

## Lab 5: fft 和 heated plate 的并行化

学号：19335109	课程：高性能计算
姓名：李雪堃	学期：Fall 2021
专业：计算机科学与技术（超算）	教师：黄聃、卢宇彤
邮箱：i@xkun.me	TAs：江嘉治、刘亚辉

### Table of Contents

#### Lab 5: fft 和 heated plate 的并行化

- (一) 实验任务
- (二) 实验环境
- (三) 实验过程和核心代码
  - (1) parallel\_for 并行化 fft
  - (2) heated\_plate\_openmp 改造为基于 MPI 的进程并行应用
  - (3) fft\_pf 和 fft\_openmp 的性能对比
- (四) 实验结果
- (五) 实验感想
- 附录：参考资料

### (一) 实验任务

- 任务 1：
  - 通过实验 4 构造的基于 Pthreads 的 `parallel_for` 函数替换 `fft_serial` 应用中的某些计算量较大的 `for` 循环，实现 `for` 循环分解、分配和线程并行执行。
- 任务 2：
  - 将 `fft_serial` 应用改造成基于 MPI 的进程并行应用（为了适合 MPI 的消息机制，可能需要对 `fft_serial` 的代码实现做一定调整）。
  - 将 `heated_plate_openmp` 应用改造成基于 MPI 的进程并行应用。
  - Bonus：使用 `MPI_Pack` / `MPI_Unpack`，或 `MPI_Type_create_struct` 实现数据重组后的消息传递。
- 任务 3：
  - 性能分析任务 1 和并行化 `fft` 应用，包括：
    - 不同问题规模的并行化 `fft` 应用并行执行时间对比，其中问题规模定义为 `N` 变化范围 2, 4, 6, 8, 16, 32, 64, 128, ....., 2097152；并行规模为 1, 2, 4, 8 进程/线程。

- 内存消耗对比，内存消耗采用 `valgrind massif` 工具采集，注意命令 `valgrind` 命令中增加 `--stacks=yes` 参数采集程序运行栈内内存消耗。

## (二) 实验环境

- Ubuntu 20.04.3 LTS x86\_64
- gcc 9.3.0
- OpenMPI 4.0.3
- GNU Make 4.2.1
- valgrind 3.15.0

## (三) 实验过程 and 核心代码

### (1) `parallel_for` 并行化 fft

代码在 `fft/fft-pf` 目录下。其中 `parallel_for.h` 和 `parallel_for.c` 是上次实验基于 `pthread` 的 `parallel_for` 函数。上次写的有些 bug，这次进行了修复和改进。

```
1 typedef struct
2 {
3     int my_rank; // these four arguments are specified by parallel_for
4     int my_start;
5     int my_end;
6     int my_increment;
7
8     void * func_arg; // thread function arguments, specified by user
9 } pf_arg_t;
```

```
1 void parallel_for(int start, int end, int increment, \
2                   void *(*functor)(void *), void *arg, int thread_count);
```

首先是 `parallel_for` 的内部参数 `pf_arg_t`，这个参数不需要用户指定，是在 `parallel_for` 函数内部传递给各个线程的参数。`pf_arg_t` 中的 `fun_arg` 才是用户自定义的结构体参数，`parallel_for` 中的 `void *arg` 参数传递的就是用户自定义的结构体。

```

1 void parallel_for(int start, int end, int increment, \
2                 void *(*functor)(void *), void *func_arg, int thread_count)
3 // functor is a function pointer, defined by user, executed by thread
4 // func_arg is the arguments of functor, defined and allocated by user
5 {
6     pthread_t *thread_handles = (pthread_t *) malloc(sizeof(pthread_t) * thread_count);
7     pf_arg_t *pf_args = (pf_arg_t *) malloc(sizeof(pf_arg_t) * thread_count);
8
9     int loop_count; // total loop times
10    if ((end - start) / increment == 0)
11        loop_count = (end - start) / increment;
12    else
13        loop_count = (end - start) / increment + 1;
14
15    if (loop_count <= thread_count) // if loop_count smaller than thread_count
16    {
17        thread_count = 1;
18    }
19
20    int my_count = loop_count / thread_count; // expected average thread loop times
21
22    for (long thread = 0; thread < thread_count; thread++)
23    {
24        if (thread == thread_count - 1) // I'm the last thread
25        {
26            pf_args[thread].my_rank = thread;
27            pf_args[thread].my_start = start + increment * my_count * thread;
28            pf_args[thread].my_end = end;
29            pf_args[thread].my_increment = increment;
30            pf_args[thread].func_arg = func_arg;
31        }
32        else // I'm not the last thread
33        {
34            pf_args[thread].my_rank = thread;
35            pf_args[thread].my_start = start + increment * my_count * thread;
36            pf_args[thread].my_end = pf_args[thread].my_start + increment * my_count;
37            pf_args[thread].my_increment = increment;
38            pf_args[thread].func_arg = func_arg;
39        }
40    }
41
42    for (long thread = 0; thread < thread_count; thread++)
43        pthread_create(&thread_handles[thread], NULL, functor, (void *)&pf_args[thread]);
44
45    for (long thread = 0; thread < thread_count; thread++)
46        pthread_join(thread_handles[thread], NULL);
47
48    free(thread_handles);
49    free(pf_args);
50 }

