

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ  
ИМЕНИ М.В. ЛОМОНОСОВА  
КУРС: «ПАРАЛЛЕЛЬНОЕ ПРОГРАММИРОВАНИЕ И  
СУПЕРКОМПЬЮТЕРЫ»

# **Рекурсивная координатная бисекция**

**Выполнил:**

**Сактаганов Нуржан (аспирант)**

**Дата подачи: 05.01.2018**

**Москва 2018**

## Table of Contents

Постановка задачи.....	3
Метод решения.....	3
Описание вычислительной системы.....	4
Анализ полученных результатов.....	5
Приложение А. Последовательная программа.....	7
Приложение Б. Параллельная программа.....	8

## Постановка задачи

Дана двумерная сетка, топологически эквивалентная индексному прямоугольнику  $n_1 \times n_2$ . Необходимо разбить сетку на  $k$  частей (доменов). Узел сетки с индексами  $(i, j)$  связан с существующими соседними по  $i$  и  $j$  узлами:  $(i-1, j)$ ,  $(i+1, j)$ ,  $(i, j-1)$ ,  $(i, j+1)$ .

Сетка представляется в виде одмерного массива структур Point.

```
Point{
```

```
    float coord[2];
```

```
    int index;
```

```
}
```

```
P[n1*n2];
```

Точки данной сетки имеют координаты  $P[i*n_2+j].coord[0] = x(i, j)$ ,  $P[i*n_2+j].coord[1] = y(i, j)$ , где  $i = 0, \dots, n_1-1$ ,  $j = 0, \dots, n_2-1$ . Индекс определяется соотношением  $P[i*n_2+j].index = i*n_2+j$ .

Для разбиения сетки нужно написать параллельную программу для вычислительной системы с распределенной памятью, обеспечивающую разбиение методом рекурсивной координатной бисекции узлов сетки на  $k$  доменов. Число вершин в доменах должно быть равно с точностью до одной вершины величине  $(n_1*n_2/k)$ .

На выходе должны получить:

1. файл, содержащий  $n_1*n_2$  строк, в каждой из которых 5 чисел:  $i$   $j$   $X_{ij}$   $Y_{ij}$   $d$  // номера вершины в сетке, координаты вершины, номер домена  $[0, \dots, k-1]$
2. число разрезанных рёбер
3. время выполнения декомпозиции (только декомпозиции, без учета формирования сетки и вывода результатов).

## Метод решения

Сначала опишем последовательный алгоритм. Далее на основе последовательного алгоритма будем описывать параллельный алгоритм.

**Последовательный алгоритм** очень простой. На вход даются массив структур в описанном выше формате и количество доменов на которое нужно разбить. Мы сначала сортируем точки по заданной координате. После этого раскладываем  $k = k_1 + k_2$ , где  $k_1 = (k + 1) / 2$ ,  $k_2 = k - k_1$ . Далее делим массив в пропорции  $k_1:k_2$ . Меняем координату, по которой нужно сортировать массив, и рекурсивно вызываем на первой половине с параметром  $k := k_1$ , и на второй половине  $k := k_2$ . После чего обратно вооставнавливаем координату для сортировки. Рекурсия останавливается в моменте, когда у нас  $k = 1$ . Это означает, что все точки массива принадлежат одному и тому же домену.

**Параллельный алгоритм** делает по сути точно то же самое что и последовательный. Будем считать, что каждый процесс запущен на отдельном узле, и что на узле находится один

процессор.

