

Bitonic Sort Report

111062117, Hsiang-Sheng Huang

May 15, 2025

Implementation Overview

Local Bitonic Sort

```
1 void compSwap(float &a, float &b, bool dir) {
2     if ((a > b) == dir) swap(a, b);
3 }
4
5 void bitonic_sort_loop(vector<float>& arr) {
6     const int n = (int)arr.size(); // n must be 2^k
7     for (int k = 2; k <= n; k <= 1)
8         for (int j = k >> 1; j > 0; j >>= 1)
9             for (int i = 0; i < n; ++i) {
10                 int ixj = i ^ j;
11                 if (ixj > i) {
12                     bool dir = ((i & k) == 0); // true=ASC, false=DESC
13                     compSwap(arr[i], arr[ixj], dir);
14                 }
15             }
16 }
```

Each MPI *rank* first performs a **local** bitonic sort on its private segment. I chose the iterative variant to avoid extra call-stack overhead incurred by the recursive implementation.

Time Complexity The three nested loops yield $k = \log_2 n$, $j = \log_2 n$, and $i = n$ iterations, giving $\mathcal{O}(n \log^2 n)$.

Parallel Bitonic Merge

After every rank holds a locally sorted block, we repeatedly compare-exchange with a partner rank so that, after $\log_2 p$ phases, the global order is established.

Merging the High / Low Halves

```
1 void merge_high(const vector<float>& a, const vector<float>& b, vector<float>& res,
2                 int n) {
3     int j = n - 1, k = n - 1;
4     for (int i = n - 1; i >= 0; --i) {
5         if (j < 0) res[i] = b[k--];
6         else if (k < 0) res[i] = a[j--];
7         else if (a[j] > b[k]) res[i] = a[j--];
8         else res[i] = b[k--];
9     }
10 }
11
12 void merge_low(const vector<float>& a, const vector<float>& b, vector<float>& res,
13                int n) {
14     int j = 0, k = 0;
15     for (int i = 0; i < n; ++i) {
```

```

14     if (j >= n) res[i] = b[k++];
15     else if (k >= n) res[i] = a[j++];
16     else if (a[j] < b[k]) res[i] = a[j++];
17     else res[i] = b[k++];
18 }
19 }

```

The `merge_high` routine keeps the larger half (descending), while `merge_low` keeps the smaller half (ascending).

Stable Tag Generation

We use the Cantor pairing function to generate unique MPI tags:

$$\text{cantor}(a, b) = \frac{(a + b)(a + b + 1)}{2} + b.$$

Compare-Exchange Kernel

```

1 void cmp(vector<float>& local, int localN, int partner, bool dir, MPI_Comm comm) {
2     int rank;
3     MPI_Comm_rank(comm, &rank);
4
5     MPI_Request req[2];
6     MPI_Status stat[2];
7
8     int cantor_arg1 = (dir == true) ? rank : partner;
9     int cantor_arg2 = (dir == true) ? partner : rank;
10    int send_tag = (1 + (dir == true)) * cantor(cantor_arg1, cantor_arg2);
11    int recv_tag = (1 + (dir == false)) * cantor(cantor_arg1, cantor_arg2);
12
13    vector<float> recv(localN);
14    if (dir == true) {
15        MPI_Isend(local.data(), localN, MPI_FLOAT, partner, send_tag, comm, &req[0]);
16        MPI_Irecv(recv.data(), localN, MPI_FLOAT, partner, recv_tag, comm, &req[1]);
17    } else {
18        MPI_Irecv(recv.data(), localN, MPI_FLOAT, partner, recv_tag, comm, &req[0]);
19        MPI_Isend(local.data(), localN, MPI_FLOAT, partner, send_tag, comm, &req[1]);
20    }
21
22    MPI_Waitall(2, req, stat);
23
24    vector<float> res(localN);
25    if (dir == true) merge_low(local, recv, res, localN);
26    else merge_high(local, recv, res, localN);
27    local.swap(res);
28 }

```

Here, `dir == true` means ascending order for the current exchange.

Main

```

1 int p = __lg(usedP);
2 for (int i = 0; i < p; ++i)
3     for (int j = i; j >= 0; --j) {
4         int partner = rank ^ (1 << j);
5         bool dir = ((rank >> (i + 1)) & 1) == ((rank >> j) & 1);
6         cmp(local, localN, partner, dir, activeComm);
7     }

```

The outer loop iterates over stages i , and the inner loop over levels j . Each rank communicates with exactly one partner per (i, j) .

