-bit field in the packet would code which operation was desired). Again, there would be adders included at the memory modules, i.e. the addition would be done at the memory end, not at the processors. When the F&A packet arrived at a memory module, our variable **X** would have 1 added to it, while the old value would be sent back in the return packet (and put into R).

Another possibility for speedup occurs if our system uses a multistage interconnection network such as a crossbar. In that situation, we can design some intelligence into the network nodes to do **packet combining**: Say more than one CPU is executing an F&A operation at about the same time for the same variable **X**.

Then more than one of the corresponding packets may arrive at the same network node at about the same time. If each one requested an incrementing of **X** by 1, the node can replace the two packets by one, with an increment of 2. Of course, this is a delicate operation, and we must make sure that different CPUs get different return values, etc.

## 2.7   Multicore Chips

A recent trend has been to put several CPUs on one chip, termed a **multicore** chip. As of March 2008, dual-core chips are common in personal computers, and quad-core machines are within reach of the budgets of many people. Just as the invention of the integrated circuit revolutionized the computer industry by making computers affordable for the average person, multicore chips will undoubtedly revolutionize the world of parallel programming.

Note that we can build an SMP or NUMA system from such chips. We could, for instance, put four dual-core chips on a bus and thus create an 8-processor SMP system.

A typical dual-core setup might have the two CPUs sharing a common L2 cache, with each CPU having its own L3 cache. The chip may interface to the bus or interconnect network of via an L1 cache.

## 2.8   Illusion of Shared-Memory through Software

### 2.8.0.1   Software Distributed Shared Memory

There are also various shared-memory software packages that run on message-passing hardware such as NOWs, called **software distributed shared memory** (SDSM) systems. Since the platforms do not have any physically shared memory, the shared-memory view which the programmer has is just an illusion. But that illusion is very useful, since the shared-memory paradigm is believed to be the easier one to program in. Thus SDSM allows us to have "the best of both worlds"—the convenience of the shared-memory world view with the inexpensive cost of some of the message-passing hardware systems, particularly networks of workstations (NOWs).

SDSM itself is divided into two main approaches, the **page-based** and **object-based** varieties. The page-based approach is generally considered clearer and easier to program in, and provides the programmer the "look and feel" of shared-memory programming better than does the object-based type.[15] We will discuss only the page-based approach here. The most popular SDSM system today is the page-based Treadmarks (Rice University). Another excellent page-based system is JIAJIA (Academy of Sciences, China).

To illustrate how page-paged SDSMs work, consider the line of JIAJIA code

---

[15]The term *object-based* is not related to the term *object-oriented programming*.

```
Prime = (int *) jia_alloc(N*sizeof(int));
```

The function **jia alloc()** is part of the JIAJIA library, **libjia.a**, which is linked to one's application program during compilation.

At first this looks a little like a call to the standard **malloc()** function, setting up an array **Prime** of size **N**. In fact, it does indeed allocate some memory. Note that each node in our JIAJIA group is executing this statement, so each node allocates some memory at that node. Behind the scenes, not visible to the programmer, each node will then have its own copy of **Prime**.

However, JIAJIA sets things up so that when one node later accesses this memory, for instance in the statement

```
Prime[I] = 1;
```

this action will eventually trigger a network transaction (not visible to the programmer) to the other JIAJIA nodes.[16] This transaction will then update the copies of **Prime** at the other nodes.[17]

How is all of this accomplished? It turns out that it relies on a clever usage of the nodes' virtual memory (VM) systems. To understand this, let's review how VM systems work.

Suppose a variable **X** has the virtual address 1200, i.e. **&X = 1200**. The actual physical address may be, say, 5000. When the CPU executes a machine instruction that specifies access to 1200, the CPU will do a lookup on the **page table**, and find that the true location is 5000, and then access 5000. On the other hand, **X** may not be **resident** in memory at all, in which case the page table will say so. If the CPU finds that **X** is nonresident, it will cause an internal interrupt, which in turn will cause a jump to the operating system (OS). The OS will then read **X** in from disk,[18] place it somewhere in memory, and then update the page table to show that **X** is now someplace in memory. The OS will then execute a return from interrupt instruction,[19], and the CPU will restart the instruction which triggered the page fault.

Here is how this is exploited to develop SDSMs on Unix systems. The SDSM will call a system function such as **mprotect()**. This allows the SDSM to deliberately mark a page as nonresident (even if the page *is* resident). Basically, anytime the SDSM knows that a node's local copy of a variable is invalid, it will mark the page containing that variable as nonresident. Then, the next time the program at this node tries to access that variable, a page fault will occur.

As mentioned in the review above, normally a page fault causes a jump to the OS. However, technically any page fault in Unix is handled as a signal, specifically SIGSEGV. Recall that Unix allows the programmer to write his/her own signal handler for any signal type. In this case, that means that the programmer—meaning

---

[16]There are a number of important issues involved with this word *eventually*, as we will see later.

[17]The update may not occur immediately. More on this later.

[18]Actually, it will read the entire page containing **X** from disk, but to simplify language we will just say **X** here.

[19]E.g. **iret** on Pentium chips.

the people who developed JIAJIA or any other page-based SDSM—writes his/her own page fault handler, which will do the necessary network transactions to obtain the latest valid value for **X**.

Note that although SDSMs are able to create an illusion of almost all aspects of shared memory, it really is not possible to create the illusion of shared pointer variables. For example on shared memory hardware we might have a variable like **P**:

```
int Y,*P;
...
...
P = &Y;
...
```

There is no simple way to have a variable like **P** in an SDSM. This is because a pointer is an address, and each node in an SDSM has its own memory separate address space. The problem is that even though the underlying SDSM system will keep the various copies of **Y** at the different nodes consistent with each other, **Y** will be at a potentially different address on each node.

All SDSM systems must deal with a software analog of the cache coherency problem. Whenever one node modifies the value of a shared variable, that node must notify the other nodes that a change has been made. The designer of the system must choose between update or invalidate protocols, just as in the hardware case.[20]  Recall that in non-bus-based shared-memory multiprocessors, one needs to maintain a directory which indicates at which processor a valid copy of a shared variable exists. Again, SDSMs must take an approach similar to this.

Similarly, each SDSM system must decide between sequential consistency, release consistency etc. More on this later.

Note that in the NOW context the internode communication at the SDSM level is typically done by TCP/IP network actions. Treadmarks uses UDP, which is faster than TCP. but still part of the slow TCP/IP protocol suite. TCP/IP was simply not designed for this kind of work. Accordingly, there have been many efforts to use more efficient network hardware and software. The most popular of these is the Virtual Interface Architecture (VIA).

Not only are coherency actions more expensive in the NOW SDSM case than in the shared-memory hardware case due to network slowness, there is also expense due to granularity. In the hardware case we are dealing with cache blocks, with a typical size being 512 bytes. In the SDSM case, we are dealing with pages, with a typical size being 4096 bytes. The overhead for a cache coherency transaction can thus be large.

---

[20]Note, though, that we are not actually dealing with a cache here. Each node in the SDSM system will have a cache, of course, but a node's cache simply stores parts of that node's set of pages. The coherency across nodes is across pages, not caches. We must insure that a change made to a given page is eventually propropagated to pages on other nodes which correspond to this one.

### 2.8.0.2 Case Study: JIAJIA

**Programmer Interface**

We will not go into detail on JIAJIA programming here. There is a short tutorial on JIAJIA at `http://heather.cs.ucdavis.edu/˜matloff/jiajia.html`, but here is an overview:

- One writes in C/C++ (or FORTRAN), making calls to the JIAJIA library, which is linked in upon compilation.

- The library calls include standard shared-memory operations for lock, unlock, barrier, processor number, etc., plus some calls aimed at improving performance.

Following is a JIAJIA example program, performing Odd/Even Transposition Sort. This is a variant on Bubble Sort, sometimes useful in parallel processing contexts.[21] The algorithm consists of n phases, in which each processor alternates between trading with its left and right neighbors.

```
1  // JIAJIA example program:  Odd-Even Tranposition Sort
2
3  // array is of size n, and we use n processors; this would be more
4  // efficient in a "chunked" versions, of course (and more suited for a
5  // message-passing context anyway)
6
7  #include <stdio.h>
8  #include <stdlib.h>
9  #include <jia.h>  // required include; also must link via -ljia
10
11 // pointer to shared variable
12 int *x;  // array to be sorted
13
14 int n,  // range to check for primeness
15     debug;  // 1 for debugging, 0 else
16
17 // if first arg is bigger, then replace it by the second
18 void cpsmaller(int *p1,int *p2)
19 {  int tmp;
20        if (*p1 > *p2)  *p1 = *p2;
21 }
22
23 // if first arg is smaller, then replace it by the second
24 void cpbigger(int *p1,int *p2)
25 {  int tmp;
26        if (*p1 < *p2)  *p1 = *p2;
27 }
28
29 // does sort of m-element array y
30 void oddeven(int *y, int m)
31 {  int i,left=jiapid-1,right=jiapid+1,newval;
```

---

[21]Though, as mentioned in the comments, it is aimed more at message-passing contexts.

```
32      for (i=0; i < m; i++)  {
33         if ((i+jiapid)%2 == 0)  {
34            if (right < m)
35               if (y[jiapid] > y[right]) newval = y[right];
36         }
37         else  {
38            if (left >= 0)
39               if (y[jiapid] < y[left]) newval = y[left];
40         }
41         jia_barrier();
42         if ((i+jiapid)%2 == 0 && right < m || (i+jiapid)%2 == 1 && left >= 0)
43               y[jiapid] = newval;
44         jia_barrier();
45      }
46   }
47
48   main(int argc, char **argv)
49   {  int i,mywait=0;
50      jia_init(argc,argv);  // required init call
51      // get command-line arguments (shifted for nodes > 0)
52      if (jiapid == 0)   {
53         n = atoi(argv[1]);
54         debug = atoi(argv[2]);
55      }
56      else  {
57         n = atoi(argv[2]);
58         debug = atoi(argv[3]);
59      }
60      jia_barrier();
61      // create a shared array x of length n
62      x = (int *) jia_alloc(n*sizeof(int));
63      // barrier recommended after allocation
64      jia_barrier();
65      // node 0 gets simple test array from command-line
66      if (jiapid == 0)   {
67         for (i = 0; i < n; i++)
68            x[i] = atoi(argv[i+3]);
69      }
70      jia_barrier();
71      if (debug && jiapid == 0)
72         while (mywait == 0)  { ; }
73      jia_barrier();
74      oddeven(x,n);
75      if (jiapid == 0)   {
76         printf("\nfinal array\n");
77         for (i = 0; i < n; i++)
78            printf("%d\n",x[i]);
79      }
80      jia_exit();
81   }
```

**System Workings**

JIAJIA's main characteristics as an SDSM are:

- page-based

- scope consistency

- home-based

- multiple writers

Let's take a look at these.

As mentioned earlier, one first calls **jia_alloc()** to set up one's shared variables. Note that this will occur at each node, so there are multiple copies of each variable; the JIAJIA system ensures that these copies are consistent with each other, though of course subject to the laxity afforded by scope consistency.

Recall that under scope consistency, a change made to a shared variable at one processor is guaranteed to be made visible to another processor if the first processor made the change between lock/unlock operations and the second processor accesses that variable between lock/unlock operations on that same lock.[22]

Each page—and thus each shared variable—has a **home** processor. If another processor writes to a page, then later when it reaches the unlock operation it must send all changes it made to the page back to the home node. In other words, the second processor calls **jia_unlock()**, which sends the changes to its sister invocation of **jia_unlock()** at the home processor.[23] Say later a third processor calls **jia_lock()** on that same lock, and then attempts to read a variable in that page. A page fault will occur at that processor, resulting in the JIAJIA system running, which will then obtain that page from the first processor.

Note that all this means the JIAJIA system at each processor must maintain a page table, listing where each home page resides.[24] At each processor, each page has one of three states: Invalid, Read-Only, Read-Write. State changes, though, are reported when lock/unlock operations occur. For example, if CPU 5 writes to a given page which had been in Read-Write state at CPU 8, the latter will not hear about CPU 5's action until some CPU does a lock. This CPU need not be CPI 8. When one CPU does a lock, it must coordinate with all other nodes, at which time state-change messages will be piggybacked onto lock-coordination messages.

Note also that JIAJIA allows the programmer to specify which node should serve as the home of a variable, via one of several forms of the **jia_alloc()** call. The programmer can then tailor his/her code accordingly. For example, in a matrix problem, the programmer may arrange for certain rows to be stored at a given node, and then write the code so that most writes to those rows are done by that processor.

The general principle here is that writes performed at one node can be made visible at other nodes on a "need to know" basis. If for instance in the above example with CPUs 5 and 8, CPU 2 does not access this

---

[22]Writes will also be propagated at barrier operations, but two successive arrivals by a processor to a barrier can be considered to be a lock/unlock pair, by considering a departure from a barrier to be a "lock," and considering reaching a barrier to be an "unlock." So, we'll usually not mention barriers separately from locks in the remainder of this subsection.

[23]The set of changes is called a **diff**, remiscent of the Unix file-compare command. A copy, called a **twin**, had been made of the original page, which now will be used to produce the diff. This has substantial overhead. The Treadmarks people found that it took 167 microseconds to make a twin, and as much as 686 microseconds to make a diff.

[24]In JIAJIA, that location is normally fixed, but JIAJIA does include advanced programmer options which allow the location to migrate.

page, it would be wasteful to send the writes to CPU 2, or for that matter to even inform CPU 2 that the page had been written to. This is basically the idea of all non-Sequential consistency protocols, even though they differ in approach and in performance for a given application.

JIAJIA allows multiple writers of a page. Suppose CPU 4 and CPU 15 are simultaneously writing to a particular page, and the programmer has relied on a subsequent barrier to make those writes visible to other processors.[25] When the barrier is reached, each will be informed of the writes of the other.[26] Allowing multiple writers helps to reduce the performance penalty due to false sharing.

---

[25]The only other option would be to use lock/unlock, but then their writing would not be simultaneous.

[26]If they are writing to the same variable, not just the same page, the programmer would use locks instead of a barrier, and the situation would not arise.

# Chapter 3

# Introduction to OpenMP

OpenMP has become the *de facto* standard for shared-memory programming.

## 3.1   Overview

OpenMP has become the environment of choice for many, if not most, practitioners of shared-memory parallel programming. It consists of a set of directives which are added to one's C/C++/FORTRAN code that manipulate threads, without the programmer him/herself having to deal with the threads directly. This way we get "the best of both worlds"—the true parallelism of (nonpreemptive) threads and the pleasure of avoiding the annoyances of threads programming.

Most OpenMP constructs are expressed via **pragmas**, i.e. directives. The syntax is

```
#pragma omp ......
```

The number sign must be the first nonblank character in the line.

## 3.2   Running Example

The following example, implementing Dijkstra's shortest-path graph algorithm, will be used throughout this tutorial, with various OpenMP constructs being illustrated later by modifying this code:

```
1  // Dijkstra.c
2
```

```
3   // OpenMP example program:  Dijkstra shortest-path finder in a
4   // bidirectional graph; finds the shortest path from vertex 0 to all
5   // others
6
7   // usage:  dijkstra nv print
8
9   // where nv is the size of the graph, and print is 1 if graph and min
10  // distances are to be printed out, 0 otherwise
11
12  #include <omp.h>
13
14  // global variables, shared by all threads by default
15
16  int nv,  // number of vertices
17      *notdone, // vertices not checked yet
18      nth,  // number of threads
19      chunk,  // number of vertices handled by each thread
20      md,  // current min over all threads
21      mv,  // vertex which achieves that min
22      largeint = -1;  // max possible unsigned int
23
24  unsigned *ohd,  // 1-hop distances between vertices; "ohd[i][j]" is
25          // ohd[i*nv+j]
26          *mind;  // min distances found so far
27
28  void init(int ac, char **av)
29  {  int i,j,tmp;
30     nv = atoi(av[1]);
31     ohd = malloc(nv*nv*sizeof(int));
32     mind = malloc(nv*sizeof(int));
33     notdone = malloc(nv*sizeof(int));
34     // random graph
35     for (i = 0; i < nv; i++)
36        for (j = i; j < nv; j++)   {
37           if (j == i) ohd[i*nv+i] = 0;
38           else  {
39              ohd[nv*i+j] = rand() % 20;
40              ohd[nv*j+i] = ohd[nv*i+j];
41           }
42        }
43     for (i = 1; i < nv; i++)   {
44        notdone[i] = 1;
45        mind[i] = ohd[i];
46     }
47  }
48
49  // finds closest to 0 among notdone, among s through e
50  void findmymin(int s, int e, unsigned *d, int *v)
51  {  int i;
52     *d = largeint;
53     for (i = s; i <= e; i++)
54        if (notdone[i] && mind[i] < *d)  {
55           *d = ohd[i];
56           *v = i;
57        }
58  }
59
60  // for each i in [s,e], ask whether a shorter path to i exists, through
```

```
61   // mv
62   void updatemind(int s, int e)
63   {  int i;
64      for (i = s; i <= e; i++)
65         if (mind[mv] + ohd[mv*nv+i] < mind[i])
66            mind[i] = mind[mv] + ohd[mv*nv+i];
67   }
68
69   void dowork()
70   {
71      #pragma omp parallel
72      {  int startv,endv,  // start, end vertices for my thread
73             step,  // whole procedure goes nv steps
74             mymv,  // vertex which attains the min value in my chunk
75             me = omp_get_thread_num();
76             unsigned mymd;  // min value found by this thread
77         #pragma omp single
78         {  nth = omp_get_num_threads();
79            if (nv % nth != 0) {
80               printf("nv must be divisible by nth\n");
81               exit(1);
82            }
83            chunk = nv/nth;
84            printf("there are %d threads\n",nth);
85         }
86         startv = me * chunk;
87         endv = startv + chunk - 1;
88         for (step = 0; step < nv; step++)  {
89            // find closest vertex to 0 among notdone; each thread finds
90            // closest in its group, then we find overall closest
91            #pragma omp single
92            {  md = largeint; mv = 0;  }
93            findmymin(startv,endv,&mymd,&mymv);
94            // update overall min if mine is smaller
95            #pragma omp critical
96            {  if (mymd < md)
97                  {  md = mymd; mv = mymv;  }
98            }
99            #pragma omp barrier
100           // mark new vertex as done
101           #pragma omp single
102           {  notdone[mv] = 0;  }
103           // now update my section of mind
104           updatemind(startv,endv);
105           #pragma omp barrier
106        }
107     }
108  }
109
110  int main(int argc, char **argv)
111  {  int i,j,print;
112     double startime,endtime;
113     init(argc,argv);
114     startime = omp_get_wtime();
115     // parallel
116     dowork();
117     // back to single thread
118     endtime = omp_get_wtime();
```

```
119     printf("elapsed time:  %f\n",endtime-startime);
120     print = atoi(argv[2]);
121     if (print)  {
122        printf("graph weights:\n");
123        for (i = 0; i < nv; i++)  {
124           for (j = 0; j < nv; j++)
125              printf("%u  ",ohd[nv*i+j]);
126           printf("\n");
127        }
128        printf("minimum distances:\n");
129        for (i = 1; i < nv; i++)
130           printf("%u\n",mind[i]);
131     }
132  }
```

The constructs will be presented in the following sections, but first the algorithm will be explained.


### 3.2.1  The Algorithm

The code implements the Dijkstra algorithm for finding the shortest paths from vertex 0 to the other vertices in an N-vertex undirected graph. Pseudocode for the algorithm is shown below, with the array G assumed to contain the one-hop distances from 0 to the other vertices.

```
1   Done = {0}  # vertices checked so far
2   NewDone = None  # currently checked vertex
3   NonDone = {1,2,...,N-1}  # vertices not checked yet
4   for J = 0 to N-1 Dist[J] = G(0,J)  # initialize shortest-path lengths
5
6   for Step = 1 to N-1
7      find J such that Dist[J] is min among all J in NonDone
8      transfer J from NonDone to Done
9      NewDone = J
10     for K = 1 to N-1
11        if K is in NonDone
12           # check if there is a shorter path from 0 to K through NewDone
13           # than our best so far
14           Dist[K] = min(Dist[K],Dist[NewDone]+G[NewDone,K])
```

At each iteration, the algorithm finds the closest vertex J to 0 among all those not yet processed, and then updates the list of minimum distances to each vertex from 0 by considering paths that go through J. Two obvious potential candidate part of the algorithm for parallelization are the "find J" and "for K" lines, and the above OpenMP code takes this approach.


### 3.2.2  The OpenMP `parallel` Pragma

As can be seen in the comments in the lines

```
      // parallel
      dowork();
      // back to single thread
```

the function **main()** is run by a **master thread**, which will then branch off into many threads running **dowork()** in parallel. The latter feat is accomplished by the directive in the lines

```
void dowork()
{
   #pragma omp parallel
   { int startv,endv,  // start, end vertices for this thread
         step,  // whole procedure goes nv steps
         mymv,  // vertex which attains that value
         me = omp_get_thread_num();
```

That directive sets up a team of threads (which includes the master), all of which execute the block following the directive in parallel.[1] Note that, unlike the **for** directive which will be discussed below, the **parallel** directive leaves it up to the programmer as to how to partition the work. In our case here, we do that by setting the range of vertices which this thread will process:

```
      startv = me * chunk;
      endv = startv + chunk - 1;
```

Again, keep in mind that *all* of the threads execute this code, but we've set things up with the variable **me** so that different threads will work on different vertices. This is due to the OpenMP call

```
      me = omp_get_thread_num();
```

which sets **me** to the thread number for this thread.

### 3.2.3   Scope Issues

Note carefully that in

```
   #pragma omp parallel
   { int startv,endv,  // start, end vertices for this thread
         step,  // whole procedure goes nv steps
         mymv,  // vertex which attains that value
         me = omp_get_thread_num();
```

---

[1]There is an issue here of thread startup time. The OMPi compiler sets up threads at the outset, so that that startup time is incurred only once. When a **parallel** construct is encountered, they are awakened. At the end of the construct, they are suspended again, until the next **parallel** construct is reached.

the pragam comes *before* the declaration of the local variables. That means that all of them are "local" to each thread, i.e. not shared by them. But if a work sharing directive comes within a function but *after* declaration of local variables, those variables are actually "global" to the code in the directive, i.e. they *are* shared in common among the threads.

This is the default, but you can change these properties, e.g. using the **shared** keyword. For instance,

```
#pragma omp parallel private(x,y)
```

would make **x** and **y** nonshared even if they were declared above the directive line.

It is crucial to keep in mind that variables which are global to the program (in the C/C++ sense) are automatically global to all threads. This is the primary means by which the threads communicate with each other.

