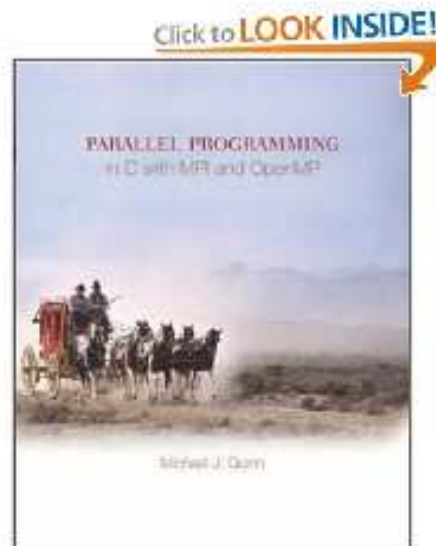


Lecture 12: Parallel quicksort algorithms

Overview

- Sequential quicksort algorithm
- Three parallel quicksort algorithms
- Chapter 14 in *Michael J. Quinn, Parallel Programming in C with MPI and OpenMP*



Recap of quicksort

- Given a list of numbers, we want to sort the numbers in an increasing order
 - The same as finding a suitable permutation
- Sequential quicksort algorithm: a recursive procedure
 - Select one of the numbers as pivot
 - Divide the list into two sublists: a “low list” containing numbers smaller than the pivot, and a “high list” containing numbers larger than the pivot
 - The low list and high list recursively repeat the procedure to sort themselves
 - The final sorted result is the concatenation of the sorted low list, the pivot, and the sorted high list

Example of quicksort

- Given a list of numbers: {79, 17, 14, 65, 89, 4, 95, 22, 63, 11}
- The first number, 79, is chosen as pivot
 - Low list contains {17, 14, 65, 4, 22, 63, 11}
 - High list contains {89, 95}
- For sublist {17, 14, 65, 4, 22, 63, 11}, choose 17 as pivot
 - Low list contains {14, 4, 11}
 - High list contains {64, 22, 63}
- ...
- {4, 11, 14, 17, 22, 63, 65} is the sorted result of sublist {17, 14, 65, 4, 22, 63, 11}
- For sublist {89, 95} choose 89 as pivot
 - Low list is empty (no need for further recursions)
 - High list contains {95} (no need for further recursions)
 - {89, 95} is the sorted result of sublist {89, 95}
- Final sorted result: {4, 11, 14, 17, 22, 63, 65, 79, 89, 95}