Handling Non-Power-of-Two Cases

Input Size $N \neq 2^k$

Pad with $+\infty$ so extra elements bubble to the end:

```
1 int paddedN = 1;
2 while (paddedN < N) paddedN <= 1;
3 // ...
4 vector<float> data;
5 data.resize(paddedN, numeric_limits<float>::infinity());
```

Processor Count $p \neq 2^k$

Dismiss ranks beyond the largest power-of-two:

```
1 int usedP = 1 << __lg(world_size);
2 int color = (world_rank < usedP) ? 0 : MPI_UNDEFINED;
3
4 MPI_Comm activeComm;
5 MPI_Comm_split(MPI_COMM_WORLD, color, world_rank, &activeComm);
6 if (color == MPI_UNDEFINED) {
7     MPI_Finalize();
8     return 0;
9 }
```

I/O Strategy

Centralised I/O on rank 0:

```
1 if (rank == 0) {
2     data.resize(paddedN, numeric_limits<float>::infinity());
3     MPI_File inFH;
4     MPI_File_open(MPI_COMM_SELF, inFile, MPI_MODE_RDONLY, MPI_INFO_NULL, &inFH);
5     MPI_File_read_at(inFH, 0, data.data(), N, MPI_FLOAT, MPI_STATUS_IGNORE);
6     MPI_File_close(&inFH);
7 }
8
9
10 MPI_Scatter(data.data(), localN, MPI_FLOAT, local.data(), localN, MPI_FLOAT, 0,
11             activeComm);
12 // ...
13 if (rank == 0) {
14     ans.resize(N);
15     MPI_File outFH;
16     MPI_File_open(MPI_COMM_SELF, outFile, MPI_MODE_CREATE | MPI_MODE_WRONLY,
17                   MPI_INFO_NULL, &outFH);
18     MPI_File_write_at(outFH, 0, ans.data(), N, MPI_FLOAT, MPI_STATUS_IGNORE);
19     MPI_File_close(&outFH);
20 }
```

Experimental Results (IPM)

Future Optimisations

1. **Fully parallel I/O:** Replace rank-0 funnel with `MPI_File_read_all` / `MPI_File_write_all`.
2. **Hybrid parallelism:** Overlap computation and communication by double-buffering `cmp` and using OpenMP within each rank.

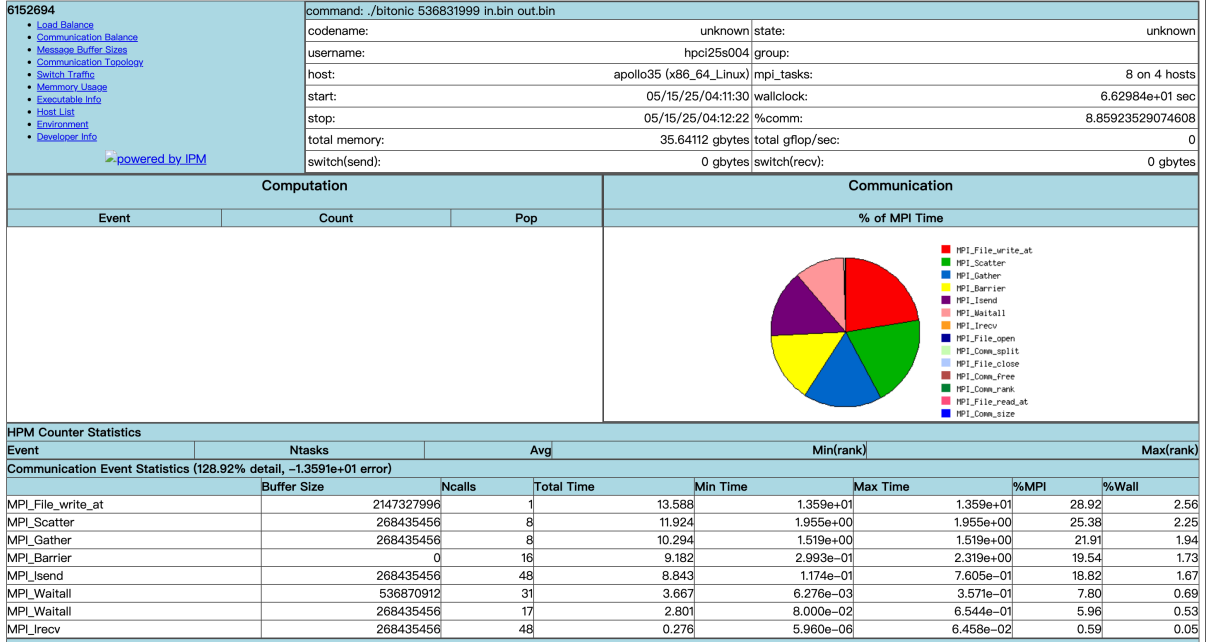


Figure 1: IPM Profile of MPI routines

MPI Routine	% of MPI	% of Wall
MPI_File_write_at	28.92%	2.56%
MPI_Scatter	25.38%	2.25%
MPI_Gather	21.91%	1.95%

Table 1: Top MPI routines by time share

3. **Custom datatypes:** Use MPI derived datatypes to pack blocks and reduce scatter/gather overhead.

Conclusion

The implementation achieves correct global sorting with per-rank complexity $\mathcal{O}(\frac{N}{p} \log^2 \frac{N}{p})$ and scales up to 32 processes. Serial I/O remains the performance bottleneck to address next.