# Multicore & GPU Programming : An Integrated Approach
# Instructor's Manual

Gerassimos Barlas

February 2, 2016

2

# Contents

# Chapter 1

# Introduction

## Exercises

1. Study one of the top 10 most powerful supercomputers in the world. Discover:

   - What kind of operating system does it run?
   - How many CPUs/GPUs is it made of?
   - What is its total memory capacity?
   - What kind of software tools can be used to program it?

   **Answer**

   Students should research the answer by visiting the Top 500 site and -if available- the site of one of the reported systems.

2. How many cores are inside the top GPU offerings from NVidia and AMD? What is the GFlop rating of these chips?

   **Answer** N/A.

3. The performance of the most powerful supercomputers in the world is usually reported as two numbers $Rpeak$ and $Rmax$, both in TFlops (tera floating point operations per second) units. Why is this done? What are the factors reducing performance from $Rpeak$ to $Rmax$? Would it be possible to ever achieve $Rpeak$?

   **Answer**

   This is done because the peak performance is unattainable. Sustained, measured performance on specific benchmarks, is a better indicator of the true machine potential.

   The reason these are different is communication overhead.

   $Rpeak$ and $Rmax$ could never be equal. Extremely compute-heavy applications, that have no inter-node communications, could asymptotically approach $Rpeak$ if they were to run for a very long time. A very long execution time is required to diminish the influence of the start-up costs.

4. A sequential application with a 20% part that must be executed sequentially, is required to be accelerated five-fold. How many CPUs are required for this task?

   **Answer**

   This requires the application of Amdahl's law. The part that can be parallelized is $\alpha = 1 - 20\% = 80\%$. The speedup predicted by Amdahl's law is $speedup = \frac{1}{1-\alpha+\frac{\alpha}{N}}$.

   Achieving a three-fold speedup requires that:

   $$\frac{1}{1-\alpha+\frac{\alpha}{N}} = 3 \Rightarrow \frac{1}{0.2+\frac{0.8}{N}} = 3 \Rightarrow \frac{0.8}{N} = \frac{1}{3} - 0.2 \Rightarrow N = \frac{0.8}{\frac{1}{3}-0.2} = 6 \tag{1.1}$$

   Achieving a 5-fold speedup requires that:

   $$\frac{1}{0.2+\frac{0.8}{N}} = 5 \Rightarrow N = \frac{0.8}{\frac{1}{5}-0.2} = \frac{0.8}{0} = \infty \tag{1.2}$$

   So, it is impossible to achieve a 5-fold speedup, according to Amdahl's law.

5. A parallel application running on 5 identical machines, has a 10% sequential part. What is the speedup relative to a sequential execution on one of the machines? If we would like to double that speedup, how many CPU would be required?

   **Answer**

   This requires the application of Gustafson-Barsis' law as the information relates to a parallel application. The parallel part is $\alpha = 1 - 10\% = 90\%$. The speedup over a single machine is $speedup = 1-\alpha+N\cdot\alpha = .1+5\cdot0.9 = 4.6$.

   Doubling the speedup would require $.1 + N \cdot 0.9 = 9.2 \Rightarrow N = \frac{9.1}{0.9} = 10.1$ machines. As $N$ has to be an integer, we have to round-up to the closest integer, i.e. $N = 11$.

6. An application with a 5% non-parallelizable part, is to be modified for parallel execution. Currently on the market there are two parallel machines available: machine X with 4 CPUs, each CPU capable of executing the application in 1hr on its own, and, machine Y with 16 CPUs, with each CPU capable of executing the application in 2hr on its own. Which is the machine you should buy, if the minimum execution time is required?

   **Answer**

   As the information provided relates to a sequential application, we have to apply Amdahl's law. The execution time for machine X is:

   $$t_X = (1-\alpha)T + \frac{\alpha T}{N} = 0.05 * 1hr + \frac{0.95 \cdot 1hr}{4} = 0.2875hr \tag{1.3}$$

   The execution time for machine Y is:

   $$t_Y = (1-\alpha)T + \frac{\alpha T}{N} = 0.05 * 2hr + \frac{0.95 \cdot 2hr}{16} = 0.21875hr \tag{1.4}$$

So we should buy machine Y.

7. Create a simple sorting application that uses the mergesort algorithm to sort a large collection (e.g. $10^7$) of 32-bit integers. The input data and output results should be stored in files, and the I/O operations should be considered a sequential part of the application. Mergesort is an algorithm that is considered appropriate for parallel execution, although it cannot be equally divided between an arbitrary number of processors, as Amdahl's and Gustafson-Barsis' laws require.

Assuming that this equal division is possible, estimate $\alpha$, i.e. the part of the program that can be parallelized, by using a profiler like `gprof` or `valgrind` to measure the duration of mergesort's execution relative to the overall execution time. Use this number to estimate the predicted speedup for your program.

Does $\alpha$ depend on the size of the input? If it does, how should you modify your predictions and their graphical illustration?

**Answer** N/A

8. A parallel application running on 10 CPUs, spends 15% of its total time, in sequential execution. What kind of CPU (how much faster) would we need to run this application completely sequentially, while keeping the same total time?

**Answer**

This is an application of Gustafson-Barsis' law. If $T$ is the parallel execution time on the 10 CPUs, the sequential execution time on a single one would be $T_s = (1-\alpha)T + N \cdot \alpha \cdot T$. The fast CPU should match the parallel time, i.e. $T_f = T$ which means if should be $\frac{T_s}{T_f} = (1-\alpha) + N \cdot \alpha = 0.15 + 10 \cdot 0.85 = 8.65$ times faster.

# Chapter 2

# Multicore and Parallel Program Design

1. Perform a 2D agglomeration step for the image convolution problem of Section 2.2. What is the resulting number of communication operations?

   **Answer**

   Assuming that we target a grid of task groups forming $K$ rows x $M$ columns and that $K$ and $M$ divide the corresponding dimensions evenly, each group will hold $\frac{IMGX}{M}\frac{IMGY}{K}$ tasks.

   The number of communication operations are:

   - Four for "internal" groups. There are $(K-1)(M-1)$ internal groups. Each group sends $\frac{IMGX}{M}$ pixel values to its top and bottom neighbors, and $\frac{IMGY}{K}$ pixels to its left and right neighbors.
   - Two for corner groups. There are four corner groups.
   - Three for "boundary" groups. There are $2(K-2)+2(M-2)$ boundary groups.

   The total data volume communicated is :

   $$totalComm = (2\frac{IMGX}{M} + 2\frac{IMGY}{K})(K-1)(M-1) + \qquad (2.1)$$

   $$(2\frac{IMGX}{M} + \frac{IMGY}{K})2(M-2) + \qquad (2.2)$$

   $$(\frac{IMGX}{M} + 2\frac{IMGY}{K})2(K-2) + \qquad (2.3)$$

   $$(\frac{IMGX}{M} + \frac{IMGY}{K})4 \qquad (2.4)$$

2. Perform the comparison between the 1D and 2D decompositions of the heat diffusion example in Section 2.3.3, by assuming that (a) half-duplex communication links are available and (b) n-port communications are possible, i.e. all communications can take place at the same time over all the links.

   **Answer**

(a) If half-duplex communication links where used, then the time spend on communication would be doubled:

$$comp_{1D} + comm_{1D} = \frac{N^2}{P} \cdot t_{comp} + 4 \cdot (t_{start} + t_{comm}N)$$

$$comp_{2D} + comm_{2D} = \frac{N^2}{P} t_{comp} + 8 \left( t_{start} + t_{comm} \cdot \frac{N}{\sqrt{P}} \right) \qquad (2.5)$$

which means that

$$comm_{1D} + comp_{1D} < comm_{2D} + comp_{2D} \Rightarrow$$

$$t_{start} > t_{comm} N \left( 1 - \frac{2}{\sqrt{P}} \right) \quad (2.6)$$

So there is no change in the condition that favors 1D over 2D.

(b) In the n-port case, the communication time per time step will be:

$$comm_{1D} = t_{start} + t_{comm}N$$

$$comm_{2D} = t_{start} + t_{comm} \cdot \frac{N}{\sqrt{P}}$$

so the comparison between the two decompositions would be based on

$$comm_{1D} + comp_{1D} < comm_{2D} + comp_{2D} \Rightarrow \sqrt{P} < 1$$

which is obviously false. So, 2D is always better than 1D.

3. How would communication costs affect the pipeline performance? Derive variations of Equations 2.16 to 2.18 that take into account a constant communication overhead between the pipeline stages.

**Answer**

Let's assume that communication between the stages costs a fixed amount of time $t_c$. Then we would have an additional overall time of $t_c(N + M)$, assuming that the last stage also sends data back to whoever is controlling the execution. Thus:

$$t_{total} = \sum_{j=0}^{l-1} t_j + N \cdot t_l + \sum_{j=l+1}^{M-1} t_j + t_c(N + M)$$

The processing rate of the pipeline is:

$$rate = \frac{N}{\sum_{j=0}^{l-1} t_j + N \cdot t_l + \sum_{j=l+1}^{M-1} t_j + t_c(N + M)}$$

The latency of the pipeline is:

$$latency = \sum_{j=0}^{M-2} t_j + t_c M$$

4. The total number of tasks calculated in Section 2.4.5 for the parallel quick-sort of Listing 2.8, is based on the best-case assumption that the input is split in equal halves by every call to the `PartitionData` function. What would be the result if the worst-case (i.e. one part gets $N-1$ elements and the other part 0) were considered?

**Answer**

If we assume that the `PartitionData` function produces a zero sized part and a part with $N-1$ elements, then:

$$T(N) = \begin{cases} 0 & \text{if } N \leq THRES \\ 1 + T(N-1) & \text{if } N > THRES \end{cases}$$

as one of the calls to `QuickSort` would always have no input.

Backward substitution can provide the answer:

$$T(N) = 1 + T(N-1) = 1 + 1 + T(N-2) = k + T(N-k)$$

$T(N-k)$ can be eliminated when $N-k = THRES \Rightarrow k = N - THRES$. Substituting this value of $k$ in the previous equation yields:

$$T(N) = N - THRES + T(THRES) = N - THRES$$

as $T(THRES) = 0$. This is obviously a poor result. What this formula does not convey is that parallelism is also eliminated as a result: there is effectively no overlap between the generated tasks.

5. Use a simple problem instance (e.g. a small array of integers) to trace the execution of the parallel quicksort of Listing 2.8. Create a Gantt graph for the tasks generated, assuming an infinite number of compute nodes is available for executing them. Can you calculate an upper bound for the speedup that can be achieved?

**Answer**

Given $N$ input elements, `PartitionData` executes between $N-1$ and $N+1$ key comparisons. For simplicity we will assume that $N$ comparisons are done and that the array is evenly split into parts equal in size to $\frac{N-1}{2}$ (taking out the pivot element).

The first initial task that starts the sorting operation, will perform $N$ comparisons, before spawning a task for $\frac{N-1}{2}$ elements and continuing with the remaining $\frac{N-1}{2}$. The comparisons are done exclusively inside the `PartitionData` function, which means we have the following sequence of task executions and corresponding comparison steps:

- 1 task, $N$ steps
- 2 tasks, $\frac{N-1}{2}$ steps each
- $2^2$ tasks, $\frac{N-2^2+1}{2^2}$ steps each
- $2^3$ tasks, $\frac{N-2^3+1}{2^3}$ steps each
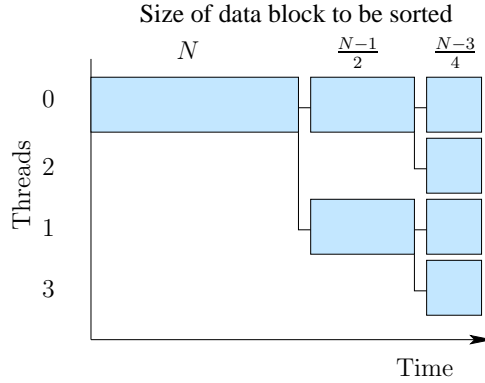- And so on...

Figure 2.1: The first three steps (assuming an ideal scenario of a perfect partition each time) in the parallel quicksort algorithm of Listing 2.8.

Figure 2.1 illustrates this process.

So, the overall duration in number of comparisons is :

$$totalComp = \sum_{i=0}^{L-1} \frac{N - 2^i + 1}{2^i} + THRES$$

where $L$ is associated with when the size become smaller or equal to $THRES$ and no more tasks are spawned.

The value of $L$ is:

$$\frac{N - 2^L + 1}{2^L} = THRES \Rightarrow N - 2^L + 1 = 2^L THRES \Rightarrow$$

$$2^L = \frac{N + 1}{THRES + 1} \Rightarrow L = lg(\frac{N + 1}{THRES + 1})$$

Thus:

$$totalComp = \sum_{i=0}^{L-1} \frac{N - 2^i + 1}{2^i} + THRES =$$

$$(N + 1) \sum_{i=0}^{L-1} (\frac{1}{2})^i - \sum_{i=0}^{L-1} 1 + THRES =$$

$$(N + 1) \frac{(\frac{1}{2})^L - 1}{\frac{1}{2} - 1} - L + THRES =$$

$$2(N + 1)(1 - \frac{THRES + 1}{N + 1}) - L + THRES =$$

$$2(N + 1) - 2(THRES + 1) - L + THRES =$$

$$2N - THRES - L \approx 2N$$

So the maximum speedup that could be ever achieved, given that the sequential quicksort performs $NlgN$ comparisons, is:

$$speedup_{max} = \frac{NlgN}{2N} = \frac{lgN}{2}$$

# Chapter 3

# Shared-memory programming : Threads

## Exercises

1. Enumerate and create the other timing diagrams that show the alternatives of Figure 3.4, when it comes to the final balance of the bank account.

   **Answer**

   All the possible permutations of the four events are allowed, as long as 1 proceeds 3 and 2 proceeds 4. So we have:

   - 1, 2, 3, 4 : produces wrong result
   - 1, 3, 2, 4 : produces correct result
   - 1, 2, 4, 3 : produces wrong result
   - 2, 1, 3, 4 : produces wrong result
   - 2, 1, 4, 3 : produces wrong result
   - 2, 4, 1, 3 : produces correct result

2. Research the term "fork bomb" and write a program that performs as such.

   **Answer**

   ```
   #include <stdlib.h>
   #include <unistd.h>
   #include <limits.h>

   int main (int argc, char **argv)
   {
     int N = atoi (argv[1]);

     for(int i=0;i<N;i++)
        pid_t cID = fork();

     sleep(INT_MAX);
     return 0;
   }
   ```

3. Modify the producer-consumer example shown in Listing 3.11, so that the threads terminate after the number 100 is generated.

**Answer**

The only modification required affects the `consume` method, where the detection of the exit condition is done. It should become:

```
bool consume(int i) {
    // to be implemented
    cout << "@"; // just to show something is happening
    if (i == 100) return true;
    else return false;
}
```

4. Suggest a modification to the program of Listing 3.12 so that the `IntegrCalc` threads can use any function that return a double and takes a double as a parameter.

**Answer**

A pointer to such a function needs to be passed to the `initClass` method, so that it can be called by the threads. The modifications are highlighted in the following code:
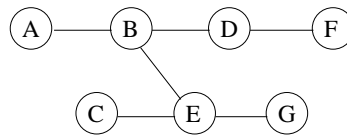
```
1   . . .
2   //————————————————————————————————————————
3   // acts as a consumer
4   class IntegrCalc : public QThread {
5   private:
6       int ID;
7       static QSemaphore *slotsAvail;
8       static QSemaphore *resAvail;
9       static QMutex l2;
10      static QMutex resLock;
11      static Slice *buffer;
12      static int out;
13      static double *result;
14      static double (*f)(double);   // <═══════
15      static QSemaphore numProducts;
16  public:
17      static void initClass(QSemaphore *s, QSemaphore *a, Slice *b,↩
               double *r, double (*f)(double));
18
19      IntegrCalc(int i) : ID(i) {
20      };
21      void run();
22  };
23  //————————————————————————————————————————
24   . . .
25  double (* IntegrCalc::f)(double);    // <═══════
26
27  //————————————————————————————————————————
28
29  void IntegrCalc::initClass(QSemaphore *s, QSemaphore *a, Slice *b↩
          , double *res, double (*g)(double)) {   // <═══════
30      slotsAvail = s;
31      resAvail = a;
32      buffer = b;
33      result = res;
34      *result = 0;
35      f = g ;   // <═══════
36  }
37  //————————————————————————————————————————
38
39  void IntegrCalc::run() {
40      while (1) {
41          resAvail->acquire(); // wait for an available item
42          l2.lock();
43          int tmpOut = out;
44          out = (out + 1) % BUFFSIZE; // update the out index
```

```cpp
45              l2.unlock();
46
47              // take the item out
48              double st = buffer[tmpOut].start;
49              double en = buffer[tmpOut].end;
50              double div = buffer[tmpOut].divisions;
51
52              slotsAvail->release(); // signal for a new empty slot
53
54              if (div == 0) break; // exit
55
56              //calculate area
57              double localRes = 0;
58              double step = (en - st) / div;
59              double x;
60              x = st;
61              localRes = f(st) + f(en);    // <=======
62              localRes /= 2;
63              for(int i=1; i< div; i++)   {
64                  x += step;
65                  localRes += f(x);    // <=======
66                }
67              localRes *= step;
68
69              // add it to result
70              resLock.lock();
71              *result += localRes;
72              resLock.unlock();
73        }
74 }
75 //------------------------------------------------
76
77 int main(int argc, char *argv[]) {
78      if (argc == 1) {
79          cerr << "Usage " << argv[0] << " #threads #jobs\n";
80          exit(1);
81      }
82      int N = atoi(argv[1]);
83      int J = atoi(argv[2]);
84      Slice *buffer = new Slice[BUFFSIZE];
85      QSemaphore avail, buffSlots(BUFFSIZE);
86      int in = 0;
87      double result;
88
89      IntegrCalc::initClass(&buffSlots, &avail, buffer, &result, ↵
               func);   // <=======
90
91   . . .
```

5. In a remote region of Siberia there are single tracks joining railroad stations. Obviously only one train can use a piece of track between two stations. The other trains can wait at the stations before they do their crossings. The following graph indicates the track and station layout:



Write a Qt program that simulates the journey of 3 trains with the following schedules:

- $A \to B \to E \to C$
- $D \to B \to E \to G$
- $C \to E \to B \to D \to F$

As each trains arrives at a station display a relative message. You can assume that a station can hold any number of trains waiting.

**Answer**

The `Train` class constructor expects two vectors : one containing the names of the stations that will be traversed, and one containing references to the binary semaphores controlling access to the lines connecting the stations.

```cpp
#include <QThread>
#include <QMutex>
#include <iostream>
#include <string>
#include <vector>

using namespace std;

class Train : public QThread
{
  private:
    int ID;
    vector<string> *stations;
    vector<QMutex *> *lineLocks;
  public:
    Train(int i, vector<string> *s, vector<QMutex *> *l) : ID(i),
          stations(s), lineLocks(l){}
    void run();
};

void Train::run()
{
    int i;
    for(i=0;i<lineLocks->size();i++)
    {
        cout << "Train " << ID << " is in station " << stations->
             at(i) << endl;
        lineLocks->at(i)->lock();
        cout << "Train " << ID << " is traveling to station " <<
             stations->at(i+1) << endl;
        lineLocks->at(i)->unlock();
    }
    cout << "Train " << ID << " arrived in station " << stations
         ->at(i) << endl;
}


int main(int argc, char *argv[])
{
    QMutex AB, BD, DF, BE, CE, EG; // mutices for the lines
    // one for each train
    vector<QMutex *> routes[3];
    vector<string> st[3];

    routes[0].push_back(&AB);
    routes[0].push_back(&BE);
    routes[0].push_back(&CE);
    st[0].push_back("A");
    st[0].push_back("B");
    st[0].push_back("E");
    st[0].push_back("C");

    routes[1].push_back(&BD);
    routes[1].push_back(&BE);
    routes[1].push_back(&EG);
    st[1].push_back("D");
    st[1].push_back("B");
    st[1].push_back("E");
    st[1].push_back("G");

    routes[2].push_back(&CE);
    routes[2].push_back(&BE);
    routes[2].push_back(&BD);
```

```
60        routes [2]. push_back(&DF);
61        st [2]. push_back("C");
62        st [2]. push_back("E");
63        st [2]. push_back("B");
64        st [2]. push_back("D");
65        st [2]. push_back("F");
66
67        // thread spawning
68        Train *t [3];
69        for(int i=0;i<3;i++)
70        {
71            t[i] = new Train(i, &st[i], &routes[i]);
72            t[i]->start();
73        }
74
75        for(int i=0;i<3;i++)
76            t[i]->wait();
77        return 0;
78  }
```

6. Modify the program of the previous exercise so that each station can hold only 2 trains. Can this lead to deadlocks?

   If you have not done so already, make sure that your program uses only one thread class.

   **Answer**

   This cannot lead to deadlocks as long as we have three or less trains. For four or more we could have pairs of trains trying to cross opposite directions of the same line, leading to deadlock.

   The only required changes involves controlling access to the stations via a set of general semaphores.

```
1   #include <QThread>
2   #include <QMutex>
3   #include <QSemaphore>
4   #include <iostream>
5   #include <string>
6   #include <vector>
7
8   using namespace std;
9
10  class Train : public QThread
11  {
12  private:
13      int ID;
14      vector<string> *stations;
15      vector<QSemaphore *> *stationsCtrl;
16      vector<QMutex *> *lineLocks;
17  public:
18      Train(int i, vector<string> *s, vector<QMutex *> *l, vector<←
            QSemaphore *> *c) : ID(i), stations(s), lineLocks(l), ←
            stationsCtrl(c){}
19      void run();
20  };
21
22  void Train::run()
23  {
24      int i;
25      stationsCtrl->at(0)->acquire();
26      for(i=0;i<lineLocks->size();i++)
27      {
28          cout << "Train " << ID << " is in station " << stations->←
                at(i) << endl;
29          stationsCtrl->at(i+1)->acquire(); // reserve station ←
                before locking the line
30
31          lineLocks->at(i)->lock();
```

```cpp
32              cout << "Train " << ID << " is traveling to station " << ↵
                    stations->at(i+1) << endl;
33              lineLocks->at(i)->unlock();

35              stationsCtrl->at(i)->release(); // release previous ↵
                    station reservation
36              cout<<ID <<" "<< stationsCtrl->at(i)->available() << endl↵
                    ;
37          }
38          cout << "Train " << ID << " arrived in station " << stations↵
                ->at(i) << endl;
39  }


int main(int argc, char *argv[])
{
    QMutex AB, BD, DF, BE, CE, EG;
    QSemaphore A(2), B(2), C(2), D(2), E(2), F(2), G(2);
    // one for each train
    vector<QMutex *> routes[3];
    vector<QSemaphore *> stCapac[3];
    vector<string> st[3];

    routes[0].push_back(&AB);
    routes[0].push_back(&BE);
    routes[0].push_back(&CE);
    st[0].push_back("A");
    st[0].push_back("B");
    st[0].push_back("E");
    st[0].push_back("C");
    stCapac[0].push_back(&A);
    stCapac[0].push_back(&B);
    stCapac[0].push_back(&E);
    stCapac[0].push_back(&C);

    routes[1].push_back(&BD);
    routes[1].push_back(&BE);
    routes[1].push_back(&EG);
    st[1].push_back("D");
    st[1].push_back("B");
    st[1].push_back("E");
    st[1].push_back("G");
    stCapac[1].push_back(&D);
    stCapac[1].push_back(&B);
    stCapac[1].push_back(&E);
    stCapac[1].push_back(&G);

    routes[2].push_back(&CE);
    routes[2].push_back(&BE);
    routes[2].push_back(&BD);
    routes[2].push_back(&DF);
    st[2].push_back("C");
    st[2].push_back("E");
    st[2].push_back("B");
    st[2].push_back("D");
    st[2].push_back("F");
    stCapac[2].push_back(&C);
    stCapac[2].push_back(&E);
    stCapac[2].push_back(&B);
    stCapac[2].push_back(&D);
    stCapac[2].push_back(&F);

    Train *t[3];
    for(int i=0;i<3;i++)
    {
        t[i] = new Train(i, &st[i], &routes[i], &stCapac[i]);
        t[i]->start();
    }

    for(int i=0;i<3;i++)
        t[i]->wait();
    return 0;
}
```

7. A desktop publishing (like PageMaker) application has two threads running: one for running the GUI and one for doing background work. Simulate this application in Qt. Your implementation should have the thread corresponding to the GUI, send requests to the other thread to run tasks on its behalf. The tasks should be (obviously just printing a message is enough for the simulation):

- Printing
- Mail merging
- PDF generation

After performing each requested task, the second thread should wait for a new request to be send to it. Make sure that the first thread does not have to wait for the second thread to finish before making new requests.

**Answer**

The following solution is a producer-consumer derivative, where just two counting semaphores are needed for each pair of interacting threads. There are two buffers, one for handling the interaction between GUI and mailmerge and one for the interaction between GUI and PDF generation.

```cpp
#include <QSemaphore>
#include <QThread>
#include <stdlib.h>
#include <unistd.h>
#include <iostream>

using namespace std;

const int TERMINFLAG = -1;
const int RUNS = 20;
const int BUFFSIZE = 5;

// the following should not be global variables
// but they are so to reduce the length of the code
QSemaphore mailReq (0), mailSpace (BUFFSIZE);
int mailJob[BUFFSIZE];              // buffer
int m_in = 0, mout = 0;

QSemaphore pdfReq (0), pdfSpace (BUFFSIZE);
int pdfJob[BUFFSIZE];               // buffer
int pin = 0, pout = 0;
//==========================================
class GUI: public QThread
{
public:
   void run ();
};
//------------------
void GUI::run ()
{
   for (int i = 0; i < RUNS; i++)
      {
         int choice = rand () % 2;
         if (choice)                      // mail merge
            {
               mailSpace.acquire ();
               mailJob[m_in] = i;
               m_in = (m_in + 1) % BUFFSIZE;
               mailReq.release ();
            }
         else                             // PDF generation
            {
               pdfSpace.acquire ();
               pdfJob[pin] = i;
               pin = (pin + 1) % BUFFSIZE;
```

```
46              pdfReq.release ();
47           }
48        }
49     // term_ination
50     mailSpace.acquire ();
51     mailJob[m_in] = TERMINFLAG;
52     mailReq.release ();
53
54     pdfSpace.acquire ();
55     pdfJob[pin] = TERMINFLAG;
56     pdfReq.release ();
57  }
58
59  //═══════════════════════════════════
60  class MailMerge:public QThread
61  {
62  public:
63     void run ();
64  };
65  //────────────────
66  void MailMerge::run ()
67  {
68     while (1)
69        {
70           mailReq.acquire ();
71           int job = mailJob[mout];
72           mout = (mout + 1) % BUFFSIZE;
73           if (job == TERMINFLAG)
74              break;
75           cout << "Mail merge job #" << job << endl;
76           mailSpace.release ();
77        }
78  }
79
80  //═══════════════════════════════════
81  class PDFGen:public QThread
82  {
83  public:
84     void run ();
85  };
86  //────────────────
87  void PDFGen::run ()
88  {
89     while (1)
90        {
91           pdfReq.acquire ();
92           int job = pdfJob[pout];
93           pout = (pout + 1) % BUFFSIZE;
94           if (job == TERMINFLAG)
95              break;
96           cout << "PDF generation job #" << job << endl;
97           pdfSpace.release ();
98        }
99  }
100
101 //═══════════════════════════════════
102 int main (int argc, char *argv[])
103 {
104    GUI g;
105    PDFGen p;
106    MailMerge m;
107    g.start ();
108    p.start ();
109    m.run ();                        // m.start ();  // can use run ↵
                 instead of start to avoid starting an extra thread
110
111    g.wait ();
112    p.wait ();
113    // m.wait (); // see above
114    return 0;
115 }
```

8. A popular bakery has a baker that cooks a loaf of bread at a time and

deposits it on a counter. Incoming customers pick up a loaf from the counter and exit the bakery. The counter can hold 20 loafs. If it is full the baker stops baking bread. If it is empty, a customer waits. Use semaphores to solve the coordination problem between the baker and the customers.

**Answer**

This is an instance of the producer-consumers problem. Because the baker thread is producing just simple integers in sequence, we can use a `QAtomicInt` variable for retrieving loaf IDs on the customers' side.

```cpp
#include <QThread>
#include <QSemaphore>
#include <QAtomicInt>
#include <iostream>
#include <stdlib.h>

using namespace std;

const int SPACE=20;

QSemaphore loafs(0), counterSpace(SPACE);
QAtomicInt loafID;
//==================================
class Baker : public QThread
{
private:
    int totalBread;
public:
    Baker(int t) : totalBread(t){}
    void run();
};
//----------------------------
void Baker::run()
{
    for(int i=0;i<totalBread;i++)
    {
        counterSpace.acquire();
        cout << "Baker baked # " << i << endl;
        loafs.release();
    }
}
//==================================
class Customer : public QThread
{
private:
    int ID;
public:
    Customer(int i) : ID(i){}
    void run();
};
//----------------------------
void Customer::run()
{
    loafs.acquire();
    int myBread = loafID.fetchAndAddOrdered(1);
    counterSpace.release();
    cout << "Customer " << ID << " got bread #" << myBread << endl
        ;
}
//==================================
int main(int argc, char *argv[])
{
    int totalCustomers = atoi(argv[1]);
    Customer *c[totalCustomers];
    for(int i=0;i<totalCustomers;i++)
    {
        c[i]=new Customer(i);
        c[i]->start();
    }
    Baker b(totalCustomers);
    b.run();
```

```
61
62          // wait for termination
63          for(int i=0;i<totalCustomers;i++)
64              c[i]->wait();
65
66          return 0;
67      }
```

9. Because of customer demand, the bakery owner is considering the following enhancements to his shop:

   (a) Increase the capacity of the counter to 1000

   (b) Hire 3 more bakers

   Modify the solution of the previous exercise to accommodate these changes. Which is the easiest to implement?

   **Answer**

   (a) In this case, the only change required is to modify line 9 in the previous listing to:

```
const int SPACE=1000;
```

   (b) The complication in this case is the termination of the baker threads. For this purpose, the `totalCustomers` variable is used to initialize a general semaphore that is tested and decremented prior to the execution of each iteration of a baker thread (line 27):

```
1   #include <QThread>
2   #include <QSemaphore>
3   #include <QAtomicInt>
4   #include <iostream>
5   #include <stdlib.h>
6
7   using namespace std;
8
9   const int SPACE=20;
10
11  QSemaphore loafs(0), counterSpace(SPACE);
12  QSemaphore moreLoafs(0);
13  QAtomicInt loafID;
14  //======================
15  class Baker : public QThread
16  {
17  private:
18      int ID;
19  public:
20      Baker(int i) : ID(i){}
21      void run();
22  };
23  //--------------------
24  void Baker::run()
25  {
26      int myProd=0;
27      while(moreLoafs.tryAcquire())
28      {
29          counterSpace.acquire();
30          loafs.release();
31          myProd++;
32      }
33      cout << "Baker #" << ID << " baked a total of " << myProd << ↩
               " breads" << endl;
34  }
35  //======================
36  class Customer : public QThread
37  {
38  private:
```

```
39        int ID;
40    public:
41        Customer(int i) : ID(i){}
42        void run();
43    };
44    //————————————————
45    void Customer::run()
46    {
47        loafs.acquire();
48        int myBread = loafID.fetchAndAddOrdered(1);
49        counterSpace.release();
50        cout << "Customer " << ID << " got bread #" << myBread << ↩
                 endl;
51    }
52    //================================
53    int main(int argc, char *argv[])
54    {
55        int totalCustomers = atoi(argv[1]);
56        Customer *c[totalCustomers];
57        moreLoafs.release(totalCustomers);
58        for(int i=0;i<totalCustomers;i++)
59        {
60            c[i]=new Customer(i);
61            c[i]->start();
62        }
63
64        Baker *b[4];
65        for(int i=0;i<4;i++)
66        {
67            b[i] = new Baker(i);
68            b[i] -> start();
69        }
70
71
72        // wait for termination
73        for(int i=0;i<totalCustomers;i++)
74            c[i]->wait();
75        for(int i=0;i<4;i++)
76            b[i]->wait();
77
78        return 0;
79    }
```

10. A bank account class is defined as follows:

```
class BankAccount {
  protected:
      double balance;
      string holderName;
  public:
      double getBalance();
      void deposit(double);
      void withdraw(double, int); // the highest the second ↩
          argument, the higher the priority of the request
};
```

Write the implementation of the three methods given above, so that withdraw operations are prioritized: if there are not enough funds in the account for all, the withdrawals must be done in order of priority regardless if there are some that can be performed with the available funds. You can assume that the priority level in the `withdraw` method is by default equal to 0, and that it can is upper bounded by a fixed constant `MAXPRIORITY`.

**Answer**

Solution is monitor based, using a separate wait condition for each level of priority.

```
1   // Monitor-based solution
2   #include <QThread>
```

```cpp
#include <QMutexLocker>
#include <QWaitCondition>
#include <cstring>
#include <unistd.h>
#include <string>
#include <iostream>

using namespace std;

const int MAXPRIORITY = 10;
const int NUMAGENTS = 20;

//———————————————————————
class BankAccount
{
protected:
  double balance;
  string holderName;
  QMutex l;
  QWaitCondition priCond[MAXPRIORITY];
  int waitingCount[MAXPRIORITY];

public:
    BankAccount (string name, double init);
  double getBalance ();
  void deposit (double);
  void withdraw (double, int);   // the highest the second ←
        argument, the higher the priority of the request
};
//———————————————————————
BankAccount::BankAccount (string name, double init)
{
  holderName = name;
  balance = init;
  memset (waitingCount, 0, MAXPRIORITY * sizeof (int));
}

//———————————————————————
double BankAccount::getBalance ()
{
  return balance;
}

//———————————————————————
void BankAccount::deposit (double x)
{
  QMutexLocker ml (&l);
  balance += x;
  cerr << "Deposited " << x << endl;

  // now determine if any thread is waiting to withdraw
  int priLvl = MAXPRIORITY - 1;
  while (priLvl >= 0)
    {
      if (waitingCount[priLvl] > 0)
        {
          cerr << "Waking level " << priLvl << endl;
          priCond[priLvl].wakeAll ();   // wake up all the threads←
                at that priority level
          return;
        }
      priLvl--;
    }
}

//———————————————————————
// the highest the second argument, the higher the priority of ←
      the request
void BankAccount::withdraw (double x, int lvl)
{
  QMutexLocker ml (&l);
  lvl = (lvl >= MAXPRIORITY) ? MAXPRIORITY - 1 : lvl;   //make ←
        sure priority is not too high
```

```cpp
73      // determine if others at the same or higher priority are ←
                waiting
74      int othersWaiting = 0;
75      int priLvl = MAXPRIORITY − 1;
76      while (priLvl >= lvl)
77        {
78          othersWaiting += waitingCount[priLvl];
79          priLvl−−;
80        }
81
82      // if they are or the funds are not enough wait also
83      if (othersWaiting > 0 || balance < x)
84        {
85          waitingCount[lvl]++;
86          while (x > balance)
87            priCond[lvl].wait (&l);
88
89          balance −= x;
90          waitingCount[lvl]−−;
91          if (waitingCount[lvl] == 0)
92            {
93              priLvl = lvl;
94              while (priLvl >= 0 && waitingCount[priLvl] == 0)
95                priLvl−−;
96              if (priLvl >= 0)
97                priCond[priLvl].wakeAll ();
98            }
99        }
100     else
101       balance −= x;
102
103     cerr << "Withdrew " << x << " with pri " << lvl << endl;
104   }
105
106   //−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−
107   // class for testing purposes
108   class Agent:public QThread
109   {
110   private:
111     BankAccount * ba;
112     double amount;
113     int priority;
114   public:
115     Agent (BankAccount * b, double x, int pri):ba (b), amount (x), ←
              priority (pri) {}
116     void run ();
117   };
118
119   //−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−
120   void Agent::run ()
121   {
122     if (amount > 0)
123       ba−>deposit (amount);
124     else
125       ba−>withdraw (−amount, priority);
126   }
127
128   //−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−
129   int main (int argc, char *argv[])
130   {
131     BankAccount b ("John Doe", 0);
132
133     Agent *ag[NUMAGENTS];
134     for (int i = 0; i < NUMAGENTS; i++)
135       {
136         ag[i] = new Agent (&b, −rand () % 20, i / 2);
137         ag[i]−>start ();
138       }
139     sleep (1);
140     b.deposit (NUMAGENTS * 20); // make sure enough funds are ←
              available
141     for (int i = 0; i < NUMAGENTS; i++)
142       ag[i]−>wait ();
143
```

```
144     return 0;
145  }
```

11. The IT department of a big corporation is equipped with 5 high speed
    printers that are used by a multitude of threads. The threads are part of
    the same accounting process. Each of the threads is supposed to perform
    the following (pseudocode) sequence in order to printout any material:

    ```
    ...
    printerID = get_available_printer();
    // print to printerID printer
    releasePrinter(printerID);
    ...
    ```

    Write an appropriate implementation for the two functions listed above
    using semaphores.  You can assume that the available printer IDs are
    stored in a shared buffer.

    **Answer**

    Solution is based on the producers-consumers pattern. A buffer with the
    available printer IDs is setup and managed. Function `get_available_printer()`
    operates as a consumer and function `releasePrinter()` operates as a pro-
    ducer. Only one counting semaphore is required as the buffer is as big as
    the overall number of printers.

    ```cpp
    1   #include <QThread>
    2   #include <QMutex>
    3   #include <QSemaphore>
    4   #include <unistd.h>
    5   #include <iostream>
    6
    7   using namespace std;
    8
    9   const int PRINTERS = 5;
    10  const int NUMAGENTS = 20;
    11
    12  int printIDbuffer[PRINTERS]={0,1,2,3,4};
    13  QSemaphore prnAvail(PRINTERS);
    14  int in=0, out=0;
    15  QMutex l1, l2;
    16
    17  //------------------------------------
    18  int get_available_printer()
    19  {
    20      prnAvail.acquire();
    21      l1.lock();
    22      int tmp = printIDbuffer[out];
    23      out = (out +1)%PRINTERS;
    24      l1.unlock();
    25      return tmp;
    26  }
    27  //------------------------------------
    28  void releasePrinter(int printerID)
    29  {
    30      l2.lock();
    31      printIDbuffer[in]=printerID;
    32      in = (in +1)%PRINTERS;
    33      l2.unlock();
    34      prnAvail.release();
    35  }
    36
    37
    38  //------------------------------------
    39  // class for testing purposes
    40  class Agent:public QThread
    41  {
    42  private:
    ```

```
43    int ID;
44  public:
45    Agent (int i) : ID(i){}
46    void run ();
47  };
48
49  //——————————————————————
50  void Agent::run ()
51  {
52    sleep(rand()%3);
53    int printerID = get_available_printer();
54
55    // print to printerID printer
56    cout << ID << " is printing to " << printerID << endl;
57    releasePrinter(printerID);
58  }
59
60  //——————————————————————
61  int main (int argc, char *argv[])
62  {
63    Agent *ag[NUMAGENTS];
64    for (int i = 0; i < NUMAGENTS; i++)
65      {
66        ag[i] = new Agent (i);
67        ag[i]->start ();
68      }
69
70    for (int i = 0; i < NUMAGENTS; i++)
71      ag[i]->wait ();
72
73    return 0;
74  }
```

12. Create 3 threads, each printing out the letters 'A', 'B' and 'C'. The printing must adhere to these rules:

    - The total number of 'B's and 'C's that have been output at any point in the output string cannot exceed the total number of 'A's that have been output at that point.

    - After a 'C' has been output, another 'C' cannot be output until one or more 'B's have been output.

    Use semaphores to solve the problem.

    **Answer**

```
1   #include <QThread>
2   #include <QMutex>
3   #include <QSemaphore>
4   #include <iostream>
5   #include <unistd.h>
6
7   using namespace std;
8
9   QSemaphore permit(0);
10  QMutex allowC;
11  volatile bool Bflag=true;
12
13  //——————————————————————
14  class ThrA : public QThread
15  {
16  public:
17      void run();
18  };
19  //————————————
20  void ThrA::run()
21  {
22    while(1)
23      {
```

```
24            cout << 'A';
25            permit.release();
26        }
27    }
28
29    //—————————————————————————————
30    class ThrB : public QThread
31    {
32    public:
33        void run();
34    };
35    //———————————————
36    void ThrB::run()
37    {
38        while(1)
39        {
40            permit.acquire();
41            cout << 'B';
42            if(Bflag==false)
43            {
44                allowC.unlock();
45                Bflag=true;
46            }
47        }
48    }
49    //—————————————————————————————
50    class ThrC : public QThread
51    {
52    public:
53        void run();
54    };
55    //———————————————
56    void ThrC::run()
57    {
58        while(1)
59        {
60            allowC.lock();
61            permit.acquire();
62            cout << 'C';
63            Bflag=false;
64        }
65    }
66    //—————————————————————————————
67    int main (int argc, char *argv[])
68    {
69        ThrA a;
70        ThrB b;
71        ThrC c;
72
73        a.start();
74        b.start();
75        c.run();
76
77        return 0;
78    }
```

13. Modify the previous exercise so that the printing is governed by this set of rules:

   - One 'C' must be output after two 'A's and three 'B's are output.
   - While there is no restriction on the order of printing 'A' and 'B', the corresponding threads must wait for a 'C' to be printed when the previous condition is met.

   Use a monitor to solve the problem.

   **Answer**

   The following solution uses the "critical section outside the monitor" approach.

```cpp
#include <QThread>
#include <QMutex>
#include <QMutexLocker>
#include <QWaitCondition>
#include <iostream>
#include <unistd.h>

using namespace std;

class Monitor
{
private:
    int cntA, cntB;
    QMutex l;
    QWaitCondition cAB, cC;
public:
    Monitor() : cntA(0), cntB(0){}
    void allowA();
    void allowB();
    void allowC();
    void doneA();
    void doneB();
    void doneC();
};
//_____
void Monitor::allowA()
{
    QMutexLocker ml(&l);
    while(cntA == 3)
        cAB.wait(&l);
}
//_____
void Monitor::allowB()
{
    QMutexLocker ml(&l);
    while(cntB == 2)
        cAB.wait(&l);
}
//_____
void Monitor::allowC()
{
    QMutexLocker ml(&l);
    while(cntA != 3 || cntB !=2)
        cC.wait(&l);
}
//_____
void Monitor::doneA()
{
    QMutexLocker ml(&l);
    cntA++;
    if(cntA==3)
        cC.wakeOne();
}
//_____
void Monitor::doneB()
{
    QMutexLocker ml(&l);
    cntB++;
    if(cntB==2)
        cC.wakeOne();
}
//_____
void Monitor::doneC()
{
    QMutexLocker ml(&l);
    cntA = cntB = 0;
    cAB.wakeAll();
}
//_____
class ThrA : public QThread
{
private:
Monitor *m;
public:
```

```
75          ThrA ( Monitor  ∗x )  :  m ( x ) { }
76          void  run ( ) ;
77     } ;
78     //————————————
79     void  ThrA : : run ( )
80     {
81          while ( 1 )
82          {
83               m−>allowA ( ) ;
84               cout<<'A' ;
85               m−>doneA ( ) ;
86          }
87     }
88
89     //————————————————————————————————
90     class  ThrB  :  public  QThread
91     {
92     private :
93     Monitor  ∗m ;
94     public :
95          ThrB ( Monitor  ∗x )  :  m ( x ) { }
96          void  run ( ) ;
97     } ;
98     //————————————
99     void  ThrB : : run ( )
100    {
101          while ( 1 )
102          {
103               m−>allowB ( ) ;
104               cout<<'B' ;
105               m−>doneB ( ) ;
106          }
107    }
108    //————————————————————————————————
109    class  ThrC  :  public  QThread
110    {
111    private :
112    Monitor  ∗m ;
113    public :
114          ThrC ( Monitor  ∗x )  :  m ( x ) { }
115          void  run ( ) ;
116    } ;
117    //————————————
118    void  ThrC : : run ( )
119    {
120          while ( 1 )
121          {
122               m−>allowC ( ) ;
123               cout<<'C' ;
124               m−>doneC ( ) ;
125          }
126    }
127    //————————————————————————————————
128    int  main  ( int  argc ,  char  ∗argv [ ] )
129    {
130       Monitor  m ;
131       ThrA  a(&m ) ;
132       ThrB  b(&m ) ;
133       ThrC  c(&m ) ;
134
135       a . start ( ) ;
136       b . start ( ) ;
137       c . run ( ) ;
138
139       return  0 ;
140    }
```

14. Address the termination problem in the previous exercise. How can the three threads terminate after e.g. a fixed number of As has been output? Or when a fixed total number of character has been output?

