

Project of Parallel Programming (2024 Spring)

PANGBO

ADI

For each computation, there is data dependency among only one dimension, so we can compute parallelly among the other two dimensions.

We will partition the whole 3-dimensional array in 3 dimensions, as shown in Fig.1, so each of the subblock will have a 3-dimensional id (i, j, k) . Assume the size of the whole array is $N \times N \times N$, and there are p processors, where p is a square number. The size of a subblock is $\frac{N^3}{\sqrt{p}}$ and each processor will be assigned with \sqrt{p} subblocks.

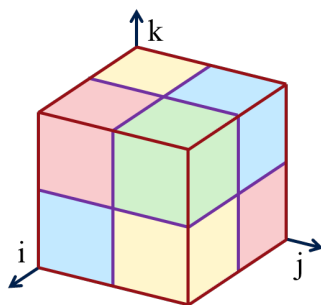


Fig. 1. Partition of the 3-dimensional array

Since we hope that all processors can run simultaneously at the same time, for a slice in any dimension, the $\frac{N^2}{p}$ sub-blocks it contains need to be allocated to p processes respectively. This requires that the subblock ids to which each processor is assigned cannot have the same i or j or k .

So here we come up with a solution to assign the subblocks to processors. We will assign the subblock (i, j, k) to processor $p_{(i,j,k)}$ where $p_{(i,j,k)} = [(k - i) \bmod \sqrt{p}] \times \sqrt{p} + [(j - i) \bmod \sqrt{p}]$. To intuitively understand this, for a single processor, everytime i dimension increased by 1, it will also increase the position on the j and k dimension by 1, so we can ensure that the subblocks with the same position in any dimension will not be assigned to the same processor.

Now suppose we are processing among the i dimension and a processor reaches the boundary of the subblock in the i dimension. Assume the processor is currently at subblock (i, j, k) , it will then calculate the owner $p_{(i+1,j,k)}$ of the next subblock $(i + 1, j, k)$. Then it will send the data of the boundary of the subblock (i, j, k) to the processor $p_{(i+1,j,k)}$, and receive the data of the boundary of the subblock $(i - 1, (j - 1) \bmod \sqrt{p}, (k - 1) \bmod \sqrt{p})$.

Author's address: pangbo.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

Since MPI communication can only send contiguous data, we need to pack the data before sending and unpack the data after receiving. The following is the code for packing and unpacking the data.

```

56 1 void pack_array(double *buffer, int i0, int j0, int k0, int i1, int j1, int k1){
57 2     int index = 0;
58 3     for (int k = k0; k <= k1; k++){
59 4         for (int j = j0; j <= j1; j++){
60 5             for (int i = i0; i <= i1; i++){
61 6                 buffer[index++] = A[k][j][i];
62 7             }
63 8         }
64 9     }
65 10 }
66 11
67 12 void unpack_array(double *buffer, int i0, int j0, int k0, int i1, int j1, int k1){
68 13     int index = 0;
69 14     for (int k = k0; k <= k1; k++){
70 15         for (int j = j0; j <= j1; j++){
71 16             for (int i = i0; i <= i1; i++){
72 17                 A[k][j][i] = buffer[index++];
73 18             }
74 19         }
75 20     }
76 21 }

```

The code of computing among the i dimension is shown below.

```

85 1 for (block_i = 0; block_i < DIM_SIZE; block_i++){
86 2     block_j = block_k = -1;
87 3     block_id(id, &block_i, &block_j, &block_k);
88 4     if (block_i != 0) {
89 5         int source = block_owner(block_i - 1, block_j, block_k);
90 6         MPI_Recv(recv_buffer, BLOCK_SIZE * BLOCK_SIZE, MPI_DOUBLE, source, 0, MPI_COMM_WORLD,
91 7             MPI_STATUS_IGNORE);
92 8         unpack_array(recv_buffer, BLOCK_LOW(block_i) - 1, BLOCK_LOW(block_j), BLOCK_LOW(block_k),
93 9             BLOCK_LOW(block_i) - 1, BLOCK_HIGH(block_j), BLOCK_HIGH(block_k));
94 10     }
95 11     for (int i = BLOCK_LOW(block_i); i <= BLOCK_HIGH(block_i); i++){
96 12         {
97 13             if (i == 0) continue;
98 14             for (int j = BLOCK_LOW(block_j); j <= BLOCK_HIGH(block_j); j++){
99 15                 {
100 16
101 17
102 18
103 19
104 20