### 3.2.4   The OpenMP `single` **Pragma**

In some cases we want just one thread to execute some code, even though that code is part of a **parallel** or other **work sharing** block.[2] We use the **single** directive to do this, e.g.:

```
#pragma omp single
{  nth = omp_get_num_threads();
   if (nv % nth != 0) {
      printf("nv must be divisible by nth\n");
      exit(1);
   }
   chunk = nv/nth;
   printf("there are %d threads\n",nth);  }
```

Since the variables **nth** and **chunk** are global and thus shared, we need not have all threads set them, hence our use of **single**.

### 3.2.5   The OpenMP `barrier` **Pragma**

As see in the example above, the **barrier** implements a standard barrier, applying to all threads.

### 3.2.6   **Implicit Barriers**

Note that there is an implicit barrier at the end of each **single** block, which is also the case for **parallel**, **for**, and **sections** blocks. This can be overridden via the **nowait** clause, e.g.

---

[2]This is an OpenMP term. The **for** directive is another example of it. More on this below.

```
#pragma omp for nowait
```

Needless to say, the latter should be used with care, and in most cases will not be usable. On the other hand, putting in a barrier where it is not needed would severely reduce performance.

### 3.2.7   The OpenMP `critical` **Pragma**

The last construct used in this example is **critical**, for critical sections.

```
#pragma omp critical
{   if (mymd < md)
      {   md = mymd; mv = mymv;   }
}
```

It means what it says, allowing entry of only one thread at a time while others wait. Here we are updating global variables **md** and **mv**, which has to be done atomically, and **critical** takes care of that for us. This is much more convenient than setting up lock variables, etc., which we would do if we were programming threads code directly.

## 3.3   The OpenMP `for` **Pragma**

This one breaks up a C/C++ **for** loop, assigning various iterations to various threads. This way the iterations are done in parallel. Of course, that means that they need to be independent iterations, i.e. one iteration cannot depend on the result of another.

### 3.3.1   Basic Example

Here's how we could use this construct in the Dijkstra program :

```
1   // Dijkstra.c
2
3   // OpenMP example program (OMPi version):  Dijkstra shortest-path finder
4   // in a bidirectional graph; finds the shortest path from vertex 0 to
5   // all others
6
7   // usage:  dijkstra nv print
8
9   // where nv is the size of the graph, and print is 1 if graph and min
10  // distances are to be printed out, 0 otherwise
11
12  #include <omp.h>
```

```
13
14   // global variables, shared by all threads by default
15
16   int nv,  // number of vertices
17       *notdone, // vertices not checked yet
18       nth,  // number of threads
19       chunk,  // number of vertices handled by each thread
20       md,  // current min over all threads
21       mv,  // vertex which achieves that min
22       largeint = -1;  // max possible unsigned int
23
24   unsigned *ohd,  // 1-hop distances between vertices; "ohd[i][j]" is
25                   // ohd[i*nv+j]
26            *mind;  // min distances found so far
27
28   void init(int ac, char **av)
29   {  int i,j,tmp;
30      nv = atoi(av[1]);
31      ohd = malloc(nv*nv*sizeof(int));
32      mind = malloc(nv*sizeof(int));
33      notdone = malloc(nv*sizeof(int));
34      // random graph
35      for (i = 0; i < nv; i++)
36         for (j = i; j < nv; j++)  {
37            if (j == i) ohd[i*nv+i] = 0;
38            else  {
39               ohd[nv*i+j] = rand() % 20;
40               ohd[nv*j+i] = ohd[nv*i+j];
41            }
42         }
43      for (i = 1; i < nv; i++)  {
44         notdone[i] = 1;
45         mind[i] = ohd[i];
46      }
47   }
48
49   void dowork()
50   {
51      #pragma omp parallel
52      {  int step,  // whole procedure goes nv steps
53             mymv,  // vertex which attains that value
54             me = omp_get_thread_num(),
55             i;
56         unsigned mymd;  // min value found by this thread
57         #pragma omp single
58         { nth = omp_get_num_threads();
59           printf("there are %d threads\n",nth);  }
60         for (step = 0; step < nv; step++)  {
61            // find closest vertex to 0 among notdone; each thread finds
62            // closest in its group, then we find overall closest
63            #pragma omp single
64            {  md = largeint; mv = 0;  }
65            mymd = largeint;
66            #pragma omp for nowait
67            for (i = 1; i < nv; i++)  {
68               if (notdone[i] && mind[i] < mymd)  {
69                  mymd = ohd[i];
70                  mymv = i;
```

```
71                    }
72               }
73               // update overall min if mine is smaller
74               #pragma omp critical
75               {   if (mymd < md)
76                     {   md = mymd; mv = mymv;   }
77               }
78               // mark new vertex as done
79               #pragma omp single
80               {   notdone[mv] = 0;   }
81               // now update ohd
82               #pragma omp for
83               for (i = 1; i < nv; i++)
84                   if (mind[mv] + ohd[mv*nv+i] < mind[i])
85                       mind[i] = mind[mv] + ohd[mv*nv+i];
86          }
87      }
88  }
89
90  int main(int argc, char **argv)
91  {   int i,j,print;
92      init(argc,argv);
93      // parallel
94      dowork();
95      // back to single thread
96      print = atoi(argv[2]);
97      if (print)  {
98          printf("graph weights:\n");
99          for (i = 0; i < nv; i++)   {
100             for (j = 0; j < nv; j++)
101                 printf("%u   ",ohd[nv*i+j]);
102             printf("\n");
103         }
104         printf("minimum distances:\n");
105         for (i = 1; i < nv; i++)
106             printf("%u\n",mind[i]);
107     }
108 }
109
```

The work which used to be done in the function **findmymin()** is now done here:

```
#pragma omp for
for (i = 1; i < nv; i++)   {
    if (notdone[i] && mind[i] < mymd)   {
        mymd = ohd[i];
        mymv = i;
    }
}
```

Each thread executes one or more of the iterations, i.e. takes responsibility for one or more values of *i*. This occurs in parallel, so as mentioned earlier, the programmer must make sure that the iterations are independent; there is no predicting which threads will do which values of **i**, in which order. By the way, for obvious reasons OpenMP treats the loop index, **i** here, as private even if by context it would be shared.

### 3.3.2   Controlling the Partitioning of Work to Threads

In this default version of the **for** construct, iterations are executed by threads *in unpredictable order*, with each thread taking on *one iteration's worth of work at a time*. Both of these can be changed by the programmer, using the **schedule** clause.

For instance, our original version of our program in Section 3.2 broke the work into chunks, with chunk size being the number vertices divided by the number of threads.

For the Dijkstra algorithm, for instance, we could have this:

```
...
        #pragma omp for schedule(static,chunk)
        for (i = 1; i < nv; i++)  {
           if (notdone[i] && mind[i] < mymd)  {
              mymd = ohd[i];
              mymv = i;
           }
        }
...
        #pragma omp for schedule(static,chunk)
        for (i = 1; i < nv; i++)
           if (mind[mv] + ohd[mv*nv+i] < mind[i])
              mind[i] = mind[mv] + ohd[mv*nv+i];
...
```

But one can enhance performance by considering other chunk sizes (in which case a thread would be responsible for more than one chunk). On the one hand, large chunks are good, due to there being less overhead—every time a thread finishes a chunk, it must go through the critical section, which serializes our parallel program and thus slows things down. On the other hand, if chunk sizes are large, then toward the end of the work, some threads may be working on their last chunks while others have finished and are now idle, thus foregoing potential speed enhancement. So it would be nice to have large chunks at the beginning of the run, to reduce the overhead, but smaller chunks at the end. This can be done using the **guided** clause.

For the Dijkstra algorithm, for instance, we could have this:

```
...
        #pragma omp for schedule(guided)
        for (i = 1; i < nv; i++)  {
           if (notdone[i] && mind[i] < mymd)  {
              mymd = ohd[i];
              mymv = i;
           }
        }
...
        #pragma omp for schedule(guided)
        for (i = 1; i < nv; i++)
           if (mind[mv] + ohd[mv*nv+i] < mind[i])
              mind[i] = mind[mv] + ohd[mv*nv+i];
...
```

### 3.3.3 The OpenMP `reduction` Clause

The name of this OpenMP clause alludes to the term **reduction** in functional programming. Many parallel programming languages include such operations, to enable the programmer to more conveniently (and often more efficiently) have threads/processors cooperate in computing sums, products, etc. OpenMP does this via the **reduction** clause.

For example, consider

```
1  int z;
2  ...
3  #pragma omp for reduction(+:z)
4  for (i = 0; i < n; i++)  z += x[i];
```

The pragma says that the threads will share the work as in our previous discussion of the **for** pragma. In addition, though, there will be independent copies of **z** maintained for each thread, each initialized to 0 before the loop begins. When the loop is entirely done, the values of **z** from the various threads will be summed, of course in an atomic manner.

Note that the **+** operator not only indicates that the values of **z** are to be summed, but also that their initial values are to be 0. If the operator were **\***, say, then the product of the values would be computed, and their initial values would be 1.

Our use of the **reduction** clause here makes our programming much easier. Indeed, if had old serial code that we wanted to parallelize, we would have to make no change to it! OpenMP is taking care of both the work splitting across values of **i**, and the atomic operations. Moreover—note this carefully—it is efficient, because by maintaining separate copies of **z** until the loop is done, we are reducing the number of serializing atomic actions, and are avoiding time-costly cache coherency transactions and the like.

Without this construct, we would have to do

```
int z,myz=0;
...
#pragma omp for private(myz)
for (i = 0; i < n; i++)  myz += x[i];
#pragma omp critical
{ z += myz; }
```

Here are the eligible operators and the corresponding initial values:

In C/C++, you can use **reduction** with +, -, \*, &, |, && and || (and the exclusive-or operator).

| operator | initial value |
|----------|---------------|
| +        | 0             |
| -        | 0             |
| *        | 1             |
| &        | bit string of 1s |
| \|       | bit string of 0s |
| ^        | 0             |
| &&       | 1             |
| \|\|     | 0             |

The lack of other operations typically found in other parallel programming languages, such as min and max, is due to the lack of these operators in C/C++. The FORTRAN version of OpenMP does have min and max.[3]

## 3.4   Other OpenMP Synchronization Issues

Earlier we saw the **critical** and **barrier** constructs. There is more to discuss, which we do here.

### 3.4.1   The OpenMP `atomic` **Clause**

The **critical** construct not only serializes your program, but also it adds a lot of overhead. If your critical section involves just a one-statement update to a shared variable, e.g.

```
x += y;
```

etc., then the OpenMP compiler can take advantage of an atomic hardware instruction, e.g. the LOCK prefix on Intel, to set up an extremely efficient critical section, e.g.

```
#pragma omp atomic
x += y;
```

Since it is a single statement rather than a block, there are no braces.

The eligible operators are:

```
++, --, +=, *=, <<=, &=, |=
```

---

[3]Note, though, that plain min and max would not help in our Dijkstra example above, as we not only need to find the minimum value, but also need the vertex which attains that value.

### 3.4.2  Memory Consistency and the `flush` Pragma

Consider a shared-memory multiprocessor system with coherent caches, and a shared, i.e. global, variable **x**. If one thread writes to **x**, you might think that the cache coherency system will ensure that the new value is visible to other threads. But it is is not quite so simple as this.

For example, the compiler may store **x** in a register, and update **x** itself at certain points. In between such updates, since the memory location for **x** is not written to, the cache will be unaware of the new value, which thus will not be visible to other threads. If the processors have write buffers etc., the same problem occurs.

In other words, we must account for the fact that our program could be run on different kinds of hardware with different memory consistency models. Thus OpenMP must have its own memory consistency model, which is then translated by the compiler to mesh with the hardware.

OpenMP takes a **relaxed consistency** approach, meaning that it forces updates to memory ("flushes") at all synchronization points, i.e. at:

- **barrier**

- entry/exit to/from **critical**

- entry/exit to/from **ordered**

- entry/exit to/from **parallel**

- exit from **parallel for**

- exit from **parallel sections**

- exit from **single**

In between synchronization points, one can force an update to **x** via the **flush** pragma:

```
#pragma omp flush (x)
```

The flush operation is obviously architecture-dependent. OpenMP compilers will typically have the proper machine instructions available for some common architectures. For the rest, it can force a flush at the hardware level by doing lock/unlock operations, though this may be costly in terms of time.

## 3.5   Compiling, Running and Debugging OpenMP Code

### 3.5.1   Compiling

There are a number of open source compilers available for OpenMP, including:

- Omni: This is available at (`http://phase.hpcc.jp/Omni/`). To compile an OpenMP program in **x.c** and create an executable file **x**, run

  ```
  omcc -g -o x x.c
  ```

- Ompi:  You can download this at `http://www.cs.uoi.gr/~ompi/index.html`. Compile **x.c** by

  ```
  ompicc -g -o x x.c
  ```

- GCC, version 4.2 or later:[4] Compile **x.c** via

  ```
  gcc -fopenmp -g -o x x.c
  ```

### 3.5.2   Running

Just run the executable as usual.

The number of threads will be the number of processors, by default. To change that value, set the OMP_NUM_THREADS environment variable. For example, to get four threads in the C shell, type

```
setenv OMP_NUM_THREADS 4
```

### 3.5.3   Debugging

OpenMP's use of pragmas makes it difficult for the compilers to maintain your original source code line numbers, and your function and variable names. But with a little care, a symbolic debugger such as GDB can still be used. Here are some tips for the compilers mentioned above, using GDB as our example debugging tool:

- Omni: The function **main()** in your executable is actually in the OpenMP library, and your function **main()** is renamed **_ompc_main()**. So, when you enter GDB, first set a breakpoint at your own code:

---

[4]You may find certain subversions of GCC 4.1 can be used too.

```
(gdb) b _ompc_main
```

Then run your program to this breakpoint, and set whatever other breakpoints you want.

You should find that your other variable and function names are unchanged.

- Ompi: During preprocessing of your file **x.c**, the compiler produces a file **x_ompi.c**, and the latter is what is actually compiled. Your function **main** is renamed to **_ompi_originalMain()**. Your other functions and variables are renamed. For example in our Dijkstra code, the function **dowork()** is renamed to **dowork_parallel_0**. And by the way, all indenting is lost! Keep these points in mind as you navigate through your code in GDB.

- GCC: GCC maintains line numbers and names well, except (as of July 2007) it does not retain names of local variables within blocks controlled by **omp parallel** at all. This is not a big problem if you have some function calls in such blocks. Pass the values of those phantom locals to those functions, and place breakpoints in them.

## 3.6 Combining Work-Sharing Constructs

In our examples of the **for** pragma above, that pragma would come within a block headed by a **parallel** pragma. The latter specifies that a team of theads is to be created, with each one executing the given block, while the former specifies that the various iterations of the loop are to be distributed among the threads. As a shortcut, we can combine the two pragmas:

```
#pragma omp parallel for
```

This also works with the **sections** pragma.

## 3.7 Performance

As is usually the case with parallel programming, merely parallelizing a program won't necessarily make it faster, even on shared-memory hardware. Operations such as critical sections, barriers and so on serialize an otherwise-parallel program, sapping much of its speed. In addition, there are issues of cache coherency transactions, false sharing etc.

### 3.7.1 The Effect of Problem Size

To illustrate this, I ran our original Dijkstra example (Section 3.2 on various graph sizes, on a quad core machine. Here are the timings:

| nv   | nth | time     |
|------|-----|----------|
| 1000 | 1   | 0.005472 |
| 1000 | 2   | 0.011143 |
| 1000 | 4   | 0.029574 |

The more parallelism we had, the *slower* the program ran! The synchronization overhead was just too much to be compensated by the parallel computation.

However, parallelization did bring benefits on larger problems:

| nv    | nth | time     |
|-------|-----|----------|
| 25000 | 1   | 2.861814 |
| 25000 | 2   | 1.710665 |
| 25000 | 4   | 1.453052 |

### 3.7.2   Some Fine Tuning

How could we make our Dijkstra code faster? One idea would be to eliminate the critical section. Recall that in each iteration, the threads compute their local minimum distance values **md** and **mv**, and then update the global values **md** and **mv**. Since the update must be atomic, this causes some serialization of the program. Instead, we could have the threads store their values **mymd** and **mymv** in a global array **mymins**, with each thread using a separate pair of locations within that array, and then at the end of the iteration we could have just one task scan through **mymins** and update **md** and **mv**.

Here is the resulting code:

```
1    // Dijkstra.c
2
3    // OpenMP example program:  Dijkstra shortest-path finder in a
4    // bidirectional graph; finds the shortest path from vertex 0 to all
5    // others
6
7    // **** in this version, instead of having a critical section in which
8    // each thread updates md and mv, the threads record their mymd and mymv
9    // values in a global array mymins, which one thread then later uses to
10   // update md and mv
11
12   // usage:  dijkstra nv print
13
14   // where nv is the size of the graph, and print is 1 if graph and min
15   // distances are to be printed out, 0 otherwise
16
17   #include <omp.h>
18
19   // global variables, shared by all threads by default
20
21   int nv,  // number of vertices
22       *notdone, // vertices not checked yet
23       nth,  // number of threads
24       chunk,  // number of vertices handled by each thread
```

```
25        md,  // current min over all threads
26        mv,  // vertex which achieves that min
27        largeint = -1;  // max possible unsigned int
28
29    int *mymins;  // (mymd,mymv) for each thread; see dowork()
30
31    unsigned *ohd,  // 1-hop distances between vertices; "ohd[i][j]" is
32             // ohd[i*nv+j]
33             *mind;  // min distances found so far
34
35    void init(int ac, char **av)
36    {  int i,j,tmp;
37       nv = atoi(av[1]);
38       ohd = malloc(nv*nv*sizeof(int));
39       mind = malloc(nv*sizeof(int));
40       notdone = malloc(nv*sizeof(int));
41       // random graph
42       for (i = 0; i < nv; i++)
43          for (j = i; j < nv; j++)  {
44             if (j == i) ohd[i*nv+i] = 0;
45             else  {
46                ohd[nv*i+j] = rand() % 20;
47                ohd[nv*j+i] = ohd[nv*i+j];
48             }
49          }
50       for (i = 1; i < nv; i++)  {
51          notdone[i] = 1;
52          mind[i] = ohd[i];
53       }
54    }
55
56    // finds closest to 0 among notdone, among s through e
57    void findmymin(int s, int e, unsigned *d, int *v)
58    {  int i;
59       *d = largeint;
60       for (i = s; i <= e; i++)
61          if (notdone[i] && mind[i] < *d)  {
62             *d = ohd[i];
63             *v = i;
64          }
65    }
66
67    // for each i in [s,e], ask whether a shorter path to i exists, through
68    // mv
69    void updatemind(int s, int e)
70    {  int i;
71       for (i = s; i <= e; i++)
72          if (mind[mv] + ohd[mv*nv+i] < mind[i])
73             mind[i] = mind[mv] + ohd[mv*nv+i];
74    }
75
76    void dowork()
77    {
78       #pragma omp parallel
79       {  int startv,endv,  // start, end vertices for my thread
80                step,  // whole procedure goes nv steps
81                me,
82                mymv;  // vertex which attains the min value in my chunk
```

```
83              unsigned mymd;  // min value found by this thread
84          int i;
85          me = omp_get_thread_num();
86          #pragma omp single
87          {  nth = omp_get_num_threads();
88             if (nv % nth != 0) {
89                printf("nv must be divisible by nth\n");
90                exit(1);
91             }
92             chunk = nv/nth;
93             mymins = malloc(2*nth*sizeof(int));
94          }
95          startv = me * chunk;
96          endv = startv + chunk - 1;
97          for (step = 0; step < nv; step++)  {
98             // find closest vertex to 0 among notdone; each thread finds
99             // closest in its group, then we find overall closest
100            findmymin(startv,endv,&mymd,&mymv);
101            mymins[2*me] = mymd;
102            mymins[2*me+1] = mymv;
103            #pragma omp barrier
104            // mark new vertex as done
105            #pragma omp single
106            {  md = largeint; mv = 0;
107               for (i = 1; i < nth; i++)
108                  if (mymins[2*i] < md) {
109                     md = mymins[2*i];
110                     mv = mymins[2*i+1];
111                  }
112               notdone[mv] = 0;
113            }
114            // now update my section of mind
115            updatemind(startv,endv);
116            #pragma omp barrier
117         }
118      }
119  }
120
121  int main(int argc, char **argv)
122  {  int i,j,print;
123     double startime,endtime;
124     init(argc,argv);
125     startime = omp_get_wtime();
126     // parallel
127     dowork();
128     // back to single thread
129     endtime = omp_get_wtime();
130     printf("elapsed time:  %f\n",endtime-startime);
131     print = atoi(argv[2]);
132     if (print)  {
133        printf("graph weights:\n");
134        for (i = 0; i < nv; i++)  {
135           for (j = 0; j < nv; j++)
136              printf("%u  ",ohd[nv*i+j]);
137           printf("\n");
138        }
139        printf("minimum distances:\n");
140        for (i = 1; i < nv; i++)
```

```
141          printf("%u\n",mind[i]);
142     }
143 }
```

Let's take a look at the latter part of the code for one iteration;

```
1           findmymin(startv,endv,&mymd,&mymv);
2           mymins[2*me] = mymd;
3           mymins[2*me+1] = mymv;
4           #pragma omp barrier
5           // mark new vertex as done
6           #pragma omp single
7           {  notdone[mv] = 0;
8              for (i = 1; i < nth; i++)
9                 if (mymins[2*i] < md) {
10                    md = mymins[2*i];
11                    mv = mymins[2*i+1];
12                 }
13          }
14          // now update my section of mind
15          updatemind(startv,endv);
16          #pragma omp barrier
```

The call to **findmymin()** is as before; this thread finds the closest vertex to 0 among this thread's range of vertices. But instead of comparing the result to **md** and possibly updating it and **mv**, the thread simply stores its **mymd** and **mymv** in the global array **mymins**. After all threads have done this and then waited at the barrier, we have just one thread update **md** and **mv**.

Let's see how well this tack worked:

| nv | nth | time |
|-------|-----|----------|
| 25000 | 1 | 2.546335 |
| 25000 | 2 | 1.449387 |
| 25000 | 4 | 1.411387 |

This brought us about a 15% speedup in the two-thread case, though less for four threads.

What else could we do? Here are a few ideas:

- False sharing could be a problem here. To address it, we could make **mymins** much longer, changing the places at which the threads write their data, leaving most of the array as padding.

- We could try the modification of our program in Section 3.3.1, in which we use the OpenMP **for** pragma, as well as the refinements stated there, such as **schedule**.

- We could try combining all of the ideas here.

### 3.7.3   OpenMP Internals

We may be able to write faster code if we know a bit about how OpenMP works inside.

You can get some idea of this from your compiler. For example, if you use the **-t** option with the Omni compiler, or **-k** with Ompi, you can inspect the result of the preprocessing of the OpenMP pragmas.