   **Answer**

The following code solves the termination problem for a fixed number of total characters. A character is reserved at the output before the termination of the `allow?()` methods, which compared to the previous exercise, return a boolean to indicate whether to continue execution or not. If the maximum number of characters has been reached, appropriate signals wake up any waiting threads.

```cpp
#include <QThread>
#include <QMutex>
#include <QMutexLocker>
#include <QWaitCondition>
#include <iostream>
#include <unistd.h>

using namespace std;

class Monitor
{
private:
    int cntA, cntB, total;
    QMutex l;
    QWaitCondition cAB, cC;
public:
    Monitor (int t):cntA (0), cntB (0), total (t) {}
    bool allowA ();
    bool allowB ();
    bool allowC ();
    void doneA ();
    void doneB ();
    void doneC ();
};

//----------------------
bool Monitor::allowA ()
{
    QMutexLocker ml (&l);
    while (cntA == 3 && total > 0)
        cAB.wait (&l);
    if (total)
        {
            total--;
            return true;
        }
    else
        {
            cAB.wakeOne ();
            cC.wakeOne ();
            return false;
        }
}

//----------------------
bool Monitor::allowB ()
{
    QMutexLocker ml (&l);
    while (cntB == 2 && total > 0)
        cAB.wait (&l);

    if (total)
        {
            total--;
            return true;
        }
    else
        {
            cAB.wakeOne ();
            cC.wakeOne ();
            return false;
        }
}

//----------------------
```

```cpp
66    bool Monitor::allowC ()
67    {
68      QMutexLocker ml (&l);
69      while ((cntA != 3 || cntB != 2) && (total > 0))
70        cC.wait (&l);
71
72      if (total)
73        {
74          total--;
75          return true;
76        }
77      else
78        {
79          cAB.wakeOne ();
80          cC.wakeOne ();
81          return false;
82        }
83    }
84
85    //_____
86    void Monitor::doneA ()
87    {
88      QMutexLocker ml (&l);
89      cntA++;
90      if (cntA == 3)
91        cC.wakeOne ();
92    }
93
94    //_____
95    void Monitor::doneB ()
96    {
97      QMutexLocker ml (&l);
98      cntB++;
99      if (cntB == 2)
100       cC.wakeOne ();
101   }
102
103   //_____
104   void Monitor::doneC ()
105   {
106     QMutexLocker ml (&l);
107     cntA = cntB = 0;
108     cAB.wakeAll ();
109   }
110
111   //_____
112   class ThrA:public QThread
113   {
114   private:
115     Monitor * m;
116   public:
117     ThrA (Monitor * x):m (x)
118     {
119     }
120     void run ();
121   };
122
123   //_____
124   void ThrA::run ()
125   {
126     while (m->allowA ())
127       {
128         cout << 'A';
129         m->doneA ();
130       }
131   }
132
133   //_____
134   class ThrB:public QThread
135   {
136   private:
137     Monitor * m;
138   public:
139     ThrB (Monitor * x):m (x)
```

```
140      {
141      }
142      void run ();
143   };
144
145   //------------------
146   void ThrB::run ()
147   {
148      while (m->allowB ())
149         {
150            cout << 'B';
151            m->doneB ();
152         }
153   }
154
155   //------------------------------------------
156   class ThrC:public QThread
157   {
158   private:
159      Monitor * m;
160   public:
161      ThrC (Monitor * x):m (x)
162         {
163         }
164      void run ();
165   };
166
167   //------------------
168   void ThrC::run ()
169   {
170      while (m->allowC ())
171         {
172            cout << 'C';
173            m->doneC ();
174         }
175   }
176
177   //------------------------------------------
178   int main (int argc, char *argv[])
179   {
180      Monitor m (100);
181      ThrA a (&m);
182      ThrB b (&m);
183      ThrC c (&m);
184
185      a.start ();
186      b.start ();
187      c.run ();
188
189      a.wait ();
190      b.wait ();
191      return 0;
192   }
```

15. Create 4 threads, each printing out the letters 'A', 'B', 'C', and 'D'. The printing must adhere to these rules:

   - The total number of 'A's and 'B's that have been output at any point in the output string cannot exceed the total number of 'C's and 'D's that have been output at that point.

   - The total number of 'A's that have been output at any point in the output string cannot exceed twice the number of 'B's that have been output at that point.

   - After a 'C' has been output, another 'C' cannot be output until one or more 'D's have been output.

   Solve the problem using (a) semaphores and (b) a monitor.

**Answer**

(a)

```cpp
1    #include <QThread>
2   #include <QMutex>
3   #include <QSemaphore>
4   #include <iostream>
5   #include <unistd.h>
6
7   using namespace std;
8
9   QSemaphore allowAB (0), permitA (0), total (100);
10  QMutex permitC;
11  bool DisOut = true;
12
13  //——————————————————————————
14  class ThrA:public QThread
15  {
16  public:
17    void run ();
18  };
19
20  //————————————
21  void ThrA::run ()
22  {
23    while (1)
24      {
25        permitA.acquire ();  // for rule #2
26        allowAB.acquire ();  // for rule #1
27        if (!total.tryAcquire ())
28          break;
29        cout << 'A';
30        usleep (1);
31      }
32  }
33
34  //——————————————————————————
35  class ThrB:public QThread
36  {
37  public:
38    void run ();
39  };
40
41  //————————————
42  void ThrB::run ()
43  {
44    while (1)
45      {
46        allowAB.acquire (); // for rule #1
47        if (!total.tryAcquire ())
48          break;
49        cout << 'B';
50        permitA.release (2); // for rule #2
51        usleep (1);
52      }
53  }
54
55  //——————————————————————————
56  class ThrC:public QThread
57  {
58  public:
59    void run ();
60  };
61
62  //————————————
63  void ThrC::run ()
64  {
65    while (1)
66      {
67        permitC.lock (); // for rule #3
68        if (!total.tryAcquire ())
69          break;
70        cout << 'C';
```

```
71          DisOut = false;
72          allowAB.release (); // for rule #1
73          usleep (1);
74        }
75    }
76
77    //————————————————————————————
78    class ThrD: public QThread
79    {
80    public:
81      void run ();
82    };
83
84    //—————————————————
85    void ThrD::run ()
86    {
87      while (1)
88        {
89          cout << 'D';
90          if (!total.tryAcquire ())
91            break;
92          allowAB.release (); // for rule #1
93          if (DisOut == false) // for rule #3
94            {
95              DisOut = true;
96              permitC.unlock ();
97            }
98          usleep (1);
99        }
100
101     // Make sure all terminate
102     // This is the most appropriate thread to do this, as ThrD does←
                  not wait for any other
103     allowAB.release (2);
104     permitA.release ();
105     if (DisOut == false)
106       {
107         permitC.unlock ();
108       }
109
110   }
111
112   //————————————————————————————
113   int main (int argc, char *argv[])
114   {
115     ThrA a;
116     ThrB b;
117     ThrC c;
118     ThrD d;
119
120     a.start ();
121     b.start ();
122     c.start ();
123     d.run ();
124
125     a.wait ();
126     b.wait ();
127     c.wait ();
128     return 0;
129   }
```

(b) In this case, the `Monitor` class could be equipped with generic `allow()` and `done()` methods, that would accept as a parameter the character that is supposed to be printed, i.e. the type of thread calling. In the following code we stick with separate pairs of `Monitor` methods for each type of thread, to make the program easier to read.

```
1   #include <QThread>
2   #include <QMutex>
3   #include <QMutexLocker>
4   #include <QWaitCondition>
5   #include <iostream>
```

```cpp
6   #include <unistd.h>
7
8   using namespace std;
9
10  class Monitor
11  {
12  private:
13    int allowAB, permitA, total;
14    bool permitC;
15    QMutex l;
16    QWaitCondition cAB, cA, cC;   // one condition for each of the ↩
            three rules
17  public:
18      Monitor (int t):allowAB (0), permitA (0), total (t), permitC ↩
            (true)
19    {
20    }
21    bool allowA ();
22    bool allowB ();
23    bool allowC ();
24    bool allowD ();
25    void doneA ();
26    void doneB ();
27    void doneC ();
28    void doneD ();
29  };
30
31  //_____
32  bool Monitor::allowA ()
33  {
34    QMutexLocker ml (&l);
35    while (permitA == 0 && total > 0)
36      cA.wait (&l);
37
38    while (allowAB == 0 && total > 0)
39      cAB.wait (&l);
40
41    allowAB --;
42    permitA --;
43    if (total)
44      {
45        total --;
46        return true;
47      }
48    else
49      {
50        cAB.wakeOne ();
51        cC.wakeOne ();
52        return false;
53      }
54  }
55
56  //_____
57  bool Monitor::allowB ()
58  {
59    QMutexLocker ml (&l);
60    while (allowAB == 0 && total > 0)
61      cAB.wait (&l);
62
63    allowAB --;
64    if (total)
65      {
66        total --;
67        return true;
68      }
69    else
70      {
71        cAB.wakeOne ();
72        cC.wakeOne ();
73        return false;
74      }
75  }
76
77  //_____
```

```cpp
78   bool Monitor::allowC ()
79   {
80     QMutexLocker ml (&l);
81     while (permitC == false && total > 0)
82       cC.wait (&l);
83
84     if (total)
85       {
86         total--;
87         return true;
88       }
89     else
90       {
91         cAB.wakeAll ();
92         cA.wakeOne ();
93         return false;
94       }
95   }
96
97   //----------------
98   bool Monitor::allowD ()
99   {
100    QMutexLocker ml (&l);
101    if (total)
102      {
103        total--;
104        return true;
105      }
106    else
107      {
108        cAB.wakeAll ();
109        cA.wakeOne ();
110        cC.wakeOne ();
111        return false;
112      }
113  }
114
115  //----------------
116  void Monitor::doneA ()
117  {
118    // nothing to be done here
119  }
120
121  //----------------
122  void Monitor::doneB ()
123  {
124    QMutexLocker ml (&l);
125    permitA += 2;
126    cA.wakeOne ();
127  }
128
129  //----------------
130  void Monitor::doneC ()
131  {
132    QMutexLocker ml (&l);
133    allowAB++;
134    permitC = false;
135    cAB.wakeAll ();
136  }
137
138  //----------------
139  void Monitor::doneD ()
140  {
141    QMutexLocker ml (&l);
142    if (permitC == false)
143      {
144        permitC = true;
145        cC.wakeOne ();
146      }
147    allowAB++;
148    cAB.wakeAll ();
149  }
150
151  //--------------------------------
```

```
152    class ThrA : public QThread
153    {
154    private :
155      Monitor * m;
156    public :
157      ThrA (Monitor * x):m (x)
158      {
159      }
160      void run ();
161    };
162
163    //————————————
164    void ThrA :: run ()
165    {
166      while (m->allowA ())
167        {
168          cout << 'A';
169          m->doneA ();
170          usleep (1);
171        }
172    }
173
174    //——————————————————————————
175    class ThrB : public QThread
176    {
177    private :
178      Monitor * m;
179    public :
180      ThrB (Monitor * x):m (x)
181      {
182      }
183      void run ();
184    };
185
186    //————————————
187    void ThrB :: run ()
188    {
189      while (m->allowB ())
190        {
191          cout << 'B';
192          m->doneB ();
193          usleep (1);
194        }
195    }
196
197    //——————————————————————————
198    class ThrC : public QThread
199    {
200    private :
201      Monitor * m;
202    public :
203      ThrC (Monitor * x):m (x)
204      {
205      }
206      void run ();
207    };
208
209    //————————————
210    void ThrC :: run ()
211    {
212      while (m->allowC ())
213        {
214          cout << 'C';
215          m->doneC ();
216          usleep (1);
217        }
218    }
219
220    //——————————————————————————
221    class ThrD : public QThread
222    {
223    private :
224      Monitor * m;
225    public :
```

```
226    ThrD (Monitor * x):m (x)
227    {
228    }
229    void run ();
230  };
231
232  //————————————
233  void ThrD::run ()
234  {
235     while (m->allowD ())
236       {
237          cout << 'D';
238          m->doneD ();
239          usleep (1);
240       }
241  }
242
243  //——————————————————————————
244  int main (int argc, char *argv[])
245  {
246     Monitor m (100);
247     ThrA a (&m);
248     ThrB b (&m);
249     ThrC c (&m);
250     ThrD d (&m);
251
252     a.start ();
253     b.start ();
254     c.start ();
255     d.run ();
256
257     a.wait ();
258     b.wait ();
259     c.wait ();
260     return 0;
261  }
```

16. Use semaphores to solve the typical cigarette smokers problem, where the agent signals directly the smoker missing the two ingredients placed on the table.

**Answer**

```
1   #include <QThread>
2   #include <QMutex>
3   #include <QSemaphore>
4   #include <iostream>
5   #include <stdlib.h>
6
7   using namespace std;
8
9   #define MAXSLEEP 1000
10  #define TOBACCO_PAPER 0
11  #define TOBACCO_MATHCES 1
12  #define PAPER_MATHCES 2
13
14  QSemaphore missingIngr[3];
15  QSemaphore   wakeAgent(0);
16  bool termFlag=false;
17  const char *msg[]={"having matches", "having paper", "having ←
           tobacco"};
18  //************************************************
19  class Smoker : public QThread
20  {
21     private:
22        int missing;
23     public:
24        Smoker(int);
25        void run();
26  };
27  //——————————————————————————————
28  Smoker::Smoker(int m) : missing(m){}
```

```
29   //————————————————————————————————————
30   void  Smoker :: run ()
31   {
32      while (1)
33      {
34         missingIngr [ missing ] . acquire () ;   // wait for agent to send ←
                    signal
35         if ( termFlag )  break ;                 // termination check
36         cout  <<  " Smoker "  <<  msg [ missing ]  <<  " is  smoking \n" ;
37         msleep ( rand () % MAXSLEEP ) ;
38         wakeAgent . release () ;                 // wake up agent
39      }
40   }
41   //****************************************************
42   class  Agent  :  public  QThread
43   {
44      private :
45         int  runs ;
46      public :
47         Agent ( int ) ;
48         void  run () ;
49   } ;
50   //————————————————————————————————————
51   Agent :: Agent ( int  r )  :  runs ( r ) { }
52   //————————————————————————————————————
53   void  Agent :: run ()
54   {
55      for ( int  i =0; i< runs ;  i++)
56      {
57         int  ingreds  =  rand () % 3 ;
58         missingIngr [ ingreds ] . release () ;
59         wakeAgent . acquire () ;
60      }
61   // set  termination  flag  and  wake  up  all  threads
62   termFlag= true ;
63   missingIngr [ 0 ] . release () ;
64   missingIngr [ 1 ] . release () ;
65   missingIngr [ 2 ] . release () ;
66   }
67   //****************************************************
68   int  main ( int  argc ,  char  ** argv )
69   {
70      Smoker  *s [ 3 ] ;
71      for ( int  i =0; i <3; i++)
72      {
73        s [ i ]  =  new  Smoker ( i ) ;
74        s [ i ]−> start () ;
75      }
76      Agent  a ( atoi ( argv [ 1 ] ) ) ;
77      a . run () ;
78
79      for ( int  i =0; i <3; i++)
80        s [ i ]−> wait () ;
81
82      return  EXIT_SUCCESS ;
83   }
```

17. Solve the cigarette smokers problem as described in Section 3.6.2, using semaphores.

   **Answer**

   In the following solution all smoker threads wake up and check the table contents.

```
1   #include  <QThread>
2   #include  <QMutex>
3   #include  <QSemaphore>
4   #include  <iostream>
5   #include  <stdlib .h>
6
7   using  namespace  std ;
```

```cpp
#define MAXSLEEP 10
#define TOBACCO_PAPER 0
#define TOBACCO_MATHCES 1
#define PAPER_MATHCES 2

QSemaphore wakeUpSmoker[3];
QSemaphore wakeAgent(0);
volatile bool termFlag=false;
int table;
const char *msg[]={"having matches", "having paper", "having
    tobacco"};
//***************************************************
class Smoker : public QThread
{
   private:
      int missing;
   public:
      Smoker(int);
      void run();
};
//——————————————————————————————————————
Smoker::Smoker(int m) : missing(m){}
//——————————————————————————————————————
void Smoker::run()
{
   while(1)
   {
      do
        {
          wakeUpSmoker[missing].acquire();
        }
      while(table != missing && !termFlag);

      if(termFlag) break;                // termination check
      cout << "Smoker " << msg[missing] << " is smoking\n";
      msleep(rand() % MAXSLEEP);
      wakeAgent.release();               // wake up agent
   }
}
//***************************************************
class Agent : public QThread
{
   private:
      int runs;
   public:
      Agent(int);
      void run();
};
//——————————————————————————————————————
Agent::Agent(int r) : runs(r){}
//——————————————————————————————————————
void Agent::run()
{
   for(int i=0;i<runs; i++)
   {
      table = rand() % 3;
      wakeUpSmoker[0].release();
      wakeUpSmoker[1].release();
      wakeUpSmoker[2].release();
      wakeAgent.acquire();
   }
  // set termination flag and wake up all threads
  termFlag=true;
  wakeUpSmoker[0].release();
  wakeUpSmoker[1].release();
  wakeUpSmoker[2].release();
}
//***************************************************
int main(int argc, char **argv)
{
   Smoker *s[3];
   for(int i=0;i<3;i++)
   {
```

```
81          s[i] = new Smoker(i);
82          s[i]->start();
83        }
84        Agent a(atoi(argv[1]));
85        a.run();
86
87        for(int i=0;i<3;i++)
88          s[i]->wait();
89
90        return EXIT_SUCCESS;
91      }
```

18. Model the movie-going process at a multiplex cinema using a monitor. Assume the following conditions:

   - There are 3 different movies played at the same time in 3 halls. The capacity of each hall is 4, 5 and 7 respectively.

   - One hundred customers are waiting to see a randomly chosen movie.

   - A cashier issues the tickets.

   - If a hall is full a movie begins to play.

   - A customer cannot enter a hall while a movie is playing or while the previous viewers are exiting the hall.

   - A movie will play for the last customers even if the corresponding hall is not full.

   **Answer**

```
1    #include <QThread>
2    #include <QMutex>
3    #include <QWaitCondition>
4    #include <iostream>
5    #include <vector>
6    #include <stdlib.h>
7
8    using namespace std;
9
10   const int NUMCUST=100;
11   //*****************************************
12   class Monitor
13   {
14   private:
15       int numHalls;
16       vector<int> &capacity;
17       int *inside;
18       int *left;
19       bool *finished;
20       int numCustomersRemain;
21       QMutex l;
22       QWaitCondition *condE; // 'enter' array of conditions. One ↩
                for each hall
23       QWaitCondition *condL; // 'leave' array of conditions
24   public:
25       Monitor(vector<int> &cap, int numCust);
26       ~Monitor();
27       void enterTheater(int movie);
28       void leaveTheater(int movie);
29   };
30   //————————————————————————————————————————
31   Monitor::Monitor(vector<int> &cap, int numCust)  : capacity(cap)
32   {
33       numHalls=cap.size();
34       numCustomersRemain = numCust;
35       inside = new int[numHalls];
36       for(int i=0;i<numHalls;i++) inside[i]=0;
37       left = new int[numHalls];
```

```cpp
        for(int i=0;i<numHalls;i++) left[i]=0;
        finished =  new bool[numHalls];
        for(int i=0;i<numHalls;i++) finished[i]=false;
        condE = new QWaitCondition[numHalls];
        condL = new QWaitCondition[numHalls];
}
//————————————————————————————
Monitor::~Monitor()
{
        delete[] inside;
        delete[] left;
        delete[] finished;
        delete[] condE;
        delete[] condL;
}
//————————————————————————————
void Monitor::enterTheater(int m)
{
        QMutexLocker ml(&l);
        while(inside[m] == capacity[m])
            condE[m].wait(&l);
        inside[m]++;
        if(inside[m] == capacity[m])    // is hall full?
        {
            cout << "Movie " << m << " started\n";
            finished[m]=true;          // set end-of-movie flag
            left[m]=0;
            condL[m].wakeAll();        // let all customers leave
        }

        numCustomersRemain--;
        if(numCustomersRemain==0) // start the remaining shows
        {
            for(int i=0;i<numHalls;i++)
                if(inside[i]>0)
                {
                    cout << "Movie " << i << " started\n";
                    finished[i]=true;
                    left[i]=0;
                    condL[i].wakeAll();
                }
        }
}
//————————————————————————————
void Monitor::leaveTheater(int m)
{
        QMutexLocker ml(&l);
        while(finished[m]==false)
            condL[m].wait(&l);
        left[m]++;
        if(left[m] == capacity[m])   // is hall empty
        {
            finished[m]=false;
            inside[m]=0;
            left[m]=0;
            condE[m].wakeAll();
        }
}
//*****************************************
class Customer : public QThread
{
private:
        int ID;
        int movieID;
        Monitor *m;
public:
        Customer(int i, Monitor *c, int mID) : ID(i), movieID(mID), m←
            (c) {};
        void run();
};
//————————————————————————————
void Customer::run()
{
        cout << "Customer " << ID << " wants to see " << movieID << "←
```

```
                    \n";
111        m−>enterTheater(movieID);
112        cout << "Customer " << ID << " entered the theater\n";
113        m−>leaveTheater(movieID);
114        cout << "Customer " << ID << " left the movies\n";
115    }
116    //******************************************
117    int main(int argc, char *argv[])
118    {
119        vector<int> c;
120        c.push_back(4);
121        c.push_back(5);
122        c.push_back(7);
123
124        Monitor m(c, NUMCUST);
125        Customer *cust[NUMCUST];
126        for(int i=0;i<NUMCUST;i++)
127        {
128            cust[i] = new Customer(i, &m, rand()%3);
129            cust[i]−>start();
130        }
131
132        for(int i=0;i<NUMCUST;i++)
133            cust[i]−>wait();
134
135        return 0;
136    }
```

19. Write a multi-threaded password cracker based on the Producer-Consumer paradigm. The producer should generate plain-text passwords according to a set of rules and the consumers should be hashing each password and checking whether it matches a target signature. All the threads should terminate upon the discovery of a matching password. You can use the MD5 cryptographic hash function for this exercise.

**Answer**

The following listing shows a monitor based solution to the problem. There are two classes whose instances act as a producer (class `PassGenerator`) and consumers (class `PassChecker`).

The main program is expecting two command-line parameters : the number of password checkers and the maximum length of the brute-force generated passwords to be examined. Once a password with the same MD5 hash value is found, the `finishUp` method is called, which in turn terminates all remaining threads upon their next interaction with the Monitor.

```
1    #include <iostream>
2    #include <string>
3    #include <string.h>
4    #include "md5.h"
5    #include <QThread>
6    #include <QMutex>
7    #include <QMutexLocker>
8    #include <QWaitCondition>
9
10   using namespace std;
11
12   char *alphabet="1234567890abcdefghijklmnopqrstuvwxyz";
13   string target("5912d7bfd10f631f1715bf85bbb72d97");  // ciphertext
14   const int PASSSTORESIZE=100;  // size of password buffer
15   //******************************************
16   class Monitor
17   {
18   private:
19     QMutex l;
20     QWaitCondition full;
21     QWaitCondition empty;
```

```cpp
22      char *buff[PASSSTORESIZE];
23      int in, out, N, Ncons;
24      int total;
25      bool found;
26   public:
27      Monitor(int consumers, int maxLen);
28      ~Monitor();
29      bool putCandidate(char *);   // both return false when it is ←
               time to stop
30      bool getCandidate(char *);
31      void finishUp();
32   };
33   //———————————————————————————
34   Monitor::Monitor(int consumers, int maxLen)
35   {
36      Ncons = consumers;
37      buff[0]= new char[(maxLen+1)*PASSSTORESIZE];
38      for(int i=1;i<PASSSTORESIZE;i++)
39        buff[i]= buff[0]+(maxLen+1)*i;
40      in=0;
41      out=0;
42      N=0;
43      total=0;
44   };
45   //———————————————————————————
46   Monitor::~Monitor()
47   {
48      delete [] buff[0];
49   }
50   //———————————————————————————
51   void Monitor::finishUp()
52   {
53      QMutexLocker ml(&l);
54      found=true;
55      empty.wakeAll();
56      full.wakeOne();
57   }
58   //———————————————————————————
59   bool Monitor::putCandidate(char *s)
60   {
61      QMutexLocker ml(&l);
62      while(N==PASSSTORESIZE && !found)
63        full.wait(&l);
64      strcpy(buff[in], s);
65      in=(in+1)%PASSSTORESIZE;
66      N++;
67      total++;
68      if(total%1000000==0)
69        cout << "Checking #" << total << " : " <<  s << endl;
70      empty.wakeOne();
71      return !found;
72   }
73   //———————————————————————————
74   bool Monitor::getCandidate(char *s)
75   {
76      QMutexLocker ml(&l);
77      while(N==0  && !found)
78        empty.wait(&l);
79      if(found) return false;
80
81      strcpy(s, buff[out]);
82      out=(out+1)%PASSSTORESIZE;
83      N--;
84      full.wakeOne();
85      return !found;
86   }
87   //*************************************************
88   class PassGenerator : public QThread
89   {
90   private:
91      char *candidate;
92      void recGen(int pos, int remain);
93      Monitor *m;
94      int maxLen;
```

```cpp
95      bool stopFlag;
96
97  public:
98      PassGenerator(Monitor *x, int l);
99      void run();
100 };
101 //————————————————————————————
102 PassGenerator::PassGenerator(Monitor *x, int l) : m(x), maxLen(l)↩
        , stopFlag(false)
103 {
104     candidate = new char[maxLen+1];
105 }
106 //————————————————————————————
107 void PassGenerator::recGen(int pos, int remain)
108  {
109     if(remain==0)
110     {
111       candidate[pos]=0;
112       stopFlag = ! m->putCandidate(candidate);
113     }
114     else
115     {
116       for(unsigned int i=0;i<strlen(alphabet) && !stopFlag;i++)
117       {
118         candidate[pos] = alphabet[i];
119         recGen(pos+1, remain-1);
120       }
121     }
122  }
123 //————————————————————————————
124 void PassGenerator::run()
125 {
126     int i=6;
127     while(stopFlag==false && i<=maxLen)
128       recGen(0, i++);
129 }
130 //**************************************************
131 class PassChecker : public QThread
132 {
133 private:
134     Monitor *m;
135     int maxLen;
136 public:
137     PassChecker(Monitor *x, int L) : m(x), maxLen(L){};
138     void run();
139 };
140 //————————————————————————————
141 void PassChecker::run()
142 {
143     char tmp[maxLen+1];
144     while( m->getCandidate(tmp))
145      {
146       if(target == md5(tmp))
147       {
148         m->finishUp();
149         cout << "Password : " << tmp << endl;
150         break;
151       }
152     }
153 }
154 //**************************************************
155  int main(int argc, char *argv[])
156  {
157     int Nchk = atoi(argv[1]);
158     int L =    atoi(argv[2]);
159     Monitor m(Nchk, L);
160     PassChecker *chk[Nchk];
161     PassGenerator gen(&m, L);
162
163     gen.start();
164     for(int i=0;i<Nchk;i++)
165     {
166       chk[i] = new PassChecker(&m,L);
167       chk[i]->start();
```

```
168        }
169
170        gen.wait();
171        for(int i=0;i<Nchk;i++)
172          chk[i]->wait();
173
174        return 0;
175    }
```

20. Write a multi-threaded program for finding the prime numbers in a user-supplied range of numbers. Compare the following design approaches:

    (a) Split the range in equal pieces and assign each one to a thread.

    (b) Have a shared `QAtomicInt` variable that holds the next number to be checked. Threads should read and increment this number before testing it.

    (c) Have a shared "monitor" object that returns upon request, a range of numbers to be tested. This can be considered a generalization of the previous design.

    Which of the designs is more efficient? Explain your findings.

    **Answer**

    (a) In the following listing, the `Repository` class is used to provide thread-safe storage for the discovered prime numbers. The `PrimeChecker` class is used to create threads that will check all the odd numbers in the range that is specified in their constructor. The main thread actually runs one of the `PrimeChecker` objects to avoid having to idle while they are doing useful work.

```cpp
1    #include <QThread>
2    #include <QMutex>
3    #include <QMutexLocker>
4    #include <math.h>
5    #include <stdlib.h>
6    #include <iostream>
7    #include <vector>
8
9    using namespace std;
10
11   //========================================
12   // used for storing the results in a thread-safe manner
13   class Repository
14   {
15   private:
16       vector<int> result;
17       QMutex l;
18   public:
19       void store(int i);
20       vector<int> *getResult();
21   };
22   //----------------------------------------
23   void Repository::store(int i)
24   {
25       QMutexLocker ml(&l);
26       result.push_back(i);
27   }
28   //----------------------------------------
29   vector<int> *Repository::getResult()
30   {
31     return &result;
32   }
33   //========================================
34   class PrimeChecker : public QThread
```

```cpp
35  {
36  private:
37      int first, last;
38      Repository *repo;
39
40      public:
41      PrimeChecker(int f, int l, Repository *r) : first(f), ←
            last(l), repo(r){}
42      void run();
43  };
44  //————————————————————————————————————
45  void PrimeChecker::run()
46  {
47      for(int i=first; i<=last; i+=2)
48      {
49          bool isPrime=true;
50          int j=3;
51          int limit = sqrt(i);
52          while(j<= limit && isPrime)
53          {
54              if(i % j == 0)
55                  isPrime=false;
56              j+=2;
57          }
58          if(isPrime)
59              repo->store(i);
60      }
61  }
62
63  //————————————————————————————————————
64  int main(int argc, char *argv[])
65  {
66      int numThreads;
67      int a, b;
68      a=atoi(argv[1]);
69      b=atoi(argv[2]);
70      numThreads=atoi(argv[3]);
71
72      // make sure a and b are odd
73      a = (a % 2 == 0) ? a+1 : a;
74      b = (b % 2 == 0) ? b-1 : b;
75
76      int rangePerThread = (b - a) * 1.0 / numThreads;
77      // make sure rangePerThread is even
78      rangePerThread = (rangePerThread % 2 == 0) ? ←
            rangePerThread : rangePerThread + 1;
79
80      Repository rep;
81      PrimeChecker *pc[numThreads];
82      int first, last=a-2;
83      for(int i=0;i<numThreads-1;i++)
84      {
85          first=last+2;
86          last = first + rangePerThread;
87          pc[i]= new PrimeChecker(first, last, & rep);
88          pc[i]->start();
89      }
90      pc[numThreads-1] = new PrimeChecker(last+2, b, & rep);
91      pc[numThreads-1]->run();
92
93      for(int i=0;i<numThreads-1;i++)
94          pc[i]->wait();
95
96      cout << "Total primes found " << rep.getResult()->size()←
            << endl;
97      return 0;
98  }
```

To run the above program in order to find all the prime numbers in the range $[10^3, 10^8]$ using eight threads, one has to use the following command-line (the same conventions are used by the next versions as well):

```
$ primeChecker_V3 1000 100000000 8
```

(b) In this version, a shared `QAtomicInt` object is used to fetch odd number to check from primality. The `Repository` class serves the same purpose as above.

```cpp
1  #include <QThread>
2  #include <QMutex>
3  #include <QMutexLocker>
4  #include <QAtomicInt>
5  #include <math.h>
6  #include <stdlib.h>
7  #include <iostream>
8  #include <vector>
9
10 using namespace std;
11
12 //===============================================
13 // used for storing the results in a thread-safe manner
14 class Repository
15 {
16 private:
17     vector<int> result;
18     QMutex l;
19 public:
20     void store(int i);
21     vector<int> *getResult();
22 };
23 //-----------------------------------------------
24 void Repository::store(int i)
25 {
26     QMutexLocker ml(&l);
27     result.push_back(i);
28 }
29 //-----------------------------------------------
30 vector<int> *Repository::getResult()
31 {
32   return &result;
33 }
34 //===============================================
35 class PrimeChecker : public QThread
36 {
37 private:
38     int last;
39     Repository *repo;
40     QAtomicInt *next;
41
42     public:
43     PrimeChecker(QAtomicInt *n, int l, Repository *r) : last↩
            (l), repo(r), next(n){}
44     void run();
45 };
46 //-----------------------------------------------
47 void PrimeChecker::run()
48 {
49     int candidate;
50     while((candidate = next->fetchAndAddOrdered(2)) <= last)
51     {
52         bool isPrime=true;
53         int j=3;
54         int limit = sqrt(candidate);
55         while(j<= limit && isPrime)
56         {
57             if(candidate % j == 0)
58                 isPrime=false;
59             j+=2;
60         }
61         if(isPrime)
62             repo->store(candidate);
63     }
64 }
65
66 //-----------------------------------------------
```

```
67   int main(int argc, char *argv[])
68   {
69       int numThreads;
70       int a, b;
71       a=atoi(argv[1]);
72       b=atoi(argv[2]);
73       numThreads=atoi(argv[3]);
74
75       // make sure a and b are odd
76       a = (a % 2 == 0) ? a+1 : a;
77       b = (b % 2 == 0) ? b-1 : b;
78
79       QAtomicInt first(a);
80
81       Repository rep;
82       PrimeChecker *pc[numThreads];
83       for(int i=0;i<numThreads-1;i++)
84       {
85           pc[i]= new PrimeChecker(&first, b, & rep);
86           pc[i]->start();
87       }
88       pc[numThreads-1] = new PrimeChecker(&first, b, & rep);
89       pc[numThreads-1]->run();
90
91       for(int i=0;i<numThreads-1;i++)
92           pc[i]->wait();
93
94       cout << "Total primes found " << rep.getResult()->size()↩
                << endl;
95       return 0;
96   }
```

The overall load balancing achieved by this version is substantially
better: As the numbers get larger, the checking procedure becomes
more time consuming. This makes the thread assigned the last por-
tion of the range in the previous version, dominate the execution
time.

(c) This version shares the load balancing efficiency of the second ver-
    sion, while reducing the cost of performing atomic operations. The
    **Monitor** effectively returns sets of 50 numbers to check, instead of
    the single one of the second version.

```
1    #include <QThread>
2    #include <QMutex>
3    #include <QMutexLocker>
4    #include <QAtomicInt>
5    #include <math.h>
6    #include <stdlib.h>
7    #include <iostream>
8    #include <vector>
9
10   using namespace std;
11
12   const int PARTLEN=100;
13   //================================================
14   // used for storing the results in a thread-safe manner
15   class Repository
16   {
17   private:
18       vector<int> result;
19       QMutex l;
20   public:
21       void store(int i);
22       vector<int> *getResult();
23   };
24   //------------------------------------------------
25   void Repository::store(int i)
26   {
27       QMutexLocker ml(&l);
28       result.push_back(i);
```

51

```
29  }
30  //———————————————————————————————————
31  vector<int> *Repository::getResult()
32  {
33      return &result;
34  }
35  //===================================================
36  class Monitor
37  {
38  private:
39      int first, last;
40      int currentAssign;
41      int assignRange;
42      QMutex l;
43  public:
44      Monitor(int a, int b, int r) : first(a), last(b), ↩
            assignRange(r){currentAssign = a;}
45      void getNextRange(int *a, int *b);
46  };
47  //———————————————————————————————————
48  void Monitor::getNextRange(int *a, int *b)
49  {
50      QMutexLocker ml(&l);
51      if(currentAssign >= last) // terminating check
52      {
53          *a = -1;
54          return;
55      }
56      *a = currentAssign;
57      *b = currentAssign + assignRange;
58      if(*b > last) *b = last;
59      currentAssign = *b + 2;
60  }
61
62  //===================================================
63  class PrimeChecker :public QThread
64  {
65  private:
66      Repository *repo;
67      Monitor *mon;
68
69  public:
70      PrimeChecker(Monitor *m, Repository *r) : repo(r), mon(m↩
            ){}
71      void run();
72  };
73  //———————————————————————————————————
74  void PrimeChecker::run()
75  {
76      int first, last;
77      while(1)
78      {
79          mon->getNextRange(&first, &last);
80          if(first == -1) break;  // if the Monitor returns -1↩
                for first, exit
81          for(int i=first; i<=last; i+=2)
82          {
83              bool isPrime=true;
84              int j=3;
85              int limit = sqrt(i);
86              while(j<= limit && isPrime)
87              {
88                  if(i % j == 0)
89                      isPrime=false;
90                  j+=2;
91              }
92              if(isPrime)
93                  repo->store(i);
94          }
95      }
96  }
97
98  //———————————————————————————————————
99  int main(int argc, char *argv[])
```

Table 3.1: Average execution times (in sec) over 10 runs, of the three prime checker variations. The best times in each case, are highlighted.

|  | Threads | | | |
| --- | --- | --- | --- | --- |
| **Version** | **1** | **2** | **4** | **8** |
| Even assignment to threads | **1.826** | 1.14 | 0.61 | 0.495 |
| `QAtomicInt` based | 2.048 | 1.04 | 0.54 | 0.426 |
| Monitor based | 1.837 | **0.92** | **0.46** | **0.418** |

```
100   {
101       int numThreads;
102       int a, b;
103       a=atoi(argv[1]);
104       b=atoi(argv[2]);
105       numThreads=atoi(argv[3]);
106
107       // make sure a and b are odd
108       a = (a % 2 == 0) ? a+1 : a;
109       b = (b % 2 == 0) ? b-1 : b;
110
111       Monitor mon(a, b, PARTLEN);   // range of PARTLEN numbers↩
                 , or PARTLEN/2 odd ones
112       Repository rep;
113       PrimeChecker *pc[numThreads];
114       for(int i=0;i<numThreads-1;i++)
115       {
116           pc[i]= new PrimeChecker(&mon, & rep);
117           pc[i]->start();
118       }
119       pc[numThreads-1] = new PrimeChecker(&mon, & rep);
120       pc[numThreads-1]->run();
121
122       for(int i=0;i<numThreads-1;i++)
123           pc[i]->wait();
124
125       cout << "Total primes found " << rep.getResult()->size()↩
                 << endl;
126       return 0;
127   }
```

The exact performance figures depend on the execution platform. On a Intel i7 3770K CPU, running under Linux and compiled with the GCC 4.8.2 compiler and -O2 optimization, we were able to measure the execution times shown in Table 3.1, for checking the range $[10^3, 10^7]$. With the exception of the single thread case, where the no-coordination-overheads version #1 wins, in all other cases the last, monitor-based version is superior.

21. Use the `QtConcurrent` functionality to implement a prime number checker. Compare it in terms of speed, efficiency and programming effort to your `QThread`-based attempt of the previous exercise.

   **Answer**

   The solution involves a `NumberRange` class for representing parts of the number range to scan for prime numbers. When an instance of this class is processed by the static `NumberRange::process` method, the prime numbers in that range are deposited in the designated repository as referenced by the `repo` pointer.

The main method uses the user-supplied desired length of the `NumberRange` instances (variable `PARTLEN`), to populate a vector of them (lines 80-90). This vector is subsequently processed by using the "mapping" functionality of the QtConcurrent namespace (line 92). The `blockingMap` function applies the `NumberRange::process` method on all elements of the input vector, effectively calculating the desired results.

```cpp
#include <QtConcurrent/QtConcurrentMap>
#include <QMutex>
#include <QMutexLocker>
#include <math.h>
#include <stdlib.h>
#include <iostream>
#include <vector>

using namespace std;

//================================================
// used for storing the results in a thread-safe manner
class Repository
{
private:
    vector<int> result;
    QMutex l;
public:
    void store(int i);
    vector<int> *getResult();
};
//------------------------------------------------
void Repository::store(int i)
{
    QMutexLocker ml(&l);
    result.push_back(i);
}
//------------------------------------------------
vector<int> *Repository::getResult()
{
    return &result;
}
//================================================
class NumberRange
{
private:
    int first, last;
    Repository *repo;
public:
    NumberRange(){}
    NumberRange(int a, int b, Repository *r) : first(a), last(b),←
            repo(r){}
    void setFirst(int f) {first=f;}
    void setLast(int l) {last=l;}
    void setRepo(Repository *r) {repo=r;}
    static void process(NumberRange &);
};
//------------------------------------------------
void NumberRange::process(NumberRange &n)
{
    for(int i=n.first; i<=n.last; i+=2)
    {
        bool isPrime=true;
        int j=3;
        int limit = sqrt(i);
        while(j<= limit && isPrime)
        {
            if(i % j == 0)
                isPrime=false;
            j+=2;
        }
        if(isPrime)
            n.repo->store(i);
    }
}
```

```
65   //======================================================
66   int main( int argc , char *argv [ ] )
67   {
68       int a , b;
69       a=atoi ( argv [ 1 ] ) ;
70       b=atoi ( argv [ 2 ] ) ;
71       int PARTLEN=atoi ( argv [ 3 ] ) ;
72
73       // make sure a and b are odd
74       a = (a % 2 == 0) ? a+1 : a ;
75       b = (b % 2 == 0) ? b-1 : b ;
76       PARTLEN = (PARTLEN % 2 == 0) ? PARTLEN : PARTLEN+1; // ↵
                 PARTLEN should be even
77
78       Repository rep ;
79       vector<NumberRange> data ;
80       int numParts = ( int ) ceil ( ( b-a ) *1.0/ PARTLEN ) ;
81       data . resize ( numParts ) ;
82       int f , l=a-2;
83       for ( int i =0; i<numParts ; i++)
84       {
85           f=l+2;
86           l=f+PARTLEN ;
87           if ( l>b) l=b;
88           data . at ( i ) . setFirst ( f ) ;
89           data . at ( i ) . setLast ( l ) ;
90           data . at ( i ) . setRepo(&rep ) ;
91       }
92
93       QtConcurrent :: blockingMap ( data , NumberRange :: process ) ;
94
95       cout << "Total primes found " << rep . getResult ( )->size ( ) << ↵
                 endl ;
96       return 0;
97   }
```

To run the above program in order to find all the prime numbers in the range $[10^3, 10^8]$, one has to use the following command-line:

```
$ primeChecker_V3 1000 100000000 100000
```

The last parameter concerns the value of `PARTLEN`.

In terms of performance, an identical set of tests as in the previous exercise, results in an average execution time of 0.42 sec, for `PARTLEN = 10000`, which is on par with the monitor-based solution. In terms of development effort, the `QtConcurrent` based solution is marginally shorter.

22. Create a big array of randomly generated 2D coordinates $(x, y)$. Each of the coordinates should be a number in the range [-1000, 1000]. Use appropriate `QtConcurrent` functions to find the points that are in a ring of distances between 100 and 200 from the point of origin. Compare the performance of your solution against a sequential implementation.

**Answer**

In the following program a `Point` structure is used to hold randomly generated data. A `vector<Point>` container holds the input data, instead of a `vector<Point*>`, in order to simplify and speedup memory management. Upon resizing the vector container in line 57, all the memory required for holding the `Point` instances is reserved.

Also, in order to avoid the use of the costly `sqrt` function, the filtering is performed on the squares of the distances.

```cpp
1   #include <QtConcurrent/QtConcurrentMap>
2   #include <QtConcurrent/QtConcurrentFilter>
3   #include <QThreadPool>
4   #include <QTime>
5   #include <boost/bind.hpp>
6   #include <vector>
7   #include <iostream>
8   #include <stdlib.h>
9
10  using namespace std;
11
12  //———————————————————————————————————
13  struct Point
14  {
15     int x, y;
16     int d2;                          // distance square
17        Point () {};
18  };
19
20  //———————————————————————————————————
21  void d2Calc (Point & p)
22  {
23     p.d2 = p.x * p.x + p.y * p.y;
24  }
25
26  //———————————————————————————————————
27  bool inRing (Point p, int minDist, int maxDist)
28  {
29     int d = p.d2;
30     if (d >= minDist && d <= maxDist)
31        return true;
32     return false;
33  }
34
35  //———————————————————————————————————
36  int main (int argc, char *argv[])
37  {
38     int numThr;
39     int numPoints;
40     vector < Point > data;
41     double tio, tfull;
42     QTime t;
43
44     if (argc < 3)
45        {
46           cerr << argv[0] << " numThreads numPoints\n";
47           exit (1);
48        }
49     // for timing
50     t.start ();
51
52     numThr = atoi (argv[1]);
53     numPoints = atoi (argv[2]);
54
55     // initialize data array
56     srand (time (0));
57     data.resize (numPoints);       // one memory allocation only
58     for (int i = 0; i < numPoints; i++)
59        {
60           data[i].x = rand () % 2000 - 1000;
61           data[i].y = rand () % 2000 - 1000;
62        }
63
64     tio = t.elapsed () / 1000.0;
65     // control the number of threads
66     QThreadPool::globalInstance ()->setMaxThreadCount (numThr);
67
68     // calculate distances first
69     QtConcurrent::blockingMap (data, d2Calc);
70
71     // filter pointes based on distance next
72     QtConcurrent::blockingFilter (data, boost::bind (inRing, _1, ↩
            100 * 100, 200 * 200));
73
```

```
74    tfull = t.elapsed () / 1000.0;
75    cout << "Total: " << tfull << " Comp.time: " << tfull - tio << ↵
         endl;
76    cout << data.size () << endl;
77    return 0;
78  }
```

23. Use the `QtConcurrent` functionality to implement a parallel bucketsort. Does the number of buckets play a significant role in your implementation's performance?.

    **Answer**

    In the following Listing, the `bucketsort` template function implements a variation of the algorithm, whereas (i) the input is scanned and distributed to buckets (lines 43-48), (ii) the buckets are sorted separately via STL's `sort` function (lines 51-54), and (iii) the resulting partial sorted arrays are copied back to the original data repository (lines 57-62).

    The number of buckets matches the number of threads used by `QThreadPool`, and for this reason the number of buckets is invariably connected with the level of concurrency.

    To avoid expensive synchronization logic, such as using mutices to control access to the buckets (STL's `vector` is not thread safe), each thread is responsible for a dedicated bucket. The bucket array address and bucket index for which a thread is responsible for, are passed as parameters to the `filterOp` template function that is invoked by `QtConcurrent::run` in line 45. The minimum value of the input range and the value range per bucket are also passed to enable `filterOp` to determine the bucket that an item belongs to.