Допустим, что распределенный массив расположен на  $n$  процессорах. На каждом процессоре находится по  $part\_size$  точек. И пусть имеется число  $k$  — количество доменов, на которое нужно разбить. Тогда, сначала сортируем глобальный массив по заданной координате (алгоритм сортировки Бэтчера). Далее, точно так же раскладываем  $k=k_1+k_2$ . Делим глобальный массив в пропорции  $k_1:k_2$ . Находим разбивающий элемент — первый элемент второй половины глобального массива. Множество процессоров надо разбить на два подмножества, первое подмножество будет обрабатывать первую половину, второе — вторую. Узел, который содержит в себе разбивающий элемент отнесем ко второму подмножеству. Поэтому, все элементы до разбивающего элемента нужно равномерно распределить по процессорам первого подмножества. Причем, некоторые процессоры первого подмножества нужно дополнить фиктивными элементами, чтобы у всех было одинаковое количество элементов (это требуется для сортировки бэтчера). А в процессорах второго подмножества нужно сделать перераспределение элементов в связи с уменьшением их количества в этом множестве. После этого меняем координату для сортировки и применяем описанный алгоритм рекурсивно для двух подмножеств процессоров с соответствующими параметрами  $k_1$  и  $k_2$ . После возврата из рекурсии надо не забыть воостановить координату для сортировки. Когда  $k$  или количество процессоров будет равно 1 — это означает конец рекурсии параллельного алгоритма, теперь уже на каждом узле выполним последовательный алгоритм с имеющим значением  $k$ .

Случай, когда разбивающий элемент будет на нулевом процессоре обрабатывается чуть по-другому. В этом случае нулевой процессор отнесем к первому подмножеству, а остальные — ко второму. Все элементы, начиная с разбивающего, распределим по процессорам из второго подмножества, пометив отправленные как фиктивные и добавив фиктивные по необходимости.

**Подсчет числа разрезанных ребер.** Возьмем произвольный домен. Для каждой его точки проверим принадлежность каждой из четырех соседних точек данному домену. Если какой-то сосед не будет принадлежать данному домену — это означает что есть разрез ребра. Если для каждого домена посчитаем количество разрезов, и потом просуммируем — получим удвоенное количество всех разрезанных ребер.

Будем считать, что  $k$  больше количества процессоров. Тогда не будет случая, когда один домен будет размещен на нескольких процессорах. Это означает, что каждый домен будет целиком помещаться внутри одного узла. Тогда мы можем для каждого домена локально посчитать количество разрезов. После чего легко получить количество всех разрезанных ребер.

## Описание вычислительной системы

Вычислительная система представляет собой комплекс IBM Blue Gene/P. В составе имеется две стойки по 1024 четырехъядерных процессоров в каждой стойке. Пиковая производительность — 27.9 терафлопс. Размер оперативной памяти — 2 Гб на каждый процессор. Каждое ядро модели PowerPC 450 с рабочей частотой 850 МГц, с 32 битной адресацией. Производительность ядра — 3.4 Гфлопс.

Описание коммуникационной сети:

- сеть общего назначения, объединяющие все вычислительные узлы; предназначена для операций типа «точка-точка»
- вычислительный узел имеет двунаправленные связи с шестью соседями
- пропускная способность каждого соединения — 425 MB/s (5,1 GB/s для всех 12 каналов)
- латентность (ближайший сосед):
  - 32-байтный пакет: 0,1  $\mu$ s
  - 256-байтный пакет: 0,8  $\mu$ s

## Анализ полученных результатов

Ниже приведена таблица времен (в секундах) запусков на различном наборе сеток и процессоров (число доменов  $k = 256$ ).

Сетка	1 процессор	8 процессоров	32 процессора	128 процессоров
4k x 2.5k	79.0919	11.055	3.56367	1.3142
4k x 5k	167.652	23.4848	7.44001	2.71177

Посчитаем ускорение.

Сетка	1 процессор	8 процессоров	32 процессора	128 процессоров
4k x 2.5k	1.0	7.1544	22.194	60.183
4k x 5k	1.0	7.1387	22.5338	61.8238

Посчитаем эффективность.

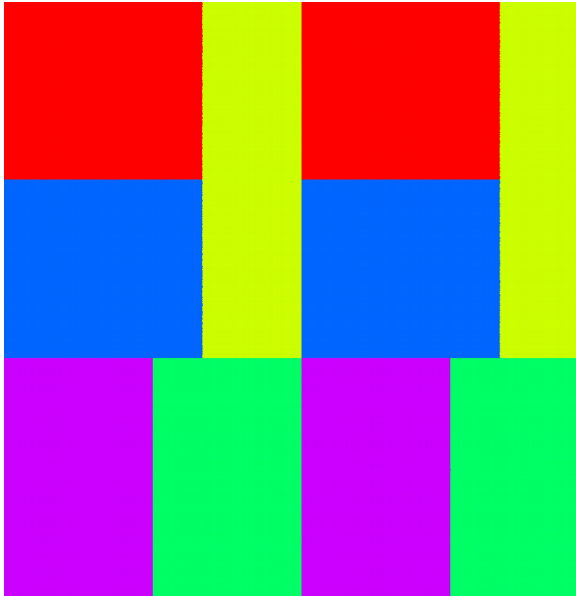
Сетка	1 процессор	8 процессоров	32 процессора	128 процессоров
4k x 2.5k	100%	89.43%	69.36%	48.02%
4k x 5k	100%	89.23%	70.42%	48.30%