Here for instance is the code produced by Omni from the call to **findmymin()** in our Dijkstra program:

```
# 93 "Dijkstra.c"
findmymin(startv,endv,&(mymd),&(mymv));{
_ompc_enter_critical(&__ompc_lock_critical);
# 96 "Dijkstra.c"
if((mymd)<(((unsigned )(md)))){

# 97 "Dijkstra.c"
(md)=(((int )(mymd)));
# 97 "Dijkstra.c"
(mv)=(mymv);
}_ompc_exit_critical(&__ompc_lock_critical);
```

Fortunately Omni saves the line numbers from our original source file, but the pragmas have been replaced by calls to OpenMP library functions.

The document, *The GNU OpenMP Implementation*, `http://pl.postech.ac.kr/~gla/cs700-07f/ref/openMp/libgomp.pdf`, includes good outline of how the pragmas are translated.

## 3.8   The Rest of OpenMP

There is much, much more to OpenMP than what we have seen here. Here are a couple of examples:

- The **sections** pragma: Suppose at a certain point in the code, there are several different code blocks which you want different threads to execute simultaneously. You can use this pragma to specify this.

- The **ordered** pragma: This is used to ensure that a subblock within a loop is executed in sequential order.

To see the details of these and other OpenMP constructs, there are many Web pages you can check, and there is also the excellent book, *Using OpenMP: Portable Shared Memory Parallel Programming*, by Barbara Chapman, Gabriele Jost and Ruud Van Der Pas, MIT Press, 2008.

# Chapter 4

# Introduction to GPU Programming with CUDA

Even if you don't play video games, you can be grateful to the game players, as their numbers have given rise to a class of highly powerful parallel processing devices—**graphics processing units** (GPUs). Yes, you program right on the video card in your computer, even though your program may have nothing to do with graphics.

## 4.1   Overview

The video game market is so lucrative that the industry has developed ever-faster GPUs, in order to handle ever-faster and ever-more visually detailed video games. These actually are parallel processing hardware devices, so around 2003 some people began to wonder if one might use them for parallel processing of nongraphics applications.

Originally this was cumbersome. One needed to figure out clever ways of mapping one's application to some kind of graphics problem. Though some high-level interfaces were developed to automate this transformation, effective coding required some understanding of graphics principles.

But the current generation of GPUs have separated out the graphics operations, and now consist of multiprocessor elements that run under the familar threads model. Thus they are easily programmable. Granted, effective coding still requires an intimate knowledge of the hardwre, but at least it's (more or less) familar hardware, not requiring knowledge of graphics.

Moreover, unlike a multicore machine, with the ability to run just a few threads at one time, e.g. four threads on a quad core machine, GPUs can run *hundreds or thousands* of threads at once. There are various

restrictions that come with this, but you can see that there is fantastic potential for speed here.

NVIDIA has developed the CUDA language as a vehicle for programming on their GPUs. It's basically just a slight extension of C, and has become very popular. More recently, the OpenCL language has been developed by Apple, AMD and others (including NVIDIA). It too is a slight extension of C, and it aims to provide a uniform interface that works with multicore machines in addition to GPUs.

Our discussion here focuses on CUDA and NVIDIA GPUs.

A CUDA program consists of code to be run on the **host**, i.e. the computer as a whole, and code to run on the **device**, i.e. the GPU. A function that is called by the host to execute on the device is called a **kernel**.

## 4.2   Hardware Structure

*Scorecards, get your scorecards here! You can't tell the players without a scorecard*—classic cry of vendors at baseball games

*Know thy enemy*—Sun Tzu, *The Art of War*

The enormous computational potential of GPUs cannot be unlocked without an intimate understanding of the hardware. This of course is a fundamental truism in the parallel processing world, but it is acutely important for GPU programming. This section presents an overview of the hardware, but true mastery of the GPU software genre requires delving into further details.

### 4.2.1   Processing Units

A GPU consists of a large set of **streaming multiprocessors** (SMs); you might say it's a multi-multiprocessor. Each SM consists of a number of **streaming processors** (SPs). It is important to understand the motivation for this hierarchy: Two threads located in different SMs cannot synchronize with each other. Though this sounds like a negative at first, it is actually a great advantage, as the independence of threads in separate SMs means that the hardware can run faster. So, if the CUDA application programmer can write his/her algorithm so as to have certain independent chunks, those chunks can be assigned to different SMs (we'll see how, shortly), then that's a "win."

The threads running within an SM *can* synchronize with each other, but there is further hierarchy: The threads within an SM are subdivided by the hardware into groups called **warps**. The key point is that *all the threads in a warp run the code in lockstep*. During the machine instruction fetch cycle, the same instruction will be fetched for all of the threads. Then in the execution cycle, each thread will either execute that particular instruction or execute nothing. This is the classical **single instruction, multiple data** (SIMD) pattern used in some early special-purpose computers such as the ILLIAC; here it is called **single instruction, multiple thread** (SIMT).

*This trait of thread execution has major implications for performance.* Consider what happens with if/then/else code. If some threads in a warp take the "then" branch and others go in the "else" direction, they cannot operate in lockstep. That means that some threads must wait while others execute. This renders the code at that point serial rather than parallel, a situation called **thread divergence**. As one CUDA Web tutorial points out, this is a "performance killer."

### 4.2.2 Memory Structure

Yet another key hierarchy—memory structure. Here's how it works:

- Registers: Each thread is allocated it's own set of registers. They are much more numerous than in a CPU, say in the hundreds, and access to them is very fast.

- Shared memory: All the threads in an SM share this memory, and use it to communicate, just as is the case with threads in CPUs. Access is said to be as fast as to registers!

- Global memory: This is shared by all the threads in an entire application, and is persistent across calls. It is very slow.

- Local memory: This is actually part of global memory, but is an area within that memory that is allocated by the compiler for a given thread. As such, it is slow, and accessible only by that thread.

### 4.2.3 Thread Management

Each SM runs the threads on a timesharing basis, just like an operating system (OS). This timesharing is implemented in the hardware, though, not in software as in the OS case.

With an OS, if a thread reaches an input/output operation, the OS suspends the thread while I/O is pending, and runs some other thread instead, so as to avoid wasting CPU cycles during the long period of time needed for the I/O. With an SM, the analogous situation is a long memory operation, so global memory; if a a warp of threads needs to access global memory (including local memory), the SM will schedule some other warp while the memory access is pending.

## 4.3 Software Structure

We'll start with a running example, then go into the details and generalizations.

### 4.3.1   Sample Program

Here's a sample program. And I've kept the sample simple: It just finds the sums of all the rows of a matrix.

```
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <cuda.h>
4
5   // CUDA example:  finds row sums of the integer matrix m, placing them
6   // in the array rs
7
8   // finds one element of rs, determined by x, of the n x n matrix m
9   // coordinate of this thread; assume grid consists of 1 block, with
10  // threads in n x 1 x 1 arrangement; matrices stored as 1-dimenstional,
11  // row-major order
12  __global__ void find1elt(int *dm, int *drs, int n)
13  {
14     int rownum = threadIdx.x;  // this thread will handle row # rownum
15     int sum=0;
16     for (int k = 0; k < n; k++)
17        sum += dm[rownum*n+k];
18     drs[rownum] = sum;
19  }
20
21  int main(int argc, char **argv)
22  {
23     // the size of the matrix is specified on the command line
24     int n = atoi(argv[1]);
25     int *hm, // host matrix
26         *dm, // device matrix
27         *hrs, // host rowsums
28         *drs; // device rowsums
29     int msize = n * n * sizeof(int);
30     hm = (int *) malloc(msize);
31     // as a test, fill matrix with consecutive integers
32     int t = 0,i,j;
33     for (i = 0; i < n; i++) {
34        for (j = 0; j < n; j++) {
35           hm[i*n+j] = t++;
36        }
37     }
38     cudaMalloc((void **)&dm,msize);
39     cudaMemcpy(dm,hm,msize,cudaMemcpyHostToDevice);
40     int rssize = n * sizeof(int);
41     hrs = (int *) malloc(rssize);
42     cudaMalloc((void **)&drs,rssize);
43     cudaMemcpy(drs,hrs,rssize,cudaMemcpyHostToDevice);
44     dim3 dimGrid(1,1);
45     dim3 dimBlock(n,1,1);
46     find1elt<<<dimGrid,dimBlock>>>(dm,drs,n);
47     cudaThreadSynchronize();
48     cudaMemcpy(hrs,drs,rssize,cudaMemcpyDeviceToHost);
49     for(int i=0; i<n; i++)  {
50        printf("%d\n",hrs[i]);
51     }
52     free(hm);
```

```
53      cudaFree(dm);
54      free(hrs);
55      cudaFree(drs);
56   }
```

Well, this is mostly C, with a bit of CUDA added here and there. Here's how the program works:

- Our **main()** runs on the host.

- Kernel functions are identified by __**global**__ **void**, are called by the host, and serve as entries to the device. We have one such function here.

- When a kernel is called, each thread runs it. Each thread receives the same arguments, though different threads may act differently based on programmer use of **threadIdx**.

- One calls **cudaMalloc()** to allocate space on the device's memory.

- Data is transferred to and from the host and device via **cudaMemcpy()**.

- Kernels return values via their arguments. However, the *call* to the kernel returns immediately. For that reason, the code above has a barrier call, to avoid copying the results back to the host from the device before they're ready:

  ```
  cudaThreadSynchronize();
  ```

  There is also a barrier available for the threads themselves, needed in many applications. The call is

  ```
  __syncthreads();
  ```

I've written the program so that each thread will handle one row of the matrix; thread i will find the sum in row i. Since I've chosen to store the matrix in one-dimensional form, and since the matrix is of size n x n, the loop

```
for (int k = 0; k < n; k++)
   sum += dm[rownum*n+k];
```

will indeed traverse the n elements of row number **rownum**, and compute their sum. That sum is then placed in the proper element of the output array:

```
drs[rownum] = sum;
```

### 4.3.2   Threads Hierarchy

Like the hardware, threads in CUDA software follow a hierarchy:

- The entirety of threads for an application is called a **grid**.

- A grid consists of one or more **blocks** of threads.

- Each block has its own ID within the grid, consisting of an "x coordinate" and a "y coordinate."

- Likewise each thread has x, y and z coordinates within whichever block it belongs to.

- Just as an ordinary CPU thread needs to be able to sense its ID, e.g. by calling **omp_get_thread_num**() in OpenMP, CUDA threads need to do the same. A CUDA thread can access its block ID via the built-in variables **blockIdx.x** and **blockIdx.y**, and can access its thread ID within its block via **threadIdx.x** and **threadIdx.y**.

- The programmer specifies the grid size (the numbers of rows and columns of blocks within a grid) and the block size (numbers of rows, columns and layers of threads within a block). In the example above, this was done by the code

  ```
  dim3 dimGrid(1,1);
  dim3 dimBlock(n,1,1);
  find1elt<<<dimGrid,dimBlock>>>(dm,drs,n);
  ```

  Here the grid is specified to consist of just one ($1 \times 1$) block, and each block consists of just n ($n \times 1 \times 1$) threads.

  That last line is of course the call to the kernel. As you can see, CUDA extends C syntax to allow specifying the grid and block sizes. CUDA will store this information in **blockDim** and **threadDim**.

- All threads in a block must run in the same SM, though more than one block might be on the same SM.

- The "coordinates" of a block within the grid, and of a thread within a block, are merely abstractions. **They do not correspond to any physical arrangment.**

The motivation for the two-dimensional block arrangment is to make coding conceptually simpler for the programmer if he/she is working an application that is two-dimensional in nature.

For example, in a matrix application one's parallel algorithm might be based on partitioning the matrix into rectangular submatrices, as we'll do in Section 7.3. In a small example there, the matrix

$$A = \left( \begin{array}{ccc} 1 & 5 & 12 \\ 0 & 3 & 6 \\ 4 & 8 & 2 \end{array} \right) \tag{4.1}$$

is partitioned as

$$A = \begin{pmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{pmatrix},$$ (4.2)

where

$$A_{00} = \begin{pmatrix} 1 & 5 \\ 0 & 3 \end{pmatrix},$$ (4.3)

$$A_{01} = \begin{pmatrix} 12 \\ 6 \end{pmatrix},$$ (4.4)

$$A_{10} = \begin{pmatrix} 4 & 8 \end{pmatrix}$$ (4.5)

and

$$A_{11} = \begin{pmatrix} 2 \end{pmatrix}.$$ (4.6)

We might then have one block of threads handle $A_{00}$, another block handle $A_{01}$ and so on. CUDA's two-dimensional ID system for blocks makes life easier for programmers in such situations. Indeed, the submatrices in such partitioning are often called *blocks* in the parallel algorithms literature.

### 4.3.3 Memory Placement

When you copy data from the host to the device, it goes into the latter's global memory. As noted, this memory is slow. If it is to be repeatedly read, copy it to shared memory, stored in a variable whose declaration is marked __**shared**.

One also uses such variables for communication between threads in the same block.

### 4.3.4 What's NOT There

*We're not in Kansas anymore, Toto*–character Dorothy Gale in *The Wizard of Oz*

It looks like C, it feels like C, and for the most part, it *is* C. But in many ways, it's quite different from what you're used to:

- You don't have access to the C library, e.g. **printf()** (the library consists of host machine language, after all). There are special versions of math functions, however, e.g. __**sin()**.

- No recursion.

- No stack. Functions are essentially inlined.

## 4.4   Synchronization

As noted earlier, a barrier for the threads in a block is available by calling __**syncthreads()**. Note carefully that if one thread writes a variable to shared memory and another then reads that variable, one must call this function in order to get the latest value.

CUDA and GPU hardware allow for certain atomic operations, such as atomic fetch-and-add, in functions such as **atomicExch()** and **atomicAdd()**. However, classic locks are NOT advised; they can easily lead to deadlock, due to thread divergence problems.

Remember, threads across blocks cannot sync with each other.

## 4.5   Performance Issues

Our brief introduction here barely scratches the surface. There are many considerations and tricks, especially in terms of memory access. The interested reader is referred to the many resources on the Web.

## 4.6   CUBLAS

CUDA includes some parallel linear algebra routines callable from straight C code. In other words, you can get the benefit of GPU in linear algebra contexts without using CUDA!

See `http://www.gsic.titech.ac.jp/~ccwww/tebiki/tesla_e/tesla5_e.html` for an example and compilation instructions.

## 4.7   Hardware Requirements, Installation, Compilation, Debugging

You do need what is currently (March 2010) a high-end NVIDIA video card. There is a list at `http://www.nvidia.com/object/cuda_gpus.html`. If you have a Linux system, run **lspci** to determine what kind you have.

Download the CUDA toolkit from NVIDIA. Just plug "CUDA download" into a Web search engine to find the site. Install as directed.

You'll need to set your search and library paths to include the CUDA **bin** and library directories.

To compile **x.cu**, type

```
$ nvcc -g -G x.cu -I/your_CUDA_include_path
```

The **-g -G** options are for setting up debugging.

Run the code as you normally would.

To determine the limits, e.g. maximum number of threads, for your device, use code like this:

```
cudaDeviceProp Props;
cudaGetDeviceProperties(Props,0);
```

The 0 is for device 0, assuming you only have one device. The return value of **cudaGetDeviceProperties()**. I recommend printing it from within GDB to see the values. One of the fields gives clock speed, which is typically slower than that of the host.

For debugging, CUDA includes a special version of GDB, **cuda-gdb**.

# Chapter 5

# Message Passing Systems

Message passing systems are probably the most common platforms for parallel processing today.

## 5.1  Overview

Traditionally, shared-memory hardware has been extremely expensive, with a typical system costing hundreds of thousands of dollars. Accordingly, the main users were for very large corporations or government agencies, with the machines being used for heavy-duty server applications, such as for large databases and World Wide Web sites. The conventional wisdom is that these applications require the efficiency that good shared-memory hardware can provide.

But the huge expense of shared-memory machines led to a quest for high-performance message-passing alternatives, first in hypercubes and then in networks of workstations (NOWs).

The situation changed radically around 2005, when "shared-memory hardware for the masses" became available in dual-core commodity PCs. Chips of higher core multiplicity are commercially available, with a decline of price being inevitable. Ordinary users will soon be able to afford shared-memory machines featuring dozens of processors.

Yet the message-passing paradigm continues to thrive. Many people believe it is more amenable to writing really fast code, and the the advent of **cloud computing** has given message-passing a big boost. In addition, many of the world's very fastest systems (see `www.top500.org` for the latest list) are in fact of the message-passing type.

In this unit, we take a closer look at this approach to parallel processing.

## 5.2   A Historical Example: Hypercubes

A popular class of parallel machines in the 1980s and early 90s was that of **hypercubes**. Intel sold them, for example, as did a subsidiary of Oracle, nCube. A hypercube would consist of some number of ordinary Intel processors, with each processor having some memory and serial I/O hardware for connection to its "neighbor" processors.

Hypercubes proved to be too expensive for the type of performance they could achieve, and the market was small anyway. Thus they are not common today, but they are still important, both for historical reasons (in the computer field, old techniques are often recycled decades later), and because the algorithms developed for them have become quite popular for use on general machines. In this section we will discuss architecture, algorithms and software for such machines.

**5.2.0.0.1   Definitions**   A **hypercube** of dimension d consists of $D = 2^d$ **processing elements** (PEs), i.e. processor-memory pairs, with fast serial I/O connections between neighboring PEs. We refer to such a cube as a **d-cube**.

The PEs in a d-cube will have numbers 0 through D-1. Let $(c_{d-1}, ..., c_0)$ be the base-2 representation of a PE's number. The PE has fast point-to-point links to d other PEs, which we will call its **neighbors**. Its i$th$ neighbor has number $(c_{d-1}, ..., 1 - c_{i-1}, ..., c_0)$.[1]

For example, consider a hypercube having D = 16, i.e. d = 4. The PE numbered 1011, for instance, would have four neighbors, 0011, 1111, 1001 and 1010.



It is sometimes helpful to build up a cube from the lower-dimensional cases. To build a (d+1)-dimensional cube from two d-dimensional cubes, just follow this recipe:

---

[1]Note that we number the digits from right to left, with the rightmost digit being digit 0.

(a) Take a d-dimensional cube and duplicate it. Call these two cubes subcube 0 and subcube 1.

(b) For each pair of same-numbered PEs in the two subcubes, add a binary digit 0 to the front of the number for the PE in subcube 0, and add a 1 in the case of subcube 1. Add a link between them.

The following figure shows how a 4-cube can be constructed in this way from two 3-cubes:



Given a PE of number $(c_{d-1}, ..., c_0)$ in a d-cube, we will discuss the i-cube to which this PE belongs, meaning all PEs whose first d-i digits match this PE's.[2] Of all these PEs, the one whose last i digits are all 0s is called the **root** of this i-cube.

For the 4-cube and PE 1011 mentioned above, for instance, the 2-cube to which that PE belongs consists of 1000, 1001, 1010 and 1011—i.e. all PEs whose first two digits are 10—and the root is 1000.

Given a PE, we can split the i-cube to which it belongs into two (i-1)-subcubes, one consisting of those PEs whose digit i-1 is 0 (to be called subcube 0), and the other consisting of those PEs whose digit i-1 is 1 (to be called subcube 1). Each given PE in subcube 0 has as its **partner** the PE in subcube 1 whose digits match those of the given PE, except for digit i-1.

To illustrate this, again consider the 4-cube and the PE 1011. As an example, let us look at how the 3-cube it belongs to will split into two 2-cubes. The 3-cube to which 1011 belongs consists of 1000, 1001, 1010, 1011, 1100, 1101, 1110 and 1111. This 3-cube can be split into two 2-cubes, one being 1000, 1001, 1010

---

[2]Note that this is indeed an i-dimensional cube, because the last i digits are free to vary.

and 1011, and the other being 1100, 1101, 1110 and 1111. Then PE 1000 is partners with PE 1100, PE 1001 is partners with PE 1101, and so on.

Each link between two PEs is a dedicated connection, much preferable to the shared link we have when we run, say, MPI, on a collection of workstations on an Ethernet. On the other hand, if one PE needs to communicate with a <u>non</u>-neighbor PE, multiple links (as many as d of them) will need to be traversed. Thus the nature of the communications costs here is much different than for a network of workstations, and this must be borne in mind when developing programs.

## 5.3   Networks of Workstations (NOWs)

The idea here is simple: Take a bunch of commodity PCs and network them for use as parallel processing systems. They are of course individual machines, capable of the usual uniprocessor, nonparallel applications, but by networking them together and using message-passing software environments such as MPI, we can form very powerful parallel systems.

The networking does result in a significant loss of performance, but the price/performance ratio in NOW can be much superior in many applications to that of shared-memory or hypercube hardware of comparable number of CPUs.

### 5.3.1   The Network Is Literally the Weakest Link

Still, one factor which can be key to the success of a NOW is to use a fast network, both in terms of hardware and network protocol. Ordinary Ethernet and TCP/IP are fine for the applications envisioned by the original designers of the Internet, e.g. e-mail and file transfer, but they are slow in the NOW context.

A popular network for a NOW today is Infiniband (IB) (`www.infinibandta.org`). It features low latency, about 1.0-3.0 microseconds, high bandwidth, about 1.0-2.0 gigaBytes per second), and uses a low amount of the CPU's cycles, around 5-10%.

The basic building block of IB is a switch, with many inputs and outputs, similar in concept to $\Omega$-net. You can build arbitrarily large and complex topologies from these switches.

A central point is that IB, as with other high-performance networks designed for NOWs, uses RDMA (Remote Direct Memory Access) read/write, which eliminates the extra copying of data between the application program's address space to that of the operating system.

IB has high performance and scalable[3] implementations of distributed locks, semaphores, collective communication operations. An atomic operation takes about 3-5 microseconds.

---

[3]The term *scalable* arises frequently in conversations on parallel processing. It means that this particular method of dealing with some aspect of parallel processing continues to work well as the system size increases. We say that the method *scales*.

IB implements true **multicast**, i.e. the simultaneous sending of messages to many nodes. Note carefully that even though MPI has its **MPI_Bcast()** function, it will send things out one at a time unless your network hardware is capable of multicast, and the MPI implementation you use is configured specifically for that hardware.

For information on network protocols, e.g. for example `www.rdmaconsortium.org`. A research paper evaluating a tuned implementation of MPI on IB is available at `nowlab.cse.ohio-state.edu/publications/journal-papers/2004/liuj-ijpp04.pdf`.

### 5.3.2   Other Issues

Increasingly today, the workstations themselves are multiprocessor machines, so a NOW really is a hybrid arrangement. They can be programmed either purely in a message-passing manner—e.g. running eight MPI processes on four dual-core machines—or in a mixed way, with a shared-memory approach being used within a workstation but message-passing used between them.