Illustration of quicksort

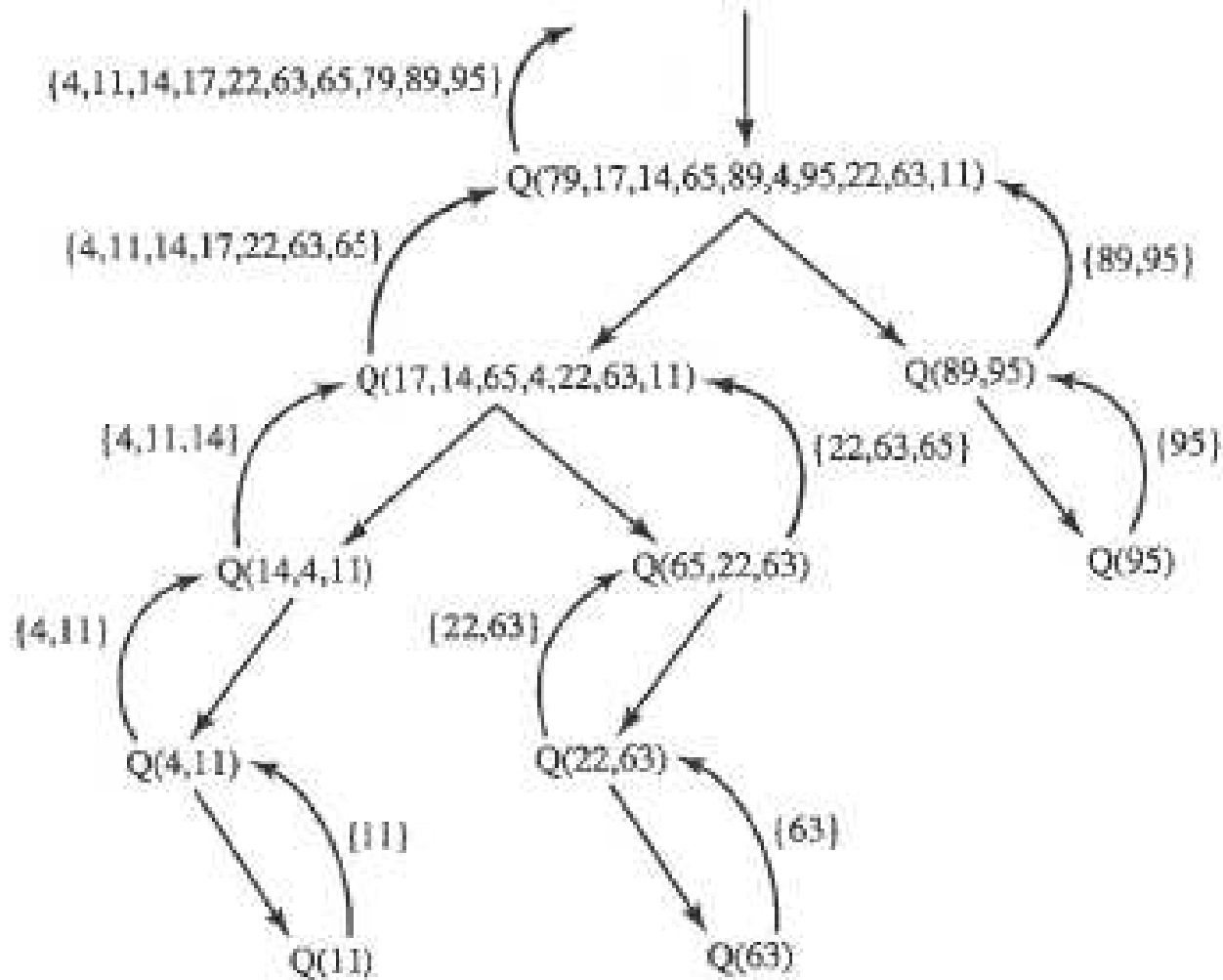


Figure 14.1 from Parallel Programming in C with MPI and OpenMP

Two observations

- Quicksort is generally recognized as the fastest sorting algorithm based on comparison of keys, in the average case
- Quicksort has some natural concurrency
 - The low list and high list can sort themselves concurrently

Parallel quicksort

- We consider the case of distributed memory
- Each process holds a segment of the unsorted list
 - The unsorted list is evenly distributed among the processes
- Desired result of a parallel quicksort algorithm:
 - The list segment stored on each process is sorted
 - The last element on process i 's list is smaller than the first element on process $i + 1$'s list

Parallel quicksort algorithm 1

- We randomly choose a pivot from one of the processes and broadcast it to every process
- Each process divides its unsorted list into two lists: those smaller than (or equal) the pivot, those greater than the pivot
- Each process in the upper half of the process list sends its “low list” to a partner process in the lower half of the process list and receives a “high list” in return
- Now, the upper-half processes have only values greater than the pivot, and the lower-half processes have only values smaller than the pivot.
- Thereafter, the processes divide themselves into two groups and the algorithm recurses.
- After $\log P$ recursions, every process has an unsorted list of values completely disjoint from the values held by the other processes.
 - The largest value on process i will be smaller than the smallest value held by process $i + 1$.
 - Each process can sort its list using sequential quicksort.