Из таблиц видим, что программа хорошо параллелится и время выполнения продолжает сокращаться с ростом числа процессоров.

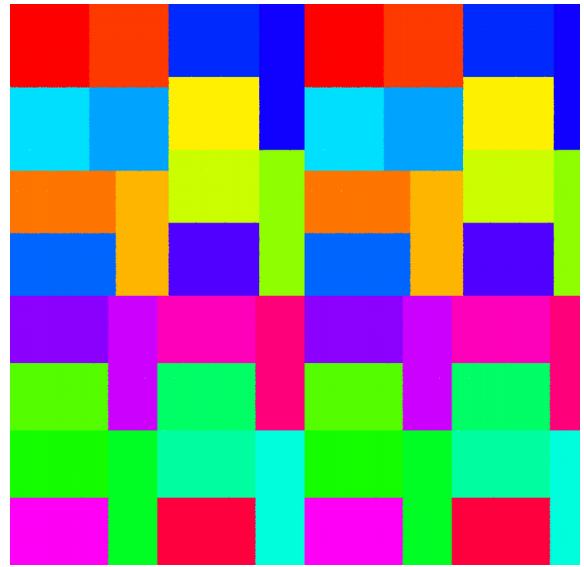
Для разнообразия сетка генерировалась псевдослучайно для каждого запуска (в точке  $(i, j)$  значение  $x$  и  $y$  — равномерное отклонение от  $i$  и  $j$  на 0.25), число разрезов тоже колеблется в некотором промежутке. Следующая таблица хорошо это демонстрирует.

Сетка	1 процессор	8 процессоров	32 процессора	128 процессоров
4k x 2.5k	118058	118124	117946	118835
4k x 5k	150875	151328	150984	154872

В качестве примера работы алгоритма на сетке 1000x1000 приведены две иллюстрации. В первом случае разбиение на 10 доменов с использованием 5 процессоров. Во втором случае разбиение на 50 доменов с использованием 20 процессоров.



*Illustration 1: 1000x1000,  $d=10$ ,  $p=5$*



*Illustration 2: 1000x1000,  $d=50$ ,  $p=20$*

## Приложение А. Последовательная программа

```
1. #pragma once
2. #include <algorithm>
3. #include "Point.hpp"
4.
5. //struct represents that point with given index
6. //belongs to given domain
7. typedef struct _point_domain_t {
8.     Point point;
9.     int domain;
10. } point_domain_t;
11.
12. /*
13. points - given array of points
14. points_domain - elements belonging
15. length - length of points array
16. k - wanted domains count
17. lowest_domain - lowest available domain number
18. */
19. void decompose(
20.     Point *points,
21.     point_domain_t *points_domain,
22.     const int length,
23.     const int k,
24.     const int lowest_domain)
25. {
26.     if (k == 1) {
27.         for (int i = 0; i < length; ++i) {
28.             points_domain[i].domain = lowest_domain;
29.             points_domain[i].point = points[i];
30.         }
31.         return;
32.     }
33.
34.     std::sort(points, points + length);
35.
36.     const int k1 = (k + 1) / 2;
37.     const int k2 = k - k1;
38.     const int l1 = 1.0 * k1 / k * length;
39.     const int l2 = length - l1;
40.
41.     Point::switch_sort_way();
42.     decompose(points, points_domain, l1, k1, lowest_domain);
43.     decompose(points + l1, points_domain + l1, l2, k2, lowest_domain + k1);
44.     Point::switch_sort_way();
45. }
```