NOWs have become so popular that there are now "recipes" on how to build them for the specific purpose of parallel processing. The term **Beowulf** come to mean a NOW, usually with a fast network connecting them, used for parallel processing. The term *NOW* itself is no longer in use, replaced by *cluster*. Software packages such as ROCKS (`http://www.rocksclusters.org/wordpress/`) have been developed to make it easy to set up and administer such systems.

## 5.4   Parallelizing Serial Code Via Hidden Message-Passing

Say you have a Python program that you wish to run on a very large data set, and suppose the application is **embarrassingly parallel**, meaning that it partitions very easily into more or less independent chunks.

How can we parallelize such a program, given that Python is fundamentally a serial language? It does include a threads capability, but due to Python's controversial Global Interpreter Lock, only one thread may run at a time, even on a multiprocessor machine. And of course there is no provision in the Python language for running in parallel on multiple machines.

This situation also arises with the statistical/data manipulation language R, as well as just about any other serial language.

### 5.4.1   One Solution: Piggyback on MPI

A solution is to somehow create an interface from the serial language to a message-passing mechanism that operates at a lower level, to varying degrees transparently to the programmer. A common approach has been

to set up an interface from the serial language to MPI, exemplified by PyMPI and Rmpi, for Python and R, respectively. For details on these, see Chapter 6.

A more recent approach, less generally applicable but more convenient, has been **MapReduce** software, described in the following section.

### 5.4.2   Another Solution: MapReduce

MapReduce was developed as part of a recently-popularized computational approach known as **cloud computing**. The idea is that a large corporation that has many computers could sell time on them, thus making profitable use of excess capacity. The typical customer would have occasional need for large-scale computing—and often large-scale data storage. The customer would submit a program to the cloud computing vendor, who would run it in parallel on the vendor's many machines (unseen, thus forming the "cloud"), then return the output to the customer.

Google, Yahoo! and Amazon, among others, have recently gotten into the cloud computing business. The open-source application of choice for this is Hadoop.

The key issue, of course, is the parallelizability of the inherently serial code. But all the user need do is provide code to break the data into chunks, code to work on a chunk, and code to collect the outputs from the chunks back into the overall output of the program.

For this to work, the program's data usage pattern must have a simple, regular structure, as in these examples:

**Example 1:** Suppose we wish to list all the words used in a file, together with the counts of the numbers of instances of the words. If we have 100000 lines in the file and 10 processors, we could divide the file into chunks of 10000 lines each, have each processor run code to do the word counts in its chunk, and then combine the results.

**Example 2:** Suppose we wish to multiply an nx1 vector X by an nxn matrix A. Say n = 100000, and again we have 10 processors. We could divide A into chunks of 10000 rows each, have each processor multiply X by its chunk, and then combine the outputs.

To illustrate this, here is a pseudocode summary of a word-count program written in Python by Michael Noll; see `http://www.michael-noll.com/wiki/Writing_An_Hadoop_MapReduce_Program_In_Python`. Actually Hadoop is really written for Java applications. However, Hadoop can work with programs in any language under Hadoop's Streaming option, by reading from STDIN and writing to STDOUT. This does cause some slowdown in numeric programs, for the conversion of strings to numbers and vice versa.[4]

**mapper.py:**

---

[4]In the case of Python, we could also run Jython, a Python interpreter that produces Java byte code. Hadoop also offers communication via Unix pipes.

```
1  for each line in STDIN
2      break line into words, placed in wordarray
3      for each word in wordarray
4          print word, '1' to STDOUT  # we have found 1 instance of the word
```

**reducer.py:**

```
1  # dictionary will consist of (word,count) pairs
2  dictionary = empty
3  for each line in STDIN
4      split line into word, thiscount
5      if word not in dictionary:
6          add (word,thiscount) to dictionary
7      else
8          change (word,count) entry to (word,count+thiscount)
9  print dictionary to STDOUT
```

Note that these two user programs have nothing in them at all regarding parallelism. Instead, the process works as follows:

- the user provides Hadoop the original data file, by copying the file to Hadoop's own file system, the Hadoop Distributed File System (HDFS)

- the user provides Hadoop with the mapper and reducer programs; Hadoop runs several instances of each

- Hadoop forms chunks by forming groups of lines in the file

- Hadoop has each instance of the mapper program work on a chunk:

  ```
  mapper.py < chunk > outputchunk
  # output is replicated and sent to the various instances of reducer
  ```

- Hadoop runs

  ```
  reducer.py < outputchunk > myfinalchunk
  # in this way final output is distributed to the nodes in HDFS
  ```

In the matrix-multiply model, the mapper program would produce chunks of X, together with the corresponding row numbers. Then the reducer program would sort the rows by row number, and place the result in X.

Note too that by having the file in HDFS, we minimize communications costs in shipping the data. "Moving computation is cheaper than moving data."

Hadoop also incorporates rather sophisticated fault tolerance mechanisms. If a node goes down, the show goes on.

Note again that this works well only on problems of a certain structure. Also, some say that the idea has been overpromoted; see for instance "MapReduce: A Major Step Backwards," *The Database Column*, by Professor David DeWitt, `http://www.databasecolumn.com/2008/01/mapreduce-a-major-step-back.html`

# Chapter 6

# Introduction to MPI

MPI is the *de facto* standard for message-passing software.

## 6.1 Overview

### 6.1.1 History

Though (small) shared-memory machines have come down radically in price, to the point at which a dual-core PC is affordable in the home, historically shared-memory machines were available only to the "very rich"—large banks, national research labs and so on.

The first "affordable" message-machine type was the Hypercube, developed by a physics professor at Cal Tech. It consisted of a number of **processing elements** (PEs) connected by fast serial I/O cards. This was in the range of university departmental research labs. It was later commercialized by Intel and NCube.

Later, the notion of **networks of workstations** (NOWs) became popular. Here the PEs were entirely independent PCs, connected via a standard network. This was refined a bit, by the use of more suitable network hardware and protocols, with the new term being **clusters**.

All of this necessitated the development of standardized software tools based on a message-passing paradigm. The first popular such tool was Parallel Virtual Machine (PVM). It still has its adherents today, but has largely been supplanted by the Message Passing Interface (MPI).

MPI itself later became MPI 2. Our document here is intended mainly for the original.

### 6.1.2   Structure and Execution

MPI is merely a set of Application Programmer Interfaces (APIs), called from user programs written in C, C++ and other languages. It has many implementations, with some being open source and generic, while others are proprietary and fine-tuned for specific commercial hardware.

Suppose we have written an MPI program **x**, and will run it on four machines in a cluster. Each machine will be running its own copy of **x**. Official MPI terminology refers to this as four **processes**. Now that multicore machines are commonplace, one might indeed run two or more cooperating MPI processes—where now we use the term *processes* in the real OS sense—on the same multicore machine. In this document, we will tend to refer to the various MPI processes as **nodes**, with an eye to the cluster setting.

Though the nodes are all running the same program, they will likely be working on different parts of the program's data. This is called the Single Program Multiple Data (SPMD) model. This is the typical approach, but there could be different programs running on different nodes. Most of the APIs involve a node sending information to, or receiving information from, other nodes.

### 6.1.3   Implementations

Two of the most popular implementations of MPI are MPICH and LAM. MPICH offers more tailoring to various networks and other platforms, while LAM runs on networks. Introductions to MPICH and LAM can be found, for example, at `http://heather.cs.ucdavis.edu/~matloff/MPI/NotesMPICH.NM.html` and `http://heather.cs.ucdavis.edu/~matloff/MPI/NotesLAM.NM.html`, respectively.

### 6.1.4   Performance Issues

Mere usage of a parallel language on a parallel platform does not guarantee a performance improvement over a serial version of your program. The central issue here is the overhead involved in internode communication.

As of 2008, Infiniband, one of the fastest cluster networks commercially available, has a **latency** of about 1.0-3.0 microseconds, meaning that it takes the first bit of a packet that long to get from one node on an Infiniband switch to another. Comparing that to the nanosecond time scale of CPU speeds, one can see that the communications overhead can destroy a program's performance. And Ethernet is quite a bit slower than Infiniband.

Note carefully that latency is a major problem even if the **bandwidth**—the number of bits per second which are sent—is high. For this reason, it is quite possible that your parallel program may actually run more slowly than its serial version.

Of course, if your platform is a shared-memory multiprocessor (especially a multicore one, where communication between cores is particularly fast) and you are running all your MPI processor on that machine, the problem is less severe.

## 6.2 Running Example

### 6.2.1 The Algorithm

The code implements the Dijkstra algorithm for finding the shortest paths in an undirected graph. Pseudocode for the algorithm is

```
1   Done = {0}
2   NonDone = {1,2,...,N-1}
3   for J = 1 to N-1 Dist[J] = infinity'
4   Dist[0] = 0
5   for Step = 1 to N-1
6      find J such that Dist[J] is min among all J in NonDone
7      transfer J from NonDone to Done
8      NewDone = J
9      for K = 1 to N-1
10         if K is in NonDone
11            Dist[K] = min(Dist[K],Dist[NewDone]+G[NewDone,K])
```

At each iteration, the algorithm finds the closest vertex J to 0 among all those not yet processed, and then updates the list of minimum distances to each vertex from 0 by considering paths that go through J. Two obvious potential candidate part of the algorithm for parallelization are the "find J" and "for K" lines, and the above OpenMP code takes this approach.

### 6.2.2 The Code

```
1   // Dijkstra.c
2
3   // MPI example program:  Dijkstra shortest-path finder in a
4   // bidirectional graph; finds the shortest path from vertex 0 to all
5   // others
6
7   // command line arguments:  nv print dbg
8
9   // where:  nv is the size of the graph; print is 1 if graph and min
10  // distances are to be printed out, 0 otherwise; and dbg is 1 or 0, 1
11  // for debug
12
13  // node 0 will both participate in the computation and serve as a
14  // "manager"
15
16  #include <stdio.h>
```

```
17   #include <mpi.h>
18
19   #define MYMIN_MSG 0
20   #define OVRLMIN_MSG 1
21   #define COLLECT_MSG 2
22
23   // global variables (but of course not shared across nodes)
24
25   int nv,  // number of vertices
26       *notdone, // vertices not checked yet
27       nnodes,  // number of MPI nodes in the computation
28       chunk,  // number of vertices handled by each node
29       startv,endv,  // start, end vertices for this node
30       me,  // my node number
31       dbg;
32   unsigned largeint,  // max possible unsigned int
33            mymin[2],  // mymin[0] is min for my chunk,
34                       // mymin[1] is vertex which achieves that min
35            othermin[2],  // othermin[0] is min over the other chunks
36                          // (used by node 0 only)
37                          // othermin[1] is vertex which achieves that min
38            overallmin[2],  // overallmin[0] is current min over all nodes,
39                            // overallmin[1] is vertex which achieves that min
40            *ohd,  // 1-hop distances between vertices; "ohd[i][j]" is
41                   // ohd[i*nv+j]
42            *mind;  // min distances found so far
43
44   double T1,T2;  // start and finish times
45
46   void init(int ac, char **av)
47   {  int i,j,tmp; unsigned u;
48      nv = atoi(av[1]);
49      dbg = atoi(av[3]);
50      MPI_Init(&ac,&av);
51      MPI_Comm_size(MPI_COMM_WORLD,&nnodes);
52      MPI_Comm_rank(MPI_COMM_WORLD,&me);
53      chunk = nv/nnodes;
54      startv = me * chunk;
55      endv = startv + chunk - 1;
56      u = -1;
57      largeint = u >> 1;
58      ohd = malloc(nv*nv*sizeof(int));
59      mind = malloc(nv*sizeof(int));
60      notdone = malloc(nv*sizeof(int));
61      // random graph
62      // note that this will be generated at all nodes; could generate just
63      // at node 0 and then send to others, but faster this way
64      for (i = 0; i < nv; i++)
65         for (j = i; j < nv; j++)  {
66            if (j == i) ohd[i*nv+i] = 0;
67            else  {
68               ohd[nv*i+j] = rand() % 20;
69               ohd[nv*j+i] = ohd[nv*i+j];
70            }
71         }
72      for (i = 0; i < nv; i++)  {
73         notdone[i] = 1;
74         mind[i] = largeint;
```

```
75      }
76      mind[0] = 0;
77      while (dbg) ;  // stalling so can attach debugger
78   }
79
80   // finds closest to 0 among notdone, among startv through endv
81   void findmymin()
82   {  int i;
83      mymin[0] = largeint;
84      for (i = startv; i <= endv; i++)
85         if (notdone[i] && mind[i] < mymin[0])  {
86            mymin[0] = mind[i];
87            mymin[1] = i;
88         }
89   }
90
91   void findoverallmin()
92   {  int i;
93      MPI_Status status;  // describes result of MPI_Recv() call
94      // nodes other than 0 report their mins to node 0, which receives
95      // them and updates its value for the global min
96      if (me > 0)
97         MPI_Send(mymin,2,MPI_INT,0,MYMIN_MSG,MPI_COMM_WORLD);
98      else  {
99         // check my own first
100        overallmin[0] = mymin[0];
101        overallmin[1] = mymin[1];
102        // check the others
103        for (i = 1; i < nnodes; i++)  {
104           MPI_Recv(othermin,2,MPI_INT,i,MYMIN_MSG,MPI_COMM_WORLD,&status);
105           if (othermin[0] < overallmin[0])  {
106              overallmin[0] = othermin[0];
107              overallmin[1] = othermin[1];
108           }
109        }
110     }
111  }
112
113  void updatemymind()  // update my mind segment
114  {  // for each i in [startv,endv], ask whether a shorter path to i
115     // exists, through mv
116     int i, mv = overallmin[1];
117     unsigned md = overallmin[0];
118     for (i = startv; i <= endv; i++)
119        if (md + ohd[mv*nv+i] < mind[i])
120           mind[i] = md + ohd[mv*nv+i];
121  }
122
123  void disseminateoverallmin()
124  {  int i;
125     MPI_Status status;
126     if (me == 0)
127        for (i = 1; i < nnodes; i++)
128           MPI_Send(overallmin,2,MPI_INT,i,OVRLMIN_MSG,MPI_COMM_WORLD);
129     else
130        MPI_Recv(overallmin,2,MPI_INT,0,OVRLMIN_MSG,MPI_COMM_WORLD,&status);
131  }
132
```

```
133  void updateallmind()  // collects all the mind segments at node 0
134  {  int i;
135     MPI_Status status;
136     if (me > 0)
137        MPI_Send(mind+startv,chunk,MPI_INT,0,COLLECT_MSG,MPI_COMM_WORLD);
138     else
139        for (i = 1; i < nnodes; i++)
140           MPI_Recv(mind+i*chunk,chunk,MPI_INT,i,COLLECT_MSG,MPI_COMM_WORLD,
141              &status);
142  }
143
144  void printmind()  // partly for debugging (call from GDB)
145  {  int i;
146     printf("minimum distances:\n");
147     for (i = 1; i < nv; i++)
148        printf("%u\n",mind[i]);
149  }
150
151  void dowork()
152  {  int step,  // index for loop of nv steps
153         i;
154     if (me == 0) T1 = MPI_Wtime();
155     for (step = 0; step < nv; step++)  {
156        findmymin();
157        findoverallmin();
158        disseminateoverallmin();
159        // mark new vertex as done
160        notdone[overallmin[1]] = 0;
161        updatemymind(startv,endv);
162     }
163     updateallmind();
164     T2 = MPI_Wtime();
165  }
166
167  int main(int ac, char **av)
168  {  int i,j,print;
169     init(ac,av);
170     dowork();
171     print = atoi(av[2]);
172     if (print && me == 0)  {
173        printf("graph weights:\n");
174        for (i = 0; i < nv; i++)  {
175           for (j = 0; j < nv; j++)
176              printf("%u  ",ohd[nv*i+j]);
177           printf("\n");
178        }
179        printmind();
180     }
181     if (me == 0) printf("time at node 0: %f\n",(float)(T2-T1));
182     MPI_Finalize();
183  }
184
```

The various MPI functions will be explained in the next section.

### 6.2.3 Introduction to MPI APIs

#### 6.2.3.1 MPI_Init() and MPI_Finalize()

These are required for starting and ending execution of an MPI program. Their actions may be implementation-dependent. For instance, if our platform is an Ethernet-based cluster , **MPI_Init()** will probably set up the TCP/IP sockets via which the various nodes communicate with each other. On an Infiniband-based cluster, connections in the special Infiniband network protocol will be established. On a shared-memory multiprocessor, an implementation of MPI that is tailored to that platform would take very different actions.

#### 6.2.3.2 MPI_Comm_size() and MPI_Comm_rank()

In our function **init()** above, note the calls

```
MPI_Comm_size(MPI_COMM_WORLD,&nnodes);
MPI_Comm_rank(MPI_COMM_WORLD,&me);
```

The first call determines how many nodes are participating in our computation, placing the result in our variable **nnodes**. Here **MPI_COMM_WORLD** is our node group, termed a **communicator** in MPI parlance. MPI allows the programmer to subdivide the nodes into groups, to facilitate performance and clarity of code. Note that for some operations, such as barriers, the only way to apply the operation to a proper subset of all nodes is to form a group. The totality of all groups is denoted by **MPI_COMM_WORLD**. In our program here, we are not subdividing into groups.

The second call determines this node's ID number, called its **rank**, within its group. As mentioned earlier, even though the nodes are all running the same program, they are typically working on different parts of the program's data. So, the program needs to be able to sense which node it is running on, so as to access the appropriate data. Here we record that information in our variable **me**.

#### 6.2.3.3 MPI_Send()

To see how MPI's basic send function works, consider our line above,

```
MPI_Send(mymin,2,MPI_INT,0,MYMIN_MSG,MPI_COMM_WORLD);
```

Let's look at the arguments:

**mymin:** We are sending a set of bytes. This argument states the address at which these bytes begin.

**2, MPI_INT:** This says that our set of bytes to be sent consists of 2 objects of type **MPI_INT**. That means 8 bytes on 32-bit machines, so why not just collapse these two arguments to one, namely the number 8? Why did the designers of MPI bother to define data types? The answer is that we want to be able to run MPI on a heterogeneous set of machines, with MPI serving as the "broker" between them in case different architectures among those machines handle data differently.

First of all, there is the issue of **endianness**. Intel machines, for instance, are **little-endian**, which means that the least significant byte of a memory word has the smallest address among bytes of the word. Sun SPARC chips, on the other hand, are **big-endian**, with the opposite storage scheme. If our set of nodes included machines of both types, straight transmission of sequences of 8 bytes might mean that some of the machines literally receive the data backwards!

Secondly, these days 64-bit machines are becoming more and more common. Again, if our set of nodes were to include both 32-bit and 64-bit words, some major problems would occur if no conversion were done.

**0:** We are sending to node 0.

**MYMIN_MSG:** This is the message type, programmer-defined in our line

```
#define MYMIN_MSG 0
```

Receive calls, described in the next section, can ask to receive only messages of a certain type.

**MPI_COMM_WORLD:** This is the node group to which the message is to be sent. Above, where we said we are sending to node 0, we technically should say we are sending to node 0 within the group **MPI_COMM_WORLD**.

### 6.2.3.4   MPI_Recv()

Let's now look at the arguments for a basic receive:

```
MPI_Recv(othermin,2,MPI_INT,i,MYMIN_MSG,MPI_COMM_WORLD,&status);
```

**othermin:** The received message is to be placed at our location **othermin**.

**2,MPI_INT:** Two objects of **MPI_INT** type are to be received.

**i:** Receive only messages of from node **i**. If we did not care what node we received a message from, we could specify the value **MPI_ANY_SOURCE**.

**MYMIN_MSG:** Receive only messages of type **MYMIN_MSG**. If we did not care what type of message we received, we would specify the value **MPI_ANY_TAG**.

**MPI_COMM_WORLD:** Group name.

**status:** Recall our line

```
MPI_Status status;  // describes result of MPI_Recv() call
```

The type is an MPI **struct** containing information about the received message. Its primary fields of interest are **MPI_SOURCE**, which contains the identity of the sending node, and **MPI_TAG**, which contains the message type. These would be useful if the receive had been done with **MPI_ANY_SOURCE** or **MPI_ANY_TAG**; the status argument would then tell us which node sent the message and what type the message was.

## 6.3 Collective Communications

MPI features a number of **collective communication** capabilities, a number of which are used in the following refinement of our Dijkstra program:

### 6.3.1 Example

```
1   // Dijkstra.coll1.c
2
3   // MPI example program:  Dijkstra shortest-path finder in a
4   // bidirectional graph; finds the shortest path from vertex 0 to all
5   // others; this version uses collective communication
6
7   // command line arguments:  nv print dbg
8
9   // where:  nv is the size of the graph; print is 1 if graph and min
10  // distances are to be printed out, 0 otherwise; and dbg is 1 or 0, 1
11  // for debug
12
13  // node 0 will both participate in the computation and serve as a
14  // "manager"
15
16  #include <stdio.h>
17  #include <mpi.h>
18
19  // global variables (but of course not shared across nodes)
20
21  int nv,  // number of vertices
22      *notdone, // vertices not checked yet
23      nnodes,  // number of MPI nodes in the computation
24      chunk,  // number of vertices handled by each node
25      startv,endv,  // start, end vertices for this node
26      me,  // my node number
27      dbg;
28  unsigned largeint,  // max possible unsigned int
29           mymin[2],  // mymin[0] is min for my chunk,
30                      // mymin[1] is vertex which achieves that min
31           overallmin[2],  // overallmin[0] is current min over all nodes,
```

```
32                             // overallmin[1] is vertex which achieves that min
33             *ohd,   // 1-hop distances between vertices; "ohd[i][j]" is
34                   // ohd[i*nv+j]
35             *mind;  // min distances found so far
36
37  double T1,T2;  // start and finish times
38
39  void init(int ac, char **av)
40  {  int i,j,tmp; unsigned u;
41     nv = atoi(av[1]);
42     dbg = atoi(av[3]);
43     MPI_Init(&ac,&av);
44     MPI_Comm_size(MPI_COMM_WORLD,&nnodes);
45     MPI_Comm_rank(MPI_COMM_WORLD,&me);
46     chunk = nv/nnodes;
47     startv = me * chunk;
48     endv = startv + chunk - 1;
49     u = -1;
50     largeint = u >> 1;
51     ohd = malloc(nv*nv*sizeof(int));
52     mind = malloc(nv*sizeof(int));
53     notdone = malloc(nv*sizeof(int));
54     // random graph
55     // note that this will be generated at all nodes; could generate just
56     // at node 0 and then send to others, but faster this way
57     for (i = 0; i < nv; i++)
58        for (j = i; j < nv; j++)  {
59           if (j == i) ohd[i*nv+i] = 0;
60           else  {
61              ohd[nv*i+j] = rand() % 20;
62              ohd[nv*j+i] = ohd[nv*i+j];
63           }
64        }
65     for (i = 0; i < nv; i++)  {
66        notdone[i] = 1;
67        mind[i] = largeint;
68     }
69     mind[0] = 0;
70     while (dbg) ;  // stalling so can attach debugger
71  }
72
73  // finds closest to 0 among notdone, among startv through endv
74  void findmymin()
75  {  int i;
76     mymin[0] = largeint;
77     for (i = startv; i <= endv; i++)
78        if (notdone[i] && mind[i] < mymin[0])  {
79           mymin[0] = mind[i];
80           mymin[1] = i;
81        }
82  }
83
84  void updatemymind()  // update my mind segment
85  {  // for each i in [startv,endv], ask whether a shorter path to i
86     // exists, through mv
87     int i, mv = overallmin[1];
88     unsigned md = overallmin[0];
89     for (i = startv; i <= endv; i++)
```

```
90         if (md + ohd[mv*nv+i] < mind[i])
91            mind[i] = md + ohd[mv*nv+i];
92  }
93
94  void printmind()  // partly for debugging (call from GDB)
95  {  int i;
96     printf("minimum distances:\n");
97     for (i = 1; i < nv; i++)
98        printf("%u\n",mind[i]);
99  }
100
101  void dowork()
102  {  int step,  // index for loop of nv steps
103         i;
104     if (me == 0) T1 = MPI_Wtime();
105     for (step = 0; step < nv; step++)  {
106        findmymin();
107        MPI_Reduce(mymin,overallmin,1,MPI_2INT,MPI_MINLOC,0,MPI_COMM_WORLD);
108        MPI_Bcast(overallmin,1,MPI_2INT,0,MPI_COMM_WORLD);
109        // mark new vertex as done
110        notdone[overallmin[1]] = 0;
111        updatemymind(startv,endv);
112     }
113     // now need to collect all the mind values from other nodes to node 0
114     MPI_Gather(mind+startv,chunk,MPI_INT,mind,chunk,MPI_INT,0,MPI_COMM_WORLD);
115     T2 = MPI_Wtime();
116  }
117
118  int main(int ac, char **av)
119  {  int i,j,print;
120     init(ac,av);
121     dowork();
122     print = atoi(av[2]);
123     if (print && me == 0)  {
124        printf("graph weights:\n");
125        for (i = 0; i < nv; i++)  {
126           for (j = 0; j < nv; j++)
127              printf("%u  ",ohd[nv*i+j]);
128           printf("\n");
129        }
130        printmind();
131     }
132     if (me == 0) printf("time at node 0: %f\n",(float)(T2-T1));
133     MPI_Finalize();
134  }
```

The new calls will be explained in the next section.

