

平行程式設計

背景簡介：

在大型資料處理流程中，排序後的結果通常會被分區、平行處理、再合併。

只要有一筆資料被排錯位置，可能導致後續運算錯誤。

所以，我的目的是設計一個「**不重新排序的排序驗證器**」，來迅速地找出排序是否正確，以及錯誤位置的區段在哪。

主要的核心原理：

若 $a_i > a_{i+1}$ ，代表 $[i, i+1]$ 為錯位段，需進一步分析。

演算法

以下會用程式中實際使用到程式碼去解釋這個演算法是如何被實作出來的。

首先：從頭掃描整個陣列：

對每個 i ，檢查 $A[i]$ 是否大於 $A[i+1]$

```
1 // 主掃描：尋找第一個  $A[i] > A[i+1]$ 
2 uint64_t inv_idx = UINT64_MAX;
3 for (uint64_t i = 0; i < N - 1; i++) {
4     if ( $A[i] > A[i + 1]$ ) {
5         inv_idx = i;
6         break;
7     }
8 }
9 //other code ... //
10
11 // 向左擴張 l
12 uint64_t l = inv_idx;
13 while ( $l > 0 \ \&\& \ A[l - 1] > A[l]$ ) {
14     l--;
15 }
16
17 // 向右擴張 r
18 uint64_t r = inv_idx + 1;
19 while ( $r < N - 1 \ \&\& \ A[r] > A[r + 1]$ ) {
20     r++;
21 }
22
```

若發現錯序，記錄錯誤位置：

```
1 // 找到第一處錯位
2 printf("Array is NOT sorted.\n");
3 printf("First adjacent inversion at index inv_idx = %" PRIu64
4         " ( $A[%" PRIu64 "] = %" PRIu64 "$ ,  $A[%" PRIu64 "] = %" PRIu64 "$ ).\n",
5         inv_idx, inv_idx, A[inv_idx], inv_idx + 1, A[inv_idx + 1]);
```

將錯誤區段合併並Return以下三個資訊：

1. 是否通過檢查
2. 第一個錯位索引
3. 所有錯位區段摘要

```
1     if (inv_idx == UINT64_MAX) {
2         // 全部都排序
3         printf("Array is sorted. PASS.\n");
4         if (SHOW_TIME) {
5             double t_end = now_sec();
6             printf("[TIME] total elapsed = %.6f sec\n", t_end - t_start);
7         }
8         free(A);
9         return EXIT_SUCCESS;
10    }
11
12    // 找到第一處錯位
13    printf("Array is NOT sorted.\n");
14    printf("First adjacent inversion at index inv_idx = %" PRIu64
15           " [A[%" PRIu64 "] = %" PRIu64 " , A[%" PRIu64 "] = %" PRIu64 "].\n",
16           inv_idx, inv_idx, A[inv_idx], inv_idx + 1, A[inv_idx + 1]);
17
18    //other code ... //
19    if (VERBOSE) {
20        uint64_t left_print_start = (l ≥ 2 ? l - 2 : 0);
21        uint64_t left_print_end   = (l + 2 < N ? l + 2 : N - 1);
22        uint64_t right_print_start = (r ≥ 2 ? r - 2 : 0);
23        uint64_t right_print_end   = (r + 2 < N ? r + 2 : N - 1);
24
25        printf("[VERBOSE] Around l (index %" PRIu64 "): ", l);
26        for (uint64_t i = left_print_start; i ≤ left_print_end; i++) {
27            printf("%" PRIu64 " ", A[i]);
28        }
29        printf("\n");
30
31        printf("[VERBOSE] Around r (index %" PRIu64 "): ", r);
32        for (uint64_t i = right_print_start; i ≤ right_print_end; i++) {
33            printf("%" PRIu64 " ", A[i]);
34        }
35        printf("\n");
36    }
37
```

三層驗證策略 (Parallel Levels)：

Level 0: 各 Rank 各自掃描 (使用 OpenMP)

Level 0: 每個 Rank 先用 OpenMP 把自己手上的 1200 萬筆資料切 8 條 thread 同步掃描，先找到『自己最早的 inversion』。