Illustration of parallel quicksort 1

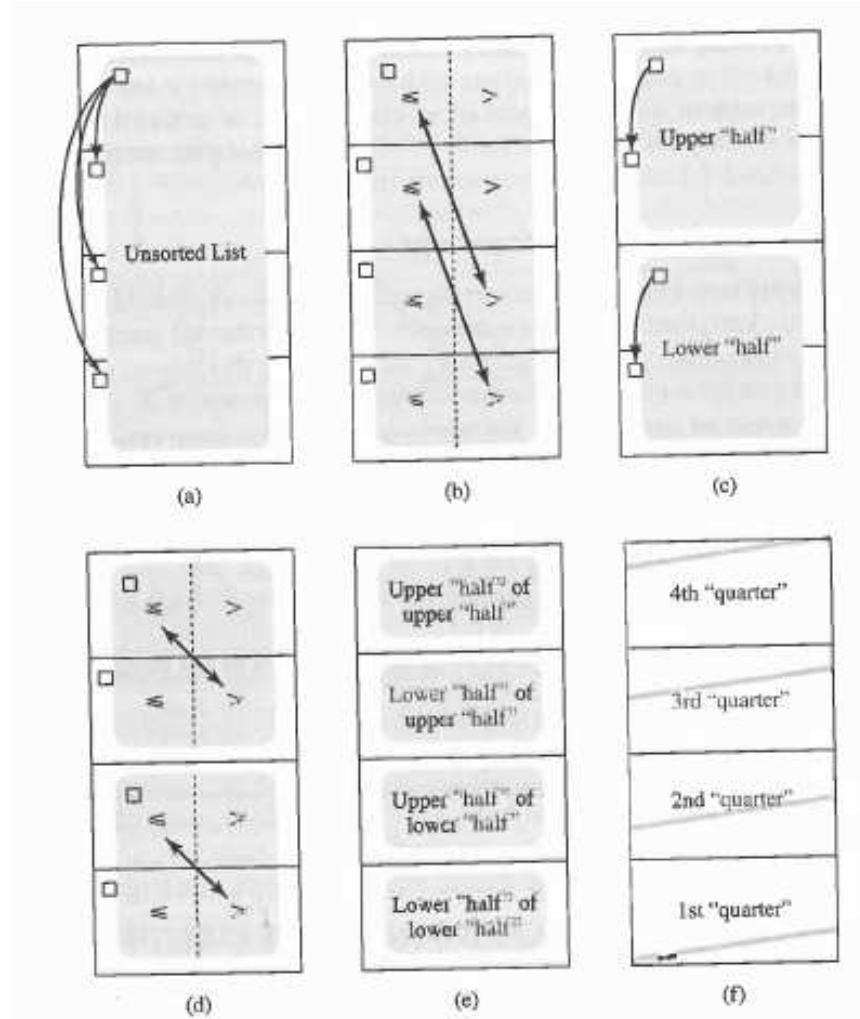


Figure 14.2 from Parallel Programming in C with MPI and OpenMP

Analysis of parallel quicksort 1

- This parallel quicksort algorithm is likely to do a poor job of load balancing
 - If the pivot value is not the median value, we will not divide the list into two equal sublists
 - Finding the median value is prohibitively expensive on a parallel computer
- The remedy is to choose the pivot value close to the true median!

Hyperquicksort – parallel quicksort 2

- Each process starts with a sequential quicksort on its local list
- Now we have a better chance to choose a pivot that is close to the true median
 - The process that is responsible for choosing the pivot can pick the median of its local list
- The three next steps of hyperquicksort are the same as in parallel quicksort 1
 - broadcast
 - division of “low list” and high list”
 - swap between partner processes
- The next step is different in hyperquicksort
 - One each process, the remaining half of local list and the received half-list are merged into a sorted local list
- Recursion within upper-half processes and lower-half processes . . .