    Finally, the sorting is performed via the unary front-end `vecSort` function that is "mapped" on all the buckets. To this end, a `vector` of `vector<T>` pointers is populated in lines 52-53 prior to the application of `QtConcurrent::blockingMap`.

```
1   #include <QtConcurrent/QtConcurrentMap>
2   #include <QtConcurrent/QtConcurrentRun>
3   #include <QFuture>
4   #include <QThreadPool>
5   #include <QTime>
6   #include <boost/bind.hpp>
7   #include <vector>
8   #include <algorithm>
9   #include <iostream>
10  #include <stdlib.h>
11
12  using namespace std;
13
14  const int MINV = -1000;
15  const int MAXV = 1000;
16
17  //————————————————————————————————
18  template <typename T> inline void filterOp (vector<T> *v, int ↵
        targetBucketIdx, vector<T> *b, T m, T bWidth)
19  {
20      for (int i=0;i< v->size (); i++)
21      {
22          T tmp = v->at (i);
23          int bucketIdx = (tmp - m) / bWidth;
24          if (bucketIdx == targetBucketIdx)
25              b->push_back (tmp);
26      }
```

```cpp
27  }
28
29  //————————————————————————————
30  template <typename T> inline void vecSort(vector<T> *v)
31  {
32      sort(v->begin(), v->end());
33  }
34
35  //————————————————————————————
36  template <typename T> void bucketSort(vector<T> &data, const T &↩
        minV, const T &maxV)
37  {
38      int numBuckets = QThreadPool::globalInstance()->↩
            maxThreadCount(); //match buckets with the number of ↩
            threads
39      vector<T> bucket[numBuckets];
40      T bwidth = (maxV - minV)/numBuckets;   // range per bucket
41
42      // launch as many requests as the number of buckets/threads
43      QFuture<void> f[numBuckets];
44      for(int i=0;i<numBuckets;i++)
45          f[i] = QtConcurrent::run(boost::bind(filterOp<T>, &data, ↩
              i, bucket + i, minV, bwidth));
46
47      for(int i=0;i<numBuckets;i++)
48          f[i].waitForFinished();
49
50      // sort each of the buckets separately
51      vector< vector<T> *> vv;
52      for(int i=0;i<numBuckets;i++)
53          vv.push_back(&(bucket[i]));
54      QtConcurrent::blockingMap(vv, vecSort<T>);
55
56      // copy back to original vector
57      int loc=0;
58      for(int i=0;i<numBuckets;i++)
59      {
60          for(int j=0;j<bucket[i].size();j++)
61              data[loc++] = bucket[i].at(j);
62      }
63  }
64
65  //————————————————————————————
66  int main (int argc, char *argv[])
67  {
68      int numThr;
69      int numItems;
70      vector < int > data;
71      double tio, tfull;
72      QTime t;
73
74      if (argc < 3)
75      {
76          cerr << argv[0] << " numThreads numItems\n";
77          exit (1);
78      }
79      // for timing
80      t.start ();
81
82      numThr = atoi (argv[1]);
83      numItems = atoi (argv[2]);
84
85      // initialize data array
86      srand (time (0));
87      data.resize (numItems);        // one memory allocation only
88      for (int i = 0; i < numItems; i++)
89      {
90          data[i] = (rand () % (MAXV-MINV)) + MINV;
91      }
92
93      tio = t.elapsed () / 1000.0;
94      // control the number of threads
95      QThreadPool::globalInstance ()->setMaxThreadCount (numThr);
96
```

```
97        bucketSort<int>(data, MINV, MAXV);
98
99        tfull = t.elapsed () / 1000.0;
100       cout << "Total: " << tfull << " Comp.time: " << tfull - tio ↵
             << endl;
101       return  0;
102   }
```

# Chapter 4

# Shared-memory programming : OpenMP

## Exercises

1. Modify the program in Listing 4.1, so that `printf` is used instead of `cout`. Do you see any difference in the output compared to the one reported in Section 4.2? Can you explain it?

   **Answer**

```
1   #include <stdlib.h>
2   #include <stdio.h>
3   #include <omp.h>
4
5   using namespace std;
6
7   int main (int argc, char **argv)
8   {
9       int numThr = atoi (argv[1]);
10  #pragma omp parallel num_threads(numThr)
11      printf("Hello from thread %i\n", omp_get_thread_num ());
12
13      return 0;
14  }
```

   The output does not appear mangled, as the `printf` strings are build and output as a unit, while the `cout` object goes through several method invocations.

2. In the matrix multiplication example of Section 4.4.2 we could get three perfectly nested loops, if we initialize the `C` matrix outside:

```
    double A[K][L];
    double B[L][M];
    double C[K][M];
. . .
#pragma omp parallel for collapse(3)
    for (int i = 0; i < K; i++)
      for (int j = 0; j < M; j++)
        for (int k = 0; k < L; k++)
          C[i][j] += A[i][k] * B[k][j];
```

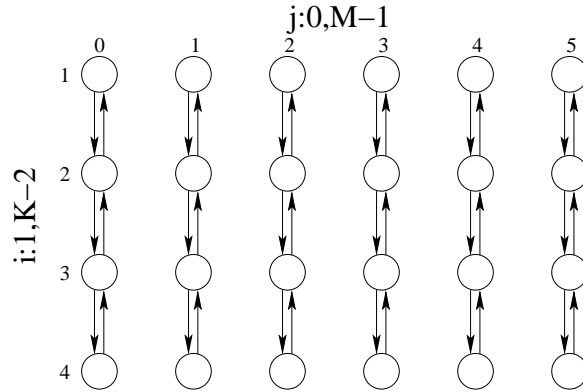   Is the above code correct? If not, what kind of modification is required for fixing it?

Figure 4.1: Dependency graph for Exercise 3.

**Answer**

The code is not correct, because for the inner loop `C[i][j]` is effectively a reduction variable. To eliminate the race condition, a critical section is needed:

```
#pragma omp parallel for collapse(3)
   for (int i = 0; i < K; i++)
     for (int j = 0; j < M; j++)
       for (int k = 0; k < L; k++)
#pragma omp critical
         C[i][j] += A[i][k] * B[k][j];
```

In terms of performance, this will effectively make all threads execute sequentially, because unnamed `critical` constructs use the same mutex.

A better alternative is to make sure that each execution of the inner loop is assigned to a single thread, by using a `schedule` clause, which is equivalent to converting the `collapse(3)` clause to `collapse(2)`:

```
#pragma omp parallel for collapse(3) schedule(static, L)
   for (int i = 0; i < K; i++)
     for (int j = 0; j < M; j++)
       for (int k = 0; k < L; k++)
         C[i][j] += A[i][k] * B[k][j];
```

3. Draw the iteration space dependency graph for the following program:

```
   for (int i = 1; i < K-1; i++)
     for (int j = 0; j < M; j++)
       a[i][j] = a[i-1][j] + a[i+1][j];
```

What kind of dependencies exist? How can you eliminate them?

**Answer**

The ISDG is shown in Figure 4.1.

We have a flow dependence from `a[i-1][j]` and an anti-dependence for a[i+1][j]. The anti-dependence can be eliminated by creating a copy `a2` of the original matrix `a`. The flow dependence can be honored (not removed) by parallelizing only the inner loop, resulting in the following code:

```
    for (int i = 0; i < K-1; i++)
      for (int j = 0; j < M; j++)
          a2[i][j] = a[i+1][j];

    for (int i = 1; i < K-1; i++)
#pragma omp parallel for
      for (int j = 0; j < M; j++)
          a[i][j] = a[i-1][j] + a2[i][j];
```

4. Create a C++ program for visualizing the thread iteration assignment performed by a `parallel for` directive, for different `schedule` schemes. Your program should have a per-thread vector that accumulates the loop control variable values assigned to each thread. Use this program, and appropriate `schedule` settings, to experiment with multiple schemes, without recompiling your program. This is a sample of the output you should be able to get for a loop of 100 iterations:

```
$ export OMP_SCHEDULE="static,5"
$ ./solution
Thread 0 : 0 1 2 3 4 40 41 42 43 44 80 81 82 83 84
Thread 1 : 5 6 7 8 9 45 46 47 48 49 85 86 87 88 89
Thread 2 : 10 11 12 13 14 50 51 52 53 54 90 91 92 93 94
Thread 3 : 15 16 17 18 19 55 56 57 58 59 95 96 97 98 99
Thread 4 : 20 21 22 23 24 60 61 62 63 64
Thread 5 : 25 26 27 28 29 65 66 67 68 69
Thread 6 : 30 31 32 33 34 70 71 72 73 74
Thread 7 : 35 36 37 38 39 75 76 77 78 79
```

**Answer**

The solution is based on having a vector per thread (line 11), to store the values of the loop control variable (line 18):

```cpp
1  #include <iostream>
2  #include <stdlib.h>
3  #include <vector>
4  #include <omp.h>
5
6  using namespace std;
7
8  int main ()
9  {
10    int N = omp_get_max_threads ();
11    vector < int >*store = new vector < int >[N];
12
13
14 #pragma omp parallel for schedule( runtime )
15    for (int k = 0; k < 100; k++)
16      {
17        int idx = omp_get_thread_num ();
18        store[idx].push_back (k);
19      }
20
21    //output
22    for (int i = 0; i < N; i++)
23      {
24        cout << "Thread " << i << " : ";
25        for (int j = 0; j < store[i].size (); j++)
26          cout << store[i][j] << " ";
27        cout << endl;
28      }
29
30    return 0;
31 }
```

5. In Listing 4.19 we examined the issue of processing a linked-list concurrently. Is there a way to improve this program for a doubly-linked list? Write the corresponding program.

**Answer**

In the following code, the list is traversed in two opposite directions, allowing the generation of two tasks per iteration of the `while` loop of lines 61-75.

```cpp
1   #include <iostream>
2   #include <math.h>
3   #include <stdlib.h>
4   #include <omp.h>
5   #include <unistd.h>
6
7   using namespace std;
8
9   // template structure for a list's node
10  template < class T > struct Node
11  {
12    T info;
13    Node *next;
14    Node *prev;
15  };
16  //———————————————————————————————————
17  // Appends a value at the end of a list pointed by the head *h ←↩
          and tail *t
18  template < class T > void append (int v, Node < T > **h, Node < T←↩
          > **t)
19  {
20    Node < T > *tmp = new Node < T > ();
21    tmp->info = v;
22    tmp->next = NULL;
23    tmp->prev = *t;
24
25    if (*t == NULL)
26       {
27         *h = tmp;
28         *t = tmp;
29       }
30    else
31       {
32         (*t)->next = tmp;
33         *t = tmp;
34       }
35  }
36
37  //———————————————————————————————————
38  // function stub for processing a node's data
39  template < class T > void process (Node < T > *p)
40  {
41  #pragma omp critical
42    cout << p->info << " by thread " << omp_get_thread_num () << ←↩
          endl;
43  }
44
45  //———————————————————————————————————
46  int main (int argc, char *argv[])
47  {
48    // build a sample list
49    int N = atoi (argv[1]);
50    Node < int >*head = NULL;
51    Node < int >*tail = NULL;
52    for (int i = 0; i < N; i++)
53      append (i, &head, &tail);
54
55  #pragma omp parallel
56    {
57  #pragma omp single
58       {
59         Node < int >*tmp1 = head;
60         Node < int >*tmp2 = tail;
61         while (tmp1 != NULL && tmp1 != tmp2)
62            {
63  #pragma omp task
64              process (tmp1);
65  #pragma omp task
```

```
66              process (tmp2);
67
68              tmp1 = tmp1->next;
69              if (tmp1 == tmp2)       // even-sized list termination ←
                    check
70                {
71                  tmp1 = NULL;
72                  break;
73                }
74              tmp2 = tmp2->prev;
75          }
76        // process the last node, in the middle of an odd-sized ←
                list
77        if (tmp1 != NULL)
78          process (tmp1);
79      }
80    }
81
82    return 0;
83 }
```

6. Write a program for traversing and processing the elements of a binary tree in parallel, using a pre-order, in-order or post-order traversal. Use the `task` construct to that effect.

   **Answer**

   In-order traversal cannot be parallelized as in enforces a total order between the tree nodes.

```
1  //————————————————————
2  template < typename T > struct TreeNode
3  {
4    T value;
5      TreeNode < T > *left;
6      TreeNode < T > *right;
7  };
8  //————————————————————
9  // empty placeholder for a processing function
10 template < typename T > void processStub (T n)
11 {
12 }
13
14 //————————————————————
15
16 template < typename T > void preOrder (TreeNode < T > *n, void *←
       process (T n))
17 {
18   if (n == NULL)
19     return;
20
21   process (n->value);
22 #pragma omp parallel
23   {
24 #pragma omp single
25     {
26 #pragma omp task
27       preOrder (n->left, process);
28
29       // keep using this thread for the other child
30       preOrder (n->right, process);
31 #pragma omp taskwait
32     }
33   }
34 }
35
36 //————————————————————
37
38 template < typename T > void postOrder (TreeNode < T > *n, void *←
       process (T n))
39 {
40   if (n == NULL)
```

```
41        return ;
42
43  #pragma omp parallel
44      {
45  #pragma omp single
46        {
47  #pragma omp task
48          postOrder (n->left , process );
49
50          // keep using this thread for the other child
51          postOrder (n->right , process );
52  #pragma omp taskwait
53        }
54      }
55      process (n->value );
56  }
57
58  //——————————————————————
59
60  // inOrder cannot be parallelized as in enforces a total order ↩
          between the tree nodes .
61  // postOrder and preOrder only have a partial ordering
62  template < typename T > void inOrder (TreeNode < T > *n, void *↩
          process (T n))
63  {
64      if (n == NULL)
65        return ;
66
67      inOrder (n->left , process );
68      process (n->value );
69      inOrder (n->right , process );
70
71  }
```

7. Modify the program of the previous exercise, so that *undeferred tasks* are generated after the traversal has moved beyond the fifth level of the tree (assume that the root sits at level 0).

   **Answer**

   The required modifications relative to the answer of the previous question, are shown below. The solution is based on the introduction of a function parameter for counting the depth of the recursion/depth of the tree nodes examined. The default value of "0" means that the caller does not need to be aware of this parameter.

```
1  template < typename T > void preOrder (TreeNode < T > *n, void *↩
          process (T n), int level=0)
2  {
3      if (n == NULL)
4        return ;
5
6      process (n->value );
7  #pragma omp parallel
8        {
9  #pragma omp single
10        {
11  #pragma omp task  if (level >=5)
12          preOrder (n->left , process , level+1);
13
14          // keep using this thread for the other child
15          preOrder (n->right , process , level+1);
16  #pragma omp taskwait
17        }
18      }
19  }
20
21  //——————————————————————
22
23  template < typename T > void postOrder (TreeNode < T > *n, void *↩
          process (T n), int level=0)
```

```
24  {
25     if  (n == NULL)
26         return ;
27
28  #pragma omp  parallel
29     {
30  #pragma omp  single
31        {
32  #pragma omp  task  if (level >=5)
33           postOrder  (n−>left ,  process ,  level+1);
34
35           // keep  using  this  thread  for  the  other  child
36           postOrder  (n−>right ,  process ,  level+1);
37  #pragma omp  taskwait
38        }
39     }
40     process  (n−>value );
41  }
```

8. Modify the Fibonacci sequence-calculating program of Listing 4.20, so that it counts the number of child tasks generated.

   **Answer**

   Use of an `atomic` pragma is required for properly incrementing the shared counter. Only changes relative to Listing 4.20 are shown.

```
1   int  numTasks = 0;
2
3   int  fib  (int  i)
4   {
5      int  t1 ,  t2 ;
6      if  (i == 0  ||  i == 1)
7         return  1;
8      else
9         {
10  #pragma omp  task  shared (t1)  if (i >25)  mergeable
11           {
12  #pragma omp  atomic
13             numTasks++;
14             t1 = fib  (i − 1);
15
16           }
17  // keep  thread  for  calling  function  again
18           t2 = fib  (i − 2);
19  #pragma omp  taskwait
20           return  t1 + t2 ;
21         }
22  }
```

9. Use the `task` directive to create a solution to the single producer, multiple consumers problem, with a variable number of consumers, as specified in the command-line.

   **Answer**

   The following code solves the integration example that is presented in Section 4.5.1.1 and solved with a fixed number of consumers in Listing 4.18. To minimize clutter, we show only the `main` function that generates the consumer threads (via a `for` loop in lines 21-25) and executes the producer thread code (lines 28-52).

```
1   int  main  (int  argc ,  char  ** argv )
2   {
3      if  (argc == 1)
4         {
5            cerr << "Usage "  <<  argv [0]  << "  #threads #jobs\n";
6            exit  (1);
```

```
7            }
8      int N = atoi (argv[1]);
9      int J = atoi (argv[2]);
10
11     Slice *buffer = new Slice[BUFFSIZE];
12     int in = 0, out = 0;
13     QSemaphore avail, buffSlots (BUFFSIZE);
14     QMutex l, integLock;
15     double integral = 0;
16 #pragma omp parallel default(none) shared(buffer, in, out, avail,←↩
          buffSlots, l, integLock, integral, J, N, cout) num_threads(←↩
          N+1)
17     {
18 #pragma omp single
19        {
20 // consumers' part
21        for (int i = 0; i < N; i++)
22           {
23 #pragma omp task
24           integrCalc (buffer, buffSlots, avail, l, out, integLock←↩
                , integral);
25           }
26
27 // producer thread, responsible for handing out 'jobs'
28        double divLen = (UPPERLIMIT − LOWERLIMIT) / J;
29        double st, end = LOWERLIMIT;
30        for (int i = 0; i < J; i++)
31           {
32             st = end;
33             end += divLen;
34             if (i == J − 1)
35               end = UPPERLIMIT;
36
37             buffSlots.acquire ();
38             buffer[in].start = st;
39             buffer[in].end = end;
40             buffer[in].divisions = 1000;
41             in = (in + 1) % BUFFSIZE;
42             avail.release ();
43           }
44
45        // put termination sentinels in buffer
46        for (int i = 0; i < N; i++)
47           {
48             buffSlots.acquire ();
49             buffer[in].divisions = 0;
50             in = (in + 1) % BUFFSIZE;
51             avail.release ();
52           }
53
54 // wait for all threads to finish
55 #pragma omp taskwait
56        }
57     }
58
59     cout << "Result is : " << integral << endl;
60     delete [] buffer;
61
62     return 0;
63 }
```

10. Use the `task` directive and its `depend` clause, to model the dependency
    graph of Figure 4.12. You can use stub functions to represent each of the
    tasks in the figure.

    **Answer**

```
1 #include <iostream>
2 #include <stdlib.h>
3 #include <omp.h>
4 #include <unistd.h>
5
```
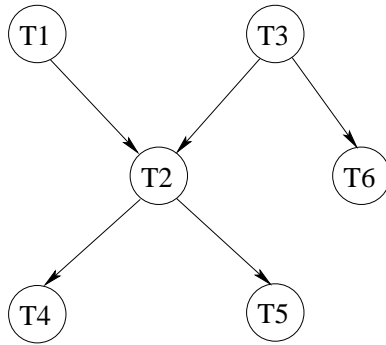
Figure 4.2: A dependency graph consisting of six tasks.

```cpp
6    using namespace std;
7
8    int main (int argc, char **argv)
9    {
10     int x, y, z;
11
12   #pragma omp parallel
13     {
14   #pragma omp single
15       {
16   #pragma omp task   shared(x) depend(out : x)
17         {   // T1
18           cout << "T1\n";
19         }
20
21   #pragma omp task shared(x) depend(out : y, z)
22         {   // T3
23           cout << "T3\n";
24         }
25
26   #pragma omp task   shared(x) depend(inout : x, y)
27         {   // T2
28           sleep (1);
29           cout << "T2\n";
30         }
31
32   #pragma omp task shared(x) depend(in : y)
33         {   // T4
34           cout << "T4\n";
35         }
36
37   #pragma omp task   shared(x) depend(in : y)
38         {   // T5
39       sleep(1);
40           cout << "T5\n";
41         }
42
43   #pragma omp task shared(x) depend(in : z)
44         {   // T6
45           cout << "T6\n";
46         }
47
48   #pragma omp taskwait
49       }
50     }
51
52   return 0;
53   }
```

11. Finding the odd integers in a array can be accomplished by the following

OpenMP code:

```
int data[N];
int oddCount=0;
#pragma omp parallel for
 for (int i = 0; i < N; i++)
       if( data[i] % 2 )
#pragma omp atomic
           oddCount ++;
```

Modify the above code, so that there is no need for a `critical` or `atomic` directive, by introducing a counter array, with one element per thread. Compare the performance of your version with and without cache false sharing. If we were to use a reduction variable, what false-sharing elimination technique, does this correspond to?

**Answer**

The following listing contains four different versions of the code snippet in question, encapsulated in appropriately named functions:

```
1   // Original code wih atomic
2   int count_atomic (int *data, int N)
3   {
4      int oddCount = 0;
5   #pragma omp parallel for
6      for (int i = 0; i < N; i++)
7        if (data[i] % 2)
8   #pragma omp atomic
9           oddCount++;
10     return oddCount;
11  }
12
13  //————————————————————————————————————
14  // Array of counters, but there is false sharing
15  int count_false_sharing (int *data, int N)
16  {
17     int numThr = omp_get_max_threads ();
18     int oddCount[numThr];
19     memset (oddCount, 0, numThr * sizeof (int));
20     int totalCount = 0;
21  #pragma omp parallel for
22     for (int i = 0; i < N; i++)
23        {
24          if (data[i] % 2)
25            oddCount[omp_get_thread_num ()]++;
26        }
27
28     for (int i = 0; i < numThr; i++)
29        totalCount += oddCount[i];
30     return totalCount;
31  }
32
33  //————————————————————————————————————
34  // Eliminates false sharing with padding
35  int count_without_false_sharing (int *data, int N)
36  {
37     int numThr = omp_get_max_threads ();
38     int oddCount[numThr * 8];
39     memset (oddCount, 0, numThr * 8 * sizeof (int));
40     int totalCount = 0;
41  #pragma omp parallel for
42     for (int i = 0; i < N; i++)
43        {
44          if (data[i] % 2)
45            oddCount[omp_get_thread_num () * 8]++;
46        }
47
48     for (int i = 0; i < numThr; i++)
49        totalCount += oddCount[i * 8];
50     return totalCount;
51  }
```

```
52
53   //————————————————————————————————
54   int count_with_reduction (int *data, int N)
55   {
56      int oddCount = 0;
57   #pragma omp parallel for reduction(+:oddCount)
58      for (int i = 0; i < N; i++)
59         {
60            if (data[i] % 2)
61               oddCount++;
62         }
63      return oddCount;
64   }
```

The use of a reduction variable is equivalent to using private variables for eliminating false sharing.

12. Quicksort rightfully holds the place of the top-performer amongst generic sorting algorithms. It's design can be considered a reflection of mergesort, although both algorithms employ a divide-and-conquer design. Write a sequential implementation of quicksort and proceed to parallelize it using OpenMP constructs. Measure the speedup that can be achieved, and compare it to the performance that can be obtained from the bottom-up mergesort of Section 4.8.1.

**Answer**

The solution involves the use of task constructs in lines and . The code incorporates several quicksort optimizations such as median-of-three pivot element selection (lines 54-62) and a switch to insertion sort once the size of the data fall below a fixed threshold (lines 92-96). To reduce the generation of tasks, a similar heuristic to the one described in Section 4.8.2 is used (_thresh_ variable in line 121).

In terms of performance, the program more closely matches the top-down mergesort, failing to fully exploit a multicore CPU because of its recursive nature. The bottom-up mergesort was about twice as fast as the quicksort on an i7 3770K CPU.

```
1    #include <stdlib.h>
2    #include <stdio.h>
3    #include <unistd.h>
4    #include <omp.h>
5    #include <QTime>
6    #include <iostream>
7
8    using namespace std;
9
10   const int INSERTION_THRESHOLD=50;
11   const int maxTasks=256;
12   int _thresh_; // used for the if clauses
13
14   //***********************************************
15   void insertionSort(int *data, int N)
16   {
17       for (int i = 1; i < N; i++)
18       {
19           int tmp = data[i];
20           int loc=i-1;
21           while(loc >=0 && data[loc] > tmp)
22           {
23               data[loc+1] = data[loc];
24               loc--;
25           }
26           data[loc+1] = tmp;
27       }
```

```c
28  }
29
30  //*************************************************
31  void numberGen (int N, int max, int *store)
32  {
33      int i;
34      srand (time (0));
35      for (i = 0; i < N; i++)
36          store[i] = rand () % max;
37  }
38
39  //*************************************************
40  void swap (int *data, int x, int y)
41  {
42      int temp = data[x];
43      data[x] = data[y];
44      data[y] = temp;
45  }
46
47  //*************************************************
48  int partition (int *data, int N)
49  {
50      int i = 0, j = N;
51
52      // median of 3
53      int a=data[0], b=data[N >> 1], c=data[N-1];
54      if((a > b && b >c) || (c > b && b >a))
55      {
56          swap(data, 0, N>>1);
57      }
58      else if((a > c && c >b) || (b > c && c>a))
59      {
60          swap(data, 0, N-1);
61      }
62
63      int pivot = data[0];
64
65      do
66      {
67          do
68          {
69              i++;
70          }
71          while (pivot > data[i] && i < N);
72
73          do
74          {
75              j--;
76          }
77          while (pivot < data[j] && j > 0);
78          swap (data, i, j);
79      }
80      while (i < j);
81
82      swap (data, i, j);
83
84      swap (data, 0, j);
85      return j;
86  }
87
88  //*************************************************
89  void QSort (int *data, int N)
90  {
91      if (N <= 1)
92          return;
93      else if(N<INSERTION_THRESHOLD)
94      {
95          insertionSort(data, N);
96          return;
97      }
98
99      int pivotPos = partition (data, N);
100 #pragma omp parallel
101     {
```

```cpp
#pragma omp single
        {
#pragma omp task  if(N > _thresh_ ) mergeable
            QSort (data, pivotPos);
#pragma omp task   if(N > _thresh_ ) mergeable
            QSort (&(data[pivotPos + 1]), N - pivotPos - 1);
#pragma omp taskwait
        }
    }
}

//————————————————————————————————————————————————————
int main (int argc, char *argv[])
{
    if (argc == 1)
    {
        fprintf (stderr, "%s N\n", argv[0]);
        exit (0);
    }
    int N = atoi (argv[1]);
    _thresh_ = 2.0 * N / (maxTasks + 1);

    int *data = new int[N];
    numberGen (N, 1000, data);
    QTime t;
    t.start();

    QSort (data, N);

    cout << "Sorted in " << t.elapsed()/1000.0 << " sec\n";

    delete [] data;
    return 0;
}
```

# Chapter 5

# Distributed memory programming

## Exercises

1. Write a MPMD version of the "Hello World" program of Figure 5.2, in effect eliminating the `if/else` structure around which the program is built. You may use the C or C++ bindings.

   **Answer**

   Two source code files are needed:

```c
// File : hello_node0.c
#include<mpi.h>
#include<string.h>
#include<stdio.h>
#define MESSTAG 0
#define MAXLEN 100
int main (int argc, char **argv)
{
  MPI_Init (&argc, &argv);

  int rank, num, i;
  MPI_Comm_rank (MPI_COMM_WORLD, &rank);
  MPI_Comm_size (MPI_COMM_WORLD, &num);
  char mess[] = "Hello World";
  for (i = 1; i < num; i++)
    MPI_Send (mess, strlen (mess) + 1, MPI_CHAR, i, MESSTAG, ↵
        MPI_COMM_WORLD);

  MPI_Finalize ();
  return 0;
}
```

```c
// File : hello_node_i.c
#include<mpi.h>
#include<string.h>
#include<stdio.h>
#define MESSTAG 0
#define MAXLEN 100
int main (int argc, char **argv)
{
  MPI_Init (&argc, &argv);

  int rank, num, i;
  MPI_Comm_rank (MPI_COMM_WORLD, &rank);

  char mess[MAXLEN];
```

```
  MPI_Status status;
  MPI_Recv (mess, MAXLEN, MPI_CHAR, 0, MESSTAG, MPI_COMM_WORLD, &↩
      status);
  printf ("%i received %s\n", rank, mess);

  MPI_Finalize ();
  return 0;
}
```

The program can run with the following appfile:

```
# File :  hello_app
−host localhost −np 1 hello_node0
−host localhost −np 3 hello_node_i
```

and via the command sequence:

```
$ mpicc hello_node0.c −o hello_node0
$ mpicc hello_node_i.c −o hello_node_i
$ mpirun −app hello_app
```

2. Write a SPMD version of the two programs shown in Listing 5.5.

   **Answer**

   The only requirement is the inclusion of an `if-else` control logic.

```
 1  #include <mpi.h>
 2  #include <stdio.h>
 3  #include <string.h>
 4  #include <stdlib.h>
 5
 6  #define MAXLEN 100
 7  char *greetings[] = { "Hello", "Hi", "Awaiting your command" };
 8
 9  char buff[MAXLEN];
10
11  int main (int argc, char **argv)
12  {
13    MPI_Status st;
14    int procNum;
15    int rank;
16
17    MPI_Init (&argc, &argv);
18    MPI_Comm_size (MPI_COMM_WORLD, &procNum);
19    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
20    if (rank == 0)
21      {
22        while (−−procNum)
23          {
24            MPI_Recv (buff, MAXLEN, MPI_CHAR, MPI_ANY_SOURCE, ↩
                  MPI_ANY_TAG, MPI_COMM_WORLD, &st);
25            int aux;
26            MPI_Get_count (&st, MPI_CHAR, &aux);
27            buff[aux] = 0;
28            printf ("%s\n", buff);
29          }
30      }
31    else
32      {
33        srand (time (0));
34        int grID = rand () % 3;
35        sprintf (buff, "Node %i says %s", rank, greetings[grID]);
36        MPI_Send (buff, strlen (buff), MPI_CHAR, 0, 0, ↩
                  MPI_COMM_WORLD);
37      }
38
39    MPI_Finalize ();
40  }
```

3. Modify the program shown in Listing 5.5 so that the master node prints out a list of the processes IDs for which the message has not been read yet. Your output should be similar to:

```
1  $ mpirun −np 1 master : −np 3 worker
2  Node 2 says Hi. Awaiting nodes : 1 3
3  Node 3 says Hi. Awaiting nodes : 1
4  Node 1 says Hi. Awaiting nodes :
```

**Answer**

An STL `set<int>` container can be used to hold the ID's of the processes still pending:

```cpp
1  #include <mpi.h>
2  #include <stdio.h>
3  #include <string.h>
4  #include <set>
5  #include <iostream>
6
7  using namespace std;
8
9  #define MAXLEN 100
10 char buff[MAXLEN];
11
12 int main (int argc, char **argv)
13 {
14   MPI_Status st;
15   int procNum;
16
17   MPI_Init (&argc, &argv);
18   MPI_Comm_size (MPI_COMM_WORLD, &procNum);
19   set<int> leftOver;
20   for(int i=1;i<procNum;i++)
21     leftOver.insert(i);
22
23   while (−−procNum)
24     {
25       MPI_Recv (buff, MAXLEN, MPI_CHAR, MPI_ANY_SOURCE, ↩
                MPI_ANY_TAG,
26         MPI_COMM_WORLD, &st);
27       int aux;
28       MPI_Get_count (&st, MPI_CHAR, &aux);
29       leftOver.erase(st.MPI_SOURCE);
30
31       buff[aux] = 0;
32       printf ("%s. Awaiting nodes", buff);
33       for (set<int>::iterator it=leftOver.begin(); it!=leftOver.↩
                end(); ++it)
34           cout << ' ' << *it;
35       cout << '\n';
36
37     }
38   MPI_Finalize ();
39 }
```

4. Create a model of the characteristics of the communication link joining two processes running on two different machines, i.e. calculate the start-up latency and communication rate, by implementing and testing a **ping-pong** benchmark program. A ping-pong program measures the time elapsed between sending a message, having it bounce at its destination, and receiving it back at its origin. By varying the message size, you can use statistical methods (least-squares) to estimate the start-up latency and rate as the intercept and slope respectively, of the line fitted to the experimental data.

**Answer**

The following code will calculate the average communication time over
REP number of tests between two processes, and for message sizes ranging
from 0 to MAX_MESG in increments of 1000.

```cpp
#include <stdio.h>
#include <stdlib.h>
#include <mpi.h>

using namespace std;

#define MESG_TAG 0
#define END_TAG 1
#define MAX_MESG 1000000

const int REP = 10;

int main (int argc, char *argv[])
{
  int size, rank;
  int namelen;
  char processor_name[MPI_MAX_PROCESSOR_NAME];
  int mesg_size;
  int tag;
  char *buffer;
  double start_time, end_time;
  MPI_Status status;

  buffer = new char[MAX_MESG];

  MPI_Init (&argc, &argv);
  MPI_Comm_size (MPI_COMM_WORLD, &size);
  MPI_Comm_rank (MPI_COMM_WORLD, &rank);
  MPI_Get_processor_name (processor_name, &namelen);

  printf ("Process %d of %d on %s\n", rank, size, processor_name)
      ;
  if (size < 2)
    {
      printf ("Need more than 1 processor to run\n");
      exit (1);
    }

  if (rank == 0)
    {
      for (int mesg_size = 0; mesg_size <= MAX_MESG; mesg_size +=
          1000)
        {
          start_time = MPI_Wtime ();
          for (int i = 0; i < REP; i++)
            {
              MPI_Send (buffer, mesg_size, MPI_CHAR, 1, tag,
                  MPI_COMM_WORLD);
              MPI_Recv (buffer, MAX_MESG, MPI_CHAR, 1,
                  MPI_ANY_TAG, MPI_COMM_WORLD, &status);
            }

          end_time = MPI_Wtime ();
          printf ("%i %lf\n", mesg_size, (end_time - start_time)
              / 2 / REP);
        }

      tag = END_TAG;
      MPI_Send (buffer, 0, MPI_CHAR, 1, tag, MPI_COMM_WORLD);
    }
  else
    {
      while (1)
        {
          MPI_Recv (buffer, MAX_MESG, MPI_CHAR, 0, MPI_ANY_TAG,
              MPI_COMM_WORLD, &status);
          if (status.MPI_TAG == END_TAG)
            break;
          MPI_Get_count (&status, MPI_CHAR, &mesg_size);
          MPI_Send (buffer, mesg_size, MPI_CHAR, 0, tag,
```

```
                    MPI_COMM_WORLD ) ;
65            }
66        }
67
68    delete [] buffer ;
69    MPI_Finalize () ;
70    return  0;
71  }
```

5. How would we need to modify the broadcasting program of Listing 5.9, if the source of the message was an arbitrary process, and not the one with rank 0?

**Answer**

The solution follows the guidelines of Listing 5.9, i.e. treating the processes as nodes in a hypercube, and proceeding to propagate the message along the hypercube's dimension one-by-one. However, the source-destination process pairings has to be done using an alternative approach. The key to the pairing calculation is the communication phase during which a node receives data. This is equal to the distance of a process/node from the root on the virtual hypercube, which can be calculated by counting the number of bits in the rank XOR root expression (variable recvPhase, line 35).

The sender node can be found by reversing the bit of a node's rank in the recvPhase - 1 position (line 51). Once a node receives the message, it can start to send it to all the nodes which are adjacent to it (ranks different by a single bit). The ranks of the receivers can be determined by reversing the bits in a node's rank from position recvPhase and onwards (lines 58-64).

```
1   #include<mpi.h>
2   #include<string.h>
3   #include<stdio.h>
4   #define MESSTAG 0
5
6   //*****************************************
7   // Returns the number of set bits in its argument
8   int bitCount (int i)
9   {
10    int count = 0;
11    while (i != 0)
12      {
13        int j = i;
14        i >>= 1;
15        if (i != j)
16          count++;
17      }
18    return count;
19  }
20
21  //*****************************************
22
23  int main (int argc, char **argv)
24  {
25    MPI_Init (&argc, &argv);
26
27    int rank, num, i;
28    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
29    MPI_Comm_size (MPI_COMM_WORLD, &num);
30    MPI_Status status;
31    int root = atoi (argv[1]);
32    double data = rank * 1.0; // something simple to enable ↩
            verification of correct receipt
33
```

```
34    printf ("%i will be served during phase %i\n", rank, bitCount (←
         rank ^ root));
35    int recvPhase = bitCount (rank ^ root);
36    if (rank == root)
37      {
38        int phase = 0;
39        int destID = rank ^ (1 << phase);
40        while (destID < num)
41          {                           // a subset of nodes gets a ←
                 message
42            printf ("#%i sending to %i\n", rank, destID);
43            MPI_Send (&data, 1, MPI_DOUBLE, destID, MESSTAG, ←
                 MPI_COMM_WORLD);
44            phase++;
45            destID = rank ^ (1 << phase);
46          }
47      }
48    else
49      {
50        int srcID;
51        srcID = rank ^ (1 << (recvPhase - 1));
52        printf ("SRC of #%i : %i\n", rank, srcID);
53
54        MPI_Recv (&data, 1, MPI_DOUBLE, srcID, MESSTAG, ←
                 MPI_COMM_WORLD, &status);
55
56        // calculate the ID of the node that will receive a copy of←
                 the message
57        int destID = rank ^ (1 << recvPhase);
58        while (destID < num)
59          {
60            printf ("#%i sendin to %i\n", rank, destID);
61            MPI_Send (&data, 1, MPI_DOUBLE, destID, MESSTAG, ←
                 MPI_COMM_WORLD);
62            recvPhase++;
63            destID = rank ^ (1 << recvPhase);
64          }
65      }
66    printf ("Node #%i has %lf\n", rank, data);
67    MPI_Finalize ();
68    return 0;
69 }
```

6. Assuming that the execution platform of your program, consists of 4 machines with identical architecture but different CPU clocks: one with 4GHz, one with 3 GHz and 2 with 2 GHz. How should you split the matrix A used in the example of Section 5.11.1, in order to solve the matrix-vector product problem in the smallest possible time?

**Answer**

The heterogeneity of the execution platform calls for an uneven partitioning of the problem's data, i.e. we should assign a bigger portion of the problem to the faster nodes. If we assume that the computing speed of the machines is linearly dependent on their clock, and we ignore communication overheads, then if $T$ is the time to complete the multiplication on a 1GHz machine, then on a $C$ times faster-clock machine we would need time $\frac{T}{C}$. If each node $i$ were assigned $part_i$ percent of the $M$ rows of the input matrix, then it would need time $T_i = part_i \frac{T}{C_i}$ to complete its part of the computation.

It can be easily proven by contradiction, that the overall execution time is minimized if $T_i = T_j$ for $\forall i \neq j$, which translates to:

$$T_i = T_j \Rightarrow part_i \frac{T}{C_i} = part_j \frac{T}{C_j} \Rightarrow part_i = part_j \frac{C_i}{C_j} \qquad (5.1)$$

As the sum of all parts should be equal to 1, we can derive the value of $part_0$ and subsequently from Eq. 5.1 all the other parts:

$$\sum_{\forall i} part_i = 1 \Rightarrow part_0 \sum_{\forall i} \frac{C_i}{C_0} = 1 \Rightarrow part_0 = (\sum_{\forall i} \frac{C_i}{C_0})^{-1} \qquad (5.2)$$

From the problem description we have $C_0 = 4GHz$, $C_1 = 3GHz$, $C_2 = 2GHz$ and $C_3 = 2GHz$, which when substituted in Eq. 5.1 and 5.2 return:

$$part_0 = (\frac{4}{4} + \frac{3}{4} + \frac{2}{4} + \frac{2}{4})^{-1} = \frac{4}{11}$$

$$part_1 = \frac{4}{11}\frac{3}{4} = \frac{3}{11}$$

$$part_2 = part_3 = \frac{4}{11}\frac{2}{4} = \frac{2}{11}$$

The modified code for setting up the scattering of the matrix to the four nodes, is shown below:

```
1    double part[] = { 4.0 / 11, 3.0 / 11, 2.0 / 11, 2.0 / 11 };
2
3    if (rank == 0)
4      {
5        int displs[N];
6        int sendcnts[N];
7        for (int i = 0; i < N; i++)
8          {
9            int rowsForProcess = M * part[i];
10           sendcnts[i] = M * rowsForProcess;
11           displs[i] = (i == 0) ? 0 : displs[i - 1] + sendcnts[i -
                 1];
12           if (i == N - 1)
13             sendcnts[i] = M * M - displs[i];
14         }
15
16       MPI_Scatterv (A, sendcnts, displs, MPI_DOUBLE, MPI_IN_PLACE
             , 0, MPI_DOUBLE, 0, MPI_COMM_WORLD);
```

7. Write a program that performs gathering as efficiently as possible, using point-to-point communications, i.e. the equivalent of MPI_Gather. What is the *time complexity* of your algorithm?

**Answer**

The answer is very closely related to the solution of Exercise 5.5. The major thing that needs to be altered, is the reversal of data flow, i.e. replacing MPI_Send with MPI_Recv and vice-versa. The XOR operator can be still used to calculate the destination and source nodes in the communications.

The other change stems from the need to accumulate the data that are collected. Because the order of collection does not necessarily (and typically does not) match the node ranking, all the parts that are accumulated are prefixed in the same buffer with the rank of the node that contributed them. In this fashion, the designated root of the gathering operation, can reshuffle the parts and produce an outcome compatible with what MPI_Gather would produce.

The following program requires as a command-line parameter, the rank of the process/node that will gather the data.

```c
#include<mpi.h>
#include<string.h>
#include<stdio.h>
#include<stdlib.h>
#include<math.h>
#define MESSTAG 0
#define MESSSIZE 10

//*******************************************
// Returns the number of set bits in its argument
int bitCount (int i)
{
   int count = 0;
   while (i != 0)
     {
        int j = i;
        i >>= 1;
        if (i != j)
           count++;
     }
   return count;
}

//*******************************************

int main (int argc, char **argv)
{
   MPI_Init (&argc, &argv);

   int rank, num, i;
   MPI_Comm_rank (MPI_COMM_WORLD, &rank);
   MPI_Comm_size (MPI_COMM_WORLD, &num);
   MPI_Status status;
   int root = atoi (argv[1]);
   int data[MESSSIZE];
   for (int i = 0; i < MESSSIZE; i++)
     data[i] = rank * 1.0;          // something simple to enable ↩
         verification of correct receipt

   int totalPhases = ceil (log2 (num)); // total collection phases

   int sendPhase = bitCount (rank ^ root); // phase when a node ↩
         forwards its part of the data
   printf ("%i will sent to root during phase %i\n", rank, ↩
         sendPhase);

   // allocate the temp. buffer for holding the data to be ↩
         communicated towards the "root"
   int buffSpace = (MESSSIZE + 1) << (totalPhases - sendPhase);
   if (rank == root)
     buffSpace = num * (MESSSIZE + 1);
   int *commBuff = new int[buffSpace];
   int inCommBuff = MESSSIZE + 1;
   // setup local part of the comm data and the relevant "header"
   commBuff[0] = rank;
   memcpy (commBuff + 1, data, sizeof (int) * MESSSIZE);
   printf ("#%i has BUFF %i\n", rank, buffSpace);

   if (rank == root)
     {
        int phase = totalPhases - 1;
        int srcID = rank ^ (1 << phase);
        while (phase >= 0)
          {
             if (srcID < num)
               {
                  printf ("#%i recv from %i  SPACE %i\n", rank, srcID↩
                        , buffSpace - inCommBuff);
                  MPI_Recv (commBuff + inCommBuff, buffSpace - ↩
                        inCommBuff, MPI_INT, srcID, MESSTAG, ↩
                        MPI_COMM_WORLD, &status);
```

```
65            int recvd;
66            MPI_Get_count (&status, MPI_INT, &recvd);
67            inCommBuff += recvd;
68          }
69        phase--;
70        srcID = rank ^ (1 << phase);
71      }
72
73    // now re-arrange the collected data
74    int gatheredData[MESSSIZE * num];
75    int pos = 0;
76    for (int i = 0; i < num; i++)
77      {
78        int destPartIdx = commBuff[pos++];
79        for (int j = 0; j < MESSSIZE; j++)
80          gatheredData[j + destPartIdx * MESSSIZE] = commBuff[↩
              pos++];
81      }
82
83    // print-out for verification purposes
84    for (int i = 0; i < MESSSIZE * num; i++)
85      printf ("%i ", gatheredData[i]);
86    printf ("\n");
87    }
88  else
89    {
90      int phase = totalPhases;
91      // calculate the ID of the node that will send its part of ↩
            the data to this node
92      int srcID = rank ^ (1 << (phase - 1));
93      // collect from other nodes first
94      while (phase > sendPhase)
95        {
96          if (srcID < num)
97            {
98              printf ("#%i recv from %i\n", rank, srcID);
99              MPI_Recv (commBuff + inCommBuff, buffSpace - ↩
                  inCommBuff, MPI_INT, srcID, MESSTAG, ↩
                  MPI_COMM_WORLD, &status);
100             int recvd;
101             MPI_Get_count (&status, MPI_INT, &recvd);
102             inCommBuff += recvd;
103           }
104         phase--;
105         srcID = rank ^ (1 << (phase - 1));
106       }
107
108     // then send towards the gathering root
109     int destID = rank ^ (1 << (sendPhase - 1));
110     printf ("DEST of #%i : %i\n", rank, destID);
111
112     MPI_Send (commBuff, inCommBuff, MPI_INT, destID, MESSTAG, ↩
            MPI_COMM_WORLD);
113   }
114
115   MPI_Finalize ();
116   return 0;
117 }
```

8. Write a program that performs scattering as efficiently as possible, using point-to-point communications, i.e. the equivalent of MPI_Scatter. What is the *volume* of data collectively communicated, if each process is to receive $K$ number of items? Express this number as a function of $K$ and the number of processes $N$.