## Приложение Б. Параллельная программа

```
1. #include <cstring>
2. #include "Point.hpp"
3. #include "decompose.hpp"
4. #include "parallel_sort.hpp"
5.
6. //represents a 'vector' with given length (i.e. size) and ptr
7. typedef struct _points_domain_t {
8.     int k;
9.     int procs_num;
10.    int size;
11.    point_domain_t *info;
12.} points_domain_t;
13.
14. points_domain_t _parallel_decompose(
15.    MPI::Intracomm &comm,
16.    MPI::Datatype &datatype,
17.    Point *points,
18.    int part_size,
19.    int k,
20.    int lowest_domain,
21.    int *sendcounts_buffer, // for using in Scatterv
22.    int *displacements_buffer // for using in Scatterv
23.);
24.
25. int calculate_max_part_size(int procs_num, int part_size, int k);
26.
27. points_domain_t parallel_decompose(
28.    MPI::Intracomm &comm,
29.    MPI::Datatype &datatype,
30.    Point *points,
31.    int part_size,
32.    int k,
33.    int lowest_domain)
34. {
35.     points_domain_t points_domain;
36.
37.     const int procs_num = comm.Get_size();
38.     const int max_part_size = calculate_max_part_size(procs_num, part_size, k);
39.
40.     Point *points_copy = new Point[max_part_size];
41.     memcpy(points_copy, points, part_size * sizeof(Point));
42.
43.     //Now we can do everyting with this copy.
44.     //However, passed user's array will remain unchanged.
45.
46.     int *sendcounts_buffer = new int [procs_num];
47.     int *displacements_buffer = new int [procs_num];
48.
49.     points_domain = _parallel_decompose(
50.         comm, datatype, points_copy,
51.         part_size, k, lowest_domain,
```



```

52.     sendcounts_buffer,
53.     displacements_buffer
54. );
55.
56. delete [] sendcounts_buffer;
57. delete [] displacements_buffer;
58.
59. delete [] points_copy;
60.
61. return points_domain;
62.};
63.
64.points_domain_t _parallel_decompose(
65. MPI::Intracomm &comm,
66. MPI::Datatype &datatype,
67. Point *points,
68. int part_size,
69. int k,
70. int lowest_domain,
71. int *sendcounts_buffer, // for using in Scatterv
72. int *displacements_buffer // for using in Scatterv
73.)
74.{
75. const int procs_num = comm.Get_size();
76. const int rank = comm.Get_rank();
77.
78. points_domain_t points_domain;
79.
80. if (procs_num == 1 || k == 1) { // there is no more processes or domains
81.     if (k > 1) {
82.         for (int i = part_size - 1; i >= 0; --i) {
83.             if (points[i].index != -1) continue;
84.             std::swap(points[i], points[part_size - 1]);
85.             --part_size;
86.         }
87.     }
88.     points_domain.k = k;
89.     points_domain.procs_num = procs_num;
90.     points_domain.size = part_size;
91.     points_domain.info = new point_domain_t[part_size];
92.     decompose(points, points_domain.info, part_size, k, lowest_domain);
93.     return points_domain;
94. }
95.
96. std::sort(points, points + part_size); // sort each local array
97. parallel_sort(comm, datatype, points, part_size);
98.
99. //now we have global sorted array of points
100.
101. // searching global middle element's coordinates
102. const long long global_length = part_size * procs_num;
103. const int k1 = (k + 1) / 2;
104. const int k2 = k - k1;