### 6.3.2  MPI_Bcast()

In our original Dijkstra example, we had a loop

```
for (i = 1; i < nnodes; i++)
   MPI_Send(overallmin,2,MPI_INT,i,OVRLMIN_MSG,MPI_COMM_WORLD);
```

in which node 0 sends to all other nodes. We can replace this by

```
MPI_Bcast(overallmin,2,MPI_INT,0,MPI_COMM_WORLD);
```

In English, this call would say,

> At this point all nodes participate in a broadcast operation, in which node 0 sends 2 objects of type **MPI_INT**. The source of the data will be located at address **overallmin** at node 0, and the other nodes will receive the data at a location of that name.

Note my word "participate" above. The name of the function is "broadcast," which makes it sound like only node 0 executes this line of code, which is not the case; all the nodes in the group (in this case that means all nodes in our entire computation) execute this line. The only difference is the action; most nodes participate by receiving, while node 0 participates by sending.

Why might this be preferable to using an explicit loop? First, it would obviously be much clearer. That makes the program easier to write, easier to debug, and easier for others (and ourselves, later) to read.

But even more importantly, using the broadcast may improve performance. We may, for instance, be using an implementation of MPI which is tailored to the platform on which we are running MPI. If for instance we are running on a network designed for parallel computing, such as Myrinet or Infiniband, an optimized broadcast may achieve a much higher performance level than would simply a loop with individual send calls. On a shared-memory multiprocessor system, special machine instructions specific to that platform's architecture can be exploited, as for instance IBM has done for its shared-memory machines. Even on an ordinary Ethernet, one could exploit Ethernet's own broadcast mechanism, as had been done for PVM, a system like MPI (G. Davies and N. Matloff, Network-Specific Performance Enhancements for PVM, *Proceedings of the Fourth IEEE International Symposium on High-Performance Distributed Computing*, 1995, 205-210).

### 6.3.2.1    MPI_Reduce()/MPI_Allreduce()

Look at our call

```
MPI_Reduce(mymin,overallmin,1,MPI_2INT,MPI_MINLOC,0,MPI_COMM_WORLD);
```

above. In English, this would say,

> At this point all nodes in this group participate in a "reduce" operation. The type of reduce operation is **MPI_MINLOC**, which means that the minimum value among the nodes will be

computed, and the index attaining that minimum will be recorded as well. Each node contributes a value to be checked, and an associated index, from a location **mymin** in their programs; the type of the pair is **MPI_2INT**. The overall min value/index will be computed by combining all of these values at node 0, where they will be placed at a location **overallmin**.

MPI also includes a function **MPI_Allreduce()**, which does the same operation, except that instead of just depositing the result at one node, it does so at all nodes. So for instance our code above,

```
MPI_Reduce(mymin,overallmin,1,MPI_2INT,MPI_MINLOC,0,MPI_COMM_WORLD);
MPI_Bcast(overallmin,1,MPI_2INT,0,MPI_COMM_WORLD);
```

could be replaced by

```
MPI_Allreduce(mymin,overallmin,1,MPI_2INT,MPI_MINLOC,MPI_COMM_WORLD);
```

Again, these can be optimized for particular platforms.

### 6.3.2.2   MPI_Gather()/MPI_Allgather()

A classical approach to parallel computation is to first break the data for the application into chunks, then have each node work on its chunk, and then gather all the processed chunks together at some node. The MPI function **MPI_Gather()** does this.

In our program above, look at the line

```
MPI_Gather(mind+startv,chunk,MPI_INT,mind,chunk,MPI_INT,0,MPI_COMM_WORLD);
```

In English, this says,

At this point all nodes participate in a gather operation, in which each node contributes data, consisting of **chunk** number of MPI integers, from a location **mind+startv** in its program. All that data is strung together and deposited at the location **mind** in the program running at node 0.

There is also **MPI_Allgather()**, which places the result at all nodes, not just one.

### 6.3.2.3   The MPI_Scatter()

This is the opposite of **MPI_Gather()**, i.e. it breaks long data into chunks which it parcels out to individual nodes.

Here is MPI code to count the number of edges in a directed graph. (A link from i to j does not necessarily imply one from j to i.) In the context here, **me** is the node's rank; **nv** is the number of vertices; **oh** is the one-hop distance matrix; and **nnodes** is the number of MPI processes. At the beginning only the process of rank 0 has a copy of **oh**, but it sends that matrix out in chunks to the other nodes, each of which stores its chunk in an array **ohchunk**.

```
1   MPI_Scatter(oh, nv*nv, MPI_INT, ohchunk, nv/nnodes, MPI_INT, 0,
2   MPI_COMM_WORLD);
3   mycount = 0;
4   for (i = 0; i < nv*nv/nnodes)
5       if (ohchunk[i] != 0) mycount++;
6   MPI_Reduce(&mycount,&numedge,1,MPI_INT,MPI_SUM,0,MPI_COMM_WORLD);
7   if (me == 0) printf("there are %d edges\n",numedge);
```

### 6.3.2.4   The MPI_Barrier()

This implements a barrier for a given communicator. The name of the communicator is the sole argument for the function.

Explicit barriers are less common in message-passing programs than in the shared-memory world.

### 6.3.3   Creating Communicators

Again, a communicator is a subset (either proper or improper) of all of our nodes. MPI includes a number of functions for use in creating communicators. Some set up a virtual "topology" among the nodes.

For instance, many physics problems consist of solving differential equations in two- or three-dimensional space, via approximation on a grid of points. In two dimensions, groups may consists of rows in the grid.

We will not pursue this further here.

## 6.4   Buffering, Synchrony and Related Issues

As noted several times so far, interprocess communication in parallel systems can be quite expensive in terms of time delay. In this section we will consider some issues which can be extremely important in this regard.

### 6.4.1 Buffering, Etc.

To understand this point, first consider situations in which MPI is running on some network, under the TCP/IP protocol. Say an MPI program at node A is sending to one at node B.

It is extremely import to keep in mind the levels of abstraction here. The OS's TCP/IP stack is running at the Session, Transport and Network layers of the network. MPI—meaning the MPI internals—is running above the TCP/IP stack, in the Application layers at A and B. And the MPI user-written application could be considered to be running at a "Super-application" layer, since it calls the MPI internals. (From here on, we will refer to the MPI internals as simply "MPI.")

MPI at node A will have set up a TCP/IP socket to B during the user program's call to **MPI_Init()**. The other end of the socket will be a corresponding one at B. This setting up of this socket pair as establishing a **connection** between A and B. When node A calls **MPI_Send()**, MPI will write to the socket, and the TCP/IP stack will transmit that data to the TCP/IP socket at B. The TCP/IP stack at B will then send whatever bytes come in to MPI at B.

Now, it is important to keep in mind that in TCP/IP the totality of bytes sent by A to B during lifetime of the connection is considered one long message. So for instance if the MPI program at A calls **MPI_Send()** five times, the MPI internals will write to the socket five times, but the bytes from those five messages will not be perceived by the TCP/IP stack at B as five messages, but rather as just one long message (in fact, only part of one long message, since more may be yet to come).

MPI at B continually reads that "long message" and breaks it back into MPI messages, keeping them ready for calls to **MPI_Recv()** from the MPI application program at B. Note carefully that phrase, *keeping them ready*; it refers to the fact that the order in which the MPI application program requests those messages may be different from the order in which they arrive.

On the other hand, looking again at the TCP/IP level, even though all the bytes sent are considered one long message, it will physically be sent out in pieces. These pieces don't correspond to the pieces written to the socket, i.e. the MPI messages. Rather, the breaking into pieces is done for the purpose of **flow control**, meaning that the TCP/IP stack at A will not send data to the one at B if the OS at B has no room for it. The **buffer** space the OS at B has set up for receiving data is limited. As A is sending to B, the TCP layer at B is telling its counterpart at A when A is allowed to send more data.

Think of what happens the MPI application at B calls **MPI_Recv()**, requesting to receive from A, with a certain tag T. Say the first argument is named **x**, i.e. the data to be received is to be deposited at **x**. If MPI sees that it already has a message of tag T, it will have its **MPI_Recv()** function return the message to the caller, i.e. to the MPI application at B. **If no such message has arrived yet, MPI won't return to the caller yet, and thus the caller blocks.**

**MPI_Send()** can block too. If the platform and MPI implementation is that of the TCP/IP network context described above, then the send call will return when its call to the OS' **write()** (or equivalent, depending on OS) returns, but that could be delayed if the OS' buffer space is full. On the other hand, another implemen-

tation could require a positive response from B before allowing the send call to return.

Note that buffering slows everything down. In our TCP scenario above, **MPI Recv()** at B must copy messages from the OS' buffer space to the MPI application program's program variables, e.g. **x** above. This is definitely a blow to performance. That in fact is why networks developed specially for parallel processing typically include mechanisms to avoid the copying. Infiniband, for example, has a Remote Direct Memory Access capability, meaning that A can write directly to **x** at B. Of course, if our implementation uses **synchronous** communication, with A's send call not returning until A gets a response from B, we must wait even longer.

Technically, the MPI standard states that **MPI Send(x,...)** will return only when it is safe for the application program to write over the array which it is using to store its message, i.e. **x**. As we have seen, there are various ways to implement this, with performance implications. Similarly, **MPI Recv(y,...)** will return only when it is safe to read **y**.

### 6.4.2   Safety

With **synchronous** communication, deadlock is a real risk. Say A wants to send two messages to B, of types U and V, but that B wants to receive V first. Then A won't even get to send V, because in preparing to send U it must wait for a notice from B that B wants to read U—a notice which will never come, because B sends such a notice for V first. This would not occur if the communication were asynchronous.

But beyond formal deadlock, programs can fail in other ways, even with buffering, as buffer space is always by nature finite. A program can fail if it runs out of buffer space, either at the sender or the receiver. See `www.llnl.gov/computing/tutorials/mpi_performance/samples/unsafe.c` for an example of a test program which demonstrates this on a certain platform, by deliberating overwhelming the buffers at the receiver.

In MPI terminology, asynchronous communication is considered **unsafe**. The program may run fine on most systems, as most systems are buffered, but fail on some systems. Of course, as long as you know your program won't be run in nonbuffered settings, it's fine, and since there is potentially such a performance penalty for doing things synchronously, most people are willing to go ahead with their "unsafe" code.

### 6.4.3   Living Dangerously

If one is sure that there will be no problems of buffer overflow and so on, one can use variant send and receive calls provided by MPI, such as **MPI Isend()** and **MPI Irecv()**. The key difference between them and **MPI Send()** and **MPI Recv()** is that they return immediately, and thus are termed **nonblocking**. Your code can go on and do other things, not having to wait.

This does mean that at A you cannot touch the data you are sending until you determine that it has either been

buffered somewhere or has reached **x** at B. Similarly, at B you can't use the data at **x** until you determine that it has arrived. Such determinations can be made via **MPI_Wait()**. In other words, you can do your send or receive, then perform some other computations for a while, and then call **MPI_Wait()** to determine whether you can go on. Or you can call **MPI_Probe()** to ask whether the operation has completed yet.

### 6.4.4 Safe Exchange Operations

In many applications A and B are swapping data, so both are sending and both are receiving. This too can lead to deadlock. An obvious solution would be, for instance, to have the lower-rank node send first and the higher-rank node receive first.

But a more convenient, safer and possibly faster alternative would be to use MPI's **MPI_Sendrecv()** function. Its prototype is

```
intMPI_Sendrecv_replace(void* buf, int count, MPI_Datatype datatype,
   int dest, int sendtag, int source, int recvtag, MPI_Comm comm,
   MPI_Status *status)
```

Note that the sent and received messages can be of different lengths and can use different tags.

## 6.5 Use of MPI from Other Languages

MPI is a vehicle for parallelizing C/C++, but some clever people have extended the concept to other languages, such as the cases of Python and R that we treat here.

### 6.5.1 Python: pyMPI

A number of interfaces of Python to MPI have been developed.[1] A well-known example is pyMPI, developed by a PhD graduate in computer science in UCD, Patrick Miller.

One writes one's pyMPI code, say in **x.py**, by calling pyMPI versions of the usual MPI routines. To run the code, one then runs MPI on the program **pyMPI** with **x.py** as a command-line argument.

Python is a very elegant language, and pyMPI does a nice job of elegantly interfacing to MPI. Following is a rendition of Quicksort in pyMPI. Don't worry if you haven't worked in Python before; the "non-C-like" Python constructs are explained in comments at the end of the code.

---

[1]If you are not familiar with Python, I have a quick tutorial at `http://heather.cs.ucdavis.edu/~matloff/python.html`.

```
1   # a type of quicksort; break array x (actually a Python "list") into
2   # p quicksort-style piles, based # on comparison with the first p-1
3   # elements of x, where p is the number # of MPI nodes; the nodes sort
4   # their piles, then return them to node 0, # which strings them all
5   # together into the final sorted array
6
7   import mpi  # load pyMPI module
8
9   # makes npls quicksort-style piles
10  def makepiles(x,npls):
11     pivot = x[:npls]  # we'll use the first npls elements of x as pivots,
12                       # i.e. we'll compare all other elements of x to these
13     pivot.sort()  # sort() is a member function of the Python list class
14     pls = []  # initialize piles list to empty
15     lp = len(pivot)  # length of the pivot array
16     # pls will be a list of lists, with the i-th list in pls storing the
17     # i-th pile; the i-th pile will start with ID i (to enable
18     # identification later on) and pivot[i]
19     for i in range(lp):  # i = 0,1,...lp-1
20        pls.append([i,pivot[i]])  # build up array via append() member function
21     pls.append([lp])
22     for xi in x[npls:]:  # now place each element in the rest of x into
23                          # its proper pile
24        for j in range(lp):  # j = 0,1,...,lp-1
25           if xi <= pivot[j]:
26              pls[j].append(xi)
27              break
28           elif j == lp-1: pls[lp].append(xi)
29     return pls
30
31  def main():
32     if mpi.rank == 0:  # analog of calling MPI_Rank()
33        x = [12,5,13,61,9,6,20,1]  # small test case
34        # divide x into piles to be disbursed to the various nodes
35        pls = makepiles(x,mpi.size)
36     else:  # all other nodes set their x and pls to empty
37        x = []
38        pls = []
39     mychunk = mpi.scatter(pls)  # node 0 (not an explicit argument) disburses
40                                 # pls to the nodes, each of which receives
41                                 # its chunk in its mychunk
42     newchunk = []  # will become sorted version of mychunk
43     for pile in mychunk:
44        # I need to sort my chunk but most remove the ID first
45        plnum = pile.pop(0)  # ID
46        pile.sort()
47        # restore ID
48        newchunk.append([plnum]+pile)  # the + is array concatenation
49     # now everyone sends their newchunk lists, which node 0 (again an
50     # implied argument) gathers together into haveitall
51     haveitall = mpi.gather(newchunk)
52     if mpi.rank == 0:
53        haveitall.sort()
54        # string all the piles together
55        sortedx = [z for q in haveitall for z in q[1:]]
56        print sortedx
57
58  # common idiom for launching a Python program
```

```
59   if __name__ == '__main__': main()
```

Some examples of use of other MPI functions:

```
mpi.send(mesgstring,destnodenumber)
(message,status) = mpi.recv()  # receive from anyone
print message
(message,status) = mpi.recv(3)  # receive only from node 3
(message,status) = mpi.recv(3,ZMSG)  # receive only message type ZMSG,
                                # only from node 3
(message,status) = mpi.recv(tag=ZMSG)  # receive from anyone, but
                                  # only message type ZMSG
```

## 6.5.2 R

### 6.5.2.1 Rmpi

The Rmpi package provides an interface from R to MPI, much like that of pyMPI.[2]

So, we run Rmpi on top of MPI. Even nicer, we can run the Snow package on top of Rmpi! Snow provides a higher-level interface, which is very convenient.

**Installation:**

Say you want to install in the directory **/a/b/c/**. The easiest way to do so is

```
> install.packages("Rmpi","/a/b/c/")
```

This will install Rmpi in the directory **/a/b/c/Rmpi**.

You'll need to arrange for the directory **/a/b/c** (not **/a/b/c/Rmpi**) to be added to your R library search path. I recommend placing a line

```
.libPaths("/a/b/c/")
```

in a file **.Rprofile** in your home directory.

**Usage:**

Fire up MPI, and then in R load in Rmpi, by typing

```
> library(Rmpi)
```

---

[2]R is a widely-used language for statistics/data. I have a programmer's tutorial for it at `http://heather.cs.ucdavis.edu/~matloff/R/RProg.pdf`.

Then start Rmpi:

```
> mpi.spawn.Rslaves()
```

This will start R on all machines in the group you started MPI on. Optionally, you can specify fewer machines via the named argument **nslaves**.

The first time you do this, try this test:

```
mpi.remote.exec(paste("I am",mpi.comm.rank(),"of",mpi.comm.size()))
```

The available functions are similar to those of pyMPI, such as

- **mpi.comm.size():**

  Returns the number of MPI processes, including the master that spawned the other processes.

- **mpi.comm.rank():**

  Returns the rank of the process that executes it.

- **mpi.send()**, mpi.recv():

  The usual send/receive operations.

- **mpi.bcast(), mpi.scatter(), mpi.gather():**

  The usual broadcast, scatter and gather operations.

- Etc.

Details are available at:

- `http://cran.r-project.org/web/packages/Rmpi/index.html`

  Site for download of package and manual.

- `http://ace.acadiau.ca/math/ACMMaC/Rmpi/sample.html`

  Nice tutorial.

But we forego details here, as Snow provides a nicer programmer interface, to be described next.

### 6.5.2.2 Snow

Snow runs on top of Rmpi (or directly via sockets), allowing the programmer to more conveniently express the parallel disposition of work.

For instance, just as the ordinary R function **apply()** applies the same function to all rows of a matrix, the Snow function **parApply()** does that in parallel, across multiple machines; different machines will work on different rows.

**Installation:**

Follow the same pattern as described above for Rmpi. If you plan to have Snow run on top of Rmpi, you'll of course need the latter too.

**Usage:**

Make sure Snow is in your library path (see material on Rmpi above).

Load Snow:

```
> library(snow)
```

One then sets up a cluster, by calling the Snow function **makeCluster()**. The named argument **type** of that function indicates the networking platform, e.g. "MPI," "PVM" or "SOCK." The last indicates that you wish Snow to run on TCP/IP sockets that it creates itself, rather than going through MPI.

It is generally preferable to use MPI for Snow. This provides more flexibility, as one's code could include calls to both Snow functions and MPI (i.e. Rmpi) functions. In the examples here, I used "SOCK," on machines named **pc48** and **pc49**, setting up the cluster this way:

```
> cls <- makeCluster(type="SOCK",c("pc48","pc49"))
```

For MPI or PVM, one specifies the number of nodes to create, rather than specifying the nodes themselves.

Note that the above R code sets up worker nodes at the machines named **pc48** and **pc49**; these are in addition to the master node, which is the machine on which that R code is executed

There are various other optional arguments. One you may find useful is **outfile**, which records the result of the call in the file **outfile**. This can be helpful if the call fails.

Let's look at a simple example of multiplication of a vector by a matrix. We set up a test matrix:

```
> a <- matrix(c(1,2,3,4,5,6,7,8,9,10,11,12),nrow=6)
> a
     [,1] [,2]
```

```
[1,]    1    7
[2,]    2    8
[3,]    3    9
[4,]    4    10
[5,]    5    11
[6,]    6    12
```

We will multiply the vector $(1,1)^T$ (T meaning transpose) by our matrix **a**, by defining a dot product function:

```
> dot <- function(x,y) {return(x%*%y)}
```

Let's test it using the ordinary **apply()**:

```
> apply(a,1,dot,c(1,1))
[1]   8 10 12 14 16 18
```

To review your R, note that this applies the function **dot()** to each row (indicated by the 1, with 2 meaning column) of **a** playing the role of the first argument to **dot()**, and with c(1,1) playing the role of the second argument.

Now let's do this in parallel, across our two machines in our cluster **cls**:

```
> parApply(cls,a,1,dot,c(1,1))
[1]   8 10 12 14 16 18
```

The function **clusterCall(cls,f,args)** applies the given function **f()** at each worker node in the cluster **cls**, using the arguments provided in **args**.