**Answer**

The following program requires as a command-line parameter, the rank of the process/node that will scatter the data. The algorithm mirrors the broadcasting algorithm used in Exercise 5.5, as far as the tree of nodes
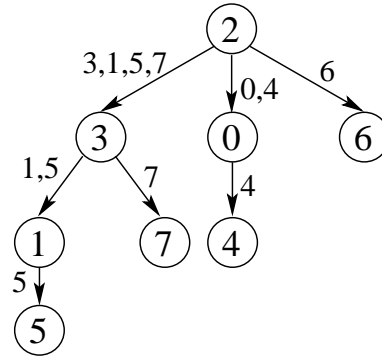
Figure 5.1: An example of a scattering operation initiated from node 2. The edges are labeled with the parts sent in a single operation.

that is implicitly formed is concerned.

The major difference is that the data to be scattered are rearranged in the root node, as the data sent-down a branch of the scatter "tree" are not related to consecutive ranking nodes. Figure 5.1 illustrates a scattering operation initiated from node 2 in an eight-node communicator.

The recursive findOrder() function of lines 31-43, calculates the order with which the parts need to be rearranged, so that parts collectively sent down a tree branch, occupy a contiguous buffer space.

For the example shown in Figure 5.1, the parts residing in the scatteredData array are rearranged accordingly to form the sequence : 2 3 1 5 7 0 4 6.

```cpp
1   #include<mpi.h>
2   #include<string.h>
3   #include<stdio.h>
4   #include<stdlib.h>
5   #include<iostream>
6   #include<vector>
7   #include<math.h>
8
9   using namespace std;
10
11  #define MESSTAG 0
12  #define MESSSIZE 10
13
14  //*****************************************
15  // Returns the number of set bits in its argument
16  int bitCount (int i)
17  {
18    int count = 0;
19    while (i != 0)
20      {
21        int j = i;
22        i >>= 1;
23        if (i != j)
24          count++;
25      }
26    return count;
27  }
28
29  //*****************************************
30  // used to calculate in what order the data to be scattered ↵
          should be arranged
31  void findOrder (int ID, int phase, int N, vector < int >&v)
32  {
```

```
33    if (ID >= N)
34      return;
35    v.push_back (ID);
36    int nextID = ID ^ (1 << phase);
37    while (nextID < N)
38      {
39        findOrder (nextID, phase + 1, N, v);
40        phase++;
41        nextID = ID ^ (1 << phase);
42      }
43  }
44
45  //*****************************************
46
47  int main (int argc, char **argv)
48  {
49    MPI_Init (&argc, &argv);
50
51    int rank, num, i;
52    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
53    MPI_Comm_size (MPI_COMM_WORLD, &num);
54    MPI_Status status;
55    int root = atoi (argv[1]);
56    int *data;  // data to be scattered
57
58    int totalPhases = ceil (log2 (num));    // total communication ←
          phases
59
60    int sendPhase = bitCount (rank ^ root); // phase when a node ←
          forwards its part of the data
61
62    // allocate the temp. buffer for holding the data to be ←
          communicated
63    int buffSpace = MESSSIZE << (totalPhases − sendPhase);
64    if (rank == root)
65      buffSpace = num * MESSSIZE;
66
67    int *scatteredData = new int[buffSpace];
68
69    if (rank == root)
70      {
71        data = new int[MESSSIZE * num];
72        for (int i = 0; i < buffSpace; i++)
73          data[i] = i;   // something simple to enable verification←
                of correct receipt
74
75        vector < int >reOrder;
76        findOrder (root, 0, num, reOrder);
77
78        // re−arrange the data to be scattered
79        int pos = 0;
80        for (int i = 0; i < num; i++)
81          {
82            int destPartIdx = reOrder[i];
83
84            for (int j = 0; j < MESSSIZE; j++)
85              scatteredData[pos ++] = data[j + destPartIdx * ←
                  MESSSIZE];
86          }
87
88        // keep the data that destined for the root of the scatter ←
              tree
89        memcpy (data, scatteredData, MESSSIZE * sizeof (int));
90
91        int phase = 0;
92        int buffPos = MESSSIZE;   // skip the part assigned to the ←
              root
93        while (phase <= totalPhases)
94          {
95            int toSend = MESSSIZE << (totalPhases − phase − 1);
96            int destID = rank ^ (1 << phase);
97
98            if (destID < num)
99              {
```

```cpp
100                     printf ("#%i sends to %i (%i) at %i\n", rank, ←
                            destID, toSend, buffPos);
101                     MPI_Send (scatteredData + buffPos, toSend, MPI_INT,←
                            destID, MESSTAG, MPI_COMM_WORLD);
102                     buffPos += toSend;
103                   }
104               phase++;
105             }
106         }
107     else
108       {
109         data = new int[MESSSIZE];
110         // calculate the ID of the node that will send the data to ←
                be scattered
111         int srcID = rank ^ (1 << (sendPhase - 1));
112         int toRecv = MESSSIZE << (totalPhases - sendPhase);
113
114         printf ("#%i getting from %i (%i)\n", rank, srcID, toRecv);
115
116         MPI_Recv (scatteredData, toRecv, MPI_INT, srcID, MESSTAG, ←
                MPI_COMM_WORLD, MPI_STATUS_IGNORE);
117
118         memcpy (data, scatteredData, MESSSIZE * sizeof (int));    ←
                // keep the data belonging to the node
119
120         // scatter to other nodes now
121         int phase = sendPhase + 1;
122         int destID = rank ^ (1 << (phase - 1));
123         int buffPos = MESSSIZE;    // skip the part assigned to the ←
                root
124         int toSend = toRecv >> 1;
125         while (phase <= totalPhases)
126           {
127             int destID = rank ^ (1 << phase - 1);
128             if (destID < num)
129               {
130                 printf ("#%i sends to %i (%i)\n", rank, destID, ←
                        toSend);
131                 MPI_Send (scatteredData + buffPos, toSend, MPI_INT,←
                        destID, MESSTAG, MPI_COMM_WORLD);
132                 buffPos += toSend;
133               }
134             phase++;
135             toSend >>= 1;
136           }
137       }
138     // verify scattering
139     printf ("%i finished : ", rank);
140     for (int i = 0; i < MESSSIZE; i++)
141       printf ("%i ", data[i]);
142     printf ("\n");
143
144     MPI_Finalize ();
145     delete [] data;
146     delete [] scatteredData;
147     return 0;
148 }
```

The scattering is performed over $\lceil lg(N) \rceil$ number of steps. In each step, the participating nodes grow exponentially, while the communicated parts shrink exponentially (by a factor of 2 in both cases).

If we number the communication phases $i = 0, \ldots, \lceil lg(N) \rceil - 1$, then we have:

- Phase 0: participating nodes $2^0$, data communicated $\frac{K \cdot N}{2}$
- Phase 1: participating nodes $2^1$, data communicated per node $\frac{K \cdot N}{2^2}$
- Phase 2: participating nodes $2^2$, data communicated per node $\frac{K \cdot N}{2^3}$
- ...

Therefore, the overall data volume is:

$$\sum_{i=0}^{\lceil lg(N)\rceil-1} 2^i \frac{K \cdot N}{2^{i+1}} = \sum_{i=0}^{\lceil lg(N)\rceil-1} \frac{K \cdot N}{2} = \frac{K \cdot N}{2}\lceil lg(N)\rceil \qquad (5.3)$$

9. The amount of data exchanged during every step of the butterfly pattern in Figure 5.9, double in relation to the previous step. If initially every process had data of size $K$ bytes to exchange, what is the total time required for the operation to complete, if we assume that each message exchange takes time $t_s + l \cdot V$, where $t_s$ is the link's start-up latency, $V$ is the volume of data to be sent and $l$ is the inverse of the communication speed.

   **Answer**

   The butterfly operation is completed in $lg(N)$ steps, with each step $i = 0, \ldots, lg(N) - 1$, requiring the exchange of $2^i K$ data moving in opposite directions. If we assume that there is no delay between the steps, and that communications are bi-directional, then the overall time would be equal to:

$$\sum_{i=0}^{lg(N)-1} t_s + l \cdot 2^i K =$$

$$lg(N)t_s + l \cdot K \sum_{i=0}^{lg(N)-1} 2^i =$$

$$lg(N)t_s + l \cdot K(2^{lg(N)} - 1) =$$

$$lg(N)t_s + l \cdot K(N - 1) \quad (5.4)$$

10. An alternative parallel bucket sort algorithm, would have the root process of a $N$-process run, scan the input data and split them into $N$ buckets before scattering the buckets to the corresponding processes. Implement this alternative design and compare its performance with the version presented in Section 5.11.5.

    **Answer**

```cpp
#include<mpi.h>
#include<stdlib.h>
#include<math.h>
#include<iostream>

using namespace std;

const int MIN = 0;
const int MAX = 10000;
//****************************************
int comp (const void *a, const void *b)
{
   return *(reinterpret_cast < const int *>(a)) -*(↩
        reinterpret_cast < const int *>(b));
}

//****************************************
void initData (int min, int max, int *d, int M)
{
   srand (time (0));
```

```cpp
20      for (int i = 0; i < M; i++)
21        d[i] = (rand () % (max - min)) + min;
22  }
23
24  //*****************************************
25  int main (int argc, char **argv)
26  {
27    MPI_Init (&argc, &argv);
28
29    int rank, N;
30    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
31    MPI_Comm_size (MPI_COMM_WORLD, &N);
32    MPI_Status status;
33
34    if (argc == 1)
35      {
36        if (rank == 0)
37          cerr << "Usage " << argv[0] << " number_of_items\n";
38        exit (1);
39      }
40
41    int M = atoi (argv[1]);
42    int maxItemsPerBucket = M;
43    int deliveredItems;
44    int bucketRange = ceil (1.0 * (MAX - MIN) / N);
45    int *data = new int[M];
46
47    if (rank == 0)
48      {
49        initData (MIN, MAX, data, M);
50        int *buckets = new int[N * maxItemsPerBucket];
51        int *bucketOffset = new int[N]; // where do buckets begin?
52        int *inBucket = new int[N];      // how many items in each ←
                one?
53
54
55        // initialize bucket counters and offsets
56        for (int i = 0; i < N; i++)
57          {
58            inBucket[i] = 0;
59            bucketOffset[i] = i * maxItemsPerBucket;
60          }
61
62        // split into buckets
63        for (int i = 0; i < M; i++)
64          {
65            int idx = (data[i] - MIN) / bucketRange;
66            int off = bucketOffset[idx] + inBucket[idx];
67            buckets[off] = data[i];
68            inBucket[idx]++;
69          }
70
71        MPI_Scatter (inBucket, 1, MPI_INT, &deliveredItems, 1, ←
                MPI_INT, 0, MPI_COMM_WORLD);
72        MPI_Scatterv (buckets, inBucket, bucketOffset, MPI_INT, ←
                data, maxItemsPerBucket, MPI_INT, 0, MPI_COMM_WORLD);
73        qsort (data, deliveredItems, sizeof (int), comp);
74
75        // calculate the offsets of the buckets again, so that they←
                end-up being contiguous after the gathering
76        for (int j = 1; j < N; j++)
77          bucketOffset[j] = bucketOffset[j - 1] + inBucket[j - 1];
78
79        MPI_Gatherv (data, deliveredItems, MPI_INT, data, inBucket,←
                bucketOffset, MPI_INT, 0, MPI_COMM_WORLD);
80
81        // print-out for verification purposes
82        for (int i = 0; i < M; i++)
83          cout << data[i] << " ";
84        cout << endl;
85
86        // release memory
87        delete [] buckets;
88        delete [] inBucket;
```

```
89            delete [] bucketOffset;
90          }
91      else
92        {
93          // get the size of the data to be sorted locally first
94          MPI_Scatter (NULL, 1, MPI_INT, &deliveredItems, 1, MPI_INT,↩
                 0, MPI_COMM_WORLD);
95          MPI_Scatterv (NULL, NULL, NULL, MPI_INT, data, ↩
                 maxItemsPerBucket, MPI_INT, 0, MPI_COMM_WORLD);
96          qsort (data, deliveredItems, sizeof (int), comp);
97          MPI_Gatherv (data, deliveredItems, MPI_INT, NULL, NULL, ↩
                 NULL, MPI_INT, 0, MPI_COMM_WORLD);
98        }
99
100     MPI_Finalize ();
101     delete [] data;
102     return 0;
103 }
```

11. Write a function that could be used for providing multicasting capabilities to a program, i.e. to be able to send a message to a subset of the processes of a communicator. Use an appropriate collective operation for the task.

    **Answer**

```
1  // Multicasts N elements of type t from address data to the list
2  // of processes included in the ranks array
3  // The first element in ranks is assumed to be the root for
4  // MPI_Bcast as it automatically get rank 0 in the new commun.
5  void multicast (void *data, int N, MPI_Datatype t, int *ranks, ↩
        int rN)
6  {
7    MPI_Group g, all;
8    MPI_Comm mc;
9
10   MPI_Comm_group (MPI_COMM_WORLD, &all);
11   MPI_Group_incl (all, rN, ranks, &g);
12   int localRank;
13   MPI_Group_rank (g, &localRank);
14
15   MPI_Comm_create (MPI_COMM_WORLD, g, &mc);
16   if (localRank != MPI_UNDEFINED)
17     {
18       MPI_Bcast (data, N, t, 0, mc);
19       MPI_Comm_free (&mc);
20     }
21
22   MPI_Group_free (&g);
23   MPI_Group_free (&all);
24 }
```

12. Write the equivalent of the ping-pong program using RMA functions, and measure the communication speed achieved versus the size of the message used. Compare your results with the data rates accomplished with point-to-point communications.

    **Answer**

    The following program prints out the average communication time over REP tests, for messages ranging in size from 0 to MAX_MESG bytes, in increments of 1000.

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include <mpi.h>
5
6  using namespace std;
7
```

```c
8   #define MAX_MESG 1000000

10  const int REP = 10;

12  int main (int argc, char *argv[])
13  {
14    int size, rank;
15    int namelen;
16    char processor_name[MPI_MAX_PROCESSOR_NAME];
17    int mesg_size;
18    int tag;
19    unsigned char *buffer;
20    double start_time, end_time;
21    MPI_Status status;
22    MPI_Win w;
23    MPI_Group otherProc, all;

25    MPI_Init (&argc, &argv);
26    MPI_Comm_size (MPI_COMM_WORLD, &size);
27    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
28    MPI_Get_processor_name (processor_name, &namelen);

30    printf ("Process %d of %d on %s\n", rank, size, processor_name)↩
          ;
31    if (size != 2)
32      {
33        printf ("Requires exactly 2 processes to run\n");
34        exit (1);
35      }

37    buffer = new unsigned char[MAX_MESG];
38    memset (buffer, rank, MAX_MESG);

40    MPI_Comm_group (MPI_COMM_WORLD, &all);
41    MPI_Group_excl (all, 1, &rank, &otherProc);
42    MPI_Win_create (buffer, MAX_MESG, 1, MPI_INFO_NULL, ↩
          MPI_COMM_WORLD, &w);

44    if (rank == 0)
45      {
46        for (int mesg_size = 0; mesg_size <= MAX_MESG; mesg_size +=↩
              1000)
47          {
48            start_time = MPI_Wtime ();
49            for (int i = 0; i < REP; i++)
50              {
51                // send
52                MPI_Win_start (otherProc, 0, w);
53                MPI_Put (buffer, mesg_size, MPI_UNSIGNED_CHAR, 1, ↩
                      0, mesg_size, MPI_CHAR, w);
54                MPI_Win_complete (w);

56                // wait for response
57                MPI_Win_post (otherProc, 0, w);
58                MPI_Win_wait (w);
59              }

61            end_time = MPI_Wtime ();
62            printf ("%i %lf\n", mesg_size, (end_time − start_time) ↩
                  / 2 / REP);
63          }
64      }
65    else
66      {
67        // other process mirrors what rank 0 is doing
68        for (int mesg_size = 0; mesg_size <= MAX_MESG; mesg_size +=↩
              1000)
69          for (int i = 0; i < REP; i++)
70            {
71              MPI_Win_post (otherProc, 0, w);
72              MPI_Win_wait (w);

74              MPI_Win_start (otherProc, 0, w);
75              MPI_Put (buffer, mesg_size, MPI_UNSIGNED_CHAR, 0, 0, ↩
```

```
                  mesg_size , MPI_CHAR , w );
76              MPI_Win_complete (w);
77            }
78        }
79
80    delete [] buffer ;
81    MPI_Win_free (&w );
82    MPI_Group_free (&all );
83    MPI_Group_free (&otherProc );
84    MPI_Finalize ();
85    return 0;
86 }
```

13. Modify the program of Listing 5.7, so that the partitioning of the range depends on the relative speed of the participating nodes. One easy approach is to make the ranges proportional to the CPU operating frequency, or the calculated bogomips. Both numbers are available in the /proc/cpuinfo pseudo-file. The master can collect the numbers and reply back to the worker nodes with the calculated ranges. If we represent as $m_i$ the $i$-node's bogomips, the percent of the range $\alpha_i$ that should be assigned to node $i$ can be calculated as $\alpha_i = \frac{m_i}{\sum_{\forall k} m_k}$

**Answer**

The master node has to gather the bogomips of the execution platform machines, and calculate the appropriate ranges before sending them back to the nodes. The bogomips can be read from the /proc/cpuinfo pseudo-file using a pipe as shown in line 23. The bogomips numbers are gathered by the master node (line 31) before calculating the parts $\alpha_i$ (lines 35-41, identified as part[] in the code) and the ranges that should be assigned to each node (lines 43-51). The remaining code is identical to Listing 5.7.

```
1  #include <mpi.h>
2  #include<stdio.h>
3  #include<stdlib.h>
4  #include<string.h>
5
6  #define RANGEMIN 0
7  #define RANGEMAX 1000
8  #define MSGTAG 0
9
10
11 int main (int argc , char **argv)
12 {
13    int rank , num , i;
14    int range [2];
15    MPI_Init (&argc , &argv );
16    MPI_Comm_rank (MPI_COMM_WORLD , &rank );
17    MPI_Comm_size (MPI_COMM_WORLD , &num );
18    MPI_Status status ;
19
20    float bogoMIPS ;
21    float *allBogo = NULL ;
22
23    FILE *f = popen ("/bin/cat /proc/cpuinfo | /bin/grep bogo | /←
          usr/bin/head −n 1  | /usr/bin/gawk −F ':'  '{print $2}'", "←
          r");
24    fscanf (f, "%f", &bogoMIPS );
25    pclose (f);
26
27    // array for gathering is only need in rank 0
28    if (rank == 0)
29      allBogo = new float [num ];
30
31    MPI_Gather (&bogoMIPS , 1, MPI_FLOAT , allBogo , 1, MPI_FLOAT , 0, ←
          MPI_COMM_WORLD );
32
```

```cpp
33     if (rank == 0)
34       {
35         float sumBogo = allBogo[0];
36         for (i = 1; i < num; i++)
37           sumBogo += allBogo[i];
38
39         float part[num];
40         for (i = 0; i < num; i++)
41           part[i] = allBogo[i] / sumBogo;
42
43         int rng[2 * num];
44         int width = RANGEMAX - RANGEMIN;
45         rng[0] = RANGEMIN;            // left limit
46         rng[1] = rng[0] + width * part[0] - 1;     // right limit
47         for (i = 1; i < num; i++)
48           {
49             rng[i * 2] = rng[i * 2 - 1] + 1;
50             rng[i * 2 + 1] = (i == num - 1) ? RANGEMAX : rng[i * 2]↩
                  + width * part[i] - 1;
51           }
52
53         MPI_Request rq[num - 1];
54         for (i = 1; i < num; i++)
55           MPI_Isend (rng + i * 2, 2, MPI_INT, i, MSGTAG, ↩
                  MPI_COMM_WORLD, &(rq[i - 1]));
56
57         for (i = 1; i < num; i++)
58           MPI_Wait (&(rq[i - 1]), &status);
59
60         range[0] = rng[0];          // master's limits
61         range[1] = rng[1];
62
63         delete [] allBogo;
64       }
65     else
66       {
67         MPI_Request rq;
68         MPI_Irecv (range, 2, MPI_INT, 0, MSGTAG, MPI_COMM_WORLD, &↩
                  rq);
69         MPI_Wait (&rq, &status);
70       }
71
72     printf ("Node %i's range : ( %i, %i )\n", rank, range[0], range↩
            [1]);
73
74     MPI_Finalize ();
75     return 0;
76   }
```

14. The butterfly communication scheme that is outlined in Section 5.11.4, is only one of the possible strategies for an all-to-all, or all-reduce data exchange. A different approach would be mandated if the underlying communication infrastructure did not provide the required links, making the procedure inefficient. An example of such an architecture is the ring, where each node is directly connected to just two others.

   (a) Write an MPI program that would implement an efficient all-to-all exchange of data on a ring of machines.

   (b) How many steps would be required in comparison to a butterfly scheme, if the number of nodes/processes were $N = 2^k$?

   (c) If we assume that time taken to send $V$ bytes over a communication link is given by $l \cdot V$, where $l$ is the (inverse of the) link speed in $sec/byte$, how does your algorithm compare against the butterfly scheme in terms of overall communication time?

   **Answer**

(a) The following program contains an `allToAll` function that is equivalent to `MPI_Alltoall`. Assuming that $N$ processes are organized in a ring, each process performs $N-1$ rounds/steps of forwarding its data to the next process (identified by `destID`) and receiving data from the previous process (identified by `srcID`). The data are rearranged so that the first block of `NperNode` items received, should reside with the process. The remaining `N - round - 1` blocks are sent to the next process and so on. The data rearrangement is performed before commencing any communications in lines 24-26. The target data layout depends on the rank of the process.

The MPI datatype of the communicated data is needed so that heterogeneous platforms can be supported, while the native type size (parameter `typeSize`) is needed so that a properly-sized data buffer is allocated for handling the communication traffic.

The communication pattern for the `allToAllGather` function is made up from the same number of steps, but with a reverse flow of communications and each communication carries `NperNode` elements.

```cpp
1   //All−to−all scattering on a ring
2   // Number of processes can be arbitrary
3   #include<mpi.h>
4   #include<string.h>
5   #include<unistd.h>
6   #include<math.h>
7   #include<iostream>
8
9   const int K = 10;
10  const int ALLGATHERTAG = 0;
11  const int ALLSCATTERTAG = 0;
12
13  using namespace std;
14
15  //−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−
16  // all−to−all scattering. outBuffer/inBuffer are supposed to↩
          have N ∗ NperNode elements/space
17  void allToAll (void ∗outBuffer, int NperNode, int typeSize, ↩
        MPI_Datatype t, void ∗inBuffer)
18  {
19     int rank, N;
20     MPI_Comm_rank (MPI_COMM_WORLD, &rank);
21     MPI_Comm_size (MPI_COMM_WORLD, &N);
22
23     // local−destined part of data copied first
24     memcpy (inBuffer + NperNode ∗ typeSize ∗ rank, outBuffer +↩
            NperNode ∗ typeSize ∗ rank, NperNode ∗ typeSize);
25
26     // rearrange local data before starting to sent
27     unsigned char ∗tempBuff = new unsigned char[(N − 1) ∗ ↩
            NperNode ∗ typeSize];
28     memcpy (tempBuff, outBuffer + NperNode ∗ typeSize ∗ (rank ↩
            + 1), (N − rank − 1) ∗ NperNode ∗ typeSize);
29     memcpy (tempBuff + (N − rank − 1) ∗ NperNode ∗ typeSize, ↩
            outBuffer, rank ∗ NperNode ∗ typeSize);
30
31     int round = 1;
32     int toSend = (N − 1) ∗ NperNode;
33     int srcID = (rank == 0) ? N − 1 : rank − 1;
34     int destID = (rank + 1) % N;
35     void ∗ptr1 = tempBuff, ∗ptr2 = outBuffer, ∗aux;
36     while (round < N)
37       {
38          MPI_Send (ptr1, toSend, t, destID, ALLSCATTERTAG, ↩
                MPI_COMM_WORLD);
39          MPI_Recv (ptr2, toSend, t, srcID, ALLSCATTERTAG, ↩
                MPI_COMM_WORLD, MPI_STATUS_IGNORE);
40
```

```cpp
41          // keep local−destined part and sent the rest
42          int offset = ((rank − round + N) % N) * NperNode * ↵
                 typeSize;
43          memcpy (inBuffer + offset, ptr2, NperNode * typeSize);

45          // switch buffer pointers
46          aux = ptr1;
47          ptr1 = ptr2;
48          ptr2 = aux;
49          ptr1 += NperNode * typeSize;      // shift beginning ↵
                 of buffers to send/receive
50          ptr2 += NperNode * typeSize;
51          toSend −= NperNode;
52          round++;
53        }
54      delete [] tempBuff;
55    }

57    //——————————————————————————————————————————————
58    // all−to−all gathering. outBuffer is supposed to have
59    // NperNode elements, while inBuffer has enough space
60    // for N * NperNode
61    void allToAllGather (void *outBuffer, int NperNode, int ↵
            typeSize, MPI_Datatype t, void *inBuffer)
62    {
63      int rank, N;
64      MPI_Comm_rank (MPI_COMM_WORLD, &rank);
65      MPI_Comm_size (MPI_COMM_WORLD, &N);

67      // local−destined part of data copied first
68      memcpy (inBuffer + NperNode * typeSize * rank, outBuffer, ↵
            NperNode * typeSize);


71      int round = 1;
72      // reverse ring traversal than allToAll
73      int srcID = (rank + 1) % N;
74      int destID = (rank == 0) ? N − 1 : rank − 1;
75      void *ptr1 = inBuffer + NperNode * typeSize * rank;
76      void *ptr2 = (rank == N − 1) ? inBuffer : ptr1 + NperNode ↵
            * typeSize;
77      void *aux;
78      while (round < N)
79        {
80          MPI_Send (ptr1, NperNode, t, destID, ALLGATHERTAG, ↵
                MPI_COMM_WORLD);
81          MPI_Recv (ptr2, NperNode, t, srcID, ALLGATHERTAG, ↵
                MPI_COMM_WORLD, MPI_STATUS_IGNORE);

83          // advance buffer pointers
84          aux = ptr2;
85          ptr1 = ptr2;
86          ptr2 = ((rank + round + 1) % N == 0) ? inBuffer : aux ↵
                + NperNode * typeSize;
87          round++;
88        }
89    }

91    //*****************************************
92    int main (int argc, char **argv)
93    {
94      MPI_Init (&argc, &argv);

96      int rank, N;
97      MPI_Comm_rank (MPI_COMM_WORLD, &rank);
98      MPI_Comm_size (MPI_COMM_WORLD, &N);
99      MPI_Status status;


102     double *localPart = new double[K * N];
103     double *allParts = new double[K * N];

105     if (rank == 0)
106       cout << "Test for alltoAll\n";
```

```
107    for (int i = 0; i < K * N; i++)
108      localPart[i] = pow (10, rank) + i;
109
110    allToAll (localPart, K, sizeof (double), MPI_DOUBLE, ←
           allParts);
111
112    sleep (rank);
113    cout << rank << " : ";
114    for (int i = 0; i < K * N; i++)
115      cout << allParts[i] << " ";
116    cout << endl;
117
118    MPI_Barrier (MPI_COMM_WORLD);
119
120    if (rank == 0)
121      cout << "Test for alltoAllGather\n";
122    for (int i = 0; i < K; i++)
123      localPart[i] = rank;
124
125    allToAllGather (localPart, K, sizeof (double), MPI_DOUBLE, ←
           allParts);
126
127    sleep (rank);
128    cout << rank << " : ";
129    for (int i = 0; i < K * N; i++)
130      cout << allParts[i] << " ";
131    cout << endl;
132
133
134    MPI_Finalize ();
135    delete [] localPart;
136    delete [] allParts;
137    return 0;
138  }
```

(b) The number of steps/rounds required is $N - 1$ which is obviously much higher than the $lgN$ of the butterfly scheme. However, the data transmitted in every step by the `allToAllGather` function is constant, while in the butterfly algorithm they grow exponentially in size by a factor of 2.

(c) Assuming that the blocks destined for each process have a size of $K$ bytes, and that the communications during a round take place simultaneously, then the total communication costs are:

$$t_{comm}^{(ring)} = \sum_{i=0}^{N-2} l \cdot K = (N - 1) \cdot l \cdot K \qquad (5.5)$$

$$t_{comm}^{(butterfly)} = \sum_{i=0}^{lgN-1} l \cdot K \cdot 2^i = l \cdot K(2^{lgN} - 1) = (N - 1) \cdot l \cdot K \quad (5.6)$$

The $t_{comm}^{(butterfly)} = t_{comm}^{(ring)}$ result may seem surprising, but there is an important factor that has been intentionally ignored in this setting : communication latency. The start-up overhead can be significant. The butterfly pattern can offer substantially lower costs if this is factored-in as well.

15. The case study on Diffusion Limited Aggregation in Section 5.20, decomposes the problem's data on a particle-wise basis. Explore the alternative of partitioning the 2D grid and assigning it to the processes involved.

What are the benefits and drawbacks of this approach? Is the communication pattern involved different in this case? Does this approach produce similar results to the sequential program and alternative parallel programs?

**Answer**

The benefit of a row-wise geometric decomposition is a more regular communication pattern that is also localized, instead of global. Changes to the crystal formation do not have to be broadcasted to all the nodes. The drawbacks are (a) a skewed load distribution, as the nodes handling the crystal region have in general fewer particles assigned to them, and (b) a potentially much higher communication cost that grows linearly with the number of columns used. Results should be identical to the sequential program, as long as the sequential program also grows the crystal at the end of each time step.

In the following listing (which is part of a project including the `dla_core.cpp` file shown later), the grid is initialized in node 0 (line 74) and scattered in a row-wise fashion to the remaining nodes (line 79).

The grid local to each node is made up of two additional columns (to simplify boundary checking) and two additional rows that contain data handled by the neighboring processes. The two top and bottom rows of each local grid are exchanged with the previous and next process respectively (the first and last processes are obviously an exception). The particles that end-up in the top/bottom rows "migrate" to the previous/next process after the message exchanges taking place in line 92-93 and 99-100. The loops of lines 113-115 and 120-122 collect those migrating particles. The `bounceUp` and `bounceDown` variables control whether this migration is possible : It is obviously impossible for the very first and very last rows in the overall grid, because that would mean that the particles are leaving completely.

The final state of the grid is gathered by process 0 (line 134) and output to the console in the form of a PBM image (lines 138-149). A schematic of the problem decomposition is shown in Figure 5.2.

```cpp
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <string.h>
4   #include <assert.h>
5   #include <math.h>
6   #include <iostream>
7   #include "dla_core.h"
8   #include <mpi.h>
9
10  using namespace std;
11
12  #define PIC_SIZE 100
13  #define PARTICLES 500
14  #define MAX_ITER 10000
15
16  /*———————————————————————————————————————*/
17  int main (int argc, char **argv)
18  {
19    int cols, rows, iter, particles, x, y;
20    int *pic, *pic2, *tmp;
21    int *nextNodeRows, *prevNodeRows;
22    int rank, num, i;
23    MPI_Init (&argc, &argv);
24    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
```

Figure 5.2: Geometric problem decomposition illustration for the solution of Exercise 15.

```
25    MPI_Comm_size (MPI_COMM_WORLD, &num);
26    MPI_Request req;
27
28    if (argc < 2)                      // use default values if user ←
          does not specify anything
29      {
30        cols = PIC_SIZE;
31        rows = PIC_SIZE;
32        iter = MAX_ITER;
33        particles = PARTICLES;
34      }
35    else
36      {
37        cols = atoi (argv[1]);
38        rows = atoi (argv[2]);
39        particles = atoi (argv[3]);
40        iter = atoi (argv[4]);
41      }
42
43    // to make code simpler and avoid special cases, rows is forced←
            to be a multiple of
44    // the number of processes
45    rows = (int) ceil (rows * 1.0 / num) * num;
46    int rowsPerNode = rows / num;
47
48    // initialize the random number generator
49    srand(time(0));
50
51    // grid has two extra rows, one above and one below the actual ←
          rows
52    // The extra rows "belong" to other nodes, so at the end of ←
          each time step
53    // there should be a pair-wise reduction of these rows
54    if (rank == 0)
55      {
56        pic = new int [(cols + 2) * (rows + 2)];
57        pic2 = new int [(cols + 2) * (rows + 2)];
58      }
59    else
60      {
```

```
61             pic = new int [( cols + 2) * ( rowsPerNode + 2)];
62             pic2 = new int [( cols + 2) * ( rowsPerNode + 2)];
63         }
64       bool bounceUp = false , bounceDown = false ;
65       if ( rank == 0)
66         bounceUp = true ;
67       else if ( rank == num - 1)
68         bounceDown = true ;
69
70       nextNodeRows = new int [2 * ( cols + 2)];
71       prevNodeRows = new int [2 * ( cols + 2)];
72
73       if ( rank == 0)
74         dla_init ( pic2 , rows + 2, cols + 2, particles , 1);
75       // clean pic buffer
76       memset ( pic , 0, ( cols + 2) * ( rowsPerNode + 2) * sizeof ( int ));
77
78       // boundary rows do not need to be communicated now .
79       MPI_Scatter ( pic2 + cols + 2, ( cols + 2) * rowsPerNode , MPI_INT↩
               , pic + cols + 2, ( cols + 2) * rowsPerNode , MPI_INT , 0, ↩
               MPI_COMM_WORLD );
80
81       int nextID = rank + 1, prevID = rank - 1;
82       while ( --iter >= 0)
83         {
84           dla_evolve ( pic , pic2 , rowsPerNode + 2, cols + 2, bounceUp ,↩
                 bounceDown );
85           tmp = pic2 ;
86           pic2 = pic ;
87           pic = tmp ;
88           // exchange information with neighboring nodes
89           // exchange with " previous "
90           if ( rank != 0)
91             {
92               MPI_Isend ( pic , 2 * ( cols + 2), MPI_INT , prevID , iter , ↩
                     MPI_COMM_WORLD , &req );
93               MPI_Recv ( prevNodeRows , 2 * ( cols + 2), MPI_INT , prevID↩
                     , iter , MPI_COMM_WORLD , MPI_STATUS_IGNORE );
94             }
95
96           // exchange with " next "
97           if ( rank != num - 1)
98             {
99               MPI_Isend ( pic + rowsPerNode * ( cols + 2), 2 * ( cols + ↩
                     2), MPI_INT , nextID , iter , MPI_COMM_WORLD , &req );
100              MPI_Recv ( nextNodeRows , 2 * ( cols + 2), MPI_INT , nextID↩
                     , iter , MPI_COMM_WORLD , MPI_STATUS_IGNORE );
101            }
102
103          // reduction is duplicated in the nodes so that they dont ↩
                 have to wait
104          // what is important in the neighbors 'rows , is the negative↩
                 cells
105          if ( rank != 0)
106            {
107              for ( int i = 0; i < cols + 2; i++)
108                if ( prevNodeRows [i] < 0)
109                  pic [i] = prevNodeRows [i];
110                else
111                  pic [i] = 0;         // reset any particles now on the↩
                         other node
112
113              for ( int i = 0; i < cols + 2; i++)
114                if ( prevNodeRows [ cols + 2 + i] > 0)
115                  pic [ cols + 2 + i] += prevNodeRows [ cols + 2 + i];   ↩
                         // add particles coming from the other node
116            }
117
118          if ( rank != num - 1)
119            {
120              for ( int i = 0; i < cols + 2; i++)
121                if ( nextNodeRows [i] > 0)
122                  pic [ rowsPerNode * ( cols + 2) + i] += nextNodeRows [ ↩
                         i];   // add particles coming from the other ↩
```

```
                          node
123
124              for  ( int  i = 0;  i < cols + 2;  i++)
125                if  (nextNodeRows[cols + 2 + i] < 0)
126                  pic[(rowsPerNode + 1) * (cols + 2) + i] = ←
                          nextNodeRows[cols + 2 + i];
127                else
128                  pic[(rowsPerNode + 1) * (cols + 2) + i] = 0;         ←
                          // reset any particles now on the other node
129          }
130
131      }
132
133
134   MPI_Gather (pic + cols + 2, rowsPerNode * (cols + 2), MPI_INT, ←
            pic + cols + 2, rowsPerNode * (cols + 2), MPI_INT, 0, ←
            MPI_COMM_WORLD);
135   /* Print to stdout a PBM picture of the simulation space */
136   if  (rank == 0)
137      {
138        printf ("P1\n%i %i\n", cols, rows);
139
140        for  (y = 1;  y <= rows;  y++)
141          {
142            for  (x = 1;  x <= cols;  x++)
143              {
144                if  (pic[y * (cols + 2) + x] < 0)
145                  printf ("1 ");
146                else
147                  printf ("0 ");
148              }
149            printf ("\n");
150          }
151      }
152
153   MPI_Finalize ();
154   delete [] pic;
155   delete [] pic2;
156   delete [] nextNodeRows;
157   delete [] prevNodeRows;
158   return 0;
159 }
```

The `dla_core.cpp` file that supplements the code above is:

```
160 #include <stdlib.h>
161 #include "dla_core.h"
162
163 /*——————————————————————————————————————————
164  * Checks the presence of a structure in neighboring cells
165  * pic points to 2D array holding data, arranged in cols
166  * columns. Because the array is two columns and two rows
167  * wider than necessary, there is no need to check for
168  * boundary values of x and y.
169  */
170 int check_proxim (int *pic, int cols, int x, int y)
171 {
172   int *row0, *row1, *row2;
173   row0 = pic + (y - 1) * cols + x - 1;
174   row1 = row0 + cols;
175   row2 = row1 + cols;
176   if (*row0 < 0 || *(row0 + 1) < 0 || *(row0 + 2) < 0 || *row1 < ←
          0 || *(row1 + 1) < 0 || *(row1 + 2) < 0 || *row2 < 0 || *(←
          row2 + 1) < 0 || *(row2 + 2) < 0)
177     return (-1);
178   else
179     return (1);
180 }
181
182 /*——————————————————————————————————————————————————*/
183 /* Returns -1,0 and 1 with equal probability */
184 inline int three_way ()
185 {
186   return (random () % 3) - 1;
```

```c
187  }
188
189  /*——————————————————————————————————————*/
190  /* Initializes the 2D array for the simulation */
191  void dla_init (int *pic, int rows, int cols, int particles, int ←
         init_seed)
192  {
193    int i, j, x, y;
194    for (i = 0; i < rows; i++)
195      for (j = 0; j < cols; j++)
196        pic[i * cols + j] = 0;
197
198
199    for (i = 0; i < particles; i++)          /* generate initial ←
           particle placement */
200      {
201        x = random () % (cols - 2) + 1;    // counting starts from 1
202        y = random () % (rows - 2) + 1;
203        if ((y == rows / 2 + 1) && (x == cols / 2 + 1))   // repeat←
               if true
204          i--;
205        else
206          pic[y * cols + x]++;
207      }
208
209    if (init_seed)
210      pic[(rows / 2) * cols + (cols / 2)] = -1;   /* place initial ←
             seed */
211
212  }
213
214  /*——————————————————————————————————————
215   * Single step evolution of the simulation.
216   *
217   * The cell values represent:
218   *   0: empty space
219   *   >0 : multiple particles
220   *   <0 : crystal
221   *
222   * Returns the address of the structure holding the last update ←
           */
223  int *dla_evolve (int *pic, int *pic2, int rows, int cols, bool ←
         bounceUp, bool bounceDown)
224  {
225    int x, y, k;
226
227    // prepare array to hold new state
228    for (y = 1; y < rows - 1; y++)
229      for (x = 1; x < cols - 1; x++)
230        pic2[y * cols + x] = pic[y * cols + x] > 0 ? 0 : pic[y * ←
               cols + x];
231
232    for (y = 1; y < rows - 1; y++)
233      for (x = 1; x < cols - 1; x++)
234        for (k = 0; k < pic[y * cols + x]; k++)
235          {
236            int new_x = x + three_way ();
237            if (new_x == 0 && bounceUp)
238              new_x = 1;
239            else if (new_x == cols - 1 && bounceDown)
240              new_x = cols - 2;
241
242            int new_y = y + three_way ();
243            if (new_y == 0)
244              new_y = 1;
245            else if (new_y == rows - 1)
246              new_y = rows - 2;
247
248            if (pic2[new_y * cols + new_x] > 0)   // steps into ←
                   empty space
249              pic2[new_y * cols + new_x]++;
250            else if (pic2[new_y * cols + new_x] == 0)     // steps ←
                   into unchecked space
251              {
```

```
252                    pic2 [ new_y  *  cols  +  new_x ]  =  check_proxim  ( pic2 ,  ↵
                           cols ,  new_x ,  new_y );
253               }
254          }
255     return  pic2 ;
256  }
```

16. Implement a 3D-space Diffusion Limited Aggregation simulation, by extending the solutions provided in Section 5.20.

**Answer**

The project reported in Section 5.20 consists of two source code files: `dla_core.c` and `dla_mpi.c`. The transition to 3D requires significant changes only to the functions in `dla_core.c`, which handle the particles as represented by the `PartStr` structure (which in turn gains a `z` coordinate).

The complete code (with a minor conversion to C++) for `dla_core.cpp` is shown below, the major differences being limited to:

(a) Lines 24-33 of the `check_proxim` function, that now has to check 27 possible neighboring positions, scattered over 3 planes along the z-axis.

(b) The calculation of the indices into the grid array (as seen in lines 69, 109 and 142), that need to take into consideration the third dimension.

```cpp
1   /* Core library for simulating 3D DLA phenomena */
2   #include <stdlib.h>
3   #include <string.h>
4   #include "dla_core.h"
5
6   /*─────────────────────────────────────────────
7    * Checks the presence of a structure in neighboring cells
8    * pic points to 3D array holding data, arranged in depth number
9    * of colsXcolumns planes. Because the array is "surrounded" by
10   * an empty cube of cells, there is no need to check for
11   * boundary values of x, y and z.
12   */
13  int check_proxim (int *pic, int rows, int cols, int x, int y, int↵
         z)
14  {
15    int *plane0, *plane1, *plane2;
16    plane0 = pic + (z − 1) * rows * cols;
17    plane1 = pic + z * rows * cols;
18    plane2 = pic + (z + 1) * rows * cols;
19
20    int row0, row1, row2;
21    row0 = (y − 1) * cols + x − 1;
22    row1 = row0 + cols;
23    row2 = row1 + cols;
24    if (plane0 [row0] < 0 || plane0 [row0 + 1] < 0 || plane0 [row0 + ↵
           2] < 0 ||
25        plane0 [row1] < 0 || plane0 [row1 + 1] < 0 || plane0 [row1 + ↵
             2] < 0 ||
26        plane0 [row2] < 0 || plane0 [row2 + 1] < 0 || plane0 [row2 + ↵
             2] < 0 ||
27        plane1 [row0] < 0 || plane1 [row0 + 1] < 0 || plane1 [row0 + ↵
             2] < 0 ||
28        plane1 [row1] < 0 || plane1 [row1 + 1] < 0 || plane1 [row1 + ↵
             2] < 0 ||
29        plane1 [row2] < 0 || plane1 [row2 + 1] < 0 || plane1 [row2 + ↵
             2] < 0 ||
30        plane2 [row0] < 0 || plane2 [row0 + 1] < 0 || plane2 [row0 + ↵
             2] < 0 ||
31        plane2 [row1] < 0 || plane2 [row1 + 1] < 0 || plane2 [row1 + ↵
             2] < 0 ||
```

```
32              plane2[row2] < 0 || plane2[row2 + 1] < 0 || plane2[row2 + ↩
                    2] < 0)
33          return (-1);
34        else
35          return (1);
36      }
37
38      /*————————————————————————————————————————————*/
39      /* Returns -1,0 and 1 with equal probability */
40      inline int three_way ()
41      {
42        return (random () % 3) - 1;
43      }
44
45      /*————————————————————————————————————————————*/
46      /* Initializes the 3D array and array of particles for the ↩
            simulation */
47      void dla_init_plist (int *pic, int rows, int cols, int depth, ↩
            PartStr * p, int particles, int init_seed)
48      {
49        int i, x, y, z, idx;
50        memset (pic, 0, rows * cols * depth * sizeof (int));
51
52
53        for (i = 0; i < particles; i++)          /* generate initial ↩
                particle placement */
54          {
55            x = random () % (cols - 2) + 1;    // counting starts from 1
56            y = random () % (rows - 2) + 1;
57            z = random () % (rows - 2) + 1;
58            if ((y == rows / 2 + 1) && (x == cols / 2 + 1) && (z == ↩
                depth / 2 + 1))   // repeat if true
59              i--;
60            else
61              {
62                p[i].x = x;
63                p[i].y = y;
64                p[i].z = z;
65              }
66          }
67        if (init_seed)
68          {
69            idx = (depth / 2) * rows * cols + (rows / 2) * cols + (cols↩
                / 2);
70            pic[idx] = -1;               /* place initial seed */
71          }
72
73      }
74
75      /*————————————————————————————————————————————————
76       * Single step evolution of the simulation.
77       *
78       * The cell values represent:
79       *  0  : empty space
80       *  <0 : crystal  -- for consistency with alternative formulation
81       *
82       * Particles are held in a separate array.
83       */
84      void dla_evolve_plist (int *pic, int rows, int cols, int depth, ↩
            PartStr * p, int *particles, PartStr * changes, int *↩
            numChanges)
85      {
86        int i;
87        *numChanges = 0;
88
89        for (i = 0; i < *particles; i++)
90          {
91            int new_x = p[i].x + three_way ();
92            if (new_x == 0)                // bounce off boundaries
93              new_x = 1;
94            else if (new_x == cols - 1)
95              new_x = cols - 2;
96
97            int new_y = p[i].y + three_way ();
```

```
98          if (new_y == 0)                  // bounce off boundaries
99            new_y = 1;
100         else if (new_y == rows - 1)
101           new_y = rows - 2;
102
103         int new_z = p[i].z + three_way ();
104         if (new_z == 0)                  // bounce off boundaries
105           new_z = 1;
106         else if (new_z == depth - 1)
107           new_z = depth - 2;
108
109         int idx = new_z * rows * cols + new_y * cols + new_x;
110         if (pic[idx] == 0)         // steps into empty space
111           {
112             int turnCrystal = check_proxim (pic, rows, cols, new_x,↩
                    new_y, new_z);
113             if (turnCrystal < 0)
114               {
115                 // record crystal change
116                 changes[*numChanges].x = new_x;
117                 changes[*numChanges].y = new_y;
118                 changes[*numChanges].z = new_z;
119                 (*numChanges)++;
120
121                 // erase particle from list
122                 p[i] = p[(*particles) - 1];
123                 i--;
124                 (*particles)--;
125               }
126             else                         // change position to particle
127               {
128                 p[i].x = new_x;
129                 p[i].y = new_y;
130                 p[i].z = new_z;
131               }
132           }
133       }
134 }
135
136 /*————————————————————————————————————————————————————*/
137 void apply_changes (int *pic, int rows, int cols, PartStr * ↩
        changes, int numChanges)
138 {
139   int i;
140   for (i = 0; i < numChanges; i++)
141     {
142       int idx = changes[i].z * rows * cols + changes[i].y * cols ↩
             + changes[i].x;
143       pic[idx] = -1;
144     }
145 }
```

On the other hand, the contents of `dla_mpi.cpp` file shown below, require only minor changes. These are limited to:

(a) Lines 206-220, where the derived MPI datatype for `PartStr` is defined.

(b) Lines 247-258 that output the resulting structure as a CSV file. This CSV file can be imported in ParaView (http://www.paraview.org/Wiki/ParaView/Data_formats) for viewing purposes.