Example of using hyperquicksort

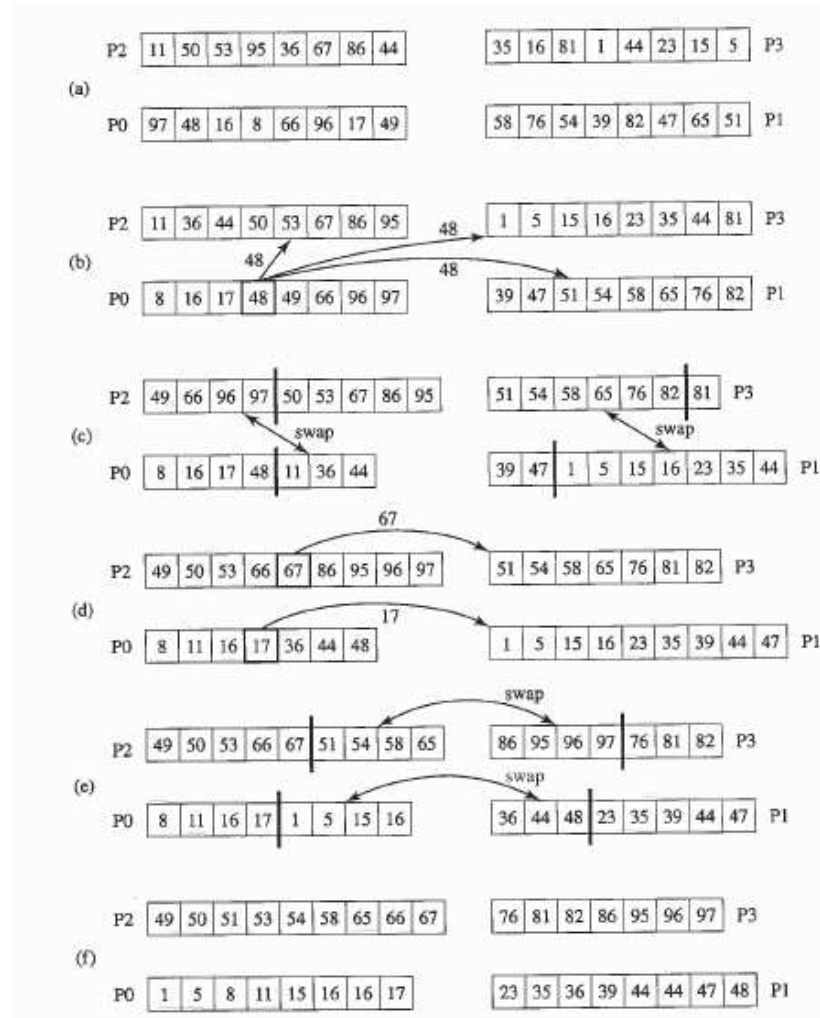


Figure 14.3 from Parallel Programming in C with MPI and OpenMP

Some observations about hyperquicksort

- $\log P$ steps are needed in the recursion
 - The expected number of times a value is passed from one process to another is $\frac{\log P}{2}$, therefore quite some communication overhead!
- The median value chosen from a local segment may still be quite different from the true median of the entire list
 - Load imbalance may still arise, although better parallel quicksort algorithm 1

Algorithm 3 – parallel sorting by regular sampling

- Parallel sorting by regular sampling (PSRS) has four phases
 1. Each process uses sequential quicksort on its local segment, and then selects data items at local indices $0, n/P^2, 2n/P^2, \dots, (P-1)(n/P^2)$ as a regular sample of its locally sorted block
 2. One process gathers and sorts the local regular samples. The process then selects $P-1$ pivot values from the sorted list of regular samples. The $P-1$ pivot values are broadcast. Each process then partitions its sorted sublist into P disjoint pieces accordingly.
 3. Each process i keeps its i th partition and sends the j th partition to process j , for all $j \neq i$
 4. Each process merges its P partitions into a single list