```

```

105. const long long l1 = 1.0 * k1 / k * global_length;
106. //const long long l2 = global_length - l1;
107.
108. const int middle_proc = l1 / part_size;
109. const int middle_index = l1 % part_size;
110.
111. MPI::Cartcomm new_comm;
112. int color; // 0 - left half, 1 - right half
113. bool i_am_sender = rank == middle_proc;
114. bool i_am_receiver;
115. int send_count;
116. int rcv_procs;
117. Point *send_points_start;
118.
119. if (middle_proc > 0) {
120.     // we have to send to left
121.     i_am_receiver = rank < middle_proc;
122.     color = rank >= middle_proc;
123.     send_count = middle_index; // elements in position 0, ..., middle_index-1
124.     send_points_start = points; // send from beginning
125.     rcv_procs = middle_proc;
126. } else { // middle_proc == 0
127.     // no way to send to left, only to right
128.     i_am_receiver = rank > middle_proc;
129.     color = rank > middle_proc;
130.     send_count = part_size - middle_index; // elements in position middle_index, ..., part_size-1
131.     send_points_start = points + middle_index; // send from middle_index and so on
132.     rcv_procs = procs_num - 1;
133. }
134.
135. // communicate only if required
136. if (send_count > 0) {
137.     memset(sendcounts_buffer, 0, procs_num * sizeof(int));
138.     memset(displacements_buffer, 0, procs_num * sizeof(int));
139.
140.     int offset = middle_proc == 0;
141.
142.     // init sendcounts
143.     for (int i = 0; i < rcv_procs; ++i) {
144.         sendcounts_buffer[i + offset] = send_count / rcv_procs;
145.     }
146.
147.     for (int i = 0; i < send_count % rcv_procs; ++i) {
148.         sendcounts_buffer[i + offset] += 1;
149.     }
150.
151.     // init displacements
152.     displacements_buffer[offset] = 0;
153.     for (int i = 1; i < rcv_procs; ++i) {
154.         displacements_buffer[i + offset] =
155.             displacements_buffer[i-1 + offset]
156.             + sendcounts_buffer[i-1 + offset];

```

```

157.     }
158.
159.     comm.Scatterv(
160.         send_points_start, // we send from here
161.         sendcounts_buffer,
162.         displacements_buffer,
163.         datatype,
164.         points + part_size, // we receive to this buffer
165.         sendcounts_buffer[rank], // we receive this amount of points
166.         datatype,
167.         middle_proc
168.     );
169.
170.     int receiver_part_size = 0;
171.     // in all receivers indentify phantom elements
172.     if (i_am_receiver) {
173.         const int rcv_n = send_count / rcv_procs + (send_count % rcv_procs > 0);
174.         receiver_part_size = part_size + rcv_n;
175.         const int phantom_elements = rcv_n - sendcounts_buffer[rank];
176.
177.         for (int i = 0; i < phantom_elements; ++i) {
178.             points[receiver_part_size-1-i].index = -1;
179.         }
180.     }
181.
182.     const bool can_compensate = middle_proc > 0 && (procs_num - rcv_procs > 1);
183.
184.     if (i_am_sender && !can_compensate) {
185.         // make sent elements phantom
186.         for (int i = 0; i < send_count; ++i) send_points_start[i].index = -1;
187.     }
188.
189.     if (can_compensate) {
190.         // compensate sent elements
191.         const int new_procs_num = procs_num - rcv_procs;
192.         const long long new_length = 1L * part_size * new_procs_num - send_count;
193.         const int new_part_size = new_length / new_procs_num + (new_length % new_procs
_num > 0);
194.         const int new_middle_proc_points = part_size - send_count;
195.         const int to_compensate = new_part_size - new_middle_proc_points;
196.         const int compensators = new_procs_num - 1;
197.
198.         int *rcv_counts = sendcounts_buffer;
199.         memset(rcv_counts, 0, procs_num * sizeof(int));
200.         memset(displacements_buffer, 0, procs_num * sizeof(int));
201.
202.         const int offset = middle_proc + 1;
203.         for (int i = 0; i < compensators; ++i) {
204.             rcv_counts[offset + i] = to_compensate / compensators;
205.         }
206.
207.         for (int i = 0; i < to_compensate % compensators; ++i) {
208.             rcv_counts[offset + i] += 1;