The function **clusterExport(cls,varlist)** copies the variables in the list **varlist** to each worker in the cluster **cls**. You can use this to avoid constant shipping of large data sets from the master to the workers; you just do so once, using **clusterExport()** on the corresponding variables, and then access those variables as global. For instance:

```
> z <- function() return(x)
> x <- 5
> y <- 12
> clusterExport(cls,list("x","y"))
> clusterCall(cls,z)
[[1]]
[1] 5

[[2]]
[1] 5
```

The function **clusterEvalQ(cls,expression)** runs **expression** at each worker node in **cls**. Continuing the above example, we have

```
> clusterEvalQ(cls,x <- x+1)
[[1]]
[1] 6

[[2]]
[1] 6

> clusterCall(cls,z)
[[1]]
[1] 6

[[2]]
[1] 6

> x
[1] 5
```

Note that **x** still has its original version back at the master.

The function **clusterApply(cls,individualargs,f,commonargsgohere)** runs **f()** at each worker node in **cls**, with arguments as follows. The first argument to **f()** for worker i is the $i^{th}$ element of the list **individualargs**, i.e. **individualargs[[i]]**, and optionally one can give additional arguments for **f()** following **f()** in the argument list for **clusterApply()**.

Here for instance is how we can assign an ID to each worker node, like MPI **rank**:[3]

```
> myid <- 0
> clusterExport(cls,"myid")
> setid <- function(i) {myid <<- i}  # note superassignment operator
> clusterApply(cls,1:2,setid)
[[1]]
[1] 1

[[2]]
[1] 2

> clusterCall(cls,function() {return(myid)})
[[1]]
[1] 1

[[2]]
[1] 2
```

Don't forget to stop your clusters before exiting R, by calling **stopCluster()clustername**.

There are various other useful Snow functions. See the user's manual for details.

---

[3] I don't see a provision in Snow itself that does this.

**To learn more about Snow:**

I recommend the following Web pages:

- `http://cran.cnr.berkeley.edu/web/packages/snow/index.html`
  CRAN page for Snow; the package and the manual are here.

- `http://www.bepress.com/cgi/viewcontent.cgi?article=1016&context=uwbiostat`
  A research paper.

- `http://www.cs.uiowa.edu/˜luke/R/cluster/cluster.html`
  Brief intro by the author.

- `http://www.sfu.ca/˜sblay/R/snow.html#clusterCall`
  Examples, short but useful.

# Chapter 7

# Introduction to Parallel Matrix Operations

## 7.1 Overview

In the early days parallel processing was mostly used in physics problems. Typical problems of interest would be grid computations such as the heat equation, matrix multiplication, matrix inversion (or equivalent operations) and so on. These matrices are not those little 3x3 things you worked with in your linear algebra class. In parallel processing applications of matrix algebra, our matrices can have thousands of rows and columns, or even larger.

The range of applications of parallel processing is of course far broader today. In many of these applications, problems which at first glance seem not to involve matrices, actually do have matrix solutions. An example in graph theory is the following.

Let n denote the number of vertices in the graph. Define the graph's **adjacency matrix** A to be the n x n matrix whose element (i,j) is equal to 1 if there is an edge connecting vertices i an j (i.e. i and j are "adjacent"), and 0 otherwise. The corresponding **reachability matrix** R has its (i,j) element equal to 1 if there is some path from i to j, and 0 otherwise.

One can prove that

$$R = b[(I + A)^{n-1}], \tag{7.1}$$

where I is the identity matrix and the function b() ('b' for "boolean") is applied elementwise to its matrix argument, replacing each nonzero element by 1 while leaving the elements which are 0 unchanged. The graph is connected if and only if all elements of R are 1s.

So, the original graph connectivity problem reduces to a matrix problem.

## 7.2   Message-Passing Setting

The algorithms presented in this introduction will mainly be written from a message-passing point of view, assuming MPI for concreteness. Adaptations to the shared-memory paradigm will be discussed in Section 7.6.

## 7.3   Partitioned Matrices

Parallel processing of course relies on finding a way to partition the work to be done. In the matrix algorithm case, this is often done by dividing a matrix into blocks.

For example, let

$$A = \begin{pmatrix} 1 & 5 & 12 \\ 0 & 3 & 6 \\ 4 & 8 & 2 \end{pmatrix} \tag{7.2}$$

and

$$B = \begin{pmatrix} 0 & 2 & 5 \\ 0 & 9 & 10 \\ 1 & 1 & 2 \end{pmatrix}, \tag{7.3}$$

so that

$$C = AB = \begin{pmatrix} 12 & 59 & 79 \\ 6 & 33 & 42 \\ 2 & 82 & 104 \end{pmatrix}. \tag{7.4}$$

We could partition A as

$$A = \begin{pmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{pmatrix}, \tag{7.5}$$

where

$$A_{00} = \begin{pmatrix} 1 & 5 \\ 0 & 3 \end{pmatrix}, \tag{7.6}$$

$$A_{01} = \begin{pmatrix} 12 \\ 6 \end{pmatrix}, \tag{7.7}$$

$$A_{10} = \begin{pmatrix} 4 & 8 \end{pmatrix} \tag{7.8}$$

and

$$A_{11} = \begin{pmatrix} 2 \end{pmatrix}. \tag{7.9}$$

Similarly we would partition B and C into blocks of the same size as in A,

$$B = \begin{pmatrix} B_{00} & B_{01} \\ B_{10} & B_{11} \end{pmatrix} \tag{7.10}$$

and

$$C = \begin{pmatrix} C_{00} & C_{01} \\ C_{10} & C_{11} \end{pmatrix}, \tag{7.11}$$

so that for example

$$B_{10} = \begin{pmatrix} 1 & 1 \end{pmatrix}. \tag{7.12}$$

The key point is that multiplication still works if we pretend that those submatrices are numbers. For example, that would give the relation

$$C_{00} = A_{00}B_{00} + A_{01}B_{10}, \tag{7.13}$$

which the reader should verify really is correct.

## 7.4 Matrix Multiplication

Let's suppose for the sake of simplicity that each of the matrices is of dimensions nxn. Let p denote the number of MPI nodes.

We assume that the matrices are **dense**, meaning that most of their entries are nonzero. This is in contrast to **sparse** matrices, such as **tridiagonal** matrices, in which the only nonzero elements are either on the diagonal

or on subdiagonals just below or above the diagonal. Clearly we would use a different type of algorithm for sparse matrices.

The obvious plan of attack here is to break the matrices into blocks, and then assign different blocks to different MPI nodes. Assume that $\sqrt{p}$ evenly divides n, and partition each matrix into submatrices of size $n/\sqrt{p}$ x $n/\sqrt{p}$. In other words, each matrix will be divided into m rows and m columns of blocks, where $m = n/\sqrt{p}$.

One of the conditions assumed here is that the matrices A and B are stored in a distributed manner across the nodes. As you'll see, the algorithms then have the nodes passing blocks among themselves.

### 7.4.1   Fox's Algorithm

Consider the node that has the responsibility of calculating block (i,j) of the product C, which it calculates as

$$A_{i0}B_{0j} + A_{i1}B_{1j} + ... + A_{ii}B_{ij} + ... + A_{i,m-1}B_{m-1,j} \tag{7.14}$$

Rearrange this as

$$A_{ii}B_{ij} + A_{i,i+1}B_{i+1j} + ... + A_{i,m-1}B_{m-1,j} + A_{i0}B_{0j} + A_{i1}B_{1j} + ... + A_{i,i-1}B_{i-1,j} \tag{7.15}$$

Written more compactly, this is

$$\sum_{k=0}^{m-1} A_{i,(i+k)mod\ m}B_{(i+k)mod\ m,j} \tag{7.16}$$

The order of summation in this rearrangement will be the actual order of computation.

The algorithm is then as follows. The node which is handling the computation of $C_{ij}$ does this (in parallel with the other nodes which are working with their own values of i and j):

```
1   iup  = i+1 mod m;
2   idown = i-1 mod m;
3   for (k = 0; k < m; k++) {
4       km = (i+k) mod m;
5       broadcast(A[i,km]) to all nodes handling row i of C;
6       C[i,j] = C[i,j] + A[i,km]*B[km,j]
7       send B[km,j] to the node handling C[idown,j]
8       receive new B[km+1 mod m,j] from the node handling C[iup,j]
9   }
```

This is Fox's algorithm. Cannon's algorithm is similar, except that it does cyclical rotation in both rows and columns, compared to Fox's rotation only in columns but broadcast within rows.

The algorithm can be adapted in the obvious way to nonsquare matrices, etc.

### 7.4.2 Performance Issues

Note that in MPI we would probably want to implement this algorithm using communicators. For example, this would make broadcasting within a block row more convenient and efficient.

Note too that there is a lot of opportunity here to overlap computation and communication, which is the best way to solve the communication problem. For instance, we can do the broadcast above at the same time as we do the computation.

Obviously this algorithm is best suited to settings in which we have PEs in a mesh topology. This includes hypercubes, though one needs to be a little more careful about communications costs there.

### 7.4.3 Finding Powers of Matrices

In Section (7.1), we saw a special case of matrix multiplication, powers, so that in our context here we have A = B. There are some small improvements that we could make in our algorithm for this case, but also there is something big.

Suppose for instance we need to find $A^{32}$. We could apply the above algorithm 31 times. But a much faster approach would be to first calculate $A^2$, then square that result to get $A^4$, then square it to get $A^8$ and so on. That would get us $A^{32}$ by applying the algorithm in Section 7.4.1 only five times, instead of 31.

## 7.5 Solving Systems of Linear Equations

Suppose we have a system of equations

$$a_{i0}x_0 + ... + a_{i,n-1}x_{n-1} = b_i, i = 0, 1, ..., n - 1, \tag{7.17}$$

where the $x_i$ are the unknowns to be solved for.

As you know, this system can be represented compactly as

$$AX = B, \tag{7.18}$$

where A is nxn and X and B are nx1.

In theory, this system could be solved by finding $A^{-1}$ and left-multiplying by it on both sides of (7.18). However, in practice, this is never done, due to potential problems with numerical stability, etc. There are many other ways (some of which amount to finding $A^{-1}$ indirectly).

### 7.5.1   Gaussian Elimination

You learned this in high school, and in your linear algebra course. Form the n x (n+1) matrix C = (A | B) by appending the column vector B to the right of A. Then we work on the rows of C, with the pseudocode for the sequential case in the most basic form being

```
1   for ii = 0 to n-1
2       divide row ii by c[i][i]
3       for r = ii+1 to n-1   // vacuous if r = n-1
4           replace row r by row r - c[r][ii] times row ii
5   set new B to be column n-1 of C
```

This transforms C to upper triangular form, i.e. all the elements $c_{ij}$ with $i > j$ are 0. Also, all diagonal elements are equal to 1. This corresponds to a new set of equations,

$$
\begin{aligned}
c_{00}x_0 + c_{11}x_1 + c_{22}x_2 + ... + c_{0,n-1}x_{n-1} &= b_0 \\
c_{11}x_1 + c_{22}x_2 + ... + c_{1,n-1}x_{n-1} &= b_1 \\
c_{22}x_2 + ... + c_{2,n-1}x_{n-1} &= b_2 \\
... \\
c_{n-1,n-1}x_{n-1} &= b_{n-1}
\end{aligned}
$$

We then find the $x_i$ via **back substitution**:

```
1   x[n-1] = b[n-1] / c[n-1,n-1]
2   for i = n-2 downto 0
3       x[i] = (b[i] - c[i][n-1] * x[n-1] - ... - c[i][i+1] * x[i+1]) / c[i][i]
```

An obvious parallelization of this algorithm would be to assign each node one contiguous group of rows. Then each node would do

```
1   for ii = 0 to n-1
2       if ii is in my group of rows
3           pivot = c[i][i]
4           divide row ii by pivot
```

```
5          broadcast row ii
6      else receive row ii
7      for r = ii+1 to n-1 in my group
8          subtract c[r][ii] times row ii from row r
9  set new B to be column n-1 of C
```

One problem with this is that in the outer loop, when ii gets past a given node's group of column indices, that node becomes idle. This can be solved by giving each node several groups of rows, in cyclic order. For example, say we have four nodes. Then node 0 could take rows 0-99, 400-499, 800-899 and so on, node 1 would take rows 100-199, 500-599 etc.

### 7.5.2   Iterative Methods

#### 7.5.2.1   The Jacobi Algorithm

One can rewrite (7.17) as

$$x_i = \frac{1}{a_{ii}}[b_i - (a_{i0}x_0 + ... + a_{i,i-1}x_{i-1} + a_{i,i+1}x_{i+1} + ... + a_{i,n-1}x_{n-1})], i = 0, 1, ..., n - 1. \quad (7.19)$$

This suggests a natural iterative algorithm for solving the equations. We start with our guess being, say, $x_i = b_i$ for all i. At our kth iteration, we find our $(k+1)^{st}$ guess by plugging in our kth guess into the right-hand side of (7.19). We keep iterating until the difference between successive guesses is small enough to indicate convergence.

Parallelization of this algorithm is easy: Just assign each node to handle a block of X. Note that this means that each node must send its portion of the new X after every iteration.

#### 7.5.2.2   The Gauss-Seidel Algorithm

This is a variant on the Jacobi algorithm, motivated by the following observation: In a sequential implementation of (7.19), when we get to $x_i$, we already know the new values of $x_0, x_1, ..., x_{i-1}$. Intuitively, we can speed up convergence of our algorithm by using those new values instead of the old ones.

In the parallel case, the easiest way to implement would be that, although we still assign each node to groups of the $x_i$, we would do this in a cyclic order as in Section 7.5.1.

## 7.6   The Shared-Memory Case

We can use all of these algorithms in the shared memory setting, with obvious modifications, a major one being that we remove the code that does send, receive and broadcast, as well as code (e.g. in the matrix multiplication algorithms) that moves data.

*However: Keep in mind that in shared-memory settings, we are in effect doing send, receive and broadcast anyway.* Every one thread does a write, it means that at some later time some other thread will try to read that item, which will mean that the latest copy of that item will be send to this second thread. That sending will be done either as a cache coherency transation in the case of shared-memory hardware, or as a similar page transation in the software distributed shared memory case.

That in turn means that we have to do our best to avoid false sharing. For instance, in the Gaussian elimination case, we have to make sure that the total number of bytes in a group of rows is a multiple of cache line size if we have shared-memory hardware, or a multiple of page size in the software case. We must also make sure that each group of rows begins at cache line/page boundaries. This is not hard, since our matrices will be stored in memory allocated by **malloc()** and the like. We may need to add some padding.

# Chapter 8

# Parallel Combinitorial Algorithms

## 8.1  Overview

In Chapter 1, we saw Dijkstra's algorithm for finding the shortest path in a graph. In Chapter 7, we saw an algorithm for finding bridges within a graph. Both of these are **combinatorial search algorithms**. Such algorithms generally have exponential time complexity, and thus are natural candidates for parallel computation. This unit will present a few more examples.

## 8.2  The 8 Queens Problem

A famous example is the 8 Queens Problem, in which one wishes to place eight queens on a standard 8x8 chessboard in such a way that no queen is attacking any other. (The generalization, of course, would involve n queens on an nxn board.) Suppose our goal is to find all possible solutions.

To start a solution to this problem, we first note in any solution will have the property that no row will contain more than one queen. This suggests building up a solution row by row: Suppose we have successfully placed queens so far in rows 0, 1, ..., k-1 (row 0 being the top row of the board). Where can we place a queen in row k? Well, since we cannot use any column already occupied by the preceding k queens, that means we have a choice of 8-k columns. But even among those k columns, there will be j of them, for some $0 \leq j \leq 8-k$ that are in the diagonal attack path of some preceding queen. Then we can extend our tentative k-row solution to 8-k-j new (k+1)-row solutions.

We will define our solution here for the shared-memory paradigm, though it would be easy to change this for the message-passing paradigm.[1] Define

---

[1] The main point would be to change linked lists and pointers to arrays and array indices.

```
struct TentSoln {
   int RowsSoFar;
   int Cols[8];
   struct TentSoln *Next;
}
```

Each such **struct** contains a partial solution, up through row number **RowsSoFar**. The array **Cols** has the interpretation that **Col[I] == J** tells us which column the queen in row **I** occupies.

Each **struct** is a task showing one partial solution. The node which obtains this task will then extend this partial solution to several new, longer partial solutions.

The tasks are all placed into a linked list. **Next** points to the next item in the work pool.

A parallel solution based on this idea would like something like this:

```
1   while (work pool nonempty or at least one nonidle processor) {
2      get a TentSoln struct from the work pool, and point P to it;
3      I = P->RowsSoFar;
4      for (J = 0; J < 8; J++) {
5         if (a queen at row I, column J would not attack the previous queens) {
6            Q = malloc(sizeof(struct TentSoln));
7            Q->RowsSoFar = I+1;
8            add the struct pointed to by Q to the work pool;
9         }
10      }
11   }
```

There of course would also be code in the case I = 8 to check and see if we have found a solution, and if so, to record it, etc.

Note that any rotation of a solution—interchanging rows and columns—is also a solution. Similarly, any reflection across one of the two main diagonals of the board is also a solution. This information could be used to speed up computation, though at the expense of additionality complexity of the code.

## 8.3   The 8-Square Puzzle Problem

This game was invented more than 100 years ago. Here is what a typical board position looks like:

| 0 | 5 | 3 |
|---|---|---|
| 1 | 4 |   |
| 7 | 2 | 6 |

(The real puzzle has numbering from 1 to 8, but we use 0 to 7.)

Each number is on a little movable square, which can be moved up, down, left and right as long as the spot in the given direction is unoccupied. In the example above, the square 3, for instance, could be moved downward, producing an empty spot at the top right of the puzzle. The object of the game is arrange the squares in ascending numerical order, with square 0 at the upper left of the puzzle (which in this example happens to be the case already).

We again solve this by setting up a work pool, in this case a pool of board positions. Each board position would be implemented in something like this:

```
struct BoardPos   {
   int Row[9];
   int Col[9];
   struct BoardPos *Next;
}
```

Here **Row[I]** and **Col[I]** would be the position of the square numbered **I**. For convenience, we also store the location of the blank position, in **Row[8]** and **Col[8]**.

Suppose a processor goes to the work pool and gets the board position depicted above. In the simplest form of the algorithm, the processor would check each of the three possible moves (4 right, 3 down, 6 up) to see if the resulting board position would duplicate one that had already been checked. All moves that lead to new positions would be added to the work pool. Each processor would loop around, pulling items from the work pool, until some processor somewhere finds a solution to the game (in which case that processor would add termination messages to the work pool, so that the other processors knew to stop). An outline of the algorithm would be as follows:

```
1   while (work pool nonempty or at least one nonidle processor) {
2      get a BoardPos struct from the work pool, and point P to it;
3      for (I = 0; I < 8; I++)   {
4         for all possible moves of square I do  {
5            Q = malloc(sizeof(struct BoardPos));
6            fill in *Q according to this move;
7            if *Q has not already been checked
8               add this board to the work pool;
9         }
10     }
11  }
```

Again, code would need to be included for checking to see if a solution has been found, whether we have found that no solution exists, and so on.

Note the operation

```
        if *Q has not already been checked
           add this board to the work pool;
```

Clearly this is needed, to avoid endless cycling. But it is not as inoccuous as it looks. If the set of all previously-checked board positions is to be made available to all processors, this may produce substantial increases in contention for memory and interprocessor interconnects. On the other hand, we could arrange the code such that only certain processors have to know about certain subsets of the set of previously-checked board positions, but this makes the code more complex and may produce load-balancing problems.

A more sophisticated version of the algorithm would use a **branch-and-bound** technique. The idea here is to reduce computation by giving priority in the work pool to those board positions which appear "promising" by some reasonable measure. For example, we could take as our measure the "distance" between a given board position and the goal board position, as defined by the sum of the distances from each numbered square to its place in the winning position. In the example above, for instance, the square numbered 5 is a distance of 2 from its ultimate place (2 meaning, one square to the right, one square down, so 1+1 = 2). The board above is a distance 15 from the winning board.

The idea, then would be that we implement the work pool as an ordered linked list (or other ordered structure), and when a board position is added to the work pool, we insert it according to its distance from the winning board. This way the processors will usually work on the more promising boards, and thus hopefully reach the solution faster.

## 8.4 Itemset Analysis in Data Mining

### 8.4.1 What Is It?

The term **data mining** is a buzzword, but all it means is the process of finding relationships among a set of variables. In other words, it would seem to simply be a good old-fashioned statistics problem.

Well, in fact it *is* simply a statistics problem—but writ large. Instead of the tiny sample sizes of 25 you likely saw in your statistics courses, typical sample sizes in the data mining arena run in the hundreds of thousands or even hundreds of millions. And there may be hundreds of variables, in constrast to the, say, half dozen you might see in a statistics course.

**Major, Major Warning:** With so many variables, the chances of picking up spurious relations between variables is large. And although many books and tutorials on data mining will at least pay lip service to this issue (referring to it as **overfitting**, they don't emphasize it enough.[2]

Putting the overfitting problem aside, though, by now the reader's reaction should be, "This calls for parallel processing," and he/she is correct. Here we'll look at parallelizing a particular problem, called **itemset analysis**, the most famous example of which is the **market basket problem**:

---

[2]Some writers recommend splitting one's data into a **training set**, which is used to discover relationships, and a **validation set**, which is used to confirm those relationships. However, overfitting can still occur even with this precaution.

### 8.4.2 The Market Basket Problem

Consider an online bookstore has records of every sale on the store's site. Those sales may be represented as a matrix S, whose (i,j)th element $S_{ij}$ is equal to either 1 or 0, depending on whether the $ith$ sale included book j, i = 0,1,...,s-1, j = 0,1,...,t-1. So each row of S represents one sale, with the 1s in that row showing which titles were bought. Each column of S represents one book title, with the 1s showing which sales transactions included that book.

Let's denote the entire line of book titles by $T_0, ..., T_{b-1}$. An **itemset** is just a subset of this. A **frequent** itemset is one which appears in many of sales transactions. But there is more to it than that. The store wants to choose some books for special ads, of the form "We see you bought books X and Y. We think you may be interested in Z."

Though we are using marketing as a running example here (which is the typical way that this subject is introduced), we will usually just refer to "items" instead of books, and to "database records" rather than sales transactions.

We have the following terminology:

- An **association rule** $I \rightarrow J$ is simply an ordered pair of disjoint itemsets I and J.

- The **support** of an an association rule $I \rightarrow J$ is the proportion of records which include both I and J.

- The **confidence** of an association rule $I \rightarrow J$ is the proportion of records which include J, *among those records which include I.*

Note that in probability terms, the support is basically P(I and J) while the confidence is P(J|I). If the confidenc the book business, it means that buyers of the books in set I also tend to buy those in J. But this information is not very useful if the support is low, because it means that the combination occurs so rarely that it's not worth our time to deal with it.