Example of using PSRS

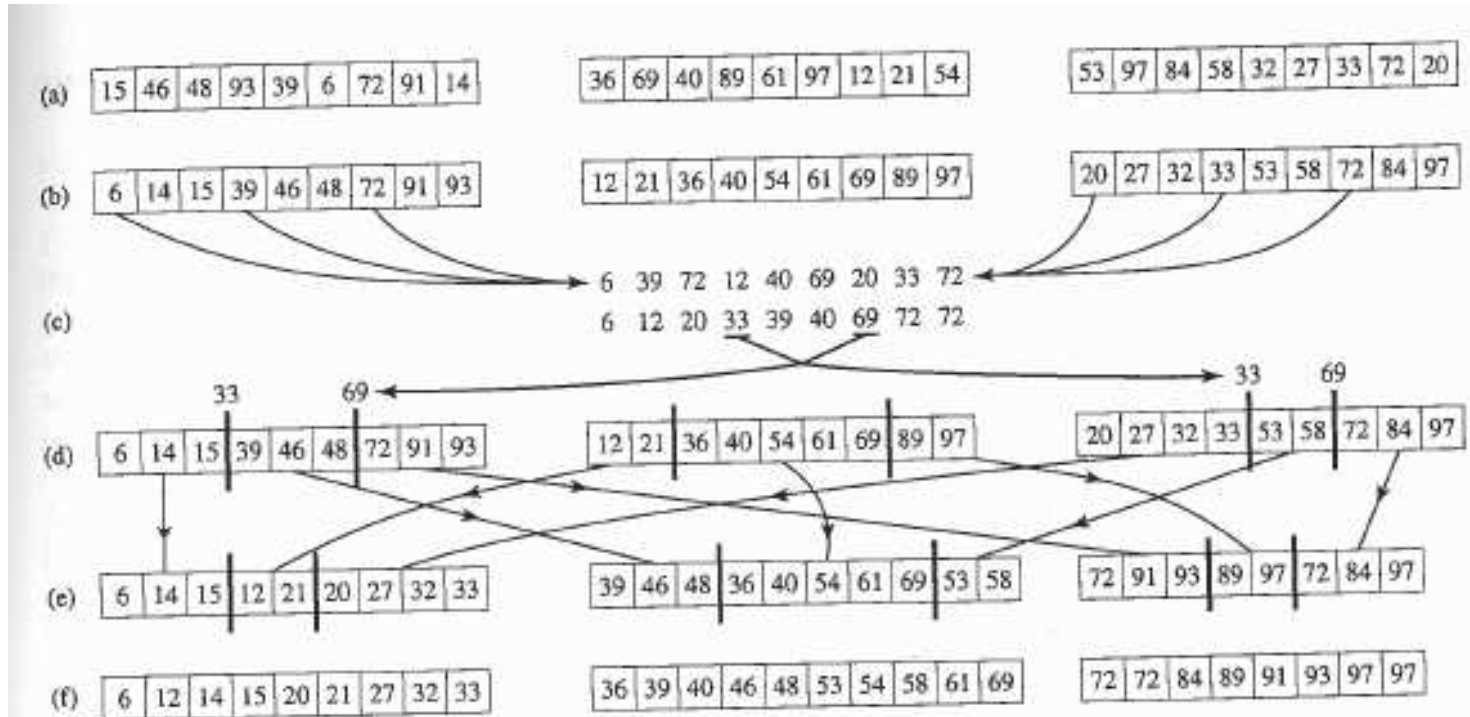


Figure 14.5 This example illustrates how three processes would sort 27 elements using the PSRS algorithm. (a) Original unsorted list of 27 elements is divided among three processes. (b) Each process sorts its share of the list using sequential quicksort. (c) Each process selects regular samples from its sorted sublist. A single process gathers these samples, sorts them, and broadcasts pivot elements from the sorted list of samples to the other processes. (d) Processes use pivot elements computed in step (c) to divide their sorted sublists into three parts. (e) Processes perform an all-to-all communication to migrate the sorted sublist parts to the correct processes. (f) Each process merges its sorted sublists.

Figure 14.5 from Parallel Programming in C with MPI and OpenMP

Three advantages of PSRS

- Better load balance (although perfect load balance can not be guaranteed)
- Repeated communications of a same value are avoided
- The number of processes does not have to be power of 2, which is required by parallel quicksort algorithm 1 and hyperquicksort

Optional exercise

- Make an MPI implementation of the parallel sorting by regular sampling algorithm

Recap; solving 2D wave equation

Mathematical model

$$\frac{\partial^2 u}{\partial t^2} = \frac{1}{2} \left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right)$$

- Spatial domain: unit square $(x, y) \in (0, 1) \times (0, 1)$
- Temporal domain: $0 < t \leq T$
- Boundary conditions: u is known on the entire boundary
- Initial conditions: $u(x, y, 0)$ is known, and “wave initially at rest”

Numerical method

- Uniform 2D mesh $(\Delta x, \Delta y)$
- Time step size Δt
- Central finite differences (same as in 1D)

$$\begin{aligned} u_{i,j}^{\ell+1} = & 2u_{i,j}^{\ell} - u_{i,j}^{\ell-1} \\ & + \frac{\Delta t^2}{2\Delta x^2} (u_{i-1,j}^{\ell} - 2u_{i,j}^{\ell} + u_{i+1,j}^{\ell}) \\ & + \frac{\Delta t^2}{2\Delta y^2} (u_{i,j-1}^{\ell} - 2u_{i,j}^{\ell} + u_{i,j+1}^{\ell}) \\ & i = 1, 2, \dots, M, \quad j = 1, 2, \dots, N \end{aligned}$$