```
146 /* DLA 3D MPI program
147  * Particles are evenly distributed among nodes/processes
148  * G. Barlas, 2014
149  */
150
151 #include <stdio.h>
152 #include <stdlib.h>
153 #include <assert.h>
154 #include "dla_core.h"
```

```
155    #include <mpi.h>
156
157    #define PIC_SIZE 100
158    #define PARTICLES 1000
159    #define MAX_ITER 10000
160
161    /*————————————————————————————————————*/
162    int main (int argc, char **argv)
163    {
164      int cols, rows, depth, iter, particles, x, y;
165      int *pic;
166      PartStr *p, *changes, *totalChanges;
167      int rank, num, i, numChanges, numTotalChanges;
168      int *changesPerNode, *buffDispl;
169      MPI_Init (&argc, &argv);
170      MPI_Comm_rank (MPI_COMM_WORLD, &rank);
171      MPI_Comm_size (MPI_COMM_WORLD, &num);
172
173      if (argc < 2)                    // use default values if user ↩
                does not specify anything
174        {
175          cols = PIC_SIZE + 2;
176          rows = PIC_SIZE + 2;
177          depth = PIC_SIZE + 2;
178          iter = MAX_ITER;
179          particles = PARTICLES;
180        }
181      else
182        {
183          cols = atoi (argv[1]) + 2;
184          rows = atoi (argv[2]) + 2;
185          depth = atoi (argv[3]) + 2;
186          particles = atoi (argv[4]);
187          iter = atoi (argv[5]);
188        }
189
190      // initialize the random number generator
191      srand(time(0));
192
193      int particlesPerNode = particles / num;
194      if (rank == num - 1)
195        particlesPerNode = particles - particlesPerNode * (num - 1); ↩
                   // in case particles cannot be split evenly
196   // printf("%i has %i\n", rank, particlesPerNode);
197      pic = (int *) malloc (sizeof (int) * cols * rows * depth);
198      p = (PartStr *) malloc (sizeof (PartStr) * particlesPerNode);
199      changes = (PartStr *) malloc (sizeof (PartStr) * ↩
              particlesPerNode);
200      totalChanges = (PartStr *) malloc (sizeof (PartStr) * ↩
              particlesPerNode);
201      changesPerNode = (int *) malloc (sizeof (int) * num);
202      buffDispl = (int *) malloc (sizeof (int) * num);
203      assert (pic != 0 && p != 0 && changes != 0 && totalChanges != 0↩
              && changesPerNode != 0);
204
205      // MPI user type declaration
206      int lengths[3] = { 1, 1, 1 };
207      MPI_Datatype types[3] = { MPI_INT, MPI_INT, MPI_INT };
208      MPI_Aint add1, add2, add3;
209      MPI_Aint displ[3];
210      MPI_Datatype Point;
211
212      MPI_Address (p, &add1);
213      MPI_Address (&(p[0].y), &add2);
214      MPI_Address (&(p[0].z), &add3);
215      displ[0] = 0;
216      displ[1] = add2 - add1;
217      displ[2] = add3 - add1;
218
219      MPI_Type_struct (3, lengths, displ, types, &Point);
220      MPI_Type_commit (&Point);
221
222      dla_init_plist (pic, rows, cols, depth, p, particlesPerNode, 1)↩
              ;
```

```
223    while (−−iter)
224      {
225        dla_evolve_plist (pic, rows, cols, depth, p, &←
                particlesPerNode, changes, &numChanges);
226
227        //exchange information with other nodes
228        MPI_Allgather (&numChanges, 1, MPI_INT, changesPerNode, 1, ←
                MPI_INT, MPI_COMM_WORLD);
229        //calculate offsets
230        numTotalChanges = 0;
231        for (i = 0; i < num; i++)
232          {
233            buffDispl[i] = numTotalChanges;
234            numTotalChanges += changesPerNode[i];
235          }
236
237        if (numTotalChanges > 0)
238          {
239            MPI_Allgatherv (changes, numChanges, Point, ←
                    totalChanges, changesPerNode, buffDispl, Point, ←
                    MPI_COMM_WORLD);
240            apply_changes (pic, rows, cols, totalChanges, ←
                    numTotalChanges);
241          }
242      }
243
244    if (rank == 0)
245      {
246        // save data points as a CSV file that can be imported in ←
                Paraview
247        FILE *f = fopen ("dla.csv", "w+t");
248        fprintf (f, "x coord, y coord, z coord, scalar\n");
249        int idx = 0;
250        for (int z = 0; z < depth; z++)
251          for (int y = 0; y < rows; y++)
252            for (int x = 0; x < cols; x++)
253              {
254                if (pic[idx] == −1)
255                  fprintf (f, "%i, %i, %i, 1\n", x, y, z);
256                idx++;
257              }
258        fclose (f);
259      }
260
261    MPI_Reduce (&particlesPerNode, &particles, 1, MPI_INT, MPI_SUM,←
            0, MPI_COMM_WORLD);
262    if (rank == 0)
263      fprintf (stderr, "Remaining particles %i\n", particles);
264
265    free (pic);
266    free (p);
267    free (changes);
268    free (changesPerNode);
269    free (buffDispl);
270    MPI_Finalize ();
271    return 0;
272  }
```

17. It is not unusual in NoWs setups, to have machines with different capabilities. Create an MPI program that would have each process read from a file a number of integers that is proportional to its CPU speed as indicated by the operating frequency (its "clock"). Use appropriate *filetypes* and *views* to perform the data distribution in a cyclic-block manner.

**Answer**

Assuming that the input file is broken into BLOCKSIZE pieces [1], the solution

---

[1]Actually, it is an implicit requirement that the number of etype elements in the file is a multiple of the filetype which is in our case equal to BLOCKSIZE in length. Otherwise, the "excess" part of the file will not be read by any process.

is assign to each node, a part of a block proportional to its clock speed.

The mathematical formulation is identical to the one used in Exercise 6 and Equations 5.1 and 5.2. Each node $i$ creates a filetype that corresponds to a part of block equal to

$$localSize_i = \lfloor part_i \cdot BLOCKSIZE \rfloor \tag{5.7}$$

in size, that starts at offset

$$localOff_i = \sum_{j=0}^{i-1} \lfloor part_j \cdot BLOCKSIZE \rfloor \tag{5.8}$$

from the beginning of a block.

In order to make sure that the sum of all parts equals `BLOCKSIZE`, the last of $N$ nodes is assigned a part equal to

$$localSize_{N-1} = BLOCKSIZE - \sum_{j=0}^{N-2} \lfloor part_j \cdot BLOCKSIZE \rfloor \tag{5.9}$$

The resulting source code is shown below. Lines 24-31 deal with the local retrieval and global all-to-all distribution of the clock frequencies. The all-to-all scattering operation is required so that the computation of the individual parts can be duplicated in all nodes, without the need to have them communicated by the master node.

Lines 33-57 implements Equations 5.1, 5.2, 5.7 and 5.8. Lines 55, 56, executed solely on the last node of the communicator, deal with the truncations issues depicted by Equation 5.9.

Lines 65 and 66 declare a different local filetype based on the characteristics of each node. This in turn is used to read the data, one filetype at a time, in lines 77-90.

Commented-out lines 67, 81 and 87 represent alternative ways for accomplishing the desired results.

```
1   // Works only if the filesize is a multiple of the calculated ↩
        actualBlock
2   #include <mpi.h>
3   #include<stdio.h>
4   #include<stdlib.h>
5   #include<string.h>
6   #include<unistd.h>
7   #include<vector>
8
9   using namespace std;
10
11  #define BLOCKSIZE 1000
12
13  int main (int argc, char **argv)
14  {
15    int rank, num, i;
16    MPI_Init (&argc, &argv);
17    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
18    MPI_Comm_size (MPI_COMM_WORLD, &num);
19    MPI_Status status;
20
```

```
21    float clockMHz;
22    float *allClock = NULL;
23
24    FILE *fp = popen ("/bin/cat /proc/cpuinfo  | /bin/grep MHz | /←
          usr/bin/head −n 1  | /usr/bin/gawk −F ':' '{print $2}'", "←
          r");
25    fscanf (fp, "%f", &clockMHz);
26    pclose (fp);
27
28    allClock = new float[num];
29
30    // all gather so calculations can be replicated across all ←
          nodes
31    MPI_Allgather (&clockMHz, 1, MPI_FLOAT, allClock, 1, MPI_FLOAT,←
           MPI_COMM_WORLD);
32
33    float sumClocks = allClock[0];
34    for (i = 1; i < num; i++)
35      sumClocks += allClock[i];
36
37    float part[num];
38    for (i = 0; i < num; i++)
39      part[i] = allClock[i] / sumClocks;
40
41    int localSize;
42    int localOff = 0;
43    int actualBlock = 0;
44    for (i = 0; i < rank; i++)
45      {
46        localOff += part[i] * BLOCKSIZE;
47        actualBlock += part[i] * BLOCKSIZE;
48      }
49    localSize = part[rank] * BLOCKSIZE;
50    actualBlock += localSize;
51    for (i = rank + 1; i < num; i++)
52      actualBlock += part[i] * BLOCKSIZE;
53
54    // last node picking up the slack, so that actualBlock == ←
          BLOCKSIZE
55    if(rank == num−1)
56        localSize += ( BLOCKSIZE − actualBlock);
57    actualBlock = BLOCKSIZE;
58
59    MPI_File f;
60    MPI_File_open (MPI_COMM_WORLD, argv[1], MPI_MODE_RDONLY, ←
          MPI_INFO_NULL, &f);
61
62    MPI_Datatype filetype;
63    int starts = 0;
64
65    MPI_Type_create_subarray (1, &actualBlock, &localSize, &starts,←
           MPI_ORDER_C, MPI_INT, &filetype);
66    MPI_Type_commit (&filetype);
67 //   MPI_File_set_view (f, 0, MPI_INT, filetype, (char *)"native←
       ", MPI_INFO_NULL);
68    MPI_File_set_view (f, localOff * sizeof (int), MPI_INT, ←
          filetype, (char *) "native", MPI_INFO_NULL);
69
70    vector < int >data;
71    int temp[BLOCKSIZE];
72
73    MPI_Offset filesize;
74    MPI_File_get_size (f, &filesize);      // get size in bytes
75    filesize /= sizeof (int);     // convert size in number of ←
          items
76    long long int pos = 0;          //localOff;         // initial file←
           position per process
77    while (pos < filesize)
78      {
79        MPI_File_read (f, temp, 1, filetype, &status);
80        int cnt;
81 //       MPI_Get_count (&status, filetype, &cnt);
82        // get the number of data read in etype units
83        MPI_Get_count (&status, MPI_INT, &cnt);
```

```
84
85
86          pos += actualBlock;
87 //          for (int i = 0; i < localSize; i++)
88          for (int i = 0; i < cnt; i++)
89            data.push_back (temp[i]);
90      }
91
92    sleep (rank);
93    cout << rank << " read " << data.size () << " numbers." << endl←
         ;
94    for (int i = 0; i < 30; i++)
95      cout << data[i] << " ";
96    cout << ".... Last one is : " << data[data.size () − 1];
97    cout << endl;
98
99    MPI_Finalize ();
100    return 0;
101 }
```

18. The details of the trapezoidal rule for computing a function integral are discussed in Section 3.5.2. Implement an MPI-based version of the trapezoidal rule using: (a) dynamic partitioning and (b) static partitioning. For part (a) you can base your solution on master-worker implementation of Section 5.22.1.

   **Answer**

   (a) The following listing is a derivation of the master-worker implementation of Section 5.22.1. Given the simple nature of the problem and the associated I/O, the resulting code is much simpler than the original. The only essential customizations are: the adaptation of the derived datatype used to communicate the description of a "job", and the elimination of the "job assignment tracking" feature (supported by the assignedPart array in Listing 5.38), as in this case the results are scalar, floating point numbers that are accumulated by the master node (line 121).

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <iostream>
4 #include <math.h>
5 #include <mpi.h>
6 using namespace std;
7
8 //************************************************************
9 // Communication tags
10 #define NULLRESULTTAG 0
11 #define RESULTTAG     1
12 #define WORKITEMTAG   2
13 #define ENDTAG        3
14
15 //************************************************************
16 struct WorkItem
17 {
18    double start;
19    double end;
20    int divisions;
21 };
22
23 //_____
24 double func (double x)
25 {
26    return fabs (sin (x));
27 }
28
29 //_____
30 double integrCalc (double st, double en, int div)
31 {
```

```cpp
32     //calculate area
33     double localRes = 0;
34     double step = (en - st) / div;
35     double x;
36     x = st;
37     localRes = func (st) + func (en);
38     localRes /= 2;
39     for (int i = 1; i < div; i++)
40       {
41         x += step;
42         localRes += func (x);
43       }
44     localRes *= step;
45     return localRes;
46  }
47
48  //*************************************************************
49  void registerWorkItem (MPI_Datatype * workItemType)
50  {
51     struct WorkItem sample;
52
53     int blklen[3] = { 1, 1, 1 };
54     MPI_Aint displ[3], off, base;
55     MPI_Datatype types[3] = { MPI_DOUBLE, MPI_DOUBLE, MPI_INT };
56
57     displ[0] = 0;
58     MPI_Get_address (&(sample.start), &base);
59     MPI_Get_address (&(sample.end), &off);
60     displ[1] = off - base;
61     MPI_Get_address (&(sample.divisions), &off);
62     displ[2] = off - base;
63
64     MPI_Type_create_struct (3, blklen, displ, types, workItemType);
65     MPI_Type_commit (workItemType);
66  }
67
68  //*************************************************************
69  int main (int argc, char *argv[])
70  {
71     int N, rank;
72     double start_time, end_time;
73     MPI_Status status;
74     MPI_Request request;
75
76     MPI_Init (&argc, &argv);
77     MPI_Comm_size (MPI_COMM_WORLD, &N);
78     MPI_Comm_rank (MPI_COMM_WORLD, &rank);
79
80     MPI_Datatype workItemType;
81     registerWorkItem (&workItemType);
82
83     if (rank == 0)                     // master code
84       {
85         double LOWERLIMIT = 0;
86         double UPPERLIMIT = 10;
87         int NDIV = 1000000;
88         int NJOBS = 10;
89
90         if (argc > 4)
91           {
92             LOWERLIMIT = atof (argv[1]);
93             UPPERLIMIT = atof (argv[2]);
94             NDIV = atoi (argv[3]);
95             NJOBS = atoi (argv[4]);
96           }
97
98         double total = 0, partialRes;
99         int divPerJob = NDIV / NJOBS;
100        WorkItem *w = new WorkItem[NJOBS];
101        double jobStep = (UPPERLIMIT - LOWERLIMIT) / NJOBS;
102        double jobSt = LOWERLIMIT;
103        double jobEnd = jobSt + jobStep;
104        for (int i = 0; i < NJOBS; i++)
105          {
```

```
106              w[i].start = jobSt;
107              w[i].end = jobEnd;
108              w[i].divisions = divPerJob;
109              jobSt = jobEnd;
110              jobEnd += jobStep;
111            }
112
113        // now distribute the work item to the worker nodes
114        for (int i = 0; i < NJOBS; i++)
115          {
116            MPI_Recv (&partialRes, 1, MPI_DOUBLE, MPI_ANY_SOURCE, ←
                    MPI_ANY_TAG, MPI_COMM_WORLD, &status);
117            int workerID = status.MPI_SOURCE;
118            int tag = status.MPI_TAG;
119            MPI_Isend (&(w[i]), 1, workItemType, workerID, ←
                    WORKITEMTAG, MPI_COMM_WORLD, &request);
120            if (tag == RESULTTAG)
121              total += partialRes;
122          }
123
124        // now send termination messages
125        for (int i = 1; i < N; i++)
126          {
127            MPI_Recv (&partialRes, 1, MPI_DOUBLE, MPI_ANY_SOURCE, ←
                    MPI_ANY_TAG, MPI_COMM_WORLD, &status);
128            int workerID = status.MPI_SOURCE;
129            int tag = status.MPI_TAG;
130
131            if (tag == RESULTTAG)
132              total += partialRes;
133            MPI_Isend (NULL, 0, workItemType, workerID, ENDTAG, ←
                    MPI_COMM_WORLD, &request);
134          }
135
136        cout << "Result is " << total << endl;
137        delete [] w;
138      }
139    else                          // worker code
140      {
141        MPI_Send (NULL, 0, MPI_DOUBLE, 0, NULLRESULTTAG, ←
                MPI_COMM_WORLD); // establish communication with ←
                master
142        while (1)
143          {
144            WorkItem w;
145            MPI_Recv (&w, 1, workItemType, 0, MPI_ANY_TAG, ←
                    MPI_COMM_WORLD, &status);        // get a new work ←
                    item
146            int tag = status.MPI_TAG;
147            if (tag == ENDTAG)
148              break;
149
150            double partRes = integrCalc (w.start, w.end, w.←
                    divisions);
151            MPI_Send (&partRes, 1, MPI_DOUBLE, 0, RESULTTAG, ←
                    MPI_COMM_WORLD);       // return the results
152          }
153      }
154
155    MPI_Finalize ();
156    return 0;
157  }
```

(b) The static partitioning is a much simpler proposition. Each node can query the command-line parameters and determine (without receiving any information from the "master") the part of the integration range that is dedicate to it (lines 211, 212). Following the partial integration of line 213, with a reduction (line 216) completes the calculation.

```
158  #include <stdio.h>
159  #include <stdlib.h>
160  #include <iostream>
```

109

```cpp
#include <math.h>
#include <mpi.h>
using namespace std;

//———————————————————————————————————————————————
double func (double x)
{
   return fabs (sin (x));
}

//———————————————————————————————————————————
double integrCalc (double st, double en, int div)
{
   //calculate area
   double localRes = 0;
   double step = (en - st) / div;
   double x;
   x = st;
   localRes = func (st) + func (en);
   localRes /= 2;
   for (int i = 1; i < div; i++)
      {
         x += step;
         localRes += func (x);
      }
   localRes *= step;
   return localRes;
}

//*************************************************************
int main (int argc, char *argv[])
{
   int N, rank;
   double LOWERLIMIT = 0;
   double UPPERLIMIT = 10;
   int NDIV = 1000000;

   MPI_Init (&argc, &argv);
   MPI_Comm_size (MPI_COMM_WORLD, &N);
   MPI_Comm_rank (MPI_COMM_WORLD, &rank);

   if (argc > 3)
      {
         LOWERLIMIT = atof (argv[1]);
         UPPERLIMIT = atof (argv[2]);
         NDIV = atoi (argv[3]);
      }

   int divPerNode = NDIV / N;
   double nodeStep = (UPPERLIMIT - LOWERLIMIT) / N;
   double localSt = LOWERLIMIT + rank * nodeStep;
   double localEnd = localSt + nodeStep;
   double localRes = integrCalc (localSt, localEnd, divPerNode);
   double total;

   MPI_Reduce (&localRes, &total, 1, MPI_DOUBLE, MPI_SUM, 0, ↩
       MPI_COMM_WORLD);
   if (rank == 0)
     cout << "Result is " << total << endl;

   MPI_Finalize ();
   return 0;
}
```

19. Use the master-worker code provided in Section 5.22, as a basis for the creation of a hierarchical master-worker configuration, where nodes are organized in a 3-level (or higher) tree instead of the 2-level tree of the simple setup. The secondary master nodes of the middle tree layer(s) should be responsible for managing the load distribution in their subtrees, while the primary master at the root of the tree should be responsible for

the overall workload distribution.

**Answer**

The calculation of the Mandelbrot fractal is by no means a suitable target application for this kind of setup. So this solution is more of a proof-of-concept, than a practical way to speed-up the particular calculation.

In the listing that follows, nodes form a 3-level hierarchy, with `SECONDARYMASTER` nodes serving in the intermediate layer. The worker nodes calculate the secondary master they are reporting to, with the simple calculation of line 347. The major points can be summarized as follows:

- The master node initially distributes groups of `BATCHSIZE` number of work items (lines 234-242). Subsequently, it sends to each secondary master a single work item for every partial result received (lines 245-255). When all the work items have been distributed, the master sends termination messages for each partial result received (lines 257-266).

- The secondary master nodes maintain a local pool of work items, that is initialized by the bulk work items sent by the primary master (lines 287-289). The work items are sent one-by-one to the workers that communicate with the secondary masters. The image parts/results that are sent back, are forwarded to the primary master (lines 296 and 316).

- The secondary master nodes -expect to- receive a new work item for each result they send back. The new work item is appended to the local pool (line 320), so it is the next one to be sent to a worker. The worker assignment is delayed, until the work item in question is received from the master node (lines 301-307).

  Obviously, a proper, circular, FIFO queue setup would be more beneficial for the handling of each local work item pool. To avoid making an already lengthy code even longer, this is not pursued.

- The worker node code is identical to the one reported in Section 5.22, with the exception of line 347.

```
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <iostream>
4   #include <math.h>
5   #include "mpi.h"
6   #include <QImage>
7   #include <QRgb>
8   using namespace std;
9
10  //*************************************************************
11  // Communication tags
12  #define NULLRESULTTAG 0
13  #define RESULTTAG       1
14  #define WORKITEMTAG     2
15  #define ENDTAG          3
16
17  #define BATCHSIZE 10              /* Jobs are initially send to ←
        secondary master nodes in groups of BATCHSIZE
18                                    * BATCHSIZE has to be greater ←
                                          than 1, otherwise the ←
                                          secondary master nodes
19                                    * terminate prematurely
20                                    */
```

```cpp
#define SECONDARYMASTER 2          /* Number of secondary master ↵
        nodes */
//*************************************************************

typedef struct WorkItem
{
   double upperX, upperY, lowerX, lowerY;
   int pixelsX, pixelsY, imageX, imageY, imgPartID;
} WorkItem;

//*************************************************************
// Class for computing a fractal set part
class MandelCompute
{
private:
   double upperX, upperY, lowerX, lowerY;
   int pixelsX, pixelsY, imgID;
   int *img;

   static int MAXITER;
   int diverge (double cx, double cy);

public:
      MandelCompute ();
   void init (WorkItem * wi);
    ~MandelCompute ();
   int *compute ();
};
int MandelCompute::MAXITER = 255;

//--------------------------------------------------

int MandelCompute::diverge (double cx, double cy)
{
   int iter = 0;
   double vx = cx, vy = cy, tx, ty;
   while (iter < MAXITER && (vx * vx + vy * vy) < 4)
     {
        tx = vx * vx - vy * vy + cx;
        ty = 2 * vx * vy + cy;
        vx = tx;
        vy = ty;
        iter++;
     }
   return iter;
}

//--------------------------------------------------
MandelCompute::~MandelCompute ()
{
   if (img != NULL)
     delete [] img;
}

//--------------------------------------------------
MandelCompute::MandelCompute ()
{
   img = NULL;
}

//--------------------------------------------------

void MandelCompute::init (WorkItem * wi)
{
   upperX = wi->upperX;
   upperY = wi->upperY;
   lowerX = wi->lowerX;
   lowerY = wi->lowerY;

   if (img == NULL || pixelsX != wi->pixelsX || pixelsY != wi->↵
        pixelsY)
     {
        if (img != NULL)
           delete [] img;
```

```
93            // img has an extra element for storing the ID of the ↩
                   picture part computed
94            img = new int[(wi->pixelsX) * (wi->pixelsY) + 1];
95          }
96      pixelsX = wi->pixelsX;
97      pixelsY = wi->pixelsY;
98      imgID = wi->imgPartID;
99    }
100
101   //———————————————————————————————————————————
102
103   int *MandelCompute::compute ()
104   {
105      double stepx = (lowerX - upperX) / pixelsX;
106      double stepy = (upperY - lowerY) / pixelsY;
107
108      for (int i = 0; i < pixelsX; i++)
109        for (int j = 0; j < pixelsY; j++)
110          {
111              double tempx, tempy;
112              tempx = upperX + i * stepx;
113              tempy = upperY - j * stepy;
114              img[j * pixelsX + i] = diverge (tempx, tempy);
115          }
116      img[pixelsX * pixelsY] = imgID;
117      return img;
118   }
119
120   //**************************************************************
121   void registerWorkItem (MPI_Datatype * workItemType)
122   {
123      struct WorkItem sample;
124
125      int blklen[2];
126      MPI_Aint displ[2], off, base;
127      MPI_Datatype types[2];
128
129      blklen[0] = 4;
130      blklen[1] = 5;                    // the part's location is ↩
              communicated regardless
131
132      types[0] = MPI_DOUBLE;
133      types[1] = MPI_INT;
134
135      displ[0] = 0;
136      MPI_Get_address (&(sample.upperX), &base);
137      MPI_Get_address (&(sample.pixelsX), &off);
138      displ[1] = off - base;
139
140      MPI_Type_create_struct (2, blklen, displ, types, workItemType);
141      MPI_Type_commit (workItemType);
142   }
143
144   //**************************************************************
145   // Uses the divergence iterations to pseudocolor the fractal set
146   void savePixels (QImage * img, int *imgPart, int imageX, int ↩
           imageY, int height, int width)
147   {
148      for (int i = 0; i < width; i++)
149        for (int j = 0; j < height; j++)
150          {
151              int color = imgPart[j * width + i];
152              img->setPixel (imageX + i, imageY + j, qRgb (256 - color,↩
                   256 - color, 256 - color));
153          }
154   }
155
156   //**************************************************************
157   int main (int argc, char *argv[])
158   {
159      int N, rank;
160      double start_time, end_time;
161      MPI_Status status;
162      MPI_Request request;
```

```cpp
163
164     start_time = MPI_Wtime ();
165
166     MPI_Init (&argc, &argv);
167     MPI_Comm_size (MPI_COMM_WORLD, &N);
168     MPI_Comm_rank (MPI_COMM_WORLD, &rank);
169
170     if (N < SECONDARYMASTER * 2 + 1)
171       {
172         cerr << argv[0] << " needs at least " << SECONDARYMASTER * ←
                  2 + 1 << " nodes to run\n";
173         MPI_Abort (MPI_COMM_WORLD, 2);
174       }
175
176     MPI_Datatype workItemType;
177     registerWorkItem (&workItemType);
178
179 //———————————————————————————————
180     if (rank == 0)                    // primary master code
181       {
182         if (argc < 6)
183           {
184             cerr << argv[0] << " upperCornerX upperCornerY ←
                      lowerCornerX lowerCornerY workItemPixelsPerSide\n"←
                      ;
185             MPI_Abort (MPI_COMM_WORLD, 1);
186           }
187
188         double upperCornerX, upperCornerY;
189         double lowerCornerX, lowerCornerY;
190         double partXSpan, partYSpan;
191         int workItemPixelsPerSide;
192         int Xparts, Yparts;
193         int imgX = 1024, imgY = 768;
194
195         upperCornerX = atof (argv[1]);
196         upperCornerY = atof (argv[2]);
197         lowerCornerX = atof (argv[3]);
198         lowerCornerY = atof (argv[4]);
199         workItemPixelsPerSide = atoi (argv[5]);
200
201         // make sure that the image size is evenly divided in work ←
                  items
202         Xparts = (int) ceil (imgX * 1.0 / workItemPixelsPerSide);
203         Yparts = (int) ceil (imgY * 1.0 / workItemPixelsPerSide);
204         imgX = Xparts * workItemPixelsPerSide;
205         imgY = Yparts * workItemPixelsPerSide;
206
207         partXSpan = (lowerCornerX - upperCornerX) / Xparts;
208         partYSpan = (upperCornerY - lowerCornerY) / Yparts;
209         QImage *img = new QImage (imgX, imgY, QImage::Format_RGB32)←
                  ;
210
211         // prepare the work items in individual structures
212         WorkItem *w = new WorkItem[Xparts * Yparts];
213         for (int i = 0; i < Xparts; i++)
214           for (int j = 0; j < Yparts; j++)
215             {
216               int idx = j * Xparts + i;
217
218               w[idx].upperX = upperCornerX + i * partXSpan;
219               w[idx].upperY = upperCornerY - j * partYSpan;
220               w[idx].lowerX = upperCornerX + (i + 1) * partXSpan;
221               w[idx].lowerY = upperCornerY - (j + 1) * partYSpan;
222
223               w[idx].imageX = i * workItemPixelsPerSide;
224               w[idx].imageY = j * workItemPixelsPerSide;
225               w[idx].pixelsX = workItemPixelsPerSide;
226               w[idx].pixelsY = workItemPixelsPerSide;
227               w[idx].imgPartID = idx;
228             }
229
230         // now distribute the work item to the worker nodes
231         int *assignedPart = new int[SECONDARYMASTER];      // keep ←
```

```
                        track of how many jobs each secondary master is ←↩
                        assigned
232          int *imgPart = new int[workItemPixelsPerSide * ←↩
                        workItemPixelsPerSide + 1];          // for collecting ←↩
                        results
233          int assignedCounter = 0;
234          for (int i = 1; i <= SECONDARYMASTER; i++)
235            {
236              MPI_Recv (NULL, 0, MPI_INT, MPI_ANY_SOURCE, ←↩
                        NULLRESULTTAG, MPI_COMM_WORLD, &status);
237              int secMasterID = status.MPI_SOURCE;
238              int toBeAssigned = min (BATCHSIZE, Xparts * Yparts − ←↩
                        assignedCounter);
239              assignedPart[secMasterID] = toBeAssigned;
240              MPI_Isend (w + assignedCounter, toBeAssigned, ←↩
                        workItemType, secMasterID, WORKITEMTAG, ←↩
                        MPI_COMM_WORLD, &request);
241              assignedCounter += toBeAssigned;
242            }
243
244          int remainPartToCollect = Xparts * Yparts;
245          while (assignedCounter < Xparts * Yparts)
246            {
247              MPI_Recv (imgPart, workItemPixelsPerSide * ←↩
                        workItemPixelsPerSide + 1, MPI_INT, MPI_ANY_SOURCE←↩
                        , RESULTTAG, MPI_COMM_WORLD, &status);
248              int secMasterID = status.MPI_SOURCE;
249              int widx = imgPart[workItemPixelsPerSide * ←↩
                        workItemPixelsPerSide];     //assignedPart[workerID←↩
                        ];
250              MPI_Isend (&(w[assignedCounter]), 1, workItemType, ←↩
                        secMasterID, WORKITEMTAG, MPI_COMM_WORLD, &request←↩
                        );
251
252              savePixels (img, imgPart, w[widx].imageX, w[widx].←↩
                        imageY, workItemPixelsPerSide, ←↩
                        workItemPixelsPerSide);
253              assignedCounter++;
254              remainPartToCollect −−;
255            }
256
257          while (remainPartToCollect > 0)
258            {
259              MPI_Recv (imgPart, workItemPixelsPerSide * ←↩
                        workItemPixelsPerSide + 1, MPI_INT, MPI_ANY_SOURCE←↩
                        , RESULTTAG, MPI_COMM_WORLD, &status);
260              int secMasterID = status.MPI_SOURCE;
261              int widx = imgPart[workItemPixelsPerSide * ←↩
                        workItemPixelsPerSide];
262              MPI_Isend (NULL, 0, workItemType, secMasterID, ENDTAG, ←↩
                        MPI_COMM_WORLD, &request);
263
264              savePixels (img, imgPart, w[widx].imageX, w[widx].←↩
                        imageY, workItemPixelsPerSide, ←↩
                        workItemPixelsPerSide);
265              remainPartToCollect −−;
266            }
267
268
269          img−>save ("mandel.png", "PNG", 0);          // save the ←↩
                        resulting image
270
271          delete [] w;
272          delete [] assignedPart;
273          delete [] imgPart;
274
275          end_time = MPI_Wtime ();
276          cout << "Total time : " << end_time − start_time << endl;
277        }
278  //−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−
279    else if (rank <= SECONDARYMASTER)      // secondary master code
280      {
281          int workItemPixelsPerSide = atoi (argv[5]);
282          WorkItem *localPool = new WorkItem[BATCHSIZE];
```

```
283         int *imgPart = new int[workItemPixelsPerSide * ←
               workItemPixelsPerSide + 1];          // for collecting ←
               results
284         int poolCount = 0;
285
286         // get batch of load from primary master
287         MPI_Send (NULL, 0, MPI_INT, 0, NULLRESULTTAG, ←
               MPI_COMM_WORLD);     // establish communication with ←
               primary master
288         MPI_Recv (localPool, BATCHSIZE, workItemType, 0, ←
               WORKITEMTAG, MPI_COMM_WORLD, &status);
289         MPI_Get_count (&status, workItemType, &poolCount);        ←
               // get the size of the assignment
290
291         MPI_Request jobReq;
292         bool waitFlag = false;
293         int pendingResults = 0;
294         while (poolCount > 0)
295           {
296             MPI_Recv (imgPart, workItemPixelsPerSide * ←
                   workItemPixelsPerSide + 1, MPI_INT, MPI_ANY_SOURCE←
                   , MPI_ANY_TAG, MPI_COMM_WORLD, &status);
297             int workerID = status.MPI_SOURCE;
298             int tag = status.MPI_TAG;
299
300             // wait for job(s) to be received from the MPI_Irecv ←
                    call following below
301             if (waitFlag)
302               {
303                 MPI_Status st2;
304                 MPI_Wait (&jobReq, &st2);
305                 if (st2.MPI_TAG == WORKITEMTAG)
306                   poolCount++;
307               }
308             poolCount--;
309             pendingResults++;
310             MPI_Isend (localPool + poolCount, 1, workItemType, ←
                   workerID, WORKITEMTAG, MPI_COMM_WORLD, &request);
311             if (tag == RESULTTAG)
312               {
313                 pendingResults--;
314
315                 // forward results to primary master
316                 MPI_Send (imgPart, workItemPixelsPerSide * ←
                       workItemPixelsPerSide + 1, MPI_INT, 0, ←
                       RESULTTAG, MPI_COMM_WORLD);
317
318                 // get new work item after the corresponding memory←
                        is free to accept new data
319                 MPI_Wait (&request, &status);
320                 MPI_Irecv (localPool + poolCount, 1, workItemType, ←
                       0, MPI_ANY_TAG, MPI_COMM_WORLD, &jobReq);
321                 waitFlag = true;
322               }
323           }
324
325         // now send termination messages
326         while (pendingResults > 0)
327           {
328             MPI_Recv (imgPart, workItemPixelsPerSide * ←
                   workItemPixelsPerSide + 1, MPI_INT, MPI_ANY_SOURCE←
                   , MPI_ANY_TAG, MPI_COMM_WORLD, &status);
329             int workerID = status.MPI_SOURCE;
330             int tag = status.MPI_TAG;
331
332             MPI_Isend (NULL, 0, workItemType, workerID, ENDTAG, ←
                   MPI_COMM_WORLD, &request);
333             if (tag == RESULTTAG)
334               {
335                 pendingResults--;
336
337                 // forward results to primary master
338                 MPI_Send (imgPart, workItemPixelsPerSide * ←
                       workItemPixelsPerSide + 1, MPI_INT, 0, ←
```

```
                          RESULTTAG ,  MPI_COMM_WORLD ) ;
339                    }
340                }
341            delete [] imgPart ;
342            delete [] localPool ;
343        }
344  //————————————————————————————
345     else                                  // worker code
346        {
347            int myMasterID = ( rank % SECONDARYMASTER ) + 1;     // which ←
                    master to report to
348            MandelCompute c ;
349            MPI_Send ( NULL ,  0 ,  MPI_INT ,  myMasterID ,  NULLRESULTTAG ,  ←
                    MPI_COMM_WORLD ) ;    // establish communication with ←
                    master
350            while  ( 1 )
351                {
352                    WorkItem w ;
353                    MPI_Recv (&w ,  1 ,  workItemType ,  myMasterID ,  MPI_ANY_TAG , ←
                            MPI_COMM_WORLD ,  &status ) ;     // get a new work ←
                            item
354                    int tag = status . MPI_TAG ;
355                    if  ( tag == ENDTAG )
356                        break ;

358                    c . init  (&w ) ;
359                    int *res = c . compute  () ;
360                    MPI_Send ( res ,  w . pixelsX * w . pixelsY + 1 ,  MPI_INT ,  ←
                            myMasterID ,  RESULTTAG ,  MPI_COMM_WORLD ) ;    // ←
                            return the results
361                }
362        }

364     MPI_Finalize  () ;
365     return  0;
366  }
```

20. Modify the multi-threaded master-worker code of Section 5.22.2 so that
    there are two separate threads in each worker process for communication
    : one for receiving work items and one for sending back results. What are
    the benefits and drawbacks of this arrangement?

    **Answer**

    The benefits are two-fold:

    (a) The decoupled "sender" and "receiver" codes are simpler to write
        and maintain.  The less complicated logic is evident, as a simple
        comparison to Listing 5.40 can attest.

    (b) The two communication operations can proceed at different speeds,
        potentially improving the utilization of the computing resources.

    The modifications required are limited to the part of the `main` function
    that deal with worker communications (lines 391-436 of Listing 5.40). The
    additional `ResultSender` class shown below, is used to spawn a thread
    that handles communication *to* the master node exclusively, leaving the
    main thread to handle communication *from* the master node.

    In terms of data items, the only change affects the `assigned` counter,
    which now needs to be changed from two threads.  Instead of using a
    mutex to control access to it, its type is changed to `QAtomicInt`.

    The following listing contains only the changes/additions to Listing 5.40,
    in order to reduce clutter.

```
1   ...
2   class ResultSender:public QThread
3   {
4   private:
5     QueueMonitor < MandelResult * >*outque;
6     int numWorkerThreads;
7     QAtomicInt *assigned;
8   public:
9     ResultSender (QueueMonitor < MandelResult * >*oq, int nw, ←
           QAtomicInt * as):outque (oq), numWorkerThreads (nw), ←
           assigned (as)  {}
10    void run ();
11  };
12
13  //**********************************************************
14  void ResultSender::run ()
15  {
16    while (numWorkerThreads)
17      {
18        MandelResult *res;
19        if (outque->availItems () > 0)
20          {
21            res = outque->deque ();
22            if (res == NULL)
23              {
24                numWorkerThreads--;
25              }
26            else
27              {
28                MPI_Request r;
29                MPI_Status s;
30
31                MPI_Isend (res->getResultAddress (), res->←
                    getResultSize (), MandelResult::type, 0, ←
                    RESULTTAG, MPI_COMM_WORLD, &r);           // ←
                    return the results
32                MPI_Wait (&r, &s);
33
34                outque->release (res);
35                assigned->fetchAndAddOrdered (-1);
36              }
37          }
38      }
39  }
40
41  //**********************************************************
42  int main (int argc, char *argv[])
43  {
44  ...
45      }
46    else                              // worker code
47      {
48  ...
49        // one loop for sending and recv messages
50        bool endOfWork = false;
51        int numWorkerThreads = numCores;
52        QAtomicInt assigned (0);
53        // spawn the thread for sending results back
54        ResultSender rs (outque, numWorkerThreads, &assigned);
55        rs.start ();
56
57        // job receiving part
58        while (!endOfWork)
59          {
60            MandelWorkItem *w = inque->reserve ();
61            MPI_Recv (w, 1, MandelWorkItem::type, 0, MPI_ANY_TAG, ←
                MPI_COMM_WORLD, &status);        // get a new work ←
                item
62            int tag = status.MPI_TAG;
63            if (tag == ENDTAG)
64              {
65                for (int i = 0; i < numCores; i++)
66                  inque->enque (NULL);
67                endOfWork = true;
```

```
68                }
69              else
70              {
71                 inque->enque (w);
72                 assigned.fetchAndAddOrdered (1);
73              }
74          }
75
76        for (int i = 0; i < numCores; i++)
77          thr[i]->wait ();
78        rs.wait ();
79  ...
```

21. Conway's Game of Life is played on a rectangular grid of cells that may or may not contain an organism. The state of the cells is updated at each time step by applying the following set of rules:

    - Every organism with 2-3 neighbors survives.
    - Every organism with 4 or more neighbors dies from overpopulation.
    - Every organism with 0 - 1 neighbors dies from isolation.
    - Every empty cell adjacent to 3 organisms gives birth to a new one.

    Create an MPI program that evolves a board of arbitrary size (dimensions could be specified at the command-line) over several iterations. The board could be randomly generated or read from a file.

    Try applying the geometric decomposition pattern to partition the work among your processes. One example could be to evenly split the board row-wise. It is clear that each process can update its part of the board only by knowing the state of the bottom board row resident in the previous process, and the state of the top board row resident in the next process (the boundary processes being an exception).