So, the user—let's call him/her the "data miner"—will first set thresholds for support and confidence, and then set out to find all association rules for which support and confidence exceed their respective thresholds.

### 8.4.3 Serial Algorithms

Various algorithms have been developed to find frequent itemsets and association rules. The most famous one for the former task is the **Apriori** algorithm. Even it has many forms. We will discuss one of the simplest forms here.

The algorithm is basically a breadth-first tree search. At the root we find the frequent 1-item itemsets. At the second level, we find the frequent 2-item itemsets, and so on. After we finish with level i, we then generate new candidate itemsets of size i+1 from the frequent itemsets we found of size i, by

The key point in the latter operation is that if an itemset is not frequent, i.e. has support less than the threshold, then adding further items to it will make it even less frequent. That itemset is then pruned from the tree, and the branch ends.

Here is the pseudocode:

set $F_1$ to the set of 1-item itemsets whose support exceeds the threshold
for i = 2 to b
    $F_i = \phi$
    for each I in $F_{i-1}$
        for each K in $F_1$
            $Q = I \cup K$
            if support(Q) exceeds support threshold
                add Q to $F_i$
    if $F_i$ is empty break
return $\cup_i F_i$

Again, there are many refinements of this, which shave off work to be done and thus increase speed. For example, we should avoid checking the same itemsets twice, e.g. first $\{1,2\}$ then $\{2,1\}$. This can be accomplished by keeping itemsets in lexicographical order. We will not pursue any refinements here.

### 8.4.4 Parallelizing the Apriori Algorithm

Clearly there is lots of opportunity for parallelizing the serial algorithm above. Both of the inner **for** loops can be parallelized in straightforward ways; they are "embarrassingly parallel." There are of course critical sections to worry about in the shared-memory setting, and in the message-passing setting one must designate a manager node in which to store the $F_i$.

However, as more and more refinements are made in the serial algorithm, then the parallelism in this algorithm become less and less "embarrassing." And things become more challenging if the storage needs of the $F_i$, and of their associated "accounting materials" such as a directory showing the current tree structure (done via hash trees), become greater than what can be stored in the memory of one node.

In other words, parallelizing the market basket problem can be very challenging. The interested reader is referred to the considerable literature which has developed on this topic.

# Chapter 9

# Introduction to Parallel Sorting

Sorting is one of the most common operations in parallel processing applications. For example, it is central to many parallel database operations, and important in areas such as image processing, statistical methodology and so on. A number of different types of parallel sorting schemes have been developed. Here we look at some of these schemes.

## 9.1 Quicksort

You are probably familiar with the idea of quicksort: First break the original array into a "small-element" pile and a "large-element" pile, by comparing to a **pivot** element. In a naive implementation, the first element of the array serves as the pivot, but better performance can be obtained by taking, say, the median of the first three elements. Then "recurse" on each of the two piles, and then string the results back together again.

This is an example of the **divide and conquer** approach seen in so many serial algorithms. It is easily parallelized (though load-balancing issues may arise). Here, for instance, we might assign one pile to one thread and the other pile to another thread.

Suppose the array to be sorted is named **x**, and consists of **n** elements.

In a naive implementation, the piles would be put into new arrays, but this is bad in two senses: It wastes memory space, and wastes time, since much copying of arrays needs to be done. A better implementation places the two piles back into the original array **x**. The following C code does that.

The function **separate()** is intended to be used in a recursive quicksort operation. It operates on **x[l]** through **x[h]**, a subarray of **x** that itself may have been formed at an earlier stage of the recursion. It forms two piles from those elements, and placing the piles back in the same region **x[l]** through **x[h]**. It also has a return value, showing where the first pile ends.

```
int separate(int l, int h)
{  int ref,i,j,k,tmp;
   ref = x[h]; i = l-1; j = h;
   do  {
      do i++; while (x[i] < ref && i < h);
      do j--; while (x[j] > ref && j > l);
      tmp = x[i];  x[i] = x[j];  x[j] = tmp;
   } while (j > i);
   x[j] = x[i];  x[i] = x[h];  x[h] = tmp;
   return i;
}
```

The function **separate()** rearranges the subarray, returning a value **m**, so that:

- **x[l]** through **x[m-1]** are less than **x[m]**,

- **x[m+1]** through **x[h]** are greater than **x[m]**, and

- **x[m]** is in its "final resting place," meaning that **x[m]** will never move again for the remainder of the sorting process. (Another way of saying this is that the current **x[m]** is the **m**-th smallest of all the original **x[i]**, **i** = 0,1,...,**n**-1.)

By the way, **x[l]** through **x[m-1]** will also be in their final resting places as a group. They may be exchanging places with each other from now on, but they will never again leave the range **i** though **m-1** within the **x** array as a whole. A similar statement holds for **x[m+1]** through **x[n-1]**.

### 9.1.1   Shared-Memory Quicksort

Here is OpenMP code which performs quicksort in the shared-memory paradigm (adapted from code in the OpenMP Source Code Repository, `http://www.pcg.ull.es/ompscr/`):

```
1   void qs(int *x, int l, int h)
2   {  int newl[2], newh[2], i, m;
3      m = separate(x,l,h);
4      newl[0] = l;   newh[0] = m-1;
5      newl[1] = m+1;   newh[1] = h;
6      #pragma omp parallel
7      {
8         #pragma omp for nowait
9         for (i = 0; i < 2; i++)
10           qs(newl[i],newh[i]);
11     }
12  }
```

Note the **nowait** clause. Since different threads are operating on different portions of the array, they need not be synchronized.

A variant on this which might achieve better load balancing would set up a **task pool**, consisting of an array of (**l**, **h**) pairs. Initially the pool consists of just [0,n-1]. The function **qs()** would then become iterative instead of recursive, with its main loop looking something like this for an array of length n:

```
fetch an (l,h) pair from the task pool
while not done
   call separate() on x[l] through x[h], yielding m
   if m < h
      add (m+1,h) to the task pool
   h = m-1
   if l == h
      fetch [l,h] from the task pool
```

This pseudocode is missing important details. For example, How does the iteration within a thread stop? The key lies in pairs of the form (i,i), which I'll call *singletons*. The sort is done when the number of singletons reaches n.

### 9.1.2 Hyperquicksort

This algorithm was originally developed for hypercubes, but can be used on any message-passing system having a power of 2 for the number of nodes.[1]

It is assumed that at the beginning each PE contains some chunk of the array to be sorted. After sorting, each PE will contain some chunk of the <u>sorted</u> array, meaning that:

- each chunk is itself in sorted form

- for all cases of $i < j$, the elements at PE i are less than the elements at PE j

If the sorted array itself were our end, rather than our means to something else, we could now collect it at some node, say node 0. If, as is more likely, the sorting is merely an intermediate step in a larger distributed computation, we may just leave the chunks at the nodes and go to the next phase of work.

Say we are on a d-cube. The intuition behind the algorithm is quite simple:

```
for i = d downto 1
   for each i-cube:
      root of the i-cube broadcasts its median to all in the i-cube,
         to serve as pivot
      consider the two (i-1)-subcubes of this i-cube
      each pair of partners in the (i-1)-subcubes exchanges data:
         low-numbered PE gives its partner its data larger than pivot
         high-numbered PE gives its partner its data smaller than pivot
```

---

[1]See Chapter 5 for definitions of hypercube terms.

To avoid deadlock, have the lower-numbered partner send then receive, and vice versa for the higher-numbered one. Better, in MPI, use **MPI_SendRcv()**.

After the first iteration, all elements in the lower (d-1)-cube are less than all elements in higher (d-1)-cube. After d such steps, the array will be sorted.

## 9.2    Mergesorts

### 9.2.1    Sequential Form

In its serial form, mergesort has the following pseudocode:

```
1  // initially called with l = 0 and h = n-1, where n is the length of the
2  // array and is assumed here to be a power of 2
3  void seqmergesort(int *x, int l, int h)
4  {  seqmergesort(x,0,h/2-1);
5     seqmergesort(x,h/2,h);
6     merge(x,l,h);
7  }
```

The function **merge()** should be done in-place, i.e. without using an auxiliary array. It basically codes the operation shown in pseudocode for the message-passing case in Section 9.2.3.

### 9.2.2    Shared-Memory Mergesort

This is similar to the patterns for shared-memory quicksort in Section 9.1.1 above.

### 9.2.3    Message Passing Mergesort on a Tree Topology

First, we organize the processing nodes into a binary tree. This is simply from the point of view of the software, rather than a physical grouping of the nodes. We will assume, though, that the number of nodes is one less than a power of 2.

To illustrate the plan, say we have seven nodes in all. We could label node 0 as the root of the tree, label nodes 1 and 2 to be its two children, label nodes 3 and 4 to be node 1's children, and finally label nodes 5 and 6 to be node 2's children.

It is assumed that the array to be sorted is initially distributed in the leaf nodes (recall a similar situation for hyperquicksort), i.e. nodes 3-6 in the above example. The algorithm works best if there are approximately the same number of array elements in the various leaves.

In the first stage of the algorithm, each leaf node applies a regular sequential sort to its current holdings. Then each node begins sending its now-sorted array elements to its parent, one at a time, in ascending numerical order.

Each nonleaf node then will merge the lists handed to it by its two children. Eventually the root node will have the entire sorted array. Specifically, each nonleaf node does the following:

```
do
   if my left-child datum < my right-child datum
      pass my left-child datum to my parent
   else
      pass my right-child datum to my parent
until receive the "no more data" signal from both children
```

Of course, due to network latency and the like, one may get better performance if each node accumulates a chunk of data before sending to the parent, rather than sending just one datum at a time.

### 9.2.4   Compare-Exchange Operations

These are key to many sorting algorithms.

A **compare-exchange**, also known as **compare-split**, simply means in English, "Let's pool our data, and then I'll take the lower half and you take the upper half." Each node executes the following pseudocode:

```
send all my data to partner
receive all my partner's data
if I have a lower id than my partner
   I keep the lower half of the pooled data
else
   I keep the upper half of the pooled data
```

### 9.2.5   Bitonic Mergesort

A sequence $(a_0, a_1, .., a_{k-1})$ is called **bitonic** if it is first nondecreasing then nonincreasing, meaning that for some r

$$(a_0 \leq a_1 \leq ... \leq a_r \geq a_{r+1} \geq a_{n-1})$$

(For convenience, from here on I will use the terms *increasing* and *decreasing* instead of *nonincreasing* and *nondecreasing*.)

This includes the cases in which the sequence is purely nondecreasing (r = n-1) or purely nonincreasing (r = 0) . By convention, it also includes sequences which can be cyclically shifted into the above form.

For instance, the sequence (3,8,12,15,14,5,1,2) can be rotated rightward by two element positions to form (1,2,3,8,12,15,14,5), so (3,8,12,15,14,5,1,2) is defined to be bitonic too.

These are the "A-type" bitonic sequences, so called because they look like the letter A (or like a carat). The "V-type" bitonic sequences consist of a nonincreasing sequence followed by a nondecreasing sequence.

Suppose we have bitonic sequence $(a_0, a_1, .., a_{k-1})$, where k is a power of 2. Rearrange the sequence by doing compare-exchange operations between $a_i$ and $a_{n/2+i}$), i = 0,1,...,n/2-1. Then it is not hard to prove that the new $(a_0, a_1, .., a_{k/2-1})$ and $(a_{k/2}, a_{k/2+1}, .., a_{k-1})$ are bitonic, and every element of that first subarray is less than or equal to every element in the second one.

So, we have set things up for yet another divide-and-conquer attack:

```
1  // x is bitonic of length n, n a power of 2
2  void sortbitonic(int *x, int n)
3  {  do the pairwise compare-exchange operations
4     if (n > 2) {
5        sortbitonic(x,n/2);
6        sortbitonic(x+n/2,n/2);
7     }
8  }
```

So much for sorting bitonic sequences. But what about general sequences? We can proceed as follows:

1. Each of the pairs $(a_i, a_{i+1})$, i = 0,2,...,n-2 is bitonic, since *any* 2-element array is bitonic!

2. For each i = 0,2,4,...,n-2:

   - Apply **sortbitonic()** to $(a_i, a_{i+1})$.
   - If i/2 is odd, reverse the pair, so that this pair and the pair immediately preceding it now form a 4-element bitonic sequence.

3. For each i = 0,4,8,...,n-4:

   - Apply **sortbitonic()** to $(a_i, a_{i+1}, a_{i+2}, a_{i+3})$.
   - If i/4 is odd, reverse the quartet, so that this quartet and the quartet immediately preceding it now form an 8-element bitonic sequence.

4. Keep building in this manner, until get to a single sorted n-element list.

There are many ways to parallelize this. In the hypercube case, the algorithm consists of doing compare-exchange operations with all neighbors, pretty much in the same pattern as hyperquicksort.

## 9.3   The Bubble Sort and Its Cousins

### 9.3.1   The Much-Maligned Bubble Sort

Recall the **bubble sort**:

```
1   void bubblesort(int *x, int n)
2   {  for i = n-1 downto 1
3         for j = 0 to i
4            compare-exchange(x,i,j,n)
5   }
```

Here the function **compare-exchange()** is as in Section 9.2.4 above. In the context here, it boils down to

```
if x[i] > x[j]
   swap x[i] and x[j]
```

In the first **i** iteration, the largest element "bubbles" all the way to the right end of the array. In the second iteration, the second-largest element bubbles to the next-to-right-end position, and so on.

You learned in your algorithms class that this is a very inefficient algorithm—when used serially. But it's actually rather usable in parallel systems.

For example, in the shared-memory setting, suppose we have one thread for each value of **i**. Then those threads can work in parallel, as long as a thread with a larger value of **i** does not overtake a thread with a smaller **i**, where "overtake" means working on a larger **j** value.

Once again, it probably pays to chunk the data. In this case, **compare-exchange()** fully takes on the meaning it had in Section 9.2.4.

### 9.3.2   A Popular Variant: Odd-Even Transposition

A popular variant of this is the **odd-even transposition sort**. The pseudocode for a shared-memory version is:

```
1   // the argument "me" is this thread's ID
2   void oddevensort(int *x, int n, int me)
3   {  for i = 1 to n
4         if i is odd
5            if me is even
6               compare-exchange(x,me,me+1,n)
7            else  // me is odd
8               compare-exchange(x,me,me-1,n)
9         else  // i is even
```

```
10          if me is even
11            compare-exchange(x,me,me-1,n)
12          else  // me is odd
13            compare-exchange(x,me,me+1,n)
```

If the second or third argument of **compare-exchange()** is less than 0 or greater than **n**-1, the function has no action.

This looks a bit complicated, but all it's saying is that, from the point of view of an even-numbered element of **x**, it trades with its right neighbor during odd phases of the procedure and with its left neighbor during even phases.

Again, this is usually much more effective if done in chunks.

## 9.4   Shearsort

In some contexts, our hardware consists of a two-dimensional mesh of PEs. A number of methods have been developed for such settings, one of the most well known being Shearsort, developed by Sen, Shamir and the eponymous Isaac Scherson of UC Irvine. Again, the data is assumed to be initially distributed among the PEs. Here is the pseudocode:

```
1  for i = 1 to ceiling(log2(n)) + 1
2      if i is odd
3          sort each even row in descending order
4          sort each odd row in ascending order
5      else
6          sort each column is ascending order
```

At the end, the numbers are sorted in a "snakelike" manner.

For example:

| 6 | 12 |
|---|----|
| 5 | 9  |

| 6 | 12 |
|---|----|
| 9 | 5  |

| 6 | 5  |
|---|----|
| 9 | 12 |

| 5  | 6 ↓  |
|----|------|
| 12 | ← 9 |

No matter what kind of system we have, a natural domain decomposition for this problem would be for each process to be responsible for a group of rows. There then is the question about what to do during the even-numbered iterations, in which column operations are done. This can be handled via a parallel matrix transpose operation. In MPI, the function **MPI_Alltoall()** may be useful.

## 9.5 Bucket Sort with Sampling

For concreteness, suppose we are using MPI on message-passing hardware, say with 10 PEs. As usual in such a setting, suppose our data is initially distributed among the PEs.

Suppose we knew that our array to be sorted is a random sample from the uniform distribution on (0,1). In other words, about 20% of our array will be in (0,0.2), 38% will be in (0.45,0.83) and so on.

What we could do is assign PE0 to the interval (0,0.1), PE1 to (0.1,0.2) etc. Each PE would look at its local data, and distribute it to the other PEs according to this interval scheme. Then each PE would do a local sort.

In general, we don't know what distribution our data comes from. We solve this problem by doing sampling. In our example here, each PE would sample some of its local data, and send the sample to PE0. From all of these samples, PE0 would find the decile values, i.e. 10th percentile, 20th percentile,..., 90th percentile. These values, called **splitters** would then be broadcast to all the PEs, and they would then distribute their local data to the other PEs according to these intervals.

# Chapter 10

# Parallel Computation of Fourier Series, with an Introduction to Parallel Imaging
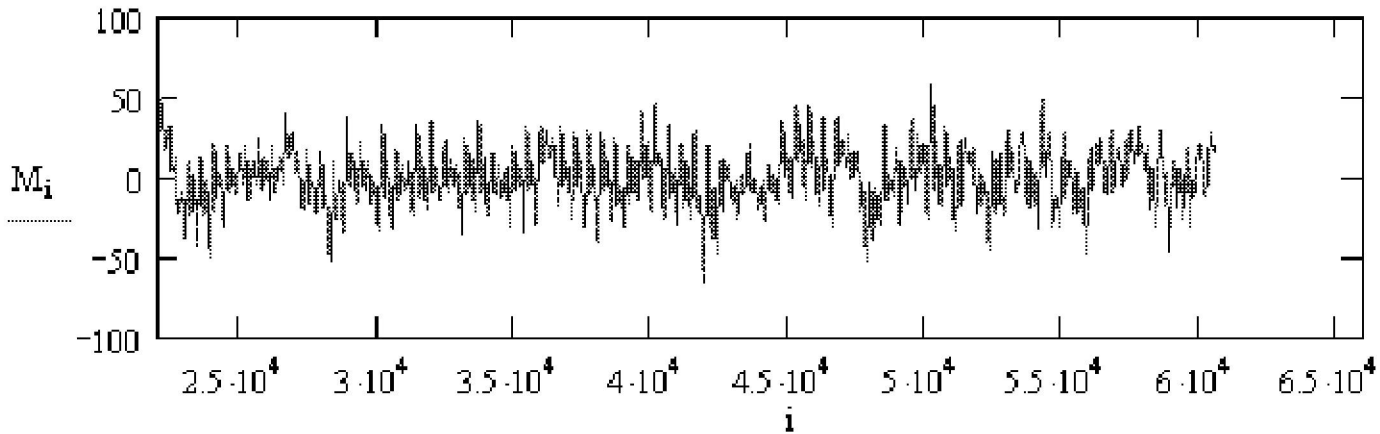
Mathematical computations involving sounds and images, for example for voice and pattern recognition are often performed using **Fourier** analysis.

## 10.1   General Principles

### 10.1.1   One-Dimensional Fourier Series

A sound **wave form** graphs volume of the sound against time. Here, for instance, is the wave form for a vibrating reed:[1]

---

[1]Reproduced here by permission of Prof. Peter Hamburger, Indiana-Purdue University, Fort Wayne. See http://www.ipfw.edu/math/Workshop/PBC.html

Recall that we say a function of time g(t) is **periodic** ("repeating," in our casual wording above) with period T if if g(u+T) = g(u) for all u. The **fundamental frequency** of g() is then defined to be the number of periods per unit time,

$$f_0 = \frac{1}{T} \tag{10.1}$$

Recall also from calculus that we can write a function g(t) (not necessarily periodic) as a Taylor series, which is an "infinite polynomial":

$$g(t) = \sum_{n=0}^{\infty} c_n t^n. \tag{10.2}$$

The specific values of the $c_n$ may be derived by differentiating both sides of (10.2) and evaluating at t = 0, yielding

$$c_n = \frac{g^{(n)}(0)}{n!}, \tag{10.3}$$

where $g^{(j)}$ denotes the ith derivative of g().

For instance, for $e^t$,

$$e^t = \sum_{n=0}^{\infty} \frac{1}{n!} t^n \tag{10.4}$$

In the case of a repeating function, it is more convenient to use another kind of series representation, an "infinite trig polynomial," called a **Fourier series**. This is just a fancy name for a weighted sum of sines and

cosines of different frequencies. More precisely, we can write any repeating function g(t) with period T and fundamental frequency $f_0$ as

$$g(t) = \sum_{n=0}^{\infty} a_n \cos(2\pi n f_0 t) + \sum_{n=1}^{\infty} b_n \sin(2\pi n f_0 t) \tag{10.5}$$

for some set of weights $a_n$ and $b_n$. Here, instead of having a weighted sum of terms

$$1, \ t, \ t^2, \ t^3, \ \ldots \tag{10.6}$$

as in a Taylor series, we have a weighted sum of terms

$$1, \ \cos(2\pi f_0 t), \ \cos(4\pi f_0 t), \ \cos(6\pi f_0 t), \ \ldots \tag{10.7}$$

and of similar sine terms. Note that the frequencies $n f_0$, in those sines and cosines are integer multiples of the fundamental frequency of x, $f_0$, called **harmonics**.