```

```

10514         for (int k = BLOCK_LOW(block_k); k <= BLOCK_HIGH(block_k); k++)
10615         {
10716             A[k][j][i] = A[k][j][i] * 0.4 - A[k][j][i-1] * 0.6;
10817         }
10918     }
11019 }
11120 if (block_i != DIM_SIZE - 1) {
11221     MPI_Barrier(MPI_COMM_WORLD);
11322     pack_array(send_buffer, BLOCK_HIGH(block_i), BLOCK_LOW(block_j), BLOCK_LOW(block_k),
11423               BLOCK_HIGH(block_i), BLOCK_HIGH(block_j), BLOCK_HIGH(block_k));
11524     int dest = block_owner(block_i + 1, block_j, block_k);
11625     MPI_Isend(send_buffer, BLOCK_SIZE * BLOCK_SIZE, MPI_DOUBLE, dest, 0, MPI_COMM_WORLD, &request);
11726 }
11827 MPI_Barrier(MPI_COMM_WORLD);
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136

```

At the beginning of the loop except for the first loop, the processor will receive the data of the boundary of the subblock $(i-1, j, k)$, and then compute the data of the subblock (i, j, k) . After the computation, the processor will pack the data of the boundary of the subblock (i, j, k) and send it to the processor $p_{(i+1,j,k)}$.

To avoid the deadlock, we use `MPI_Isend` instead of `MPI_Send` to send the data. The `MPI_Isend` will not block the processor, so the processor can continue to the next loop without waiting for the data to be received by the destination processor. The `MPI_Barrier` is used to ensure that all processors have finished the computation of the current subblock.

The code of computing among the j and k dimension is similar to the code above, so we will not show it here.

Finally, all data will be gathered to the processor 0, and the result will be printed out.

```

1371 if (id == 0) {
1382     double *buffer = (double *) malloc(sizeof(double) * BLOCK_SIZE * BLOCK_SIZE * BLOCK_SIZE);
1393     for (int source = 1; source < p; source++)
1404     {
1415         for (block_i = 0; block_i < DIM_SIZE; block_i++) {
1426             block_j = block_k = -1;
1437             block_id(source, &block_i, &block_j, &block_k);
1448             MPI_Recv(buffer, BLOCK_SIZE * BLOCK_SIZE * BLOCK_SIZE, MPI_DOUBLE, source, 0,
1459                    MPI_COMM_WORLD, MPI_STATUS_IGNORE);
14610             unpack_array(buffer, BLOCK_LOW(block_i), BLOCK_LOW(block_j), BLOCK_LOW(block_k),
14711                        BLOCK_HIGH(block_i), BLOCK_HIGH(block_j), BLOCK_HIGH(block_k));
14812         }
14913     }
15014     free(buffer);
15115 } else {
15216 }
153
154
155
156

```

```

15714 double *buffer = (double *)malloc( sizeof(double) * BLOCK_SIZE * BLOCK_SIZE * BLOCK_SIZE);
15815 for (block_i = 0; block_i < DIM_SIZE; block_i++)
15916 {
16017     block_j = block_k = -1;
16117     block_id(id, &block_i, &block_j, &block_k);
16218     pack_array(buffer, BLOCK_LOW(block_i), BLOCK_LOW(block_j), BLOCK_LOW(block_k), BLOCK_HIGH(
16319         block_i), BLOCK_HIGH(block_j), BLOCK_HIGH(block_k));
164
165     MPI_Send(buffer, BLOCK_SIZE * BLOCK_SIZE * BLOCK_SIZE, MPI_DOUBLE, 0, 0, MPI_COMM_WORLD);
16620
16721 }
16822 free( buffer );
169
17023 }

```

If we ignore the final gathering, there are $3\sqrt{p} - 3$ transmits in total, and each transmit will send a boundary of size N^2 , so the total communication cost is $O(N^2\sqrt{p})$ and $O(\frac{N^2}{\sqrt{p}})$ for each processor.