    **Answer**

    The following solution uses geometric decomposition to distribute the board evenly among the MPI processes, in a row-wise fashion. At the end of each time step, the processes communicate their top and bottom rows with their neighbors, to allow proper calculation of neighbor cell counts. The arrangement is similar to the one used in Exercise 15, for the DLA problem. The difference with the depiction shown in Figure 5.2 is that only one instead of two rows is communicated per message.

```
1  // Game of Life
2  // G. Barlas, Dec. 2014
3  #include <stdio.h>
4  #include <stdlib.h>
5  #include <string.h>
6  #include <math.h>
7  #include <mpi.h>
8
9  using namespace std;
10
11 #define PIC_SIZE 100
12 #define PARTICLES 500
13 #define MAX_ITER 10
14
15 #define DUMPIMAGES
16 /*--------------------------------------------------*/
17 int countNeighbors (unsigned char *pic, int cols, int x, int y)
18 {
19   int count = - pic[cols * y + x];
```

```
20     for  ( int  i  = −1;  i  < 2;  i++)
21       for  ( int  j  = −1;  j  < 2;  j++)
22         count  += pic [ cols  *  (y  +  j )  +  (x  +  i ) ] ;
23
24     return  count ;
25   }
26
27   /*−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−*/
28   /* Initializes  the  2D array  for  the  simulation  */
29   void  board_init  ( unsigned  char  *pic ,  int  rows ,  int  cols ,  int  ←
          particles )
30   {
31     int  i ,  j ,  x ,  y ;
32     for  ( i  = 0;  i  < rows ;  i++)
33       for  ( j  = 0;  j  < cols ;  j++)
34         pic [ i  *  cols  +  j ]  =  false ;
35
36
37     for  ( i  = 0;  i  < particles ;  i++)          /* generate  initial  ←
            particle  placement  */
38       {
39         x  = random  ()  % ( cols  −  2)  + 1;     // counting  starts  from  1
40         y  = random  ()  % ( rows  −  2)  + 1;
41         pic [ y  *  cols  +  x ]  =  true ;
42       }
43   }
44
45   /*−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−
46    * Single  step  evolution  of  the  simulation .
47    */
48   void  board_evolve  ( unsigned  char  *pic ,  unsigned  char  *pic2 ,  int  ←
          rows ,  int  cols )
49   {
50     // prepare  array  to  hold  new  state
51     memset  ( pic2 ,  0,  sizeof  ( unsigned  char )  *  rows  *  cols ) ;
52
53     for  ( int  y  = 1;  y  < rows  −  1;  y++)
54       for  ( int  x  = 1;  x  < cols  −  1;  x++)
55         {
56           int  idx  = y  *  cols  +  x ;
57           int  c  = countNeighbors  ( pic ,  cols ,  x ,  y ) ;
58           if  ( pic [ idx ])             // cell  in−place
59             {
60               switch  ( c )
61                 {
62   // Cases  are  commented  out  because  pic2  is  reset  already
63   //           case  0:
64   //           case  1:   // death  −  loniless
65   //             pic2 [ idx ]= false ;
66   //             break ;
67                   case  2:
68                   case  3:              // survival
69                     pic2 [ idx ]  =  true ;
70                     break ;
71                   default :             // death  −  overcrowded
72   //             pic2 [ idx ]= false ;
73                     break ;
74                 }
75             }
76           else  if  ( c  == 3)          // birth  in  empty  cell
77             pic2 [ idx ]  =  true ;
78         }
79   }
80
81   /*−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−*/
82   void  saveImage  ( unsigned  char  *  pic ,  int  rows ,  int  cols ,  char  *←
          fname )
83   {
84     FILE  *f  = fopen  ( fname ,  ”w+t” ) ;
85     fprintf  ( f ,  ”P1\n%i  %i\n” ,  cols ,  rows ) ;
86
87     for  ( int  y  = 1;  y  <= rows ;  y++)
88       {
89         for  ( int  x  = 1;  x  <= cols ;  x++)
```

```
 90              {
 91                if (pic[y * (cols + 2) + x])
 92                  fprintf (f, "1 ");
 93                else
 94                  fprintf (f, "0 ");
 95              }
 96            fprintf (f, "\n");
 97          }
 98      fclose (f);
 99    }
100
101    /*———————————————————————————————————*/
102    int main (int argc, char **argv)
103    {
104      int cols, rows, iter, particles, x, y;
105      unsigned char *pic, *pic2, *tmp;
106      int rank, num, i;
107      MPI_Init (&argc, &argv);
108      MPI_Comm_rank (MPI_COMM_WORLD, &rank);
109      MPI_Comm_size (MPI_COMM_WORLD, &num);
110      MPI_Request req, req2, req3;
111
112      if (argc < 2)        // use default values if user does not ←
              specify anything
113        {
114          cols = PIC_SIZE;
115          rows = PIC_SIZE;
116          iter = MAX_ITER;
117          particles = PARTICLES;
118        }
119      else
120        {
121          cols = atoi (argv[1]);
122          rows = atoi (argv[2]);
123          particles = atoi (argv[3]);
124          iter = atoi (argv[4]);
125        }
126
127      // to make code simpler and avoid special cases, rows is forced←
              to be a multiple of
128      // the number of processes
129      rows = (int) ceil (rows * 1.0 / num) * num;
130      int rowsPerNode = rows / num;
131
132      srand (time (0)); // initialize the random number generator
133
134      // grid has two extra rows, one above and one below the actual ←
              rows
135      // The extra rows "belong" to other nodes, so at the end of ←
              each time step
136      // they should be exchanged
137      if (rank == 0)
138        {
139          pic = new unsigned char[(cols + 2) * (rows + 2)];
140          pic2 = new unsigned char[(cols + 2) * (rows + 2)];
141        }
142      else
143        {
144          pic = new unsigned char[(cols + 2) * (rowsPerNode + 2)];
145          pic2 = new unsigned char[(cols + 2) * (rowsPerNode + 2)];
146        }
147
148      if (rank == 0)
149        board_init (pic2, rows + 2, cols + 2, particles);
150      // clean pic buffer
151      memset (pic, 0, (cols + 2) * (rowsPerNode + 2) * sizeof (←
              unsigned char));
152
153      // boundary rows do not need to be communicated now.
154      MPI_Scatter (pic2 + cols + 2, (cols + 2) * rowsPerNode, ←
              MPI_UNSIGNED_CHAR, pic + cols + 2, (cols + 2) * ←
              rowsPerNode, MPI_UNSIGNED_CHAR, 0, MPI_COMM_WORLD);
155
156      int nextID = rank + 1, prevID = rank - 1;
```

```cpp
157     int step = 0;
158     while (++step < iter)
159       {
160 #ifdef DUMPIMAGES
161         MPI_Gather (pic + cols + 2, rowsPerNode * (cols + 2), ↩
                MPI_UNSIGNED_CHAR, pic + cols + 2, rowsPerNode * (cols↩
                + 2), MPI_UNSIGNED_CHAR, 0, MPI_COMM_WORLD);
162         if (rank == 0)
163           {
164             char fname[20];
165             sprintf (fname, "%05i.pbm", step);
166             saveImage (pic, rows, cols, fname);
167           }
168 #endif
169
170         board_evolve (pic, pic2, rowsPerNode + 2, cols + 2);
171         tmp = pic2;
172         pic2 = pic;
173         pic = tmp;
174         //exchange information with "previous" node
175         if (rank != 0)
176           {
177             MPI_Isend (pic + cols + 2, cols + 2, MPI_UNSIGNED_CHAR,↩
                    prevID, step, MPI_COMM_WORLD, &req);
178             MPI_Irecv (pic, cols + 2, MPI_UNSIGNED_CHAR, prevID, ↩
                    step, MPI_COMM_WORLD, &req2);
179           }
180
181         // exchange with "next"
182         if (rank != num - 1)
183           {
184             MPI_Isend (pic + rowsPerNode * (cols + 2), cols + 2, ↩
                    MPI_UNSIGNED_CHAR, nextID, step, MPI_COMM_WORLD, &↩
                    req);
185             MPI_Irecv (pic + (rowsPerNode + 1) * (cols + 2), cols +↩
                    2, MPI_UNSIGNED_CHAR, nextID, step, ↩
                    MPI_COMM_WORLD, &req3);
186           }
187
188         // block until messages are received
189         if (rank != 0)
190           MPI_Wait(&req2, MPI_STATUS_IGNORE);
191         if (rank != num - 1)
192           MPI_Wait(&req3, MPI_STATUS_IGNORE);
193       }
194
195   MPI_Gather (pic + cols + 2, rowsPerNode * (cols + 2), ↩
          MPI_UNSIGNED_CHAR, pic + cols + 2, rowsPerNode * (cols + ↩
          2), MPI_UNSIGNED_CHAR, 0, MPI_COMM_WORLD);
196   if (rank == 0)
197     saveImage (pic, rows, cols, "out.pbm");
198
199   MPI_Finalize ();
200   delete [] pic;
201   delete [] pic2;
202   return 0;
203 }
```

22. Radix sort is a linear complexity non-comparison-based sorting algorithm that is susceptible to concurrent execution. Radix sort sorts data by separating them into groups based on their digits (for integers) or characters (for strings). The data must be of fixed range, i.e. the number of bits or characters used must be known a priori.

Radix sort comes in two forms: Least-Significant Digit radix sort (LSD) or Most-Significant Digit radix sort (MSD). The latter is suitable for parallel execution, as data which are partitioned in groups can be operated independently in subsequent phases of the algorithm. The MSD algorithm,

which is very close to bucket-sort [2], can be implemented recursively as shown in the following pseudocode, for binary data. An extension for data with non-binary digits (strings) is straightforward. The use of the auxiliary array B allows the sorting to be *stable*.

```
1   // Input: array A, with N elements, each D bits long.
2   // Output : A holds sorted data
3   radix_sort(A, N)
4      allocate memory B equal in size to A
5      tmp <- radix_aux(A, N, B, D-1)
6      if tmp <> A   // if sorted data end up in B array
7         copy B to A
8
9   // Auxiliary recursive function
10  // Use of temporary array B for a stable sort
11  // Returns location of sorted data
12  radix_aux( A, N, B, k)
13     if k = -1  Or N<2     // base case for termination
14         return A
15     else
16        let r be the number of items in A, with k-th bit set to 0
17        resetIdx <- 0
18        setIdx <- r
19        for i <- 0 to N-1     // separate data in two bins
20           if k-th bit of A[i] is set
21               store A[i] in B[ setIdx ]
22               setIdx <- setIdx + 1
23           else
24               store A[i] in B[ resetIdx ]
25               resetIdx <- resetIdx + 1
26
27        tmp1 <- radix_aux(B, r, A, k-1)        // sort bin with a 0 ↵
                  k-th bit, using the (k-1)-th bit
28        tmp2 <- radix_aux(B+r, N-r, A+r, k-1) // sort bin with a 1 ↵
                  k-th bit, using the (k-1)-th bit
29        if tmp1 + r <> tmp2      // pointer comparison
30           if r > N - r    // make the smallest copy
31               copy N-r elements from tmp2 to tmp1+r
32               return tmp1
33           else
34               copy r elements from tmp1 to tmp2-r
35               return tmp2-r
36        else
37            return tmp1
```

Use divide & conquer decomposition as your decomposition pattern (see Section 2.3.2) to design and implement an MPI radix sort.

**Answer**

The algorithm design is influenced by the initial distribution of data, in the same way that bucketsort does. For example, it could be possible that all the nodes get a portion of the initial data, sort it according to one of the bits and exchange the resulting two parts. In this solution we consider the alternative where the data are originally resident at the 0-ranked node, which also performs the first sorting phase according to the most significant bit of the data values.

The following listing implements the divide & conquer decomposition pattern, whereas the data after each pass/phase of the algorithm that splits them into two groups, are split between two nodes. The overall communication pattern resembles a scattering operation. Node 0 first splits the

---

[2]The major difference between radix sort and bucket sort is in how the keys are examined: piece-wise in radix sort and as a whole in bucket sort. In the latter, an arbitrary number of buckets can be prescribed.

input into two pieces according to the value of the most significant bit, keeps one part for itself and sends the other part to node 1. During the second pass/phase, node 0 splits the remaining data and sends one of the parts to node 2, while node 1 performs the same operation for its part of the data and sends one of the parts to node 3. The process is repeated as many times as the number of bits used by the data values, but the migration of data to other nodes is limited by the number of MPI processes (check of line 120).

The key points of the following listing are:

- The data are initialized in node 0, and are sorted by having all the nodes call the `radixsort` function in line 195. The `radixsort` function sets-up the auxiliary memory needed for a stable sort (line 171) and calls the recursive function `radixSort_aux` which performs the actual sorting.

- In `radixSort_aux` each node determines the phase at which it will join the sorting operation (line 84), as well as the node from which it will receive its assignment (line 108), and the nodes to which it will subsequently assign parts of its workload (line 118).

- As with the quicksort algorithm, the split of parts is not even. In the event that a node does not have data to operate on (base case `M<=0` of line 88), it terminates, but not before informing all its sub-ordinate nodes that they should terminate as well (loop of lines 94 to 99).

- The `radixSort_aux` behaves differently during its first invocation, which is flagged by setting the `phase` parameter to -1. Namely, data are received from another node (line 109), and sent back to it following the completion of the "local" sort operation (line 153).

- The auxiliary array pointed to by `aux` is used to hold the data during the split operation carried out by function `radixPhaseK`. Data are shifted between the two memory locations (`data` and `aux`) during the different phases, with `radixSort_aux` maintaining a pointer (`tmp1`) to the actual location. The data should eventually return to their original location. If for any reason (e.g. an odd number of bits is used in the sorting) this is not done, the check in line 174 restores them.

- The `radixSort_aux` function also makes sure that the two data portions produced by the application of `radixPhaseK` in line 117, end up in the same location after they are independently sorted (lines 120-128 for the "left" part and line 131 for the "right" part). The necessary move is minimized by the logic of lines 134-148.

```cpp
1   #include<mpi.h>
2   #include<string.h>
3   #include<stdio.h>
4   #include<stdlib.h>
5   #include<iostream>
6   #include<math.h>
7   #include<unistd.h>
8
9   using namespace std;
10
11  #define ASSIGNTAG 0
```

```c
12  #define SORTEDRESTAG 1
13
14  const int MINV = 0;
15  const int MAXV = 10000;
16
17  //*****************************************
18  // first bit is at position 1
19  int msb (int i)
20  {
21    int res = 0;
22    while (i)
23      {
24        i >>= 1;
25        res++;
26      }
27    return res;
28  }
29
30  //*****************************************
31  // Auxiliary function that performs a phase
32  // of a radixsort algorithm, based on the
33  // k-th bit. Temporary array B is used for
34  // a stable sort.
35  // Returns location the location where the
36  // sorted data "split apart". So B[r] is the
37  // first element with a 1-bit at position k
38  int radixPhaseK (int *A, int N, int *B, int k)
39  {
40    if ((k == -1) || (N < 2))       // base case for termination
41      {
42        memcpy (B, A, N * sizeof (int)); // data are expected to ↩
              move to B
43        return 0;
44      }
45    else
46      {
47        int r = 0;
48        int bitMask = (1 << k);
49        for (int i = 0; i < N; i++)
50          if ((A[i] & bitMask) == 0)
51            r++;
52
53        int resetIdx = 0, setIdx = r;
54        for (int i = 0; i < N; i++)
55          {
56            if ((A[i] & bitMask) > 0)
57              {
58                B[setIdx] = A[i];
59                setIdx++;
60              }
61            else
62              {
63                B[resetIdx] = A[i];
64                resetIdx++;
65              }
66          }
67        return r;
68      }
69  }
70
71  //*****************************************
72  // Recursive radisort function. Based on the
73  // phase, data maybe sent to other nodes for
74  // the subsequent phases.
75  // totalBits counts the bits left for comparison
76  // phase==-1 flags the first call to this function
77  int *radixSort_aux (int *data, int M, int *aux, int totalBits, ↩
        int rank, int num, int phase = -1)
78  {
79    int newPhase;
80    int sourceID;
81    int *tmp1, *tmp2;
82    if (phase == -1)                 // if first call to recursive ↩
            function
```

```
83        {
84          newPhase = msb (rank);
85          totalBits -= newPhase;    // corresponds to when the node ←
                  will join the sorting
86        }
87
88      if (totalBits <= 0 || M <= 0) // base cases, no more sorting to←
            be done
89        {
90          // have to terminate any nodes waiting for data
91          if (phase != -1)
92            {
93              int partnerID = rank ^ (1 << phase);
94              while (partnerID < num)
95                {
96                  MPI_Send (NULL, 0, MPI_INT, partnerID, ASSIGNTAG, ←
                        MPI_COMM_WORLD);
97                  phase++;
98                  partnerID = rank ^ (1 << phase);
99                }
100           }
101         return data;
102       }
103
104     // first get data assignment
105     MPI_Status st;
106     if (phase == -1 && rank != 0)
107       {
108         sourceID = rank ^ (1 << (newPhase - 1));
109         MPI_Recv (data, M, MPI_INT, sourceID, ASSIGNTAG, ←
                MPI_COMM_WORLD, &st);
110         MPI_Get_count (&st, MPI_INT, &M); // get how many where ←
                actually received
111       }
112
113     if (phase != -1)
114       newPhase = phase;
115
116     // now split the data in two groups based on a bit value
117     int split = radixPhaseK (data, M, aux, totalBits - 1);
118     int partnerID = rank ^ (1 << newPhase);
119     //either send one part to another node
120     if (partnerID < num)
121       {
122         MPI_Send (aux, split, MPI_INT, partnerID, ASSIGNTAG, ←
                MPI_COMM_WORLD);
123       }
124     // or sort it locally
125     else
126       {
127         tmp1 = radixSort_aux (aux, split, data, totalBits - 1, rank←
                , num, newPhase + 1);
128       }
129
130
131     tmp2 = radixSort_aux (aux + split, M - split, data + split, ←
            totalBits - 1, rank, num, newPhase + 1);
132
133     // make sure all data end-up at the same place, pointed by tmp1
134     if (partnerID < num)      // no need to move data if the other ←
            part has been sent off
135       tmp1 = tmp2 - split;
136     else if (tmp2 != tmp1 + split)
137       {
138         // minimize movement by moving the smaller part
139         if (split > M / 2)
140           {
141             memcpy (tmp1 + split, tmp2, (M - split) * sizeof (int))←
                    ;
142           }
143         else
144           {
145             memcpy (tmp2 - split, tmp1, split * sizeof (int));
146             tmp1 = tmp2;
```

```cpp
147               }
148           }
149
150       // collect if data were sent away
151       if (partnerID < num)
152         {
153           MPI_Recv (tmp1, split, MPI_INT, partnerID, SORTEDRESTAG, ↩
                  MPI_COMM_WORLD, &st);
154         }
155
156       // before recursion ends, send back to source
157       if (phase == -1 && rank != 0)
158         MPI_Send (tmp1, M, MPI_INT, sourceID, SORTEDRESTAG, ↩
                MPI_COMM_WORLD);
159
160       return tmp1;
161   }
162
163   //*****************************************
164   void radixSort (int *data, int M, int totalBits)
165   {
166       int rank, num;
167       MPI_Comm_rank (MPI_COMM_WORLD, &rank);
168       MPI_Comm_size (MPI_COMM_WORLD, &num);
169
170       int myPhase = msb (rank);
171       int *aux = new int[M];
172
173       int *tmp = radixSort_aux (data, M, aux, totalBits, rank, num);
174       if (tmp != data)
175         memcpy (data, tmp, M * sizeof (int));
176       delete [] aux;
177   }
178
179
180   //*****************************************
181
182   int main (int argc, char **argv)
183   {
184       MPI_Init (&argc, &argv);
185
186       int rank, num, totalBits;
187       MPI_Comm_rank (MPI_COMM_WORLD, &rank);
188       MPI_Comm_size (MPI_COMM_WORLD, &num);
189       MPI_Status status;
190       int M = atoi (argv[1]);          // get size of data to sort
191       int *data;                       // data buffer
192
193       data = new int[M];
194
195       totalBits = ceil (log2 (MAXV));        // total number of phases↩
                  for radix sort
196
197       // initialize data array in rank 0
198       if (rank == 0)
199         {
200           srand (time (0));
201           for (int i = 0; i < M; i++)
202             {
203               data[i] = (rand () % (MAXV - MINV)) + MINV;
204             }
205         }
206
207       radixSort (data, M, totalBits);
208
209       // correctness check
210       if (rank == 0)
211         {
212           for (int i = 1; i < M; i++)
213             if (data[i] < data[i - 1])
214               cout << "ERROR\n";
215         }
216
217       MPI_Finalize ();
```

```
218        delete [] data;
219        return  0;
220    }
```

# Chapter 6

# GPU Programming

## Exercises

1. An array of type `float` elements is to be processed in a one-element-per-thread fashion by a GPU. Suggest an execution configuration for the following scenarios:

   (a) The array is 1-D and of size $N$. The target GPU has 8 SMs, each with 16 SPs.

   (b) The array is 2-D and of size $NxN$. The target GPU has 5 SMs, each with 48 SPs.

   For each of the above scenario, calculate what is the minimum size that $N$ should satisfy to make the GPU computation a desirable alternative to CPU computation.

   **Answer**

   The number of data items should be high enough to provide work for all SPs.

   (a) A total of $8 \cdot 16 = 128$ SPs require at least $N \geq 128$. Ideally, $N$ should be a multiple of 128. An execution configuration can be obtained with the following code:

   ```
   int block = 16 * 4;   // 64 threads, a multiple of 16
   int grid = ( ( N - 1 ) / block ) + 1;
   foo<<< grid, block >>>(...);
   ```

   (b) A total of $5 \cdot 48 = 240$ SPs require at least $N^2 \geq 240 \Rightarrow N \geq 15.5$. Ideally, $N^2$ should be a multiple of 240. An execution configuration can be obtained with the following code:

   ```
   dim3 block (16, 12); // 192 threads, a multiple of 48
   dim3 grid ( ( N - 1 ) / 16 + 1,  ( N - 1 ) / 12 + 1 );
   foo<<< grid, block >>>(...);
   ```

2. A reduction is an operation frequently encountered in a many algorithms: summing-up the elements of an array, finding the minimum, maximum, etc. One possible solution to a CUDA kernel that calculates the sum of an array would be:

```
__global__ void sum(float *in, float *out)
{
    __shared__ float localStore[];   // to speedup data access

    int globalID = threadIdx.x + blockIdx.x * blockDim.x;
    int localID = threadIdx.x;

    localStore[localID] = in[globalID];    // copy to shared memory
    for(int i=1; i< blockDim.x ; i*=2)
      {
          if(localID % (2*i) == 0)
            localStore[localID] += localStore[localID + i];
          __syncthreads ();
      }
    if(localID == 0)
       out[blockIdx.x] = localStore[0];
}
```

The above needs to be called multiple times, each time reducing the size of
the data by the number of threads in a block. Data at each kernel launch
should be multiples of the block size.

Analyze the above code in relation to the following criteria and how they
reflect on the execution speed: thread divergence, memory coalescing,
use of SPs within a warp. Suggest modifications that would improve the
performance.

**Answer**

The shown kernel scores low in all the mentioned areas, i.e. thread di-
vergence, memory coalescing, and use of SPs within a warp. Analytically,
the issues can be identified as follows:

- Thread divergence : The if condition inside the `for`-loop causes
  threads to diverge. During the first iteration only the threads with
  an ID which is a multiple of 2 can execute the `if` code body. During
  the second iteration only the ones with an ID which is a multiple of
  4 can do this, and so on.

- Memory coalescing : There is no regard to memory coalescing inside
  the `for` loop. In fact, as the loop progresses, the number of bank
  conflicts increases.

- SP utilization : As soon as the `for` loop starts executing, an exponen-
  tially growing number of SPs with each iteration, become inactive,
  as they do not satisfy the `if` condition.

The following program can produce the same result with only one ker-
nel invocation, with maximum SP utilization and minimum thread diver-
gence:

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <cuda.h>
4
5  #define MAXVALUE 10000
6
7  //————————————————————————————————————
8  void numberGen (int N, int max, int *store)
9  {
10   int i;
11   srand (time (0));
12   for (i = 0; i < N; i++)
13     rand () % max;
```

```cpp
14  }
15
16  //—————————————————————————————————
17  __global__ void sum (int *d, int N, int *odds)
18  {
19    extern __shared__ int count [];
20
21    int myID = blockIdx.x * blockDim.x + threadIdx.x;
22    int numThr = blockDim.x * gridDim.x;
23    int localID = threadIdx.x;
24    count[localID] = 0;
25    if (myID < N)
26      {
27        count[localID] = d[myID];
28        int idx = myID + numThr;
29        // partial sum
30        while (idx < N)
31        {
32          count[localID] += d[idx];
33          idx += numThr;
34        }
35
36        __syncthreads ();
37        // block−wide partial sum calculation
38        int step = 1;
39        while (((localID | step) < blockDim.x) && ((localID & step)↩
              == 0))
40          {
41            count[localID] += count[localID | step];
42            step *= 2;
43            __syncthreads ();
44          }
45
46        // now the block−wide partial sum is in count[0]
47        // add to global counter
48        if (localID == 0)
49            atomicAdd (odds, count[0]);
50      }
51  }
52
53  //—————————————————————————————————
54  int main (int argc, char **argv)
55  {
56    int N = atoi (argv[1]);
57
58    int *ha, *hres, *da, *dres;   // host (h*) and device (d*) ↩
          pointers
59
60    ha = new int [N];
61    hres = new int [1];
62
63    cudaMalloc ((void **) &da, sizeof (int) * N);
64    cudaMalloc ((void **) &dres, sizeof (int) * 1);
65
66    numberGen (N, MAXVALUE, ha);
67
68    cudaMemcpy (da, ha, sizeof (int) * N, cudaMemcpyHostToDevice);
69    cudaMemset (dres, 0, sizeof (int));
70
71    int blockSize, gridSize;
72    blockSize = 256;
73    gridSize = 16;
74    sum <<< gridSize, blockSize, blockSize * sizeof (int) >>> (da, ↩
          N, dres);
75
76    cudaMemcpy (hres, dres, sizeof (int), cudaMemcpyDeviceToHost);
77
78    printf ("%i \n", *hres);
79
80    cudaFree ((void *) da);
81    cudaFree ((void *) dres);
82    delete [] ha;
83    delete [] hres;
84    cudaDeviceReset ();
```

```
85
86      return  0;
87    }
```

The kernel of lines 17-51 allocates shared memory in the form of one integer per block thread. This memory is dynamically allocated during the kernel invocation (line 74). Each thread initially spends its time calculating a partial sum by going over the input data with a stride equal to the total number of threads (`while` loop of lines 30-34). Subsequently, a block-wide partial sum is calculated by reducing the shared memory counters with the loop of lines 39-44. The resulting sum in `count[0]` is added to the global memory counter with an atomic operation.

3. The reduction operation discussed in the previous exercise, is a special case of the "scan" or prefix-sum operation that can be applied to the elements of a vector or list. In general, the operator applied can be any of summation, subtraction, minimum, maximum, etc.. Implement a CUDA kernel capable of performing a prefix-sum operation. The "Prefix-Sums and Their Applications" paper by Guy Blelloch, available at `http://www.cs.cmu.edu/~guyb/papers/Ble93.pdf`, is a wonderful resource for learning more on the topic.

**Answer**

The solution is an extension of the one given in the previous exercise. The only differences lie in how the initial partial result for each thread is initialized (lines 34-45), in the operator used to perform the reduction (lines 56-67 and 78-89), and in the atomic operation used to consolidate the block-wide partial results with the global memory result (lines 98-109).

This program can be easily extended to handle any prefix-sum operation by introducing another symbolic constant representing it (following the ones defined in lines 8-10), and appropriately augmenting the `switch` statements.

```
1    #include <stdio.h>
2    #include <stdlib.h>
3    #include <limits.h>
4    #include <cuda.h>
5
6    #define MAXVALUE 10000
7
8    #define SUM_OP  0
9    #define MIN_OP  1
10   #define MAX_OP  2
11
12   #define min(A,B)  ((A>B)  ?  B  :  A)
13   #define max(A,B)  ((A<B)  ?  B  :  A)
14   //————————————————————————————————————
15   void numberGen (int N, int max, int *store)
16   {
17     int  i;
18     srand (time (0));
19     for (i = 0; i < N; i++)
20       rand () % max;
21   }
22
23   //————————————————————————————————————
24
25   __global__ void reduce (int *d, int N, int *res, int OP)
26   {
27     extern __shared__ int partRes [];
28
```

```c
29    int myID = blockIdx.x * blockDim.x + threadIdx.x;
30    int numThr = blockDim.x * gridDim.x;
31    int localID = threadIdx.x;
32
33    // initialize all the threads' storage regardless
34    switch (OP)
35      {
36      case SUM_OP:
37        partRes[localID] = 0;
38        break;
39      case MIN_OP:
40        partRes[localID] = INT_MAX;
41        break;
42      case MAX_OP:
43        partRes[localID] = INT_MIN;
44        break;
45      };
46
47    if (myID < N)
48      {
49        partRes[localID] = d[myID];
50        int idx = myID + numThr;
51
52        // partial result
53        while (idx < N)
54          {
55            int tmp = d[idx];
56            switch (OP)
57              {
58              case SUM_OP:
59                partRes[localID] += tmp;
60                break;
61              case MIN_OP:
62                partRes[localID] = min (partRes[localID], tmp);
63                break;
64              case MAX_OP:
65                partRes[localID] = max (partRes[localID], tmp);
66                break;
67              };
68            idx += numThr;
69          }
70
71        __syncthreads ();
72        // block-wide partial result calculation
73        int step = 1;
74        int otherIdx = localID | step;
75        while ((otherIdx < blockDim.x) && ((localID & step) == 0))
76          {
77            int tmp = partRes[otherIdx];
78            switch (OP)
79              {
80              case SUM_OP:
81                partRes[localID] += tmp;
82                break;
83              case MIN_OP:
84                partRes[localID] = min (partRes[localID], tmp);
85                break;
86              case MAX_OP:
87                partRes[localID] = max (partRes[localID], tmp);
88                break;
89              };
90            step <<= 1;
91            otherIdx = localID | step;
92            __syncthreads ();
93          }
94
95        // now the block-wide partial sum is in partRes[0]
96        // add to global counter
97        if (localID == 0)
98          switch (OP)
99            {
100            case SUM_OP:
101              atomicAdd (res, partRes[0]);
102              break;
```

```
103              case MIN_OP:
104                atomicMin (res, partRes[0]);
105                break;
106              case MAX_OP:
107                atomicMax (res, partRes[0]);
108                break;
109              };
110      }
111  }
112
113  //——————————————————————————————
114  int sharedSize (int b)
115  {
116    return b * sizeof (int);
117  }
118
119  //——————————————————————————————
120
121  int main (int argc, char **argv)
122  {
123    int N = atoi (argv[1]);
124
125    int *ha, *hres, *da, *dres;   // host (h*) and device (d*) ↩
            pointers
126
127    ha = new int[N];
128    hres = new int[1];
129
130    cudaMalloc ((void **) &da, sizeof (int) * N);
131    cudaMalloc ((void **) &dres, sizeof (int) * 1);
132
133    numberGen (N, MAXVALUE, ha);
134
135    cudaMemcpy (da, ha, sizeof (int) * N, cudaMemcpyHostToDevice);
136    cudaMemset (dres, 0, sizeof (int));
137
138    int blockSize, gridSize;
139    blockSize = 256;
140    gridSize = 16;
141    reduce <<< gridSize, blockSize, blockSize * sizeof (int) >>> (↩
            da, N, dres, MAX_OP);
142
143    cudaMemcpy (hres, dres, sizeof (int), cudaMemcpyDeviceToHost);
144
145    printf ("%i \n", *hres);
146
147    cudaFree ((void *) da);
148    cudaFree ((void *) dres);
149    delete [] ha;
150    delete [] hres;
151    cudaDeviceReset ();
152
153    return 0;
154  }
```

The code shown here is not nearly as optimum as it could be; it is a trade-off between readability and performance. A more thorough derivation of an optimum CUDA reduction kernel is covered in http://developer.download.nvidia.com/assets/cuda/files/reduction.pdf, although the material in question targets a Compute Capability 1.0 device, which exhibits many more quirks than contemporary architectures.

4. Create CUDA implementations of the gsl_stats_mean() and gsl_stats_variance() functions offered by the GNU Scientific Library, that produce the mean and variance statistics of an array of type double data. Their signatures are:

```
double gsl_stats_mean (const double DATA[], // Pointer to input ↩
    data
```

135

```
                    size_t STRIDE,          // Step used to read ↵
                        the input. Normally this should be ↵
                        set to 1.
                    size_t N);              // Size of DATA array
double gsl_stats_variance (const double DATA[], // Same as above.
                        size_t STRIDE,
                        size_t N);
```

Assuming that the STRIDE is 1, create a memory access pattern that utilizes coalescing. Suggest ways to deal with the problem if the stride is not 1.

**Answer**

The `gsl_stats_mean` and `gsl_stats_variance` functions are implemented as host front-ends, to allow for the hiding of the implementation details (such as memory management and grid/block design) from the application that calls them.

The solution is based on the reduction approach that is used in the two previous exercises. Two kernels are employed for calculating the sum of values and the sum of the square of the values respectively, in order for the variance to be calculated from the formula:

$$var = \frac{\sum_{i=0}^{N-1} v_i^2}{N} - \left(\frac{\sum_{i=0}^{N-1} v_i}{N}\right)^2 \tag{6.1}$$

where $N$ is the number of values.

In order to deal with a STRIDE parameter which is different from 1, the `gsl_stats_mean` and `gsl_stats_variance` front-end functions, copy the data to be actually processed in another array (lines 127-130 and 169-172). The benefits are two-fold:

- A considerably smaller volume of data cross the PCIe bus, significantly reducing the communication cost.
- The data items reside in contiguous memory locations allowing for memory coalescing to take place, at least to the extent allowed by the `double` data type.

The last major point of the following listing is the `atomicAdd` device function of lines 22-34, that is an operation for `double`-type operands that is missing from the arsenal of CUDA primitives.

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <cuda.h>
4
5  #define MAXVALUE 10000
6  #define min(A,B) ((A>B) ? B : A)
7  //————————————————————————————————
8  void numberGen (int N, int max, double *store)
9  {
10    int i;
11    srand (time (0));
12    for (i = 0; i < N; i++)
13      store[i] = i;              //rand () % max;
14  }
15
16  //————————————————————————————————
17  int sharedSize (int b)
```

```
18  {
19     return b * sizeof (double);
20  }
21  //———————————————————————————————
22  __device__ double atomicAdd (double* address, double val)
23  {
24      unsigned long long int* address_as_ull =
25                                      (unsigned long long int↩
                                            *)address;
26      unsigned long long int old = *address_as_ull, assumed;
27      do {
28          assumed = old;
29          old = atomicCAS (address_as_ull, assumed,
30                          __double_as_longlong (val +
31                          __longlong_as_double (assumed)));
32      } while (assumed != old);
33      return __longlong_as_double (old);
34  }
35  //———————————————————————————————

37  __global__ void sum (double *d, int N, double *totalSum)
38  {
39     extern __shared__ double partialSum [];

41     int myID = blockIdx.x * blockDim.x + threadIdx.x;
42     int numThr = blockDim.x * gridDim.x;
43     int localID = threadIdx.x;
44     partialSum [localID] = 0;
45     if (myID < N)
46       {
47         partialSum [localID] = d[myID];
48         int idx = myID + numThr;
49         // partial sum
50         while (idx < N)
51           {
52             partialSum [localID] += d[idx];
53             idx += numThr;
54           }

56         __syncthreads ();

58         // reduction
59         int step = 1;
60         int otherIdx = localID | step;
61         while ((otherIdx < blockDim.x) && ((localID & step) == 0))
62           {
63             partialSum [localID] += partialSum [otherIdx];
64             step <<= 1;
65             otherIdx = localID | step;
66             __syncthreads ();
67           }

69         // now the block−wide partial sum is in partialSum [0]
70         // add to global counter
71         if (localID == 0)
72           atomicAdd (totalSum, partialSum [0]);
73       }
74  }

76  //———————————————————————————————

78  __global__ void sum2 (double *d, int N, double *totalSum)
79  {
80     extern __shared__ double partialSum [];

82     int myID = blockIdx.x * blockDim.x + threadIdx.x;
83     int numThr = blockDim.x * gridDim.x;
84     int localID = threadIdx.x;
85     partialSum [localID] = 0;
86     if (myID < N)
87       {
88         partialSum [localID] = d[myID]*d[myID];
89         int idx = myID + numThr;
90         // partial sum
```

```
 91          while  ( idx  <  N )
 92            {
 93               partialSum [ localID ]  +=  d [ idx ] * d [ idx ] ;
 94               idx  +=  numThr ;
 95            }
 96
 97          __syncthreads  ( ) ;
 98
 99          // reduction
100          int  step  =  1 ;
101          int  otherIdx  =  localID  |  step ;
102          while  (( otherIdx  <  blockDim . x )  &&  (( localID  &  step )  ==  0 ) )
103            {
104               partialSum [ localID ]  +=  partialSum [ otherIdx ] ;
105               step  <<=  1 ;
106               otherIdx  =  localID  |  step ;
107               __syncthreads  ( ) ;
108            }
109
110          // now  the  block−wide  partial  sum  is  in  partialSum [ 0 ]
111          // add  to  global  counter
112          if  ( localID  ==  0 )
113            atomicAdd  ( totalSum ,  partialSum [ 0 ] ) ;
114       }
115  }
116
117  //−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−
118
119  double  gsl_stats_mean  ( const  double  DATA [ ] ,  size_t  STRIDE ,  size_t↩
          N )
120  {
121     double  *ha ,  *da ,  *dres ;    // host  ( h∗)  and  device  ( d∗)  pointers
122     double  hres ;
123
124     int  Neffective ;
125     if ( STRIDE  !=1 )
126     {
127        Neffective  =  ( N−1 ) / STRIDE  +  1 ;
128        ha  =  new  double [ Neffective ] ;
129        for ( int  i=0 ,  j=0 ; i<N ; i+=STRIDE ,  j++)
130          ha [ j ]  =  DATA [ i ] ;
131     }
132     else
133     {
134       ha  =  ( double  ∗) DATA ;
135       Neffective  =  N ;
136     }
137
138     cudaMalloc  (( void  ∗∗)  &da ,  sizeof  ( double )  ∗  Neffective ) ;
139     cudaMalloc  (( void  ∗∗)  &dres ,  sizeof  ( double )  ∗  1 ) ;
140
141     cudaMemcpy  ( da ,  ha ,  sizeof  ( double )  ∗  Neffective ,  ↩
          cudaMemcpyHostToDevice ) ;
142     cudaMemset  ( dres ,  0 ,  sizeof  ( double ) ) ;
143
144     int  blockSize ,  gridSize ;
145     cudaOccupancyMaxPotentialBlockSizeVariableSMem  (& gridSize ,  &↩
          blockSize ,  ( void  ∗)  sum ,  sharedSize ,  Neffective ) ;
146     gridSize  =  min ( gridSize ,  Neffective / blockSize ) ;
147
148     sum  <<<  gridSize ,  blockSize ,  blockSize  ∗  sizeof  ( double )  >>>  (↩
          da ,  Neffective ,  dres ) ;
149
150     cudaMemcpy  (& hres ,  dres ,  sizeof  ( double ) ,  ↩
          cudaMemcpyDeviceToHost ) ;
151
152     if ( STRIDE  !=1 )
153        delete [ ] ha ;
154
155     cudaFree  (( void  ∗)  da ) ;
156     cudaFree  (( void  ∗)  dres ) ;
157     return  hres / Neffective ;
158  }
159  //−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−−
```

```cpp
160
161   double gsl_stats_variance (const double DATA[], size_t STRIDE, ←↩
          size_t N)
162   {
163     double *ha, *da, *dres;    // host (h*) and device (d*) pointers
164     double hsum, hsum2;
165
166     int Neffective;
167     if(STRIDE !=1)
168     {
169       Neffective = (N−1)/STRIDE + 1;
170       ha = new double[Neffective];
171       for(int i=0, j=0;i<N;i+=STRIDE, j++)
172         ha[j] = DATA[i];
173     }
174     else
175     {
176       ha = (double *)DATA;
177       Neffective = N;
178     }
179
180     cudaMalloc ((void **) &da, sizeof (double) * Neffective);
181     cudaMalloc ((void **) &dres, sizeof (double) * 1);
182
183     cudaMemcpy (da, ha, sizeof (double) * Neffective, ←↩
          cudaMemcpyHostToDevice);
184     cudaMemset (dres, 0, sizeof (double));
185
186     // first calculate the sum of values
187     int blockSize, gridSize;
188     cudaOccupancyMaxPotentialBlockSizeVariableSMem (&gridSize, &←↩
          blockSize, (void *) sum, sharedSize, Neffective);
189
190     gridSize = min(gridSize, Neffective/blockSize);
191     sum <<< gridSize, blockSize, blockSize * sizeof (double) >>> (←↩
          da, Neffective, dres);
192
193     cudaMemcpy (&hsum, dres, sizeof (double), ←↩
          cudaMemcpyDeviceToHost);
194
195     // then calculate the sum of squares
196     cudaOccupancyMaxPotentialBlockSizeVariableSMem (&gridSize, &←↩
          blockSize, (void *) sum2, sharedSize, Neffective);
197     gridSize = min(gridSize, Neffective/blockSize);
198
199     // reset the result variable
200     cudaMemset (dres, 0, sizeof (double));
201     sum2 <<< gridSize, blockSize, blockSize * sizeof (double) >>> (←↩
          da, Neffective, dres);
202
203     cudaMemcpy (&hsum2, dres, sizeof (double), ←↩
          cudaMemcpyDeviceToHost);
204
205     if(STRIDE !=1)
206       delete [] ha;
207
208     cudaFree ((void *) da);
209     cudaFree ((void *) dres);
210     double avg = hsum/Neffective;
211     return hsum2/Neffective − avg*avg;
212   }
213
214   //———————————————————————————————————————
215
216   int main (int argc, char **argv)
217   {
218     int N = atoi (argv[1]);
219
220     double *data;
221
222     data = new double[N];
223
224     numberGen (N, MAXVALUE, data);
225
```

```
226    printf ("Avg : %lf \n",gsl_stats_mean(data, 1, N));
227    printf ("Var : %lf \n",gsl_stats_variance(data, 1, N) );
228
229    delete [] data;
230    cudaDeviceReset ();
231
232    return 0;
233  }
```

5. Design and implement a CUDA program for calculating the histogram of a 24-bit color image. In this case, three separate histograms will be produced, one for each color component of the image.

   **Answer**

   The histogram-calculating kernels of Section 6.7.3 can be used as a basis for the solution. The only thing that needs to be explicitly addressed, is the handling of the three color components. Sending the raw RGB data to the GPU could result in having to launch just one kernel only. However, the memory layout in that case would be equivalent to an array of structures, which as discussed in Section 6.7.4 should be avoided.

   The alternative, is to rearrange the pixel data in color planes, as shown in lines 84-110, packing four pixel values per integer (line 106). This arrangement also permits the use of the histogram kernel presented in Listing 6.22, without any modifications.

```
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include <unistd.h>
4   #include <math.h>
5   #include <cuda.h>
6   #include <QImage>
7   #include <QRgb>
8
9   const int BINS = 256;
10  const int BINS4ALL = BINS * 32;
11
12  //←
        ********************************************************************←

13  __global__ void GPU_histogram_atomic (int *in, int N, int *h)
14  {
15    int gloID = blockIdx.x * blockDim.x + threadIdx.x;
16    int locID = threadIdx.x;
17    int GRIDSIZE = gridDim.x * blockDim.x;
18    __shared__ int localH[BINS4ALL];
19    int bankID = locID % warpSize;
20    int i;
21
22    // initialize the local, shared-memory bins
23    for (i = locID; i < BINS4ALL; i += blockDim.x)
24      localH[i] = 0;
25
26    // wait for all warps to complete the previous step
27    __syncthreads ();
28
29    //start processing the image data
30    int *mySharedBank = localH + bankID;
31    if (blockDim.x > warpSize) // if the blocksize exceeds the ←
          warpSize, it is possible multiple warps run at the same ←
          time
32      for (i = gloID; i < N; i += GRIDSIZE)
33        {
34
35          int temp = in[i];
36          int v = temp & 0xFF;
37          int v2 = (temp >> 8) & 0xFF;
```

```cpp
38              int v3 = (temp >> 16) & 0xFF;
39              int v4 = (temp >> 24) & 0xFF;
40              atomicAdd (mySharedBank + (v << 5), 1);
41              atomicAdd (mySharedBank + (v2 << 5), 1);
42              atomicAdd (mySharedBank + (v3 << 5), 1);
43              atomicAdd (mySharedBank + (v4 << 5), 1);
44          }
45     else
46        for (i = gloID; i < N; i += GRIDSIZE)
47          {
48
49              int temp = in[i];
50              int v = temp & 0xFF;
51              int v2 = (temp >> 8) & 0xFF;
52              int v3 = (temp >> 16) & 0xFF;
53              int v4 = (temp >> 24) & 0xFF;
54              mySharedBank[v << 5]++;  // Optimized version of localH[←
                     bankID + v * warpSize]++
55              mySharedBank[v2 << 5]++;
56              mySharedBank[v3 << 5]++;
57              mySharedBank[v4 << 5]++;
58          }
59
60     // wait for all warps to complete the local calculations, ←
            before updating the global counts
61     __syncthreads ();
62
63     // use atomic operations to add the local findings to the ←
            global memory bins
64     for (i = locID; i < BINS4ALL; i += blockDim.x)
65        atomicAdd (h + (i >> 5), localH[i]); // Optimized version of ←
               atomicAdd (h + (i/warpSize), localH[i]);
66  }
67
68  //←
        ********************************************************************←

69  void histogramRGB_FE (QImage &img, int *hist)
70  {
71      unsigned int *d_in, *h_in;
72      int *d_hist;
73      int N, Nx, Ny;
74
75      Nx = img.width ();
76      Ny = img.height ();
77      N = ceil ((Nx * Ny) / 4.0);
78
79      h_in = new unsigned int[N];
80      cudaMalloc ((void **) &d_in, sizeof (int) * N);
81      cudaMalloc ((void **) &d_hist, sizeof (int) * BINS);
82
83      int gridSize=16, blockSize=256;
84      for(int color=0;color<3;color++)
85      {
86          int idx=0;
87          for(int y=0;y<Ny;y++)
88            for(int x=0;x<Nx;x+=4)
89            {
90                unsigned int fourPixels=0;
91                for(int i=0;i<4;i++)
92                {
93                    unsigned int tmp;
94                    switch(color)
95                    {
96                    case 0:
97                        tmp = qRed(img.pixel(x+i,y));
98                        break;
99                    case 1:
100                       tmp = qGreen(img.pixel(x+i,y));
101                       break;
102                   default:
103                       tmp = qBlue(img.pixel(x+i,y));
104                       break;
105                   }
```

```
106                    fourPixels = (fourPixels <<8) | tmp;
107                }
108             h_in[idx] = fourPixels;
109             idx++;
110           }
111
112        cudaMemcpy (d_in, h_in, sizeof (int) * N, ←
                   cudaMemcpyHostToDevice);
113        cudaMemset (d_hist, 0, BINS * sizeof (int));
114
115        GPU_histogram_atomic <<< gridSize, blockSize >>> ((int *)←
                   d_in, N, d_hist);
116
117        cudaMemcpy (hist+color*BINS, d_hist, sizeof (int) * BINS,←
                   cudaMemcpyDeviceToHost);
118     }
119
120    cudaFree ((void *) d_in);
121    cudaFree ((void *) d_hist);
122    delete [] h_in;
123 }
124
125 //←
       ******************************************************************←
126 int main (int argc, char **argv)
127 {
128
129    QImage pic;
130    pic.load(argv[1]);
131
132    int *hist = new int[3*BINS];
133    histogramRGB_FE(pic, hist);
134
135    for(int i=0;i<3;i++)
136    {
137        for(int j=0;j<10;j++)
138            printf("%i ", hist[BINS*i+j]);
139        printf(" ...\n");
140    }
141    cudaDeviceReset ();
142    delete [] hist;
143    return 0;
144 }
```

The above program, which relies on Qt for image I/O, can be compiled with:

```
$ nvcc -Xcompiler -fPIC -m64 -O2 -D_REENTRANT \
    -DQT_NO_DEBUG -DQT_GUI_LIB -DQT_CORE_LIB
    -I/usr/lib/x86_64-linux-gnu/qt5/mkspecs/linux-g++-64
    -I.
    -I/usr/include/qt5
    -I/usr/include/qt5/QtGui
    -I/usr/include/qt5/QtCore
    -lQt5Core -lQt5Gui
    -arch=sm_21
    -o histogram_atomic histogram_atomic.cu
```

6. Create a variation of the program in Listing 6.24 to discover and plot the memory copy speed for host-to-device and device-to-host operations and for all the available types of host memory allocations : pageable, pinned and mapped. In order to test the last type of allocation, you will have to call a kernel that will try to access the memory, triggering the transfer. Compare your results with the ones returned by the $CUDA/samples/1_Utilities/bandwidthTest sample program.