The weights $a_n$ and $b_n$, n = 0, 1, 2, ... are called the **frequency spectrum** of g(). The coefficients are calculated as follows:[2]

$$a_0 = \frac{1}{T} \int_0^T g(t) \, dt \tag{10.8}$$

$$a_n = \frac{2}{T} \int_0^T g(t) \, cos(2\pi n f_0 t) \, dt \tag{10.9}$$

$$b_n = \frac{2}{T} \int_0^T g(t) \, sin(2\pi n f_0 t) \, dt \tag{10.10}$$

By analyzing these weights, we can do things like machine-based voice recognition (distinguishing one person's voice from another) and speech recognition (determining what a person is saying). If for example one person's voice is higher-pitched than that of another, the first person's weights will be concentrated more on the higher-frequency sines and cosines than will the weights of the second.
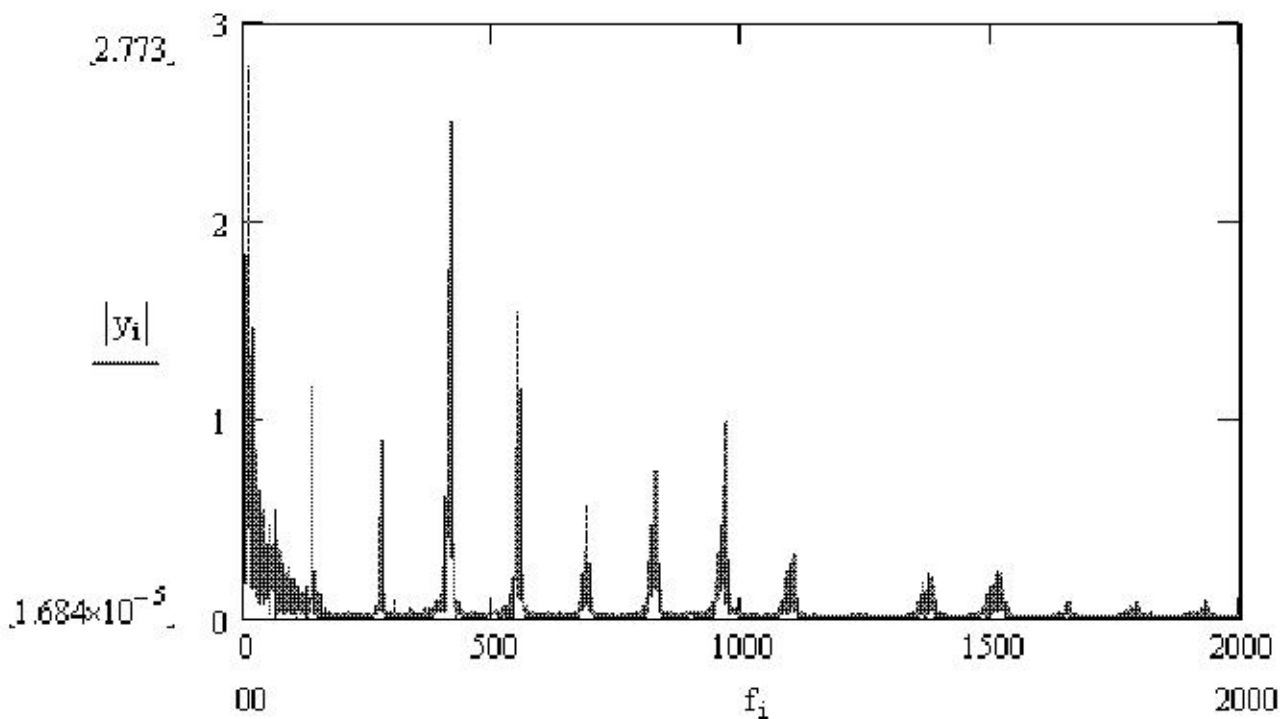
Since g(t) is a graph of loudness against time, this representation of the sound is called the **time domain**. When we find the Fourier series of the sound, the set of weights $a_n$ and $b_n$ is said to be a representation of

---

[2]The get an idea as to how these formulas arise, see Section 10.8. But for now, if you integrate both sides of (10.5), you will at least verify that the formulas below do work.

the sound in the **frequency domain**. One can recover the original time-domain representation from that of the frequency domain, and vice versa, as seen in Equations (10.8), (10.9), (10.10) and (10.5).

In other words, the transformations between the two domains are inverses of each other, and there is a one-to-one correspondence between them. Every g() corresponds to a unique set of weights and vice versa.

Now here is the frequency-domain version of the reed sound:



Note that this graph is very "spiky." In other words, even though the reed's waveform includes all frequencies, most of the power of the signal is at a few frequencies which arise from the physical properties of the reed.

Fourier series are often expressed in terms of complex numbers, making use of the relation

$$e^{i\theta} = \cos(\theta) + i\ \sin(\theta), \tag{10.11}$$

where $i = \sqrt{-1}$.[3]

---

[3]There is basically no physical interpretation of complex numbers. Instead, they are just mathematical abstractions. However, they are highly useful abstractions, with the complex form of Fourier series, beginning with (10.12), being a case in point.

The complex form of (10.5) is

$$g(t) = \sum_{j=-\infty}^{\infty} c_j e^{2\pi i j \frac{t}{T}}. \tag{10.12}$$

The $c_j$ are now generally complex numbers. They are functions of the $a_j$ and $b_j$, and thus comprise the frequency spectrum.

Equation (10.12) has a simpler, more compact form than (10.5). Do you now see why I referred to Fourier series as trig polynomials? The series (10.12) involves the $j^{th}$ powers of $e^{2\pi \frac{t}{T}}$.

### 10.1.2  Two-Dimensional Fourier Series

Let's now move from sounds images. Here g() is a function of two variables, g(u,v), where u and v are the horizontal and vertical coordinates of a pixel in the image; g(u,v) is the intensity of the image at that pixel. If it is a gray-scale image, the intensity is whiteness of the image at that pixel, typically with 0 being pure black and 255 being pure white. If it is a color image, a typical graphics format is to store three intensity values at a point, one for each of red, green and blue. The various colors come from combining three colors at various intensities.

Since images are two-dimensional instead of one-dimensional like a sound wave form, the Fourier series for an image is a sum of sines and cosines in two variables, i.e. a double sum $\Sigma_j \Sigma_k$... instead of $\Sigma_j$....

The terminology changes a bit. Our original data is now referred to as being in the **spatial domain**, rather than the time domain. But the Fourier series coefficients are still said to be in the frequency domain.

## 10.2  Discrete Fourier Transforms

In sound and image applications, we seldom if ever know the exact form of the repeating function g(). All we have is a **sampling** from g(), i.e. we only have values of g(t) for a set of discrete values of t.

In the sound example above, a typical sampling rate is 8000 samples per second.[4] So, we may have g(0), g(0.000125), g(0.000250), g(0.000375), and so on. In the image case, we sample the image pixel by pixel.

Thus we can't calculate integrals like (10.8). So, how do we approximate the Fourier transform based on the sample data?

---

[4]See Section 10.9 for the reasons behind this.

### 10.2.1 One-Dimensional Data

Let $X = (x_0, ..., x_{n-1})$ denote the sampled values, i.e. the time-domain representation of g() based on our sample data. These are interpreted as data from one period of g(), with the period being n and the fundamental frequency being 1/n. The frequency-domain representation will also consist of n numbers, $c_0, ..., c_{n-1}$, defined as follows:[5]

$$c_k = \frac{1}{n} \sum_{j=0}^{n-1} x_j e^{-2\pi i j k / n} = \frac{1}{n} \sum_{j=0}^{n-1} x_j q^{jk} \tag{10.13}$$

where

$$q = e^{-2\pi i / n} \tag{10.14}$$

again with $i = \sqrt{-1}$. The array C of complex numbers $c_k$ is called the **discrete Fourier transform** (DFT) of X.

Note that (10.13) is basically a discrete analog of (10.9) and (10.10).

As in the continuous case, we can recover each domain from the other. So, while (10.13) shows how to go to the frequency domain from the time domain, we can go from the frequency domain to the time domain via the inverse transform, whose equation is

$$x_k = \sum_{j=0}^{n-1} c_j e^{2\pi i j k / n} = \sum_{j=0}^{n-1} c_j q^{-jk} \tag{10.15}$$

Note that (10.15) is basically a discrete analog of (10.5).

Note too that instead of having infinitely many harmonics, we can only have n of them: 1, 1/n, 2/n, ..., (n-1)/n. It would be impossible to have more than n, as can be seen by reasoning as follows: The $x_k$ are given, q is a constant, and we are solving for the $c_j$. So, we have n equations in n unknowns. If we had more than n unknowns, the system would be indeterminate.

---

[5]It should be noted that there are many variant definitions of these transforms. One common variation is to include/exclude a scale factor, such as our 1/n in (10.13). Another type of variations involve changing only $c_0$, in order to make certain matrices have more convenient forms.

### 10.2.2 Two-Dimensional Data

The spectrum numbers $c_{rs}$ are double-subscripted, like the original data $x_{uv}$, the latter being the pixel intensity in row u, column v of the image, u = 0,1,...,n-1, v = 0,1,...,m-1. Equation (10.13) becomes

$$c_{rs} = \frac{1}{n}\frac{1}{m}\sum_{j=0}^{n-1}\sum_{k=0}^{m-1} x_{jk}e^{-2\pi i(\frac{jr}{n}+\frac{ks}{m})} \tag{10.16}$$

Its inverse is

$$x_{rs} = \sum_{j=0}^{n-1}\sum_{k=0}^{m-1} c_{jk}e^{2\pi i(\frac{jr}{n}+\frac{ks}{m})} \tag{10.17}$$

## 10.3 Parallel Computation of Discrete Fourier Transforms

### 10.3.1 The Fast Fourier Transform

Speedy computation of a discrete Fourier transform was developed by Cooley and Tukey in their famous Fast Fourier Transform (FFT), which takes a "divide and conquer" approach:

Equation (10.13) can be rewritten as

$$c_k = \frac{1}{n}\left[\sum_{j=0}^{m-1} x_{2j}q^{2jk} + \sum_{j=0}^{m-1} x_{2j+1}q^{(2j+1)k},\right] \tag{10.18}$$

where $m = n/2$.

After some algebraic manipulation, this becomes

$$c_k = \frac{1}{2}\left[\frac{1}{m}\sum_{j=0}^{m-1} x_{2j}z^{jk} + q^k\frac{1}{m}\sum_{j=0}^{m-1} x_{2j+1}z^{jk}\right] \tag{10.19}$$

where $z = e^{-2\pi i/m}$.

A look at Equation (10.19) shows that the two sums within the brackets have the same form as Equation (10.13). In other words, Equation (10.19) shows how we can compute an n-point FFT from two $\frac{n}{2}$-point

FFTs. That means that a DFT can be computed recursively, cutting the sample size in half at each recursive step.

In a shared-memory setting such as OpenMP, we could implement this recursive algorithm in the manners of Quicksort in Chapter 9.

In a message-passing setting, one can use the butterfly algorithm, explained for implementation of barriers in Chapter 1. Some digital signal processing chips implement this in hardware, with a special interconnection network to implement this algorithm.

### 10.3.2  A Matrix Approach

The matrix form of (10.13) is

$$C = \frac{1}{n}AX \tag{10.20}$$

where A is n x n. Element (j,k) of A is $q^{jk}$, while element j of X is $x_j$. This formulation of the problem then naturally leads one to use parallel methods for matrix multiplication; see Chapter 7.

### 10.3.3  Parallelizing Computation of the Inverse Transform

The form of the DFT (10.13) and its inverse (10.15) are very similar. For example, the inverse transform is again of a matrix form as in (10.20); even the new matrix looks a lot like the old one.[6]

Thus the methods mentioned above, e.g. FFT and the matrix approach, apply to calculation of the inverse transforms too.

### 10.3.4  Parallelizing Computation of the Two-Dimensional Transform

Regroup (10.16) as:

$$c_{rs} = \frac{1}{n}\sum_{j=0}^{n-1}\left(\frac{1}{m}\sum_{k=0}^{m-1}x_{jk}e^{-2\pi i(\frac{ks}{m})}\right)e^{-2\pi i(\frac{jr}{n})} \tag{10.21}$$

$$= \frac{1}{n}\sum_{j=0}^{n-1}y_{js}e^{-2\pi i(\frac{jr}{n})} \tag{10.22}$$

---

[6]In fact, one can obtain the new matrix easily from the old, as explained in Section 10.8.

Note that $y_{js}$, i.e. the expression between the large parentheses, is the s$^{th}$ component of the DFT of the j$^{th}$ row of our data. And hey, the last expression (10.22) above is in the same form as (10.13)! Of course, this means we are taking the DFT of the spectral coefficients rather than observed data, but numbers are numbers.

In other words: To get the two-dimensional DFT of our data, we first get the one-dimensional DFTs of each row of the data, place these in rows, and then find the DFTs of each column. This property is called **separability**.

This certainly opens possibilities for parallelization. Each thread (shared memory case) or node (message passing case) could handle groups of rows of the original data, and in the second stage each thread could handle columns.

Or, we could interchange rows and columns in this process, i.e. put the j sum inside and k sum outside in the above derivation.

## 10.4 Applications to Image Processing

In image processing, there are a number of different operations which we wish to perform. We will consider two of them here.

### 10.4.1 Smoothing

An image may be too "rough." There may be some pixels which are noise, accidental values that don't fit smoothly with the neighboring points in the image.

One way to smooth things out would be to replace each pixel intensity value[7] by the mean or median among the pixels neighbors. These could be the four immediate neighbors if just a little smoothing is needed, or we could go further out for a higher amount of smoothing. There are many variants of this.

But another way would be to apply a **low-pass filter** to the DFT of our image. This means that after we compute the DFT, we simply delete the higher harmonics, i.e. set $c_{rs}$ to 0 for the larger values of r and s. We then take the inverse transform back to the spatial domain. Remember, the sine and cosine functions of higher harmonics are "wigglier," so you can see that all this will have the effect of removing some of the wiggliness in our image—exactly what we wanted.

We can control the amount of smoothing by the number of harmonics we remove.

The term *low-pass filter* obviously alludes to the fact that the low frequencies "pass" through the filter but the high frequencies are blocked. Since we've removed the high-oscillatory components, the effect is a

---

[7]Remember, there may be three intensity values per pixel, for red, green and blue.

smoother image.[8]

To do smoothing in parallel, if we just average neighbors, this is easily parallelized. If we try a low-pass filter, then we use the parallelization methods shown here earlier.

## 10.4.2 Edge Detection

In computer vision applications, we need to have a machine-automated way to deduce which pixels in an image form an edge of an object.

Again, edge-detection can be done in primitive ways. Since an edge is a place in the image in which there is a sharp change in the intensities at the pixels, we can calculate slopes of the intensities, in the horizontal and vertical directions. (This is really calculating the approximate values of the partial derivatives in those directions.)

But the Fourier approach would be to apply a high-pass filter. Since an edge is a set of pixels which are abruptly different from their neighbors, we want to keep the high-frequency components and block out the low ones.

Below we have "before and after" pictures, first of original data and then the picture after an edge-detection process has been applied.[9]



---

[8]Note that we may do more smoothing in some parts of the image than in others.

[9]These pictures are courtesy of Bill Green of the Robotics Laboratory at Drexel University. In this case he is using a Sobel process instead of Fourier analysis, but the result would have been similar for the latter. See his Web tutorial at `www.pages.drexel.edu/~weg22/edge.html`.

The second picture looks like a charcoal sketch! But it was derived mathematically from the original picture, using edge-detection methods.

Note that edge detection methods also may be used to determine where sounds ("ah," "ee") begin and end in speech-recognition applications. In the image case, edge detection is useful for face recognition, etc.

Parallelization here is similar to that of the smoothing case.

## 10.5 The Cosine Transform

It's inconvenient, to say the least, to work with all those complex numbers. But an alternative exists in the form of the **cosine transform**, which is a linear combination of cosines in the one-dimensional case, and of products of cosines in the two-dimensional case.

$$d_{uv} = \frac{2}{\sqrt{mn}} Y(u)Y(v) \sum_{j=0}^{n-1}\sum_{k=0}^{m-1} x_{jk}\, cos\frac{(2j+1)u\pi}{2n} cos\frac{(2k+1)v\pi}{2m}, \tag{10.23}$$

where $Y(0) = 1/\sqrt{2}$ and $Y(t) = 1$ for $t > 0$.

$$x_{jk} = \frac{2}{\sqrt{mn}} \sum_{u=0}^{n-1}\sum_{v=0}^{m-1} Y(u)Y(v)d_{uv}\, cos\frac{(2j+1)u\pi}{2n} cos\frac{(2k+1)v\pi}{2m}, \tag{10.24}$$

## 10.6 Keeping the Pixel Intensities in the Proper Range

Normally pixel intensities are stored as integers between 0 and 255, inclusive. With many of the operations mentioned above, both Fourier-based and otherwise, we can get negative intensity values, or values higher than 255. We may wish to discard the negative values and scale down the positive ones so that most or all are smaller than 256.

Furthermore, even if most or all of our values are in the range 0 to 255, they may be near 0, i.e. too faint. If so, we may wish to multiply them by a constant.

## 10.7   Does the Function g() Really Have to Be Repeating?

It is clear that in the case of a vibrating reed, our loudness function g(t) really is periodic. What about other cases?

A graph of your voice would look "locally periodic." One difference would be that the graph would exhibit more change through time as you make various sounds in speaking, compared to the one repeating sound for the reed. Even in this case, though, your voice *is* repeating within short time intervals, each interval corresponding to a different sound. If you say the word *eye*, for instance, you make an "ah" sound and then an "ee" sound. The graph of your voice would show one repeating pattern during the time you are saying "ah," and another repeating pattern during the time you are saying "ee." So, even for voices, we do have repeating patterns over short time intervals.

On the other hand, in the image case, the function may be nearly constant for long distances (horizontally or vertically), so a local periodicity argument doesn't seem to work there.

The fact is, though, that it really doesn't matter in the applications we are considering here. Even though mathematically our work here has tacitly assumed that our image is duplicated infinitely times (horizontally and vertically),[10] we don't care about this. We just want to get a measure of "wiggliness," and fitting linear combinations of trig functions does this for us.

## 10.8   Vector Space Issues (optional section)

The theory of Fourier series (and of other similar transforms), relies on vector spaces. It actually is helpful to look at some of that here. Let's first discuss the derivation of (10.13).

Define X and C as in Section 10.2. X's components are real, but it is also a member of the vector space V of all n-component arrays of complex numbers.

For any complex number a+bi, define its **conjugate**, $\overline{a + bi} = a - bi$. Note that

$$\overline{e^{i\theta}} = \cos\theta - i\sin\theta == \cos(-\theta) + i\sin(-\theta) = e^{-i\theta} \tag{10.25}$$

---

[10]And in the case of the cosine transform, implicitly we are assuming that the image flips itself on every adjacent copy of the image, first right-side up, then upside-own, then right-side up again, etc.

Define an inner product ("dot product"),

$$[u, w] = \frac{1}{n} \sum_{j=0}^{n-1} u_j \bar{w}_j. \tag{10.26}$$

Define

$$v_h = (1, q^{-h}, q^{-2h}, ..., q^{-(n-1)h}), h = 0, 1, ..., n - 1. \tag{10.27}$$

Then it turns out that the $v_h$ form an orthonormal basis for V.[11] For example, to show orthnogonality, observe that for $r \neq s$

$$[v_r, v_s] = \frac{1}{n} \sum_{j=0}^{n-1} v_{rj} \overline{v_s}_j \tag{10.28}$$

$$= \frac{1}{n} \sum_{j=0}^{n-1} q^{j(-r+s)} \tag{10.29}$$

$$= \frac{1 - q^{(-r+s)n}}{n(1 - q)} \tag{10.30}$$

$$= 0, \tag{10.31}$$

due to the identity $1 + y + y^2 + .... + y^k = \frac{1-y^{k+1}}{1-y}$ and the fact that $q^n = 1$. In the case r = s, the above computation shows that $[v_r, v_s] = 1$.

The DFT of X, which we called C, can be considered the "coordinates" of X in V, relative to this orthonormal basis. The kth coordinate is then $[X, v_k]$, which by definition is (10.13).

The fact that we have an orthonormal basis for V here means that the matrix A/n in (10.20) is an orthogonal matrix. For real numbers, this means that this matrix's inverse is its transpose. In the complex case, instead of a straight transpose, we do a conjugate transpose, $B = \overline{A/n}^t$, where t means transpose. So, B is the inverse of A/n. In other words, in (10.20), we can easily get back to X from C, via

$$X = BC = \frac{1}{n} \bar{A}^t C. \tag{10.32}$$

It's really the same for the nondiscrete case. Here the vector space consists of all the possible periodic functions g() (with reasonable conditions placed regarding continuity etc.) forms the vector space, and the

---

[11]Recall that this means that these vectors are orthogonal to each other, and have length 1, and that they span V.

sine and cosine functions form an orthonormal basis. The $a_n$ and $b_n$ are then the "coordinates" of g() when the latter is viewed as an element of that space.

## 10.9 Bandwidth: How to Read the *San Francisco Chronicle* Business Page (optional section)

The popular press, especially business or technical sections, often uses the term **bandwidth**. What does this mean?

Any transmission medium has a natural range $[f_{min}, f_{max}]$ of frequencies that it can handle well. For example, an ordinary voice-grade telephone line can do a good job of transmitting signals of frequencies in the range 0 Hz to 4000 Hz, where "Hz" means cycles per second. Signals of frequencies outside this range suffer fade in strength, i.e are **attenuated**, as they pass through the phone line.[12]

We call the frequency interval [0,4000] the **effective bandwidth** (or just the **bandwidth**) of the phone line.

In addition to the bandwidth of a **medium**, we also speak of the bandwidth of a **signal**. For instance, although your voice is a mixture of many different frequencies, represented in the Fourier series for your voice's waveform, the really low and really high frequency components, outside the range [340,3400], have very low power, i.e. their $a_n$ and $b_n$ coefficients are small. Most of the power of your voice signal is in that range of frequencies, which we would call the effective bandwidth of your voice waveform. This is also the reason why digitized speech is sampled at the rate of 8,000 samples per second. A famous theorem, due to Nyquist, shows that the sampling rate should be double the maximum frequency. Here the number 3,400 is "rounded up" to 4,000, and after doubling we get 8,000.

Obviously, in order for your voice to be heard well on the other end of your phone connection, the bandwidth of the phone line must be at least as broad as that of your voice signal, and that is the case.

However, the phone line's bandwidth is not much broader than that of your voice signal. So, some of the frequencies in your voice will fade out before they reach the other person, and thus some degree of distortion will occur. It is common, for example, for the letter 'f' spoken on one end to be mis-heard as 's'on the other end. This also explains why your voice sounds a little different on the phone than in person. Still, most frequencies are reproduced well and phone conversations work well.

We often use the term "bandwidth" to literally refer to width, i.e. the width of the interval $[f_{min}, f_{max}]$.

There is huge variation in bandwidth among transmission media. As we have seen, phone lines have bandwidth intervals covering values on the order of $10^3$. For optical fibers, these numbers are more on the order of $10^{15}$.

The radio and TV frequency ranges are large also, which is why, for example, we can have many AM radio

---

[12]And in fact will probably be deliberately filtered out.

stations in a given city. The AM frequency range is divided into subranges, called **channels**. The width of these channels is on the order of the 4000 we need for a voice conversation. That means that the transmitter at a station needs to shift its content, which is something like in the [0,4000] range, to its channel range. It does that by multiplying its content times a sine wave of frequency equal to the center of the channel. If one applies a few trig identities, one finds that the product signal falls into the proper channel!

Accordingly, an optical fiber could also carry many simultaneous phone conversations.

Bandwidth also determines how fast we can set digital bits. Think of sending the sequence 10101010... If we graph this over time, we get a "squarewave" shape. Since it is repeating, it has a Fourier series. What happends if we double the bit rate? We get the same graph, only horizontally compressed by a factor of two. The effect of this on this graph's Fourier series is that, for example, our former $a_3$ will now be our new $a_6$, i.e. the $2\pi \cdot 3 f_0$ frequency cosine wave component of the graph now has the double the old frequency, i.e. is now $2\pi \cdot 6 f_0$. That in turn means that the effective bandwidth of our 10101010... signal has doubled too.

In other words: To send high bit rates, we need media with large bandwidths.