```

```

209.     }
210.
211.     displacements_buffer[offset] = 0;
212.     for (int i = 1; i < compensators; ++i) {
213.         displacements_buffer[offset + i] =
214.             displacements_buffer[offset + i-1]
215.             + recv_counts[offset + i-1];
216.     }
217.
218.     if (rank == middle_proc) {
219.         memmove(points, points + send_count, new_middle_proc_points * sizeof(Point));
220.     }
221.
222.     comm.Gatherv(
223.         points + part_size - recv_counts[rank], // from the tail
224.         recv_counts[rank], // how many
225.         datatype,
226.         points + new_middle_proc_points,
227.         recv_counts,
228.         displacements_buffer,
229.         datatype,
230.         middle_proc
231.     );
232.
233.     // make redistributed as phantom
234.     for (int i = 0; i < recv_counts[rank]; ++i){
235.         points[part_size - 1 - i].index = -1;
236.     }
237.
238.     if (rank >= middle_proc) {
239.         part_size = new_part_size;
240.     }
241. }
242.
243. if (i_am_receiver) {
244.     part_size = receiver_part_size;
245. }
246. }
247.
248. if (color == 0) {
249.     k = k1;
250. } else { // color == 1
251.     k = k2;
252.     lowest_domain += k1;
253. }
254.
255. new_comm = comm.Split(color, rank);
256.
257. Point::switch_sort_way();
258. points_domain = _parallel_decompose(
259.     new_comm, datatype, points,
260.     part_size, k, lowest_domain,

```

```

261.     sendcounts_buffer,
262.     displacements_buffer
263. );
264. Point::switch_sort_way();
265.
266. new_comm.Free();
267.
268. return points_domain;
269.}
270.
271./*    Trying predict max memory amount to prevent array reallocation.
272.    So, algorithm is recursive. In worst case we have to redistribute
273.    send_n := part_size-1 elements from sender. There is always only one sender
274.    in each recursion step. Let's consider that the amount of receivers is recv_p.
275.    It means, every receiver has to add at most ceil(send_n / recv_p) elements.
276.    There are at least recv_p := n_procs / 2 receivers, we finally get new part_size
277.    for each receiver: part_size := part_size + ceil(send_n / (n_procs / 2)).
278.    Let's consider that we already have max recursion depth in recursion_depth variable.
279.    Then we can describe algorithm in pseudocode.
280.
281.    def max_part_size(part_size, recursion_depth, n_procs):
282.        while recursion_depth > 0:
283.            send_n = part_size - 1
284.            recv_p = n_procs / 2
285.            part_size += ceil(send_n / recv_p)
286.            n_procs = ceil(n_procs / 2)
287.        return part_size
288.
289.    Recursion depth depends from k and n_procs.
290.    Exactly, it equals to rounded up value of log(min(k, procs_num)):
291.
292.    recursion_depth := ceil(log(min(k, n_procs), base=2)).
293.
294.    But it's too roughly. So, by modeling part_size changing,
295.    we can get exact value of max part size.
296.    */
297.int calculate_max_part_size(int procs_num, int part_size, int k)
298.{
299.    if (procs_num == 1 || k == 1) return part_size;
300.    // searching global middle element's coordinates
301.    const long long global_length = part_size * procs_num;
302.    const int k1 = (k + 1) / 2;
303.    const int k2 = k - k1;
304.    const long long l1 = 1.0 * k1 / k * global_length;
305.    //const long long l2 = global_length - l1;
306.
307.    const int middle_proc = l1 / part_size;
308.    const int middle_index = l1 % part_size;
309.
310.    int procs1, procs2;
311.    int p1_size, p2_size;
312.    p1_size = p2_size = part_size;
313.

```

```
314. if (middle_proc > 0) {
315.     const int send_n = middle_index;
316.     const int recv_p = middle_proc;
317.     const int recv_n = send_n / recv_p + (send_n % recv_p > 0);
318.     p1_size += recv_n;
319.
320.     procs1 = middle_proc;
321.     procs2 = procs_num - procs1;
322. } else { // middle_proc == 0
323.     const int send_n = part_size - middle_index;
324.     const int recv_p = procs_num - 1;
325.     const int recv_n = send_n / recv_p + (send_n % recv_p > 0);
326.     p2_size += recv_n;
327.
328.     procs1 = 1;
329.     procs2 = procs_num - 1;
330. }
331.
332. p1_size = calculate_max_part_size(procs1, p1_size, k1);
333. p2_size = calculate_max_part_size(procs2, p2_size, k2);
334.
335. part_size = p1_size > p2_size ? p1_size : p2_size;
336. return part_size;
337.}
```