**Answer**

This is an exercise that reveals a great deal about the underpinnings of the zero-copy memory mechanism. The program shown below transfers increasing amounts of data from the host to the device, using pageable, pinned and mapped memory allocation. In the last case, a dummy kernel is invoked to trigger the transfer of memory. The program parameters are:

(a) The number of iterations for calculating averages.

(b) The increment (in bytes) to use in testing increasing data chunks.

(c) The maximum data size (in bytes) to test for.

The maximum device memory is allocated in one go, as dictated by the third command-line parameter, in lines 59-66. The for loop of lines 68 to 72, calls the `timeCopyH2D` function once for every type of device memory and for increasing data sizes.

The `timeCopyH2D` function uses a stream and events in the same fashion as Listing 6.24, to calculate the average time taken to transfer data from the host to the device. When a zero-copy memory arrangement is used, the `dummyKernel` function is also invoked (line 33) in order to trigger the memory transfer instead of the explicit call of line 36.

```
1   #include <stdio.h>
2   #include <cuda.h>
3
4   //————————————————————————————————————————————————
5   // a dummy kernel just for initiating memory transfers when
6   // mapped memory is used
7   __global__ void dummyKernel (unsigned char *b, int N)
8   {
9     int idx = blockIdx.x * blockDim.x + threadIdx.x;
10    if (idx < N)
11      b[idx] = idx % 256;
12  }
13
14  //————————————————————————————————————————————————
15  float timeCopyH2D (unsigned char *h, unsigned char *d, int N, int↩
          iter, bool ZEROCOPY = false)
16  {
17    cudaEvent_t startT, endT;
18    cudaStream_t str;
19    float duration;
20    int i;
21
22    cudaEventCreate (&startT);
23    cudaEventCreate (&endT);
24    cudaStreamCreate (&str);
25
26    cudaEventRecord (startT, str);
27    for (i = 0; i < iter; i++)
28      {
29        if (ZEROCOPY)
30          {
31            int block = 1024;
32            int grid = (N - 1) / block + 1;
33            dummyKernel <<< 1, 1, 0, str >>> (h, N);
34          }
35        else
36          cudaMemcpyAsync (d, h, N, cudaMemcpyHostToDevice, str);
37      }
38    cudaEventRecord (endT, str);
39    cudaEventSynchronize (endT);
40    cudaEventElapsedTime (&duration, startT, endT);
41
42    cudaStreamDestroy (str);
43    cudaEventDestroy (startT);
44    cudaEventDestroy (endT);
```

```
45    return (duration / iter)/1000;  // convert to sec
46  }
47
48  //——————————————————————————————————————————————
49  int main (int argc, char **argv)
50  {
51    int iter = atoi (argv[1]);
52    int step = atoi (argv[2]);
53    int MAXDATASIZE = atoi (argv[3]);
54    unsigned char *h_data, *d_data;
55    unsigned char *h_pinned;
56    unsigned char *h_zerocopy;
57
58    // pageable data allocation
59    h_data = (unsigned char *) malloc (MAXDATASIZE);
60    cudaMalloc ((void **) &d_data, MAXDATASIZE);
61
62    //pinned host data allocation
63    cudaMallocHost ((void **) &h_pinned, MAXDATASIZE);
64
65    // zero-copy data allocation
66    cudaHostAlloc ((void **) &h_zerocopy, MAXDATASIZE, ←
           cudaHostAllocMapped);
67    printf("SIZE PAGEABLE PINNED ZERO-COPY\n");
68    for (int dataSize = 0; dataSize <= MAXDATASIZE; dataSize += ←
           step)
69        printf ("%i %f %f %f\n", dataSize,
70                   dataSize / timeCopyH2D (h_data, d_data, dataSize, ←
                        iter) / (1<<20),
71                   dataSize / timeCopyH2D (h_pinned, d_data, dataSize←
                        , iter)/ (1<<20),
72                   dataSize / timeCopyH2D (h_zerocopy, d_data, ←
                        dataSize, iter, true)/ (1<<20));
73
74    cudaFreeHost (h_pinned);
75    cudaFreeHost (h_zerocopy);
76    free (h_data);
77    cudaFree (d_data);
78    cudaDeviceReset ();
79    return 1;
80  }
```

A sample of the results (in MB/sec) as tested on a GTX 870M GPU using the CUDA 6.5.14 SDK, are shown below:

```
$ ./memcpyTest 100 1000000 10000000
SIZE PAGEABLE PINNED ZERO-COPY
1000000  8721.193359  11407.150391  251645.046875
2000000  9870.787109  11431.939453  477716.156250
3000000  8816.846680  11335.143555  729376.500000
4000000  8428.890625  10872.275391  954590.687500
5000000  8996.895508  10910.654297  1273603.500000
6000000  9064.430664  11688.812500  1586918.125000
7000000  9636.149414  11648.382812  1916372.125000
8000000  9613.311523  11710.164062  2148302.250000
9000000  9656.351562  11678.595703  2497866.500000
10000000  9543.507812  11627.116211  2602141.250000
```

The results for pageable and pinned host memory are consistent with the speeds reported from the `$CUDA/samples/1_Utilities/bandwidthTest` utility. However, the results for the zero-copy memory are clearly bogus!

It seems that the issue is that the single thread launched in line 33 does not trigger the transfer of the entire designated data block. Forcing the invocation of the kernel for all the tested memory arrangements, results in the following `timeCopyH2D` function:

```
1  float timeCopyH2D (unsigned char *h, unsigned char *d, int N, int←
        iter, bool ZEROCOPY = false)
2  {
```

```
3     cudaEvent_t startT, endT;
4     cudaStream_t str;
5     float duration;
6     int i;
7
8     cudaEventCreate (&startT);
9     cudaEventCreate (&endT);
10    cudaStreamCreate (&str);
11
12    cudaEventRecord (startT, str);
13    for (i = 0; i < iter; i++)
14      {
15        int block = 1024;
16        int grid = (N - 1) / block + 1;
17        if (!ZEROCOPY)
18          {
19            cudaMemcpyAsync (d, h, N, cudaMemcpyHostToDevice, str);
20            dummyKernel <<< grid, block, 0, str >>> (d, N);
21          }
22        else
23          dummyKernel <<< grid, block, 0, str >>> (h, N);
24      }
25    cudaEventRecord (endT, str);
26    cudaEventSynchronize (endT);
27    cudaEventElapsedTime (&duration, startT, endT);
28
29    cudaStreamDestroy (str);
30    cudaEventDestroy (startT);
31    cudaEventDestroy (endT);
32    return (duration / iter)/1000;
33  }
```

The results of the modified program are shown below:

```
$ ./memcpyTest_withKernel 100 1000000 10000000
SIZE  PAGEABLE  PINNED  ZERO-COPY
1000000  7139.290039  7236.488770  4724.393066
2000000  7622.321777  7693.185547  5189.917480
3000000  7627.702148  8057.083984  5141.410156
4000000  6630.308105  7898.557129  4420.788086
5000000  7482.463379  8091.921387  4954.079102
6000000  7526.963379  8168.914062  5034.622070
7000000  7304.780273  8191.183105  4979.723633
8000000  7429.434570  8235.093750  4981.714844
9000000  7413.144531  8203.048828  4856.591797
10000000  7374.308594  8240.675781  4991.039062
```

The above clearly illustrate that zero-copy memory is a recommended
solution only if just part of the host data are required for the computation.
Otherwise, even pageable memory is preferable in terms of performance.
The delays introduced by the suspension of the threads that access the
memory accumulate to a substantial performance hit.

7. The Mandelbrot set calculators of Section 6.12.1 are limited to a maximum
of 255 iterations per pixel. However, the beauty of the Mandelbrot set is
revealed for thousands or millions of iterations. Modify one or more of the
solution of Section 6.12.1 so that up to $2^{16} - 1$ iterations can be performed
for each pixel. Profile your program and analyze its performance. What
is the grid/block arrangement that yields the best performance?

**Answer**

The modifications required are minimum as they are concentrated on the
memory management part of the code. Using the #1 version of the solu-
tion described in Section 6.12.1.1 as the basis of the answer, we need to
modify only the `kernel.cu` file of Listing 6.29, so that the pixel values are

stored as short integers (instead of unsigned characters). The changes to Listing 6.29 are documented below:

```
// Line 5, changing the upper iteration limit
static const int MAXITER = 0xFFFF;

// Line 26, changing the signature of the kernel to accommodate ↩
    the different type of pixel storage
__global__ void mandelKernel (unsigned short int *d_res, double ↩
    upperX, double upperY, double stepX, double stepY, int resX, ↩
    int resY, int pitch)

// Lines 52, 53
unsigned short *h_res;
unsigned short *d_res;

// Lines 58, 59
CUDA_CHECK_RETURN (cudaMalloc ((void **) &d_res, resX * resY * ↩
    sizeof(unsigned short int)));
h_res = (unsigned short int *) malloc (resY * pitch * sizeof(↩
    unsigned short int));

// Line 71, getting the different-type results
CUDA_CHECK_RETURN (cudaMemcpy (h_res, d_res, resY * pitch * ↩
    sizeof(unsigned short int), cudaMemcpyDeviceToHost));
```

8. The stand-alone CUDA AES implementations of Section 6.12.2 suffer from a dominating data-transfer overhead, that exceeds the computational cost of the en-/de-cryption. Which of the following modifications will offer -if any- the biggest performance improvement?

   (a) *Use pinned host memory* : are there any concerns about the use of pinned memory for holding the whole of the data to be processed?

   (b) *Move the tables from constant to shared memory* : the constant memory cache is faster than global memory, but it is still slower than shared memory. Will the speed improvement offset the need to copy the tables to constant memory for every block of threads?

   (c) *Process multiple 16-byte blocks per thread* : turning `rijndaelGPUEncrypt()` into a `__device__` function and introducing another `__global__` function as a front-end to it, to be called by `rijndaelEncryptFE()`, should require the smallest possible effort.

   Modify the source code of the Version #2 program in order to introduce your chosen changes, and measure the performance improvement obtained.

   **Answer**

   (a) *Using pinned host memory* The problem with pinned memory is that it is a limited resource, and as such it would be prohibitive to store the entirety of the input data this way. Moving the input data to pinned memory in chunks prior to their transfer to the device would have no beneficial impact on the performance, as Version #2 already processes data in chunks of size `DEVICEMEMSIZE`, hence the piecewise transfer to pinned memory is already taking place implicitly.

   (b) *Moving the tables from constant to shared memory* The following prefix code is required in the `rijndaelGPUEncrypt` function contained in the `rijndael_device.cu` file, to copy the five encoding tables residing in constant memory to shared memory:

```cuda
__global__ void rijndaelGPUEncrypt (int nrounds, u32 * data,←
      int N)
{
  __shared__ u32 sTe0[256];
  __shared__ u32 sTe1[256];
  __shared__ u32 sTe2[256];
  __shared__ u32 sTe3[256];
  __shared__ u32 sTe4[256];

  for(int i=threadIdx.x; i< 256; i+= blockDim.x)
  {
    sTe0[i] = Te0[i];
    sTe1[i] = Te1[i];
    sTe2[i] = Te2[i];
    sTe3[i] = Te3[i];
    sTe4[i] = Te4[i];
  }
  __syncthreads();
...
```

The performance improvement is dramatic : as tested on a GTX
870M GPU using the CUDA 6.5.14 SDK, the encoding time of a
512MB file is reduced from 2.53s to 0.6s.

(c) *Processing multiple 16-byte blocks per thread* This modification can
be achieved with fewer changes, by reducing the size of the grid and
introducing a loop in the `rijndaelGPUEncrypt` function. The re-
sulting code, combining the changes of the previous step, is shown
below:

```cuda
// changes to file rijndael_host_streams.cu
const int BLOCKSPERTHREAD=4;

void rijndaelEncryptFE (const u32 * rk, int keybits, ←
    unsigned char *plaintext, unsigned char *ciphertext, int←
    N, int thrPerBlock = 256)
{
 . . .
      // grid calculation
      int grid = ceil((numDataBlocks*1.0) / (thrPerBlock * ←
          BLOCKSPERTHREAD));

      rijndaelGPUEncrypt <<< grid, thrPerBlock, 0, str[←
          whichStream] >>> (nrounds, d_buffer[whichStream], ←
          toSend);
 . . .



// changes to file rijndael_device.cu
__global__ void rijndaelGPUEncrypt (int nrounds, u32 * data,←
      int N)
{
  __shared__ u32 sTe0[256];
  __shared__ u32 sTe1[256];
  __shared__ u32 sTe2[256];
  __shared__ u32 sTe3[256];
  __shared__ u32 sTe4[256];

  for (int i = threadIdx.x; i < 256; i += blockDim.x)
    {
      sTe0[i] = Te0[i];
      sTe1[i] = Te1[i];
      sTe2[i] = Te2[i];
      sTe3[i] = Te3[i];
      sTe4[i] = Te4[i];
    }
  __syncthreads ();

  int myID = blockIdx.x * blockDim.x + threadIdx.x;
```

```
    int totalThr = gridDim.x * blockDim.x;
    for (; myID < (N >> 4); myID += totalThr)
      {
        u32 s0, s1, s2, s3, t0, t1, t2, t3;
        const u32 *rk = d_rk;   // to avoid changing d_rk
        u32 aux;

        int myDataIdx = myID << 2; // *4 to offset 16 bytes

// code below this line is shared with the original ↩
      rijndaelGPUEncrypt function
#ifndef FULL_UNROLL
        int r;
. . .
        data[myDataIdx + 3] = RESHUFFLE (s3);
      }
}
```

In terms of performance, there is virtually no change as the encryption/decryption of a single block already constitutes a sufficiently heavy computational load.

9. The MPI cluster AES implementations of Section 6.12.2 does not provide overlapping of communication and computation. This issue could be addressed if a new "work item" was downloaded by the worker nodes while the GPU was processing an already downloaded part. Modify the MPI solution to provide this functionality. Do you expect any problems with the load-balancing of the modified solution?

**Answer**

A solution to this problem could be that each worker sets-up buffer space to hold three work items, and starts processing one while using a non-blocking send and a non-blocking receive operations to communicate the completed and next pending work items respectively.

The only modifications required, relate to the worker code as shown below for the case of a GPU worker. An identical arrangement is possible for a CPU worker. The communication/computation overlap is achieved by using immediate communication primitives to exchange data with the master (lines 31, 36) while computation is taking place (line 38). There is no need to use immediate functions in lines 29 and 34, as scalar primitive types can be buffered by MPI :

```
1   // changes to file AES_MPI/main.cpp
2     if (rank == 0)
3       . . .
4       else                              // GPU worker
5         {
6           int workItemSize = atoi (argv[3]);
7           int thrPerBlock = atoi (argv[4]);
8           int pos[3] = { -1, -1, -1 };
9           int totalWork = 0;
10          unsigned char *workBuff[3];
11          workBuff[0] = new unsigned char[workItemSize];
12          workBuff[1] = new unsigned char[workItemSize];
13          workBuff[2] = new unsigned char[workItemSize];
14          int actualSize[3];
15          MPI_Request getRq, sendRq;
16
17          int activeBuff = 0;          // indices for work item buffer ↩
                management
18          int completedBuff = 2;
19          int nextBuff = 1;
20
```

```
21          MPI_Send (&(pos[completedBuff]), 1, MPI_INT, 0, TAG_RES, ←
                MPI_COMM_WORLD);
22          MPI_Recv (&(pos[activeBuff]), 1, MPI_INT, 0, TAG_WORK, ←
                MPI_COMM_WORLD, &stat);
23          MPI_Recv (workBuff[activeBuff], workItemSize, ←
                MPI_UNSIGNED_CHAR, 0, TAG_DATA, MPI_COMM_WORLD, &stat)←
                ;
24          MPI_Get_count (&stat, MPI_UNSIGNED_CHAR, &(actualSize[←
                activeBuff]));
25          totalWork += actualSize[activeBuff];
26          while (pos[activeBuff] >= 0)
27            {
28              // send completed item with immediate send
29              MPI_Send (&(pos[completedBuff]), 1, MPI_INT, 0, TAG_RES←
                  , MPI_COMM_WORLD);
30              if (pos[completedBuff] > 0)
31                MPI_Isend (workBuff[completedBuff], actualSize[←
                    completedBuff], MPI_UNSIGNED_CHAR, 0, TAG_DATA, ←
                    MPI_COMM_WORLD, &sendRq);

33              // get next item while processing
34              MPI_Recv (&(pos[nextBuff]), 1, MPI_INT, 0, TAG_WORK, ←
                  MPI_COMM_WORLD, &stat);
35              if (pos[nextBuff] > 0)
36                MPI_Irecv (workBuff[nextBuff], workItemSize, ←
                    MPI_UNSIGNED_CHAR, 0, TAG_DATA, MPI_COMM_WORLD, ←
                    &getRq);

38              rijndaelEncryptFE (rk, keybits, workBuff[activeBuff], ←
                  workBuff[activeBuff], actualSize[activeBuff], ←
                  thrPerBlock);

40          // make sure communications are complete before updating ←
                the indices and comntinuing
41              if (pos[completedBuff] > 0)
42                MPI_Wait (&sendRq, MPI_STATUS_IGNORE);
43              if (pos[nextBuff] > 0)
44                {
45                  MPI_Wait (&getRq, &stat);
46                  MPI_Get_count (&stat, MPI_UNSIGNED_CHAR, &(←
                      actualSize[nextBuff]));
47                  totalWork += actualSize[nextBuff];
48                }
49              // update indices
50              completedBuff = activeBuff;
51              activeBuff = nextBuff;
52              nextBuff = (nextBuff + 1) % 3;
53            }
54          rijndaelShutdown ();
55          cout << "Worker " << rank << " processed " << totalWork << ←
                endl;
56          delete [] workBuff[0];
57          delete [] workBuff[1];
58          delete [] workBuff[2];
59        }
```

The three work item buffers are addressed individually by the `completedBuff`, `activeBuff` and `nextBuff` indices. Each buffer is accompanied by a dedicated position (`pos` array) and length (`actualSize` array) variables. While one buffer is being processed (pointed to by `activeBuff`), another holding previous results (pointed to by `completedBuff`) is being send to the master node (line 31) and the third one (pointed to by `nextBuff`) is being filled with new input data (line 36).

In all the communication transactions, the `pos` array members guard against sending dummy data to the master (line 30) and against waiting for new data when the input is exhausted (lines 26 and 35).

A problem may arise with the load balancing of this solution if a very slow

worker is participating in the proceedings. It is possible that it can delay the termination of the overall execution by receiving a pending work item that could be processed faster at another node.

10. Modify the MPI cluster AES implementations of Section 6.12.2 so that only two types of messages are needed for data exchange, instead of the current three. How can this be combined with the modification of the previous exercise?

   **Answer**

   The messages tagged by the TAG_WORK label can be eliminated from the communication protocol if the master node maintains a record of which part of the data is assigned to each worker node respectively. When a single work item is assigned to each worker, a simple integer array with as many element as the number of workers, is sufficient for handling this requirement.

   ```
   int assignedWorkItem[comm_size];
   ```

   If more than a single work item is assigned at a time to a worker, then the assignedWorkItem array has to become an array of containers, such as a vector<int>:

   ```
   vector<int> assignedWorkItem[comm_size];
   ```

   The termination of the workers can be achieved by sending a work item of length zero.

   There is also another alternative. The whole point of this exercise is to *reduce the number of messages sent*, while at the same time being able to identify the origin/destination in the input/output data of the data parts being exchanged. Towards this end, we could utilize the tag of the messages to carry this information (it is of type int). The only drawback of collapsing the two messages into one is that there is a need for an extra data copy from the generic message buffer used to collect the incoming results (line 33) to the output data repository (line 38):

   ```
 1    if (rank == 0)
 2      {
 3        unsigned char *iobuf;
 4        unsigned char *recvbuf;
 5        if ((f = fopen (argv[1], "r")) == NULL)
 6          {
 7            fprintf (stderr, "Can't open %s\n", argv[1]);
 8            exit (EXIT_FAILURE);
 9          }
10
11        int workItemSize = atoi (argv[3]);
12
13        fseek (f, 0, SEEK_END);
14        lSize = ftell (f);
15        rewind (f);
16
17        iobuf = new unsigned char[lSize];
18        recvbuf = new unsigned char[workItemSize];
19        assert (iobuf != NULL);
20        assert (recvbuf != NULL);
21        fread (iobuf, 1, lSize, f);
22        fclose (f);
23
24        timeval tim;
25        gettimeofday (&tim, NULL);
   ```

```
26            double tm2 = tim.tv_sec + (tim.tv_usec / 1000000.0);
27
28            // master main loop
29            int pos = 0;
30            while (pos < lSize)
31              {
32                int retPos, recvSize;
33                MPI_Recv (recvbuf, workItemSize, MPI_UNSIGNED_CHAR, ←↩
                      MPI_ANY_SOURCE, MPI_ANY_TAG, MPI_COMM_WORLD, &stat←↩
                      );
34                MPI_Get_count (&stat, MPI_UNSIGNED_CHAR, &recvSize);
35                if (recvSize > 0)
36                  {
37                    retPos = stat.MPI_TAG;
38                    memcpy (iobuf + retPos, recvbuf, recvSize);
39                  }
40
41                // assign next work item
42                int actualSize = (workItemSize < lSize - pos) ? ←↩
                      workItemSize : (lSize - pos);
43                MPI_Send (iobuf + pos, actualSize, MPI_UNSIGNED_CHAR, ←↩
                      stat.MPI_SOURCE, pos, MPI_COMM_WORLD);
44                pos += actualSize;
45              }
46
47            // wait for last results
48            for (int i = 1; i < comm_size; i++)
49              {
50                int retPos, recvSize;
51                MPI_Recv (recvbuf, workItemSize, MPI_UNSIGNED_CHAR, ←↩
                      MPI_ANY_SOURCE, MPI_ANY_TAG, MPI_COMM_WORLD, &stat←↩
                      );
52                MPI_Get_count (&stat, MPI_UNSIGNED_CHAR, &recvSize);
53                if (recvSize > 0)
54                  {
55                    retPos = stat.MPI_TAG;
56                    memcpy (iobuf + retPos, recvbuf, recvSize);
57                  }
58
59                // send termination signal
60                MPI_Send (recvbuf, 0, MPI_UNSIGNED_CHAR, stat.←↩
                      MPI_SOURCE, 0, MPI_COMM_WORLD);
61              }
62            gettimeofday (&tim, NULL);
63            double tm3 = tim.tv_sec + (tim.tv_usec / 1000000.0);
64
65            FILE *fout;
66            if ((fout = fopen (argv[2], "w")) == NULL)
67              {
68                fprintf (stderr, "Can't open %s\n", argv[2]);
69                exit (EXIT_FAILURE);
70              }
71            fwrite (iobuf, 1, lSize, fout);
72            fclose (fout);
73            delete [] iobuf;
74            delete [] recvbuf;
75
76            gettimeofday (&timeMain, NULL);
77            double tm4 = timeMain.tv_sec + (timeMain.tv_usec / ←↩
                  1000000.0);
78
79            // print-out some timing information
80            printf ("%.9lf \t %.9lf \n", tm4 - tm1, tm3 - tm2);
81        }
82     else                              // GPU worker
83        {
84            int workItemSize = atoi (argv[3]);
85            int thrPerBlock = atoi (argv[4]);
86            int pos[3] = { 0, 0, 0 };
87            int totalWork = 0;
88            unsigned char *workBuff[3];
89            workBuff[0] = new unsigned char[workItemSize];
90            workBuff[1] = new unsigned char[workItemSize];
91            workBuff[2] = new unsigned char[workItemSize];
```

```
92          int actualSize[3] = { 0, 0, 0 };
93          MPI_Request getRq, sendRq;
94
95          int activeBuff = 0;          // indices for work item buffer ←
                 management
96          int completedBuff = 2;
97          int nextBuff = 1;
98
99          MPI_Send (workBuff[completedBuff], 0, MPI_UNSIGNED_CHAR, 0,←
                 pos[completedBuff], MPI_COMM_WORLD);
100         MPI_Recv (workBuff[activeBuff], workItemSize, ←
                 MPI_UNSIGNED_CHAR, 0, MPI_ANY_TAG, MPI_COMM_WORLD, &←
                 stat);
101         pos[activeBuff] = stat.MPI_TAG;
102         MPI_Get_count (&stat, MPI_UNSIGNED_CHAR, &(actualSize[←
                 activeBuff]));
103         totalWork += actualSize[activeBuff];
104         while (actualSize[activeBuff] > 0)
105           {
106             // send completed item with immediate send
107             MPI_Isend (workBuff[completedBuff], actualSize[←
                   completedBuff], MPI_UNSIGNED_CHAR, 0, pos[←
                   completedBuff], MPI_COMM_WORLD, &sendRq);
108
109             // get next item while processing
110             MPI_Irecv (workBuff[nextBuff], workItemSize, ←
                   MPI_UNSIGNED_CHAR, 0, MPI_ANY_TAG, MPI_COMM_WORLD,←
                   &getRq);
111
112             rijndaelEncryptFE (rk, keybits, workBuff[activeBuff], ←
                   workBuff[activeBuff], actualSize[activeBuff], ←
                   thrPerBlock);
113
114             // make sure communications are complete before ←
                   updating the indices and comntinuing
115             MPI_Wait (&sendRq, MPI_STATUS_IGNORE);
116             MPI_Wait (&getRq, &stat);
117             MPI_Get_count (&stat, MPI_UNSIGNED_CHAR, &(actualSize[←
                   nextBuff]));
118             pos[nextBuff] = stat.MPI_TAG;
119             totalWork += actualSize[nextBuff];
120
121             completedBuff = activeBuff;
122             activeBuff = nextBuff;
123             nextBuff = (nextBuff + 1) % 3;
124           }
125         rijndaelShutdown ();
126         cout << "Worker " << rank << " processed " << totalWork << ←
                 endl;
127         delete [] workBuff[0];
128         delete [] workBuff[1];
129         delete [] workBuff[2];
130       }
```

In the above code the explicit communication of the buffer positions is eliminated. Termination and first job request conditions are detected by checking the length of the input (line 35) and output data (line 104) communicated respectively.

11. The whole point of the multi-core "adventure", is to accelerate our programs. This should be our sole focus, beyond any mis- or pre-conceptions. The evaluation of the different AES parallel implementations conducted in Section 6.12.2.4 considered only the encryption process, disregarding any I/O costs incurred. Perform your own experiment where the overall execution time is considered and not just the encryption time. Make sure that the file cache provided by the operating system is not utilized by:

   - Either, calling the following from the command-line (root permissions

are required):

```
$ sync ; echo 3 > /proc/sys/vm/drop_caches
```

- Or, calling the posix_fadvice() function from within your program, prior to any I/O:

```
#include <unistd.h>
#include <fcntl.h>
int main(int argc, char *argv[]) {
  int fd;
  fd = open(argv[1], O_RDONLY); // Open the file holding the↩
      input data
  fdatasync(fd);
  posix_fadvise(fd, 0,0,POSIX_FADV_DONTNEED); // clear cache
  close(fd);
  . . .
}
```

Analyze your findings.

**Answer**

This is an activity that is to be carried out by the students. The following command line can be used for generating a random 32MB input file to be used in testing:

```
$ dd if=/dev/urandom of=in32MBytes bs=32M count=1
```

# Chapter 7

# The Thrust Template Library

## Exercises

1. Develop a Thrust program for calculating the inner product of two vectors by using the `thrust::transform` algorithm and an appropriate functor.

    **Answer**

    The `thrust::transform` algorithm can be used to multiple the elements of two vectors, but it has to be accompanied by the `thrust::reduce` algorithm in order to sum-up the partial products. In the following code, two randomly-initialized vectors, of user-supplied size (as specified by a single command-line parameter), are multiplied accordingly. The result is computed on both the CPU and GPU in order to compare the results:

```
#include <iostream>
#include <stdlib.h>
#include <thrust/device_vector.h>
#include <thrust/host_vector.h>
#include <thrust/transform.h>
#include <thrust/reduce.h>

using namespace std;
struct innerProdFunct
{
  __host__ __device__ float operator () (float &x, float &y)
  {
    return x * y;
  }
};

int main (int argc, char **argv)
{
  int N = atoi (argv[1]);
  thrust::host_vector < float >h_x (N);
  thrust::host_vector < float >h_y (N);
  thrust::host_vector < float >h_prod (N);
  thrust::device_vector < float >d_x;
  thrust::device_vector < float >d_y;
  thrust::device_vector < float >d_prod (N);

  srand (time (0));
  for (int i = 0; i < N; i++)
    {
      h_x[i] = rand () % 1000;
```

```cpp
      h_y[i] = rand () % 1000;
    }

  // copy data to device
  d_x = h_x;
  d_y = h_y;
  // device-based calculation
  thrust::transform (d_x.begin (), d_x.end (), d_y.begin (), ←
      d_prod.begin (), innerProdFunct ());
  float res = thrust::reduce (d_prod.begin (), d_prod.end ());
  cout << "Inner product on GPU is " << res << endl;

  // host calculation
  thrust::transform (h_x.begin (), h_x.end (), h_y.begin (), ←
      h_prod.begin (), innerProdFunct ());
  res = thrust::reduce (h_prod.begin (), h_prod.end ());
  cout << "Inner product on CPU is " << res << endl;

  return 0;
}
```

2. Calculate the inner product of two vectors by using the `thrust::inner_product` algorithm. Compare the performance of this version with the performance of a CUDA-based solution.

   **Answer** The use of the `thrust::inner_product` can replace the sequence of `thrust::transform` and `thrust::reduce` of the previous exercise. The resulting program does not require a functor for the calculation to take place:

```cpp
#include <iostream>
#include <stdlib.h>
#include <thrust/device_vector.h>
#include <thrust/host_vector.h>
#include <thrust/inner_product.h>

using namespace std;

int main (int argc, char **argv)
{
  int N = atoi (argv[1]);
  thrust::host_vector < float >h_x (N);
  thrust::host_vector < float >h_y (N);
  thrust::device_vector < float >d_x;
  thrust::device_vector < float >d_y;

  srand (time (0));
  for (int i = 0; i < N; i++)
    {
      h_x[i] = rand () % 1000;
      h_y[i] = rand () % 1000;
    }

  d_x = h_x;
  d_y = h_y;

  float res = thrust::inner_product (d_x.begin (), d_x.end (),
                       d_y.begin (),
                       0,                 // initial reduction value
                       thrust::plus < float >(),  // reduction
                       thrust::multiplies < float >()); // ←
                             transformation
  cout << "Inner product function on GPU produces " << res << ←
      endl;

  return 0;
}
```

3. Use Thrust algorithms to find the absolute maximum of an array of values.

**Answer**

The solution can be based on the `thrust::transform_reduce` algorithm, which permits kernel fusion. A custom unary functor (lines 8-14) can be used to calculate the absolute values during the transformation phase, with the built-in `thrust::maximum` functor used in the reduction phase.

```cpp
#include <iostream>
#include <stdlib.h>
#include <thrust/device_vector.h>
#include <thrust/host_vector.h>
#include <thrust/transform_reduce.h>

using namespace std;
struct absFunct
{
  __host__ __device__ float operator () (float &x)
  {
    return fabs (x);
  }
};

int main (int argc, char **argv)
{
  int N = atoi (argv[1]);
  thrust::host_vector < float >h_x (N);
  thrust::device_vector < float >d_x;

  srand (time (0));
  for (int i = 0; i < N; i++)
    {
      h_x[i] = rand () % 1000;
    }

  d_x = h_x;
  float max = thrust::transform_reduce (d_x.begin (), d_x.end (),
                                        absFunct (),
                                        0,
                                        thrust::maximum < float ↩
                                           >());
  cout << "Max by GPU: " << max << endl;


  max = 0;
  for (int i = 0; i < N; i++)
    {
      float temp = fabs (h_x[i]);
      if (temp > max)
        max = temp;
    }
  cout << "Max by CPU: " << max << endl;

  return 0;
}
```

4. Write a Thrust program for calculating the sum of two matrices.

**Answer**

The solution can be based on the `thrust::transform` algorithm, combined with the built-in `thrust::plus` functor. The matrices can be simply treated as vectors, as their dimensions have no influence on the computation (as they do in multiplication). In the following code the user supplies the dimensions of the matrices to be randomly generated for testing the code.

A `thrust::host_vector` is used for initializing the two matrices (lines 39-40 and 43-44) before passing them to the device (lines 41 and 45) and for retrieving and printing out the result (lines 49-50).

```
1    #include <iostream>
2    #include <thrust/host_vector.h>
3    #include <thrust/device_vector.h>
4    #include <thrust/transform.h>
5    #include <thrust/random.h>
6    #include <math.h>
7
8    using namespace std;
9
10   void printMatrix (thrust::host_vector < float >&m, int N, int M)
11   {
12     for (int i = 0; i < N; i++)
13       {
14         cout << "| ";
15         for (int j = 0; j < M; j++)
16           {
17             cout << m[i * M + j] << " ";
18           }
19         cout << "|" << endl;
20       }
21     cout << "═══════════════════════════\n";
22   }
23
24   //****************************************************
25   int main (int argc, char **argv)
26   {
27     // initialize the RNG
28     thrust::default_random_engine rng (time (0));
29     thrust::uniform_int_distribution < int >uniDistr (−10000, ↩
            10000);
30
31     int N = atoi (argv[1]);
32     int M = atoi (argv[2]);
33
34     // generate the data on the host and move them to the device
35     thrust::device_vector < float >x (N * M);
36     thrust::device_vector < float >y (N * M);
37     thrust::device_vector < float >z (N * M);
38     thrust::host_vector < float >aux (N * M);
39     for (int i = 0; i < x.size (); i++)
40       aux[i] = uniDistr (rng);
41     x = aux;
42     printMatrix (aux, N, M);
43     for (int i = 0; i < x.size (); i++)
44       aux[i] = uniDistr (rng);
45     y = aux;
46     printMatrix (aux, N, M);
47
48     thrust::transform (x.begin (), x.end (), y.begin (), z.begin ()↩
            , thrust::plus < float >());
49     aux = z;
50     printMatrix (aux, N, M);
51
52     return 0;
53   }
```

5. Write a Thrust program for calculating the definite integral of a function $f(x)$ over a range $[a, b]$, using the trapezoidal rule. Consider a solution that avoids the need to create a vector for all the x-values for which a trapezoid is calculated. The details of the trapezoidal rule can be found in Section 3.5.2.

   **Answer**

   The answer is to employ a `thrust::counting_iterator` to generate a sequence of indices, and a functor which scales them to generate $x$ values for the calculation of the desired function.

   The functor shown in the listing below (lines 13-28) is initialized with the appropriate offset (data member `a`) and scale (data member `h`) to convert

an index $i$ to an appropriate $x_i$ value (line 25), as dictated by the equation that calculates the integral of a function $f()$ in the range $[a, b]$ :

$$h \left( \frac{f(a) + f(b)}{2} + \sum_{i=1}^{n-1} f(x_i) \right)$$

where $x_0 = a$, $x_n = b$, $x_i = a + i \cdot h$ and $h = \frac{b-a}{n}$.

Line 26 calculates the value of the function. The `thrust::transform_reduce` algorithm of lines 39-42, sums-up the produced function values, effectively calculating the summation term of the above equation.

By having the initial reduction value set to ( `f(0.0f) + f(N * 1.0f)` ) `/2`, (`f(0.0f)` evaluates to `f(a)` and `f(N * 1.0f)` to `f(b)`)the reduction calculates the complete parenthesis expression. The multiplication by $h$ in line 43 completes the result.

```
1   #include <iostream>
2   #include <stdlib.h>
3   #include <math.h>
4   #include <thrust/iterator/counting_iterator.h>
5   #include <thrust/transform_reduce.h>
6   #include <thrust/functional.h>
7
8   const float LOWERLIMIT = 0;   // stands for range limit a
9   const float UPPERLIMIT = 10;  // stands for range limit b
10
11  using namespace std;
12  //*****************************************
13  struct integrFunct
14  {
15    float h, a;
16
17      integrFunct (float st, float width)
18    {
19      a = st;
20      h = width;
21    }
22
23    __host__ __device__ float operator () (float i)
24    {
25      float x = a + h * i;
26      return fabs (sin (x));
27    }
28  };
29  //*****************************************
30  int main (int argc, char **argv)
31  {
32    int N = atoi (argv[1]);
33
34    float h = (UPPERLIMIT - LOWERLIMIT) / N;;
35    integrFunct f (LOWERLIMIT, h);
36
37    thrust::counting_iterator < float >x (1);
38
39    float res = thrust::transform_reduce (x, x + N - 1,
40                                          f,
41                                          ( f(0.0f) + f(N * 1.0f) ) ↵
                                             /2,
42                                          thrust::plus < float >());
43    res *= h;
44    cout << "Definite integral is " << res << endl;
45
46    return 0;
47  }
```

6. Use Thrust to calculate the mean and variance of a data set $X$ of cardinality $N$. For convenience, the corresponding formulas are:

$$E[X] = \frac{\sum_{i=0}^{N-1} x_i}{N} \tag{7.1}$$

$$\sigma^2 = E[X^2] - (E[X])^2 \tag{7.2}$$

**Answer**

This is a straight-forward application of the `thrust::reduce` algorithm for calculating the sum of the data set values, and the `thrust::transform_reduce` algorithm for calculating the sum of the square data set values. The remaining calculations are very simple.

```cpp
#include <iostream>
#include <stdlib.h>
#include <math.h>
#include <thrust/host_vector.h>
#include <thrust/device_vector.h>
#include <thrust/reduce.h>
#include <thrust/random.h>
#include <thrust/transform_reduce.h>
#include <thrust/functional.h>

using namespace std;
//*****************************************
struct sqrFunct
{
  __host__ __device__ float operator () (float x)
  {
    return x*x;
  }
};
//*****************************************
int main (int argc, char **argv)
{
  int N = atoi (argv[1]);  // initialize the RNG
  thrust::default_random_engine rng (time(0));
  thrust::uniform_int_distribution<int> uniDistr(-10000,10000);

  // generate the data on the host and move them to the device
  thrust::device_vector < float >d_x (N);
  thrust::host_vector<float> h_x(N);

  for (int i = 0; i < N; i++) h_x[i] = uniDistr(rng);
  d_x = h_x;

  // calculate the sum of values
  float avg = thrust::reduce(d_x.begin(), d_x.end());
  avg /= N;
  //calculate the sum of squared values
  float sum2 = thrust::transform_reduce (d_x.begin(), d_x.end(),
                                          sqrFunct(),
                                          0,
                                          thrust::plus < float >())
                                          ;

  cout << "Mean : " << avg << endl;
  cout << "Variance : " << sum2/N - avg*avg << endl;

  return 0;
}
```

7. Measure the performance of the `thrust::sort` algorithm in Thrust by sorting varying volumes of data. Compare the achieved times with the STL version running on the host. For this purpose create a big array (make sure it is not too big to fit in the GPU's memory) and populate it with random data.

**Answer**

The following listing employs the POSIX high-resolution timing functions (see Appendix C.2) for reporting the time it takes to sort a variable-size array of integers using STL on the host and Thrust on the device respectively. The program requires three command-line parameters corresponding to the minimum (`Nmin`) and maximum array size (`Nmax`), and the iteration step used in testing (`Nstep`).

Once the arrays are initialized and placed in their respective memory spaces, the loop of lines 40-52 proceeds to sort an increasing part of the arrays and report the time and processing rate in million integers/sec units.

```cpp
1  #include <iostream>
2  #include <algorithm>
3  #include <stdlib.h>
4  #include <time.h>
5  #include <thrust/device_vector.h>
6  #include <thrust/host_vector.h>
7  #include <thrust/sort.h>
8
9  using namespace std;
10
11 //*************************************************
12
13 double hrclock_sec ()
14 {
15   timespec ts;
16   clock_gettime (CLOCK_REALTIME, &ts);
17   double aux = ts.tv_sec + ts.tv_nsec / 1000000000.0;
18   return aux;
19 }
20
21 //*************************************************
22
23 int main (int argc, char **argv)
24 {
25   int Nmin = atoi (argv[1]);
26   int Nmax = atoi (argv[2]);
27   int Nstep = atoi (argv[3]);
28   thrust::host_vector < float >h_x (Nmax);
29   thrust::device_vector < float >d_x;
30
31   srand (time (0));
32   for (int i = 0; i < Nmax; i++)
33     {
34       h_x[i] = rand () % 1000000;
35     }
36
37   d_x = h_x;
38
39   double cpu_time, gpu_time, t0, t1, t2;
40   for (int items = Nmin; items <= Nmax; items += Nstep)
41     {
42       t0 = hrclock_sec ();
43       std::sort (h_x.begin (), h_x.begin () + items);
44       t1 = hrclock_sec ();
45       cpu_time = t1 - t0;
46
47       thrust::sort (d_x.begin (), d_x.begin () + items);
48       t2 = hrclock_sec ();
49       gpu_time = t2 - t1;
50
51       cout << items << " Times CPU,GPU (sec) : " << cpu_time << "↩
                " << gpu_time << "\t Rates (million integers/sec) : "↩
                << items / cpu_time / 1000000 << " " << items / ↩
                gpu_time / 1000000 << "\n";
52     }
53
54   return 0;
```

```
55  }
```

The results of a sample run on a 3.5GHz i7-3770K CPU and a GTX 560Ti
GPU are shown below:

```
$ nvcc −O2 sortThrust.cu − o sortThrust
$ ./sortThrust 1000000 5000000 500000
1000000 Times CPU,GPU (sec) : 0.0548186 0.00260711        Rates (←↩
    million integers/sec) : 18.242 383.567
1500000 Times CPU,GPU (sec) : 0.0649869 0.00337315        Rates (←↩
    million integers/sec) : 23.0816 444.689
2000000 Times CPU,GPU (sec) : 0.0791855 0.00413704        Rates (←↩
    million integers/sec) : 25.2572 483.438
2500000 Times CPU,GPU (sec) : 0.0953443 0.00491595        Rates (←↩
    million integers/sec) : 26.2208 508.548
3000000 Times CPU,GPU (sec) : 0.103896 0.00588679         Rates (←↩
    million integers/sec) : 28.8751 509.615
3500000 Times CPU,GPU (sec) : 0.124521 0.0065577          Rates (←↩
    million integers/sec) : 28.1077 533.723
4000000 Times CPU,GPU (sec) : 0.127966 0.00732589         Rates (←↩
    million integers/sec) : 31.2583 546.009
4500000 Times CPU,GPU (sec) : 0.136482 0.00801635         Rates (←↩
    million integers/sec) : 32.9713 561.353
5000000 Times CPU,GPU (sec) : 0.151194 0.00879908         Rates (←↩
    million integers/sec) : 33.0701 568.241
```

8. Measure the performance hit that an Array-of-Structures design approach
   will have in device-based sorting of $10^8$ randomly generated instances of
   the following structure:

```
struct TestStr {
  int key;
  float value;

  __host__ __device__
  bool operator<(const TestStr &o) const
  {
    return this.key < o.key;
  }
};
```

To measure the performance deterioration, you will have to implement
and time a Structure-of-Arrays alternative.

**Answer**

The answer is given in the listing below. While the `thrust::sort` algo-
rithm can used for sorting the array of structures, the `thrust::sort_by_key`
algorithm needs to be employed for the structure-of-arrays alternative.
The data are randomly generated and are common in both arrangements.

The POSIX high-resolution timing functions (see Appendix C.2) are used
for reporting the time.

```
#include <iostream>
#include <algorithm>
#include <stdlib.h>
#include <time.h>
#include <thrust/device_vector.h>
#include <thrust/host_vector.h>
#include <thrust/sort.h>

using namespace std;

//***********************************************

double hrclock_sec ()
{
```

```cpp
  timespec ts;
  clock_gettime (CLOCK_REALTIME , &ts);
  double aux = ts.tv_sec + ts.tv_nsec / 1000000000.0;
  return aux;
}

//*************************************************
struct TestStr
{
  int key;
  float value;

  __host__ __device__ bool operator< (const TestStr & o) const
  {
    return this->key < o.key;
  }
};

//*************************************************

int main (int argc , char **argv)
{
  int N = atoi (argv[1]);
  thrust::host_vector < TestStr > h_str (N);
  thrust::device_vector < TestStr > d_str;

  thrust::host_vector < int >h_keys (N);
  thrust::host_vector < float >h_values (N);
  thrust::device_vector < int >d_keys (N);
  thrust::device_vector < float >d_values (N);

  srand (time (0));
  for (int i = 0; i < N; i++)
    {
      h_str[i].key = rand () % 1000000;
      h_str[i].value = rand () % 1000000;

      h_keys[i] = h_str[i].key;
      h_values[i] = h_str[i].value;
    }

  d_str = h_str;
  d_keys = h_keys;
  d_values = d_values;

  double t0, t1, t2;
  t0 = hrclock_sec ();
  thrust::sort (d_str.begin (), d_str.end ());
  t1 = hrclock_sec ();
  thrust::sort_by_key (d_keys.begin (), d_keys.end (), d_values.↩
      begin ());
  t2 = hrclock_sec ();
  cout << N << " " << t1 - t0 << " " << t2 - t1 << endl;

  return 0;
}
```

9. Create a Thrust program for computing the Mandelbrot set. Section 3.8.2 covers the mathematical details of how the set is calculated.