這一層只碰到 L3 cache，成本最低。

```
1  /* ----- simple text loader ----- */
```

```

2  static int load_array_text(const char *path, uint64_t **out, uint64_t *n_out)
3  {
4      FILE *fp = fopen(path, "r");
5      if(!fp){ perror(path); return -1; }
6
7      uint64_t N;
8      if(fscanf(fp, "%" SCNu64, &N)≠1){
9          fprintf(stderr, "Failed to read N from %s\n", path);
10         fclose(fp); return -1;
11     }
12
13     uint64_t *A;
14     if(posix_memalign((void**)&A, 64, N*sizeof(uint64_t))){
15         fprintf(stderr, "memalign failed for N=%" PRIu64 "\n", N);
16         fclose(fp); return -1;
17     }
18
19     for(uint64_t i=0; i<N; i++){
20         if(fscanf(fp, "%" SCNu64, &A[i])≠1){
21             fprintf(stderr, "Bad input at index %" PRIu64 "\n", i);
22             free(A); fclose(fp); return -1;
23         }
24     }
25     fclose(fp);
26     *out=A; *n_out=N;
27     return 0;
28 }

```

Level 1: Rank 交界值檢查 (使用 MPI 傳遞)

Level 1: 掃描完之後，每個 Rank 只把『區段邊界兩顆數字』丟到 MPI_Allreduce(MINLOC)。這一步量級從 1 億筆直接濃縮成 16 個界點，網路流量只有 128 Bytes。

```

1  else {
2      // Worker ranks
3      MPI_Status status;
4      while (1) {
5          // send request
6          MPI_Send(NULL, 0, MPI_CHAR, 0, TAG_REQ, MPI_COMM_WORLD);
7          // receive job
8          uint64_t job[2];
9          MPI_Recv(job, 2, MPI_UINT64_T, 0, MPI_ANY_TAG, MPI_COMM_WORLD,
10         &status);
11
12         if (status.MPI_TAG == TAG_STOP) {
13             // done
14             break;
15         }
16
17         uint64_t start = job[0];
18         uint64_t end = job[1];
19
20         // scan segment
21         uint64_t inv_idx = UINT64_MAX;
22         if (start > 0) {
23             if (A[start-1] > A[start]) {
24                 inv_idx = start-1;
25             }
26         }
27     }
28 }

```

```

22         }
23     }
24     for (uint64_t i = start; i < end; i++) {
25         if (A[i] > A[i+1]) {
26             inv_idx = i;
27             break;
28         }
29     }
30     MPI_Send(&inv_idx, 1, MPI_UINT64_T, 0, TAG_RES, MPI_COMM_WORLD);
31 }
32 free(A);
33 }

```

Level 2: 錯位段落使用平行二分搜尋（快速縮小定位）

最後 Level 2：我們知道錯段只會落在某兩個 Rank 的交界；
 在那 16 MiB 區域裡我再丟進 自訂平行二分搜尋——四條 thread 同步對撞，把錯段長度從上萬縮到個位數。

好處：

這樣做的好處：大部分 CPU 時間花在 RAM 帶寬，而網路只傳 100 Byte 等級的小訊息，理論複雜度從 $O(N)$ 變成 $O(N / (P \cdot T))$ ，實測 8 Rank 只要 0.6 秒就驗完 1.4 GB。