The snapshot of the result is shown below. Since the result is too large, we use sha256 to verify the correctness of the result. It can be figured out that our parallel implementation gives the same result as the original serial implementation.

```

• $gcc 1_adi.c && ./a.out 2>&1 | sha256sum
43960c3e10433b13c1a9459a74ba4c1c7f9eed9d2e10c96276cc87dbbfced82a -
• $mpicc 1_adi_mpi.c && mpirun -np 4 ./a.out 2>&1 | sha256sum
43960c3e10433b13c1a9459a74ba4c1c7f9eed9d2e10c96276cc87dbbfced82a -

```

Fig. 2. Snapshot of the result ($N = 64$)

```

• $gcc 1_adi.c && ./a.out 2>&1 | sha256sum
0bdbb8bf5377e2f276320214f1ce1c774b36ff71ee786a02dfd0a0ca10967db2 -
• $mpicc 1_adi_mpi.c && mpirun -np 4 ./a.out 2>&1 | sha256sum
0bdbb8bf5377e2f276320214f1ce1c774b36ff71ee786a02dfd0a0ca10967db2 -

```

Fig. 3. Snapshot of the result ($N = 1024$)

FLOYD

The Floyd algorithm is a classic algorithm to solve the all-pairs shortest path problem. For each loop k , the algorithm will update the shortest path between i and j by checking whether the path from i to j through k is shorter than the current shortest path:

$$a_{ij} = \min(a_{ij}, a_{ik} + a_{kj}).$$

We will use rowwise block striped partitioning to partition the matrix, as shown in Fig.4, so each processor will be assigned with $\frac{N}{p}$ rows. Each processor will need to access the whole k -th row and part of the k -th column of the matrix. Since the part of the k -th column is already stored in each processor's local memory, we only need to broadcast the k -th row to other processors.

The code of the Floyd algorithm is shown below.

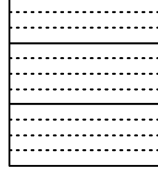


Fig. 4. Partition of the matrix

```

1  for (int k = 0; k < N; k++) {
2      MPI_Bcast(A[k], N, MPI_DOUBLE, BLOCK_OWNER(k, p, N), MPI_COMM_WORLD);
3      for (int i = BLOCK_LOW(id, p, N); i <= BLOCK_HIGH(id, p, N); i++) {
4          for (int j = 0; j < N; j++) {
5              A[i][j] = min(A[i][j], A[i][k] + A[k][j]);
6          }
7      }
8  }

```

At the beginning of each loop, the processor responsible for the k -th row will broadcast the k -th row to other processors. Then each processor will update the shortest path between i and j by checking whether the path from i to j through k is shorter than the current shortest path.

Finally, all data will be gathered to the processor 0, and the result will be printed out.

The snapshot of the result is shown below. Since the result is too large, we use sha256 to verify the correctness of the result. It can be figured out that our parallel implementation gives the same result as the original serial implementation.

```

• $gcc 3_floyd.c && ./a.out 2>&1 | sha256sum
19998e389180c17c12155149d7f8fe73317bcf6e8173ceb8768f0ae81cb823ee -
• $mpicc 3_floyd_mpi.c && mpirun -np 4 ./a.out 2>&1 | sha256sum
19998e389180c17c12155149d7f8fe73317bcf6e8173ceb8768f0ae81cb823ee -
• $mpicc 3_floyd_mpi.c && mpirun -np 16 ./a.out 2>&1 | sha256sum
19998e389180c17c12155149d7f8fe73317bcf6e8173ceb8768f0ae81cb823ee -

```

Fig. 5. Snapshot of the result ($N = 1024$)