Sequential code (part 1)

```
#include <stdio.h>
#include <malloc.h>
#include <math.h>

int main (int nargs, char** args)
{
    int M, N, L, i, j;
    double T, dx, dy, dt, x, y, t;
    double **up, **u, **um, **tmp, *up_data, *u_data, *um_data;

    M = 149; N = 99; L = 200;
    T = 1.0; dx = 1./(M+1); dy = 1./(N+1); dt = T/L;

    up_data = (double*)malloc((M+2)*(N+2)*sizeof(double));
    u_data = (double*)malloc((M+2)*(N+2)*sizeof(double));
    um_data = (double*)malloc((M+2)*(N+2)*sizeof(double));

    up = (double**)malloc((N+2)*sizeof(double*));
    u = (double**)malloc((N+2)*sizeof(double*));
    um = (double**)malloc((N+2)*sizeof(double*));

    for (j=0; j<=N+1; j++) {
        up[j] = &(up_data[j*(M+2)]);
        u[j] = &(u_data[j*(M+2)]);
        um[j] = &(um_data[j*(M+2)]);
    }
}
```

Sequential code (part 2)

```
for (j=0; j<=N+1; j++) { /* Enforcing initial condition 1 */
    y = j*dy;
    for (i=0; i<=M+1; i++) {
        x = i*dx;
        um[j][i] = sin(2*M_PI*(x+y));
    }
}

for (j=1; j<=N; j++) /* Enforcing initial condition 2 */
    for (i=1; i<=M; i++)
        u[j][i] = um[j][i]+((dt*dt)/(4*dx*dx))*(um[j][i-1]-2*um[j][i]+um[j][i+1])
        +((dt*dt)/(4*dy*dy))*(um[j-1][i]-2*um[j][i]+um[j+1][i]);

/* enforcing boundary conditions for t = dt */
t = dt;
for (j=0; j<=N+1; j++) {
    y = j*dy;
    u[j][0] = sin(2*M_PI*(y+t)); /* x = 0 */
    u[j][M+1] = sin(2*M_PI*(1+y+t)); /* x = 1 */
}
for (i=1; i<=M; i++) {
    x = i*dx;
    u[0][i] = sin(2*M_PI*(x+t)); /* y = 0 */
    u[N+1][i] = sin(2*M_PI*(x+1+t)); /* y = 1 */
}
```

Sequential code (part 3)

```
while (t<T) {
    t += dt;

    for (j=1; j<=N; j++)
        for (i=1; i<=M; i++) /* interior points */
            up[j][i] = 2*u[j][i]-um[j][i]
                +((dt*dt)/(2*dx*dx))*(u[j][i-1]-2*u[j][i]+u[j][i+1])
                +((dt*dt)/(2*dy*dy))*(u[j-1][i]-2*u[j][i]+u[j+1][i]);

    /* enforcing boundary conditions */
    for (j=0; j<=N+1; j++) {
        y = j*dy;
        up[j][0] = sin(2*M_PI*(y+t)); /* x = 0 */
        up[j][M+1] = sin(2*M_PI*(1+y+t)); /* x = 1 */
    }
    for (i=1; i<=M; i++) {
        x = i*dx;
        up[0][i] = sin(2*M_PI*(x+t)); /* y = 0 */
        up[N+1][i] = sin(2*M_PI*(x+1+t)); /* y = 1 */
    }

    /* data shuffle */
    tmp = um;
    um = u;
    u = up;
    up = tmp;
}
```

Sequential code (part 4)

```
free (up_data); free (up);  
free (u_data); free (u);  
free (um_data); free (um);
```


Parallelization

Parallelization

- Each subdomain is responsible for a rectangular region of the $M \times N$ interior points
- One layer of ghost points is needed in the local data structure
- Serial local computation + exchange of values for the ghost points