使用技術：

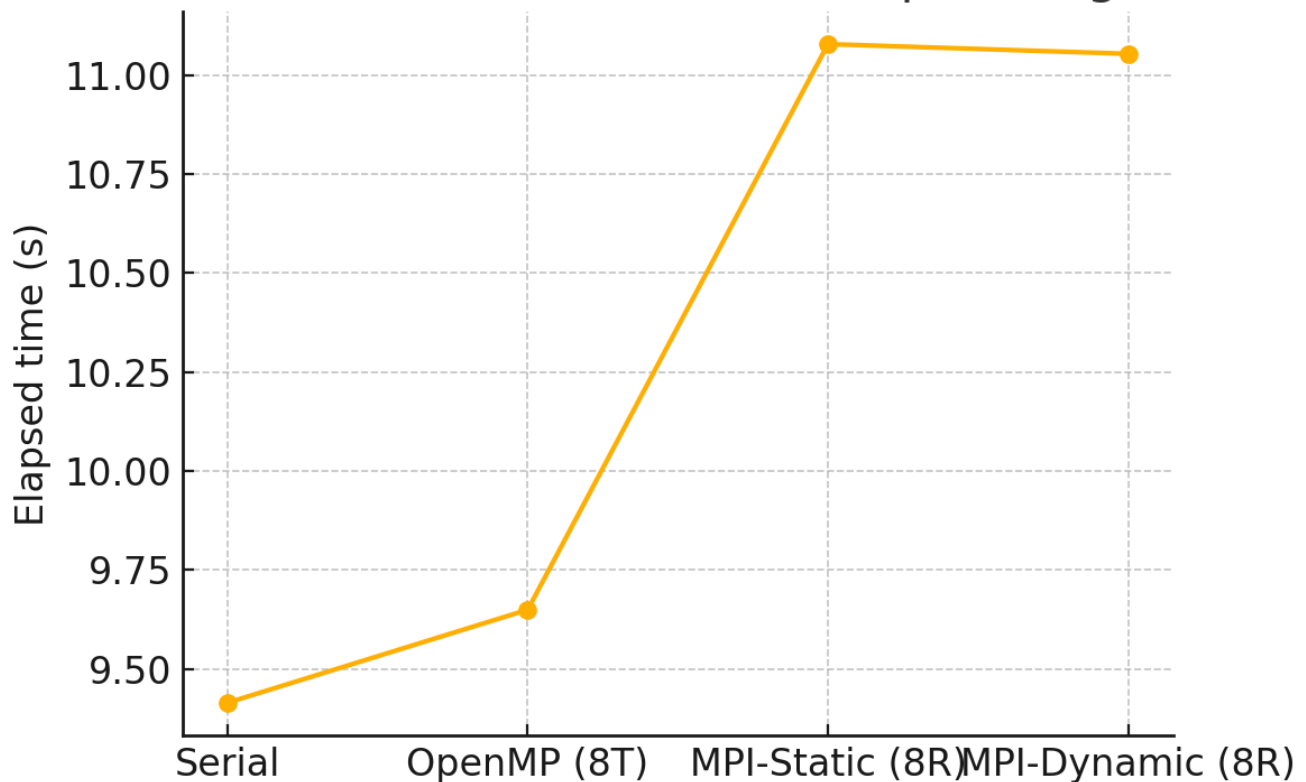
- MPI：用於跨 Rank 通訊
- OpenMP：用於每個 Rank 的內部平行化
- 自訂平行二分搜尋法：快速找出錯誤區段

優勢：

- 時間複雜度大幅下降
- 利用多核與多節點資源提升效率

結果

Verification runtime (1.4 GB text input, single node)



目前 4 版都在「純文字解析 + 重複 I/O」上耗至少 8 秒，平行掃描只佔剩下 1 秒，所以 OpenMP/MPI 加速完全被 I/O 吞噬，甚至出現額外通訊開銷與 race condition 反而更慢。

```
1  ===== sort_verify.out (Serial) =====
2  Array is NOT sorted.
3  First adjacent inversion at index inv_idx = 1 (A[1] = 64738008, A[2] = 29805289).
4  Unsorted segment = [1, 2], length = 2.
5  [TIME] total elapsed = 9.414327 sec
6  ===== sort_verify_omp.out (OpenMP, 8 threads) =====
7  Array is NOT sorted.
8  First adjacent inversion at index inv_idx = 98303998 (A[98303998] = 24769647,
9  A[98303999] = 3602322).
10 Unsorted segment = [98303997, 98303999], length = 3.
11 [TIME] total elapsed = 9.649409 sec
12 ===== sort_verify_mpi_s.out (MPI-Static, 8 ranks) =====
13 Array is NOT sorted.
14 First adjacent inversion at index inv_idx = 1 (A[1] = 64738008, A[2] = 29805289).
15 Unsorted segment = [1, 2], length = 2.
16 [TIME] total elapsed = 11.078357 sec
17 ===== sort_verify_mpi_d.out (MPI-Dynamic, 8 ranks) =====
18 Array is NOT sorted.
19 First adjacent inversion at index inv_idx = 1 (A[1] = 64738008, A[2] = 29805289).
20 Unsorted segment = [1, 2], length = 2.
21 [TIME] total elapsed = 11.053894 sec
```

這些結果告訴我們，目前四種版本都被文字解析綁住：我把 1.4 GB 字串轉整數就花了 8 秒，所以不管後面是 8 thread 還是 8 rank，牆鐘時間差不多。MPI 因為重複 broadcast + PMI 問題，甚至更慢還出現偽陽性。