**Answer**

Thrust has the limitation of not providing a 2D counting iterator, that would have been handy in this problem. We can still use the `thrust::counting_iterator` facility as long as we can couple it with a transformation that maps each generated index to a coordinate on the complex plane. In the functor of lines 19-57 in the listing that follows, this is exactly what is taking place.

The functor is initialized with all the information necessary to map what is essentially a pixel counter running from top-to-bottom and left-to-right,

into a pair of complex plane coordinates. The conversion is taking place in lines 40-43. The remaining lines in the `operator()` method is just a copy of the `MandelCompute::diverge` method shown in Listing 3.29 in page 127.

The `thrust::transform` call in line 85 calculates the number of iterations required for each examined complex plane point to diverge. These data are subsequently copied to the host (line 88), prior to using them for painting the pixels of a *QImage* object (lines 91-96).

```
1   // To compile :
2   // nvcc −m64 −O2 −Xcompiler '−fPIE ' −D_REENTRANT −DQT_NO_DEBUG −↩
        DQT_GUI_LIB −DQT_CORE_LIB −I/usr/lib/x86_64−linux−gnu/qt5/↩
        mkspecs/linux−g++−64 −I. −I/usr/include/qt5 −I/usr/include/↩
        qt5/QtGui −I/usr/include/qt5/QtCore −I. −o mandelThrust ↩
        mandelThrust.cu −L /opt/cuda/lib64/ −lcudart −rdc=true −arch↩
        =sm_21 −lm  −lQt5Core −lQt5Gui −L /usr/lib/x86_64−linux−gnu/
3
4   #include <iostream>
5   #include <stdlib.h>
6   #include <math.h>
7   #include <QImage>
8   #include <QRgb>
9   #include <thrust/device_vector.h>
10  #include <thrust/host_vector.h>
11  #include <thrust/iterator/counting_iterator.h>
12  #include <thrust/transform.h>
13
14  using namespace std;
15
16  const int MAXITER = 255;
17
18  //******************************************
19  struct mandelFunct
20  {
21    double upperX , upperY ;
22    double xStep , yStep ;
23    int imageWidth ;
24
25    mandelFunct (double uX, double uY, double xS, double yS, int w)
26    {
27      upperX = uX ;
28      upperY = uY ;
29      xStep = xS ;
30      yStep = yS ;
31      imageWidth = w ;
32    }
33
34    __host__ __device__ int operator  () (int i)
35    {
36      int xIDX , yIDX ;
37      double cx , cy ;
38      int iter = 0;
39
40      xIDX = i % imageWidth ;
41      yIDX = i / imageWidth ;
42      cx = upperX + xIDX * xStep ;
43      cy = upperY − yIDX * yStep ;
44
45      double vx = cx , vy = cy , tx , ty ;
46      while (iter < MAXITER && (vx * vx + vy * vy) < 4)
47        {
48          tx = vx * vx − vy * vy + cx ;
49          ty = 2 * vx * vy + cy ;
50          vx = tx ;
51          vy = ty ;
52          iter++;
53        }
54      return iter ;
55    }
56
```

```
57  };
58
59  //*****************************************
60  int main (int argc, char **argv)
61  {
62    int imgX = 1024, imgY = 768;
63    double upperCornerX, upperCornerY;
64    double lowerCornerX, lowerCornerY;
65    double xStep, yStep;
66
67    upperCornerX = atof (argv[1]);
68    upperCornerY = atof (argv[2]);
69    lowerCornerX = atof (argv[3]);
70    lowerCornerY = atof (argv[4]);
71
72    xStep = (lowerCornerX - upperCornerX) / imgX;
73    yStep = (upperCornerY - lowerCornerY) / imgY;
74
75    QImage *img = new QImage (imgX, imgY, QImage::Format_RGB32);
76
77    mandelFunct f (upperCornerX, upperCornerY, xStep, yStep, imgX);
78
79    thrust::host_vector < int >h_iter;
80    thrust::device_vector < int >d_iter (imgX * imgY);
81
82    thrust::counting_iterator < float >x (0);
83
84    // calculate the pixel values on the device
85    thrust::transform (x, x + imgX * imgY, d_iter.begin (), f);
86
87    // copy results to host
88    h_iter = d_iter;
89
90    // draw the image pixels accordingly
91    for (int j = 0; j < imgY; j++)
92      for (int i = 0; i < imgX; i++)
93        {
94          int color = h_iter[imgX * j + i];
95          img->setPixel (i, j, qRgb (256 - color, 256 - color, 256 ↩
              - color));
96        }
97    img->save ("mandel.png", "PNG", 0);
98    return 0;
99  }
```

10. A problem related to the furthest-distance point solved in Section 7.5, is the problem of finding the pair of points which are the furthest apart from each other. Create a brute-force solution to this problem, i.e. by examining all pairs of points, using Thrust. Consider the case where the number of points is too big to allow storage of all the calculated distances. You should avoid duplicate calculations since the distance from point A to B is the same as the distance from B to A.

**Answer**

Not being able to store the distances creates both opportunities and difficulties:

- The number $N$ of points that can be used as input is not limited by the need to allocate $\Theta(\frac{N^2-N}{2})$ storage for the pair distances.
- The distances might have to be calculated multiple times.

Each pair of points can be identified by the index of the corresponding cell in the distances table. We can compare the "indices" of two pairs of points, by calculating and comparing the corresponding pair-wise distances. In order to avoid duplicate calculations, we can limit the indices to either of

Figure 7.1: Illustration of how a pair of points and their distance can be identified by the index of the corresponding cell in the distances matrix. The indices are calculated by numbering the cells in a top-to-bottom and left-to-right fashion, skipping the gray cells. The matrix is not actually allocated.

the two table halves, as split by the main diagonal. In the following code we use the lower half. The indices are calculated as illustrated in Figure 7.1, i.e. to identify the pair of point $(i, j)$, with $i > j$, we use the formula:

$$idx_{(i,j)} = \sum_{k=1}^{i-1} k + j = \frac{i(i-1)}{2} + j \qquad (7.3)$$

The opposite mapping, i.e. from an index to a pair of points can be also performed. Given an initial approximation for $i$ :

$$idx_{(i,j)} = \frac{i^2 - i)}{2} + j \Rightarrow i^2 = 2 \cdot idx_{(i,j)} - 2j + i \Rightarrow i \approx \sqrt{2 \cdot idx_{(i,j)}} \quad (7.4)$$

we can calculate an exact value for $i$ and $j$ by utilizing the knowledge that $i > j$:

```
i = (int) floor (sqrt (2.0 * idx));
while ((j = idx - i * (i - 1) / 2) >= i)
    i++;
```

The solution given below is based on the application of a `thrust::reduce` algorithm on a `thrust::counting_iterator` that enumerates all the possible cell indices/point pairs. Upon comparing two indices, the supplied functor (lines 20-49) reverses the mapping to convert the two pair indices into two pairs of points, pair `A` and pair `B`. The point indices `Ai`, `Aj` and `Bi`, `Bj` are subsequently used to calculate and compare the square of the distances $\|\overrightarrow{AiAj}\|^2 \ \|\overrightarrow{BiBj}\|^2$ (indices used liberally here to represent the actual points).

The only restriction of the following solution is that the type `thrust::counting_iterator<int>` iterator is appropriate for up to $N < 46340$. For bigger inputs, a switch to `thrust::counting_iterator<long>` would be necessary, as a 32-bit signed integer is not sufficient for holding the generated indices.

```cpp
1   #include <iostream>
2   #include <thrust/host_vector.h>
3   #include <thrust/device_vector.h>
4   #include <thrust/iterator/counting_iterator.h>
5   #include <thrust/reduce.h>
6   #include <thrust/random.h>
7   #include <math.h>
8
9   using namespace std;
10
11  //*******************************************************
12  __host__ __device__ void reverseMap (int idx, int &i, int &j)
13  {
14    i = (int) floor (sqrt (2.0 * idx));
15    while ((j = idx - i * (i - 1) / 2) >= i)
16      i++;
17  }
18
19  //*******************************************************
20  struct maxFunct
21  {
22    thrust::device_ptr < int >x, y, z;
23
24    // comparing distances of pairs of points A and B as indexed by↩
             idxA, idxB
25    __host__ __device__ int operator () (int idxA, int idxB)
26    {
27      int Ai, Aj, Bi, Bj;
28      int distA, distB;
29
30      reverseMap (idxA, Ai, Aj);
31      reverseMap (idxB, Bi, Bj);
32
33      int xDiff, yDiff, zDiff;
34      xDiff = x[Ai] - x[Aj];
35      yDiff = y[Ai] - y[Aj];
36      zDiff = z[Ai] - z[Aj];
37      distA = xDiff * xDiff + yDiff * yDiff + zDiff * zDiff;
38
39      xDiff = x[Bi] - x[Bj];
40      yDiff = y[Bi] - y[Bj];
41      zDiff = z[Bi] - z[Bj];
42      distB = xDiff * xDiff + yDiff * yDiff + zDiff * zDiff;
43
44      if (distA > distB)
45          return idxA;
46      else
47          return idxB;
48    }
49  };
50
51  //*******************************************************
52  int main (int argc, char **argv)
53  {
54    // initialize the RNG
55    thrust::default_random_engine rng (time (0));
56    thrust::uniform_int_distribution < int >uniDistr (-10000, ↩
          10000);
57
58    int N = atoi (argv[1]);
59
60    // generate the data on the host and move them to the device
61    thrust::device_vector < int >d_x (N);
62    thrust::device_vector < int >d_y (N);
63    thrust::device_vector < int >d_z (N);
64    thrust::host_vector < int >h_x (N);
65    thrust::host_vector < int >h_y (N);
66    thrust::host_vector < int >h_z (N);
67    for (int i = 0; i < N; i++) h_x[i] = uniDistr (rng);
68    d_x = h_x;
69    for (int i = 0; i < N; i++) h_y[i] = uniDistr (rng);
70    d_y = h_y;
71    for (int i = 0; i < N; i++) h_z[i] = uniDistr (rng);
72    d_z = h_z;
```

```
73
74     // initialize the functor that will find the maximum distance , ←
           so that it has access to the coordinate arrays
75     maxFunct f;
76     f.x = d_x.data ();
77     f.y = d_y.data ();
78     f.z = d_z.data ();
79
80     // reduce the index of the most distant point
81     int furthest = thrust::reduce (thrust::counting_iterator < int ←
           >(0) ,
82                                    thrust::counting_iterator < int ←
                                          >((N * N - N) / 2) ,
83                                    0 ,
84                                    f );
85
86     cout << "The most distant pair is " << furthest << endl;
87     int Ai , Aj ;
88     reverseMap (furthest , Ai , Aj );
89     cout << "First point is " << Ai << "(" << h_x[Ai] << " " << h_y←
           [Ai] << " " << h_z[Ai] << ")" << endl;
90     cout << "Second point is " << Aj << "(" << h_x[Aj] << " " << ←
           h_y[Aj] << " " << h_z[Aj] << ")" << endl;
91
92     // CPU verification
93     double max_dist = 0, maxI , maxJ ;
94     for (int i = 0; i < N - 1; i++)
95       for (int j = i + 1; j < N; j++)
96         {
97            int xDiff , yDiff , zDiff ;
98            xDiff = h_x[i] - h_x[j];
99            yDiff = h_y[i] - h_y[j];
100           zDiff = h_z[i] - h_z[j];
101           int dist = xDiff * xDiff + yDiff * yDiff + zDiff * zDiff;
102           if (dist > max_dist)
103             {
104                max_dist = dist;
105                maxI = i;
106                maxJ = j;
107             }
108        }
109    cout << endl;
110    cout << "Distance by CPU : " << max_dist << endl;
111    cout << "First point is " << maxI << "(" << h_x[maxI] << " " <<←
           h_y[maxI] << " " << h_z[maxI] << ")" << endl;
112    cout << "Second point is " << maxJ << "(" << h_x[maxJ] << " " ←
           << h_y[maxJ] << " " << h_z[maxJ] << ")" << endl;
113
114    return 0;
115 }
```

11. The all-pairs-shortest paths graph problem can be solved by the dynamic
    programming algorithm by Floyd and Warshall. Assuming that the graph
    is described by a $VxV$ adjacency matrix where $V$ is the number of vertices,
    then the pseudocode of this algorithm below, shows how the solution is
    obtain in $V$ stages, where each stage involves the update of the distance
    matrix via the involvement of an intermediate vertex:

```
// initialization phase
Allocate a VxV distance matrix , and set all elements to infinity
for each vertex v
   distance[v][v] <- 0
for each edge (u,v)
   distance[u][v] <- adj(u,v)   // the weight of the edge (u,v)

// computation phase
for k from 1 to V    // for each k intermediate vertex
   for i from 1 to V  // reconsider each pair of vertices
      for j from 1 to V
         if distance[i][j] > distance[i][k] + distance[k][j]
            distance[i][j] <- distance[i][k] + distance[k][j]
```

```
            end  if
return  distance
```

Create a Thrust program that could perform the $V$ individual stages of the algorithm in parallel, i.e. parallelize the two inner loops of the algorithm.

**Answer**

Similarly to the previous exercise for finding the pair of most distant points in a set, we can use a `thrust::counting_iterator` to go over all the elements of the adjacency matrix. Fortunately, as all elements are to be considered[1], the mapping between an index and the adjacency matrix coordinates is straightforward as shown in lines 31,32.

The functor that is used to update the distance of two vertices (lines 22-41), is setup so that it has access to the adjacency matrix data and the number of vertices $V$ (lines 64,65). Before each phase of the algorithm executed by the `for` loop of lines 68-76, the index of the vertex to be used as an intermediate, is also set (line 71).

```
1   #include <thrust/host_vector.h>
2   #include <thrust/device_vector.h>
3   #include <thrust/iterator/counting_iterator.h>
4   #include <thrust/transform.h>
5   #include <thrust/random.h>
6   #include <stdio.h>
7
8   using namespace std;
9
10  //****************************************************
11  void printAdj (thrust::host_vector < int >&a, int V)
12  {
13      for (int i = 0; i < V; i++)
14          {
15              for (int j = 0; j < V; j++)
16                  printf ("%3i ", a[i * V + j]);
17              printf ("\n");
18          }
19  }
20
21  //****************************************************
22  struct floydFunct
23  {
24      thrust::device_ptr < int >adj;
25      int V;
26      int k;
27
28      // comparing distances of pairs of points A and B as indexed by←
                i,j
29      __host__ __device__ int operator () (int idx)
30      {
31          int row = idx / V;
32          int col = idx % V;
33
34          int currDist = adj[idx];
35          int otherRouteDist = adj[row * V + k] + adj[k * V + col];
36          if (otherRouteDist < currDist)
37              return otherRouteDist;
38          else
39              return currDist;
40      }
41  };
42
```

---

[1]Actually the $k^{th}$ column and row are unchanged during stage $k$ but we do not exclude them from the computation as this would complicate the code. Not only do we need to change the mapping from the `thrust::counting_iterator` to the table cells, but also the results of each stage's computation cannot be stored in-place.

```cpp
43  //****************************************************
44  int main (int argc, char **argv)
45  {
46    // initialize the RNG
47    thrust::default_random_engine rng (time (0));
48    thrust::uniform_int_distribution < int >uniDistr (1, 100);
49
50    int V = atoi (argv[1]);
51
52    // generate the data on the host and move them to the device
53    thrust::device_vector < int >d_adj (V * V);
54    thrust::host_vector < int >h_adj (V * V);
55    for (int i = 0; i < V * V; i++)
56      h_adj[i] = uniDistr (rng);
57    for (int i = 0; i < V; i++)
58      h_adj[i * V + i] = 0;
59
60    d_adj = h_adj;
61
62    // initialize the functor so that it has access to the ↩
63        coordinate arrays
63    floydFunct f;
64    f.V = V;
65    f.adj = d_adj.data ();
66
67    printAdj (h_adj, V);
68    for (int k = 0; k < V; k++)
69      {
70        // set the functor to use the k−th node as intermediate
71        f.k = k;
72        thrust::transform (thrust::counting_iterator <int >(0),
73                           thrust::counting_iterator <int >(V * V),
74                           d_adj.begin (),
75                           f);
76      }
77
78    printf ("————————————————\n");
79    h_adj = d_adj;
80    printAdj (h_adj, V);
81
82    return 0;
83  }
```

# Chapter 8

# Load Balancing

## Exercises

1. Design a brute-force/exhaustive algorithm for determining the optimum subset of $N$ nodes/processors in the case of static load balancing. What is the complexity of your algorithm?

   **Answer**

   In general, finding the optimum subset of nodes requires enumerating all the possible subsets, i.e. $2^N - 1$ if we ignore the empty set. The enumeration can be implemented in a variety of ways, the following recursive algorithm being a generic approach that does not impose any restrictions on $N$.

   The following description is based on the assumption that a function to evaluate the execution time on a subset $S$ exists (symbolic name `eval()` used to represent this function).

   The recursive auxiliary function OptimumSetAux, calls itself by including (line 15) or excluding (line 16) the $i$-th element of the processor set $P$, until the parameter $i$ becomes equal to $N$ (base case of line 8), which indicates that a subset is ready to be evaluated.

   1: **procedure** OPTIMUMSET($P$)
   2:     $bestS \leftarrow \emptyset$
   3:     $bestT \leftarrow \inf$
   4:     OptimumSetAux($P$,$\emptyset$, 0, $bestS$, $bestT$ )
   5:     Return $bestS$
   6: **end procedure**
   7: **procedure** OPTIMUMSETAUX($P$, $S$, $i$, $bestS$, $bestT$ )
   8:     **if** $i = N$ **then**
   9:         $T \leftarrow eval(S)$
   10:         **if** $T < bestT$ **then**
   11:             $bestS \leftarrow S$
   12:             $bestT \leftarrow T$
   13:         **end if**
   14:     **else**
   15:         OptimumSetAux($P$, $S \cup P_i$, $i + 1$, $bestS$, $bestT$ )

16:            OptimumSetAux($P$, $S$, $i + 1$, $bestS$, $bestT$ )

17:       **end if**

18: **end procedure**

The complexity of this algorithm depends on the cost of the `eval()` function. Given that there are $2^N$ subsets to be evaluated, we can conclude that the complexity $\in \Omega(2^N)$.

2. Design a greedy algorithm for determining the optimum subset of $N$ nodes/processors in the case of static load balancing. Will you always get the same solution as the one provided by an exhaustive algorithm?

   **Answer**

   A greedy algorithm can be employed that initially sorts the nodes in descending order of their computational power, i.e. ascending of their $p_i$ parameter as identified in Table 8.1.

   The following algorithm iterates over the sorted nodes, adding one node at a time for as long as the execution time decreases (`while` loop terminates otherwise via line 14) or until all nodes are examined (line 6). The complex condition of line 9 ensures that sets leading to invalid solutions (i.e. with one or more $part_k \leq 0$) are not considered.

   Actually, the condition in line 9 can be simply "`if` $T < bestT$ `then`", as a negative part to a node enlarges the parts assigned to other nodes, inevitable making the execution time higher.

   1: **procedure** GREEDYSET($P$)
   2:       $P' \leftarrow sort(P)$
   3:       $bestS \leftarrow P'_0$
   4:       $bestT \leftarrow eval(bestS)$
   5:       $i \leftarrow 1$
   6:       **while** $i < N$ **do**
   7:            $S \leftarrow bestS \cup P'_i$
   8:            $T \leftarrow eval(S)$
   9:            **if** $T < bestT$ AND $\forall P_k \in S, part_k > 0$ **then**
   10:                 $bestS \leftarrow S$
   11:                 $bestT \leftarrow T$
   12:                 $i \leftarrow i + 1$
   13:            **else**
   14:                 Break
   15:            **end if**
   16:       **end while**
   17:       Return $bestS$
   18: **end procedure**

   A heuristic algorithm cannot produce an identical solution to an exhaustive algorithm, unless this is done by chance, or special circumstances/problem settings allow the generation of an optimum solution via a greedy approach.

3. The closed-form solutions in Section 8.3.3.1 for the N-port communication setup, are based on the assumption that the master participates in the processing of the load. Derive the equations that would govern the solution

if the master abstained from this task, facilitating only I/O functionality instead.

**Answer**

The answer to this question is already provided in Section 8.3.3.1, albeit in a condensed form, without a length explanation. This exercise serves the purpose of testing student understanding and their ability to follow the proper steps towards the derivation of a solution.

Once we remove the master node from computation duties, Equation 8.9 becomes:

$$t_{distr}^{(i)} + t_{comp}^{(i)} + t_{coll}^{(i)} = t_{distr}^{(j)} + t_{comp}^{(j)} + t_{coll}^{(j)} \Rightarrow$$
$$l_i(a\ part_i L + b) + p_i(part_i L + e_i) + l_i(c\ part_i L + d) =$$
$$l_j(a\ part_j L + b) + p_j(part_j L + e_j) + l_j(c\ part_j L + d) \quad (8.1)$$

To solve the problem, we need to find one of the assigned parts, e.g. $part_1$. So from the above equation we can express all $part_i$ with $i \neq 1$ as a function of $part_1$:

$$part_i = part_1 \frac{l_1(a+c) + p_1}{l_i(a+c) + p_i} + \frac{p_1 e_1 - p_i e_i + (l_1 - l_i)(b+d)}{L(l_i(a+c) + p_i)} \quad (8.2)$$

Then, the normalization equation allows us to get a solution for $part_1$:

$$\sum_{j=1}^{N-1} part_j = 1 \Rightarrow$$

$$part_1 + part_1 \frac{l_1(a+c) + p_1}{l_2(a+c) + p_2} + \frac{p_1 e_1 - p_2 e_2 + (l_1 - l_2)(b+d)}{L(l_2(a+c) + p_2)} + \ldots$$

$$+ part_1 \frac{l_1(a+c) + p_1}{l_{N-1}(a+c) + p_{N-1}} + \frac{p_1 e_1 - p_{N-1} e_{N-1} + (l_1 - l_{N-1})(b+d)}{L(l_{N-1}(a+c) + p_{N-1})} = 1 \Rightarrow$$

$$part_1 \left( 1 + \sum_{j=2}^{N-1} \frac{l_1(a+c) + p_1}{l_j(a+c) + p_j} \right) = 1 + \sum_{j=2}^{N-1} \frac{p_j e_j - p_1 e_1 + (l_j - l_1)(b+d)}{L(l_j(a+c) + p_j)} \Rightarrow$$

$$part_1 \sum_{j=1}^{N-1} \frac{l_1(a+c) + p_1}{l_j(a+c) + p_j} = 1 + \sum_{j=2}^{N-1} \frac{p_j e_j - p_1 e_1 + (l_j - l_1)(b+d)}{L(l_j(a+c) + p_j)} \Rightarrow$$

$$part_1 = \frac{1 + L^{-1} \sum_{j=2}^{N-1} \frac{p_j e_j - p_1 e_1 + (l_j - l_1)(b+d)}{l_j(a+c) + p_j}}{(l_1(a+c) + p_1) \sum_{j=1}^{N-1} (l_j(a+c) + p_j)^{-1}} \quad (8.3)$$

4. Write a program that calculates $part_0$ from Equations 8.11 and 8.39 in a linear number of steps.

**Answer**

The following listing contains two functions that calculate $part_0$ and subsequently all $part_i$, for the two problem settings that Equations 8.11 and 8.39 correspond to. The "trick" (especially for Equation 8.39 that involves double summations) is to pre-compute and reuse terms, such as e.g. the

$\sum_{j=1}^{k} b_j$ term of Equation. 8.39, via lines 51-53 in `OneportBlockSingleInst`. The two functions also return the predicted total execution time.

It is important to note that `OneportBlockSingleInst` just implements Equations 8.38 and 8.39, without ordering the nodes as mandated by Theorem 8.3.3.

```cpp
1   #include <iostream>
2
3   using namespace std;
4   //***********************************************
5   // Stores the computed part_i in array part and the returns the ←
           total execution time
6   double NportBlockSingleInst (double *p, double *e, double *l, ←
           double a, double b, double c, double d, int N, long L, ←
           double *part)
7   {
8     double lacp[N];   // should be replaced by dynamic allocation if←
             N is big
9     double sumTerm[N];
10    double totalTime;
11    double nomin, denom;
12
13    for (int i = 0; i < N; i++)
14      lacp[i] = l[i] * (a + c) + p[i];
15
16    for (int i = 1; i < N; i++)
17      sumTerm[i] = p[0] * e[0] - p[i] * e[i] - l[i] * (b + d);
18
19    nomin = 0;
20    for (int i = 1; i < N; i++)
21      nomin += sumTerm[i] / lacp[i];
22    nomin = 1 - nomin / L;
23
24    denom = 1;
25    for (int i = 1; i < N; i++)
26      denom += 1 / lacp[i];
27    denom *= p[0];
28
29    part[0] = nomin / denom;
30
31    for (int i = 1; i < N; i++)
32      part[i] = part[0] * p[0] / lacp[i] + sumTerm[i] / (L * lacp[i←
             ]);
33
34    totalTime = p[0] * (part[0] * L + e[0]);
35    return totalTime;
36  }
37
38  //***********************************************
39  // Stores the computed part_i in array part and the returns the ←
           total execution time
40  double OneportBlockSingleInst (double *p, double *e, double l, ←
           double *b, int N, long L, double *part)
41  {
42    double pl[N]; // should be replaced by dynamic allocation if N ←
             is big
43    double bSum[N];
44
45    double totalTime;
46    double nomin, denom;
47
48    for (int i = 0; i < N; i++)
49      pl[i] = p[i] + l;
50
51    bSum[0] = 0;
52    for (int i = 1; i < N; i++)
53      bSum[i] = b[i] + bSum[i - 1];
54
55    nomin = 0;
56    for (int i = 1; i < N; i++)
57      nomin += l / pl[i] * bSum[i];
58    nomin /= L;
```

```
59     nomin++;
60     for (int i = 1; i < N; i++)
61       nomin += l / pl[i];
62
63     denom = 0;
64     for (int i = 0; i < N; i++)
65       denom += l / pl[i];
66     denom *= pl[0];
67
68     part[0] = nomin / denom;
69
70     for (int i = 1; i < N; i++)
71       part[i] = part[0] * pl[0] / pl[i] - l / (L * pl[i]) * bSum[i]↩
                  - l / pl[i];
72
73     totalTime = p[0] * part[0] * L;
74     return totalTime;
75   }
76
77   //*********************************************
78   int main ()
79   {
80     //small tester data
81     double p[4] = { 1, 2, 3, 4 };
82     double e[4] = { 1, 1, 1, 1 };
83     double b[4] = { 1000, 1000, 1000, 1000 };
84     double l[4] = { 1, 1, 1, 1 };
85     double part[4];
86
87     cout << "t : " << NportBlockSingleInst (p, e, l, 1, 0, 1, 0, 4,↩
                1000000, part) << endl;
88     for (int i = 0; i < 4; i++)
89       cout << part[i] << " ";
90     cout << endl;
91
92     cout << "t : " << OneportBlockSingleInst (p, e, 1, b, 4, ↩
                1000000, part) << endl;
93     for (int i = 0; i < 4; i++)
94       cout << part[i] << " ";
95     cout << endl;
96
97     return 0;
98   }
```

5. Solve the example at the end of Section 8.3.3.2, by reversing the computing power of the nodes: $p_0 = 4 \cdot 1.631 \cdot 10^{-07}$, $p_1 = 3 \cdot (1.631 \cdot 10^{-07})$, $p_2 = 2 \cdot (1.631 \cdot 10^{-07})$ and $p_3 = (1.631 \cdot 10^{-07})$.

   **Answer**

   Before applying the equations, the nodes (excluding the master) should be sorted according to their $b_i(p_i + l)$ property. This presents a challenge, as while we can switch the place of a node, the $b_i$ overheads are node *order* specific and not node specific. Hence. changing the order of a node, potentially causes a change in its $e_i$ parameter. Fortunately, the closed-form solutions of Section 8.3.3.2 allow the easy calculation of the predicted time for the different permutations of the nodes (surely not advisable for large $N$, but quite doable for the particular problem).

   The results show that load should be distributed in the order $P_3$, $P_2$, $P_1$. The corresponding parts are $part_0 = 0.13516$, $part_3 = 0.465104$, $part_2 = 0.238957$ and $part_1 = 0.160779$, with a total execution time of 2.19sec.

6. An alternative to having the master node distribute the load, is to have all the nodes access a network filesystem and retrieve the data from there.

What kind of communication configuration would correspond to such an arrangement? Calculate the partitioning that would be produced in this case, for the same problem setting as the previous question.

**Answer**

Such a design corresponds to a N-port communication configuration. If we assume that all nodes (including $P_0$ which was previously considered as having the load locally) perform the same sequence of getting the load, processing it and sending the results back to the filesystem server, then the minimum execution time is achieved by having all nodes finish at the same time instance, i.e. for every pair of $P_i$ and $P_j$, we have:

$$t_{distr}^{(i)} + t_{comp}^{(i)} + t_{coll}^{(i)} = t_{distr}^{(j)} + t_{comp}^{(j)} + t_{coll}^{(j)} \Rightarrow$$
$$l(part_i L + b_i) + p_i\, part_i L + l\, part_i L = l(part_j L + b_j) + p_j\, part_j L + l\, part_j L \Rightarrow$$
$$part_i L(p_i + 2\,l) = part_j L(p_j + 2\,l) + l(b_j - b_i) \Rightarrow$$
$$part_i = part_j \frac{p_j + 2\,l}{p_i + 2\,l} + \frac{l(b_j - b_i)}{L(p_i + 2\,l)} \quad (8.4)$$

Thus, we can express all $part_i$ as a function of $part_0$:

$$part_i = part_0 \frac{p_0 + 2\,l}{p_i + 2\,l} + \frac{l(b_0 - b_i)}{L(p_i + 2\,l)} \quad (8.5)$$

The normalization equation can then yield a closed-form solution for $part_0$:

$$\sum_{i=0}^{N-1} part_i = 1 \Rightarrow$$

$$part_0 \sum_{i=0}^{N-1} \frac{p_0 + 2\,l}{p_i + 2\,l} + \frac{l}{L} \sum_{i=1}^{N-1} \frac{b_0 - b_i}{p_i + 2\,l} = 1 \Rightarrow$$

$$part_0 = \frac{1 + \frac{l}{L} \sum_{i=1}^{N-1} \frac{b_i - b_0}{p_i + 2\,l}}{(p_0 + 2\,l) \sum_{i=0}^{N-1}(p_i + 2\,l)^{-1}} \quad (8.6)$$

If we assume that the nodes are ordered as specified by the problem statement, then we have $b_1 = b_2 = 2 \cdot 3 \cdot 2160B$ and $b_0 = b_3 = 3 \cdot 2160B$. Using the above equations produces $part_0 = 0.125605$, $part_1 = 0.165952$, $part_2 = 0.24447$ and $part_1 = 0.463997$, with a total execution time of 2.097sec.

7. The partitioning performed using DLT in the kernel convolution example, returns the percent of the input bytes that should be assigned to each node? How can we convert this number to pixels? One may argue that we should instead use image rows. How can we convert the result of our analysis to rows?

**Answer**

As long as there is a linear relationship between the data volume and the desired load unit, there is no need for a conversion. The $part_i$ results are

percentages and can be used without conversion to split the number of rows as well as the number of pixels.

8. Use the `tiobench` utility available at http://sourceforge.net/projects/tiobench/ to measure the performance you can get by concurrent access to a network filesystem (e.g. a NFS volume). Calculate the collective throughput of the server versus the number of threads used.

   **Answer**

   No model answer is available for this exercise. The students should experiment with the `tiobench` utility and plot the results. A plot similar to Figure 8.6 should be produced.

9. One way that can be used to improve the distribution cost of input data is to compress them. What kind of compression algorithms could be used in this case? How could we adapt the cost models to reflect this change?

   **Answer**

   There are two distinct possibilities:

   (a) The input data are already in compressed form. This imposed restrictions on the compression algorithm and data encoding format, as they should permit random access to the input data. This immediately rules out dynamic compression schemes (e.g. dictionary based) and enforces some sort of partitioning of the compressed data, in the form of packets or blocks. A static Huffman/arithmetic code compression algorithm would be suitable for the task.

   (b) The data are compressed before the distribution phase. Any compression algorithm could be used in this case, as long as the compression-decompression overhead does not exceed the gains from the reduced communication time.

   Let's assume that a node decompresses the input upon receiving it, and compresses the results prior to sending them back, spending $compr_i$ time per load unit. If on average, the compression ratio is $cr = \frac{originalDataSize}{compressedDataSize}$, then, the cost models would have to be modified as follows:

$$t^{(i)}_{comp} = p_i(part_i L + e_i) + compr_i part_i L =$$
$$(p_i + compr_i)(part_i L + \frac{p_i e_i}{p_i + compr_i}) \quad (8.7)$$

$$t^{(i)}_{distr} = l_i(a\ part_i \frac{L}{cr} + b) = l_i(\frac{a}{cr} part_i L + b) \quad (8.8)$$

$$t^{(i)}_{coll} = l_i(c\ part_i \frac{L}{cr} + d) = l_i(\frac{c}{cr} part_i L + d) \quad (8.9)$$

   Equations 8.7-8.9 adhere to the formulation of the book's Equations 8.5, 8.7 and 8.8(see pages 583, 585 and 586). This is a clear indication that existing analysis results could be applicable for such a design, by adapting the model parameters.

10. Use the DLTlib library to calculate the optimum partitioning for a "query processing" application, where the communication cost is independent of the workload, consisting of a query during distribution and the result during the collection. You can assume that $b = d = 1$, the workload consists of $L = 10^6$ items, and the parallel platform is made of 10 compute nodes connected in a single level tree, with $p_i = (i+1) \cdot 0.01 sec/item \ \ \forall i \in [0, 9]$ and $l = 0.001 sec/item$.

What would happen if the communication was 10 times slower?

**Answer** The following listing illustrates how the appropriate DLTlib function can be called and the results printed-out:

```cpp
#include <time.h>
#include <stdio.h>
#include <iostream>

using namespace std;

//——————————————————————————————————————————————
// DLTlib specific definitions that need to be used by the ↩
       library.
long global_random_seed;
#define MAX_NODE_DEGREE 10
#include "dltlib.cpp"
//——————————————————————————————————————————————

int main ()
{
   double p = 0.01;
   double l = 0.001;
   long int L = 1000000;
   double b = 1, d = 1;
   // STEP 1
   Network platform;                    // object representing parallel ↩
          platform

   // STEP 2
   // insert one−by−one the nodes that make up the machine. P0 by ↩
          default participates in
   // the computation
   platform.InsertNode ((char *) "P0", p, 0, (char *) NULL, l, ↩
          true);
   for (int i = 1; i < 10; i++)
     {
        char buff[10];
        sprintf (buff, "P%i", i);
        platform.InsertNode (buff, (i + 1) * p, 0, (char *) "P0", l↩
              , true);
     }

   // STEP 3
   // Solve the partitioning problem for 1−port, block−type ↩
          computation
   platform.SolveQuery (L, b, d);

   // print out the results, if the solution is valid
   if (platform.valid == 1)
     {
        cout << "Predicted execution time: " << platform.↩
              SimulateQuery (b, d, 0) << endl;

        // STEP 4
        // Compute nodes are stored in a public linked−list that ↩
              allows
        // rearrangement to the order of distribution and ↩
              collection
        cout << "Solution in terms of load percent :\n";
        Node *h = platform.head;
        while (h != NULL)
          {
```

```
50              // Single installment case...
51              cout << h->name << "\t" << h->part << endl;
52              h = h->next_n;
53            }
54          }
55       else
56          cout << "Solution could not be found\n";
57       return 0;
58  }
```

A noteworthy addition to the sample provided in the book, is line 10 that changes the maximum number of child nodes of a machine. The default of the library is only four, far less than the problem requirement of nine, which if left unchanged would cause the program to fail during execution with a protection fault.

This program can be compiled and run as follows:

```
$ g++ DLTexercise.cpp -o DLTexercise -I ../../DLTlib -lglpk
$ ./DLTexercise
Predicted execution time: 3414.18
Solution in terms of load percent :
P0      0.341418
P1      0.170709
P2      0.113806
P3      0.0853543
P4      0.0682834
P5      0.0569028
P6      0.0487738
P7      0.042677
P8      0.0379351
P9      0.0341416
```

Finally, testing for a 10 times slower communication medium is as simple as changing the `l` variable value of line 17 to `0.01`. The reported results are almost identical to the previous ones, as the communication cost is negligible in comparison to the computation:

```
Predicted execution time: 3414.22
Solution in terms of load percent :
P0      0.341422
P1      0.17071
P2      0.113806
P3      0.085354
P4      0.0682828
P5      0.056902
P6      0.0487729
P7      0.042676
P8      0.037934
P9      0.0341404
```

11. Derive the equivalent of Equations 8.47 and 8.48 for arbitrary $a$ and $c$, and implement them as part of a partitioning function.

**Answer**

The question is essentially to solve the N-port, block-type, single-installment distribution problem (i.e. what Section 8.3.3.1 covers) for the case when $P_0$ is not the load originating node. As such, $P_0$ has to incur communication cost both for the distribution and the result collection phases.

The solution is identical to the one in Exercise 3. The only difference is that while in Exercise 3 we assume that $P_0$ is not participating in the load processing, in this exercise we assume that $P_0$ does not hold the load originally, i.e. there is another master node that is only distributing the

data. We can reuse the solution of Exercise 3, with a simple change of the indices, so that all $part_i$ are expressed as functions of $part_0$, i.e.:

$$part_i = part_0 \frac{l_0(a+c) + p_0}{l_i(a+c) + p_i} + \frac{p_0 e_0 - p_i e_i + (l_0 - l_i)(b+d)}{L(l_i(a+c) + p_i)} \qquad (8.10)$$

and

$$part_0 = \frac{1 + L^{-1} \sum_{j=1}^{N-1} \frac{p_j e_j - p_0 e_0 + (l_j - l_0)(b+d)}{l_j(a+c) + p_j}}{(l_0(a+c) + p_0) \sum_{j=0}^{N-1} (l_j(a+c) + p_j)^{-1}} \qquad (8.11)$$

The following C++ code implements these two Equations, with function `NportBlockSingleInst` returning the $part_i$ and total execution time:

```cpp
1   #include <iostream>
2
3   using namespace std;
4   //*********************************************
5   double NportBlockSingleInst (double *p, double *e, double *l, ←
            double a, double b, double c, double d, int N, long L, ←
            double *part)
6   {
7     double lacp[N];   // should be replaced by dynamic allocation if←
             N is big
8     double sumTerm[N];
9     double totalTime;
10    double nomin, denom;
11
12    for (int i = 0; i < N; i++)
13      lacp[i] = l[i] * (a + c) + p[i];
14
15    for (int i = 1; i < N; i++)
16      sumTerm[i] = p[i] * e[i] - p[0] * e[0] + (l[i] - l[0]) * (b +←
              d);
17
18    nomin = 0;
19    for (int i = 1; i < N; i++)
20      nomin += sumTerm[i] / lacp[i];
21    nomin = 1 + nomin / L;
22
23    denom = 0;
24    for (int i = 0; i < N; i++)
25      denom += (1 / lacp[i]);
26    denom *= lacp[0];
27
28    part[0] = nomin / denom;
29
30    for (int i = 1; i < N; i++)
31      part[i] = part[0] * lacp[0] / lacp[i] - sumTerm[i] / (L * ←
              lacp[i]);
32
33    totalTime = p[0] * (part[0] * L + e[0]) + l[0] * (a * part[0] *←
              L + b) + l[0] * (c * part[0] * L + d);
34    return totalTime;
35  }
36
37
38  //*********************************************
39  int main ()
40  {
41    double p[4]  = { 1, 2, 3, 4 };
42    double e[4]  = { 1, 1, 1, 1 };
43    double b[4]  = { 1000, 1000, 1000, 1000 };
44    double l[4]  = { 1, 1, 1, 1 };
45    double part[4];
46
47    cout << "t : " << NportBlockSingleInst (p, e, l, 1, 0, 1, 0, 4,←
              1000000, part) << endl;
48    for (int i = 0; i < 4; i++)
```

```
49        cout << part[i] << " ";
50      cout << endl;
51
52      return 0;
53  }
```

12. The DLT examples presented in Sections 8.3.3.1 and 8.3.3.2, do not address an aspect of the problem which is significant : what is the optimum subset of nodes to use to process a load? Implement the heuristic algorithm you came up with in Exercise 2, to derive such a set.

Hint: nodes that should not be part of this set, get a negative load assignment $part_i$.

**Answer**

The effects of node choice and order are more significant when a one-port configuration is used. When N-port configurations are used, one needs to deal with the node selection problem, only when the communication cost to a node can exceed the total execution time on the other nodes.

In the following listing we implement the heuristic outlined in Exercise 2. The program tests for a maximum number of nodes, as specified in the command line, and reports the optimum number of nodes to be used and associated partitioning. The nodes are sorted (actually they are generated in this fashion) in descending order of their computing power.

```
1   #include <iostream>
2   #include <stdlib.h>
3
4   using namespace std;
5
6   //*********************************************
7   double OneportBlockSingleInst (double *p, double *e, double l, ↩
          double *b, int N, long L, double *part)
8   {
9     double pl[N]; // should be replaced by dynamic allocation if N ↩
            is big
10    double bSum[N];
11
12    double totalTime;
13    double nomin, denom;
14
15    for (int i = 0; i < N; i++)
16      pl[i] = p[i] + l;
17
18    bSum[0] = 0;
19    for (int i = 1; i < N; i++)
20      bSum[i] = b[i] + bSum[i - 1];
21
22    nomin = 0;
23    for (int i = 1; i < N; i++)
24      nomin += l / pl[i] * bSum[i];
25    nomin /= L;
26    nomin++;
27    for (int i = 1; i < N; i++)
28      nomin += l / pl[i];
29
30
31    denom = 0;
32    for (int i = 0; i < N; i++)
33      denom += 1 / pl[i];
34    denom *= pl[0];
35
36    part[0] = nomin / denom;
37
38    for (int i = 1; i < N; i++)
```

```
39        part[i] = part[0] * pl[0] / pl[i] - l / (L * pl[i]) * bSum[i]↩
               - l / pl[i];

40
41    totalTime = p[0] * part[0] * L;
42    return totalTime;
43  }

44
45  //***********************************************
46  int main (int argc, char **argv)
47  {
48    int N = atoi (argv[1]);

49
50    double p[N];
51    double e[N];
52    double b[N];
53    double l = .1;
54    double part[N];
55    long L = 100;

56
57    for (int i = 0; i < N; i++)
58      {
59        p[i] = (i + 1) * 1;
60        e[i] = 0.1;
61        b[i] = 1;
62      }

63
64    double bestT = p[0] * L + e[0], bestN = 1;
65    for (int i = 1; i < N; i++)
66      {
67        double t = OneportBlockSingleInst (p, e, l, b, i + 1, L, ↩
               part);
68        bool valid = true;
69        for (int j = 0; j < i + 1 && valid; j++)
70          if (part[j] < 0)
71            valid = false;
72        cout << i + 1 << " " << t << " " << valid << endl;        ↩
               // debugging output
73        if (t < bestT)  // valid does not need to be checked as per↩
               the comments in Exercise 2
74          {
75            bestN = i + 1;
76            bestT = t;
77          }
78      }

79
80    cout << "Best time " << bestT << ", best N " << bestN << endl;
81    OneportBlockSingleInst (p, e, l, b, bestN, L, part);  // solve ↩
               again to obtain part_i
82    for (int i = 0; i < bestN; i++)
83      cout << part[i] << " ";
84    cout << endl;

85
86    return 0;
87  }
```

A sample of the generated output is shown below:

```
$ ./heuristicsDLT 50
2  80.2 1
3  73.4615 1
4  70.039 1
5  67.9655 1
6  66.5785 1
7  65.5904 1
8  64.8559 1
9  64.2931 1
10 63.8524 1
...
23 62.1486 1
24 62.1394 1
25 62.1375 1
26 62.1422 0
27 62.1526 0
28 62.1683 0
```

```
...
50 63.2093 0
Best time 62.1375, best N 25
0.621375 0.0775835 0.0556876 0.0425501 0.0337917 0.0275358 ←
    0.0228438 0.0191945 0.016275 0.0138864 0.0118959 0.0102116 ←
    0.00876789 0.0075167 0.0064219 0.00545591 0.00459725 ←
    0.00382897 0.00313752 0.00251193 0.0019432 0.00142393 ←
    0.000947934 0.000510017 0.000105786
```