```

在 `parallel_for` 函数中，首先计算总的循环次数 `loop_count`，然后判断 `loop_count` 是否小于等于线程数 `thread_count`，如果是则只开启一个线程，相当于串行计算，这样做是为了不损失性能。

接着是初始化线程函数的参数，将 `pf_arg_t` 中的 `func_arg` 指针指向用户传递来的参数结构体指针。最后是线程创建和销毁。

在 `fft_serial` 中，经过观察分析，发现 `cfft_i` 和 `step` 函数可以并行，并且测试后的确带来了性能的提升；`ccopy` 也可以并行，但并行后性能反而下降。

首先是 `cfft_i`，它的作用是准备 FFT 计算时需要的 `sin` 和 `cos` 表。首先创建参数结构体 `cfft_i_arg_t`，然后编写 `cfft_i_pf` 函数，从 `pf_arg_t` 中获得 `my_start`、`my_end` 和 `my_increment`，以及用户传递的参数结构体。然后获取这些参数，用一个循环进行计算即可。

```

1  typedef struct
2  {
3      double aw;
4      double *w;
5  } cffti_arg_t;

```

```

1  void *cffti_pf(void *arg)
2  {
3      pf_arg_t * my_arg = (pf_arg_t *)arg;
4      int my_start = my_arg->my_start;
5      int my_end = my_arg->my_end;
6      int my_increment = my_arg->my_increment;
7
8      cffti_arg_t * cffti_arg = (cffti_arg_t *)my_arg->func_arg;
9      double aw = cffti_arg->aw;
10     double *w = cffti_arg->w;
11     double x;
12
13     // #ifdef DEBUG
14     //     int my_rank = my_arg->my_rank;
15     //     printf("I'm thread %d, my_start = %d, my_end = %d\n", my_rank, my_start, my_end);
16     // #endif
17
18     for (int i = my_start; i < my_end; i += my_increment)
19     {
20         x = aw * ((double)i);
21         w[i * 2 + 0] = cos(x);
22         w[i * 2 + 1] = sin(x);
23     }
24
25     return NULL;
26 }

```

在 `cffti` 函数中，将原有的循环替换为 `parallel_for` 即可。

```

1  void cffti(int n, double w[])
2  /*
3      Purpose:
4          CFFTI sets up sine and cosine tables needed for the FFT calculation.
5      Parameters:
6          Input, int N, the size of the array to be transformed.
7          Output, double W[N], a table of sines and cosines.
8  */
9  {
10     double aw;
11     int n2;
12     const double pi = 3.141592653589793;
13
14     n2 = n / 2;
15     aw = 2.0 * pi / ((double)n); // aw = 2pi / n
16
17     cffti_arg_t cffti_arg;
18     cffti_arg.aw = aw;
19     cffti_arg.w = w;
20
21     parallel_for(0, n2, 1, cffti_pf, (void *)&cffti_arg, thread_count);
22 }

```

然后是 `step` 函数的并行，该函数作用是执行 FFT 中的一次迭代。同样地，我们需要创建 `step_arg_t` 参数结构体。然后编写 `step_pf` 函数，解析参数，执行循环。

```

1  typedef struct
2  {
3      int mj;
4      int mj2;
5      double *a;
6      double *b;
7      double *c;
8      double *d;
9      double *w;
10     double sgn;
11 } step_arg_t;

```

```

1  void *step_pf(void *arg)
2  {
3      pf_arg_t *my_arg = (pf_arg_t *)arg;
4      int my_start = my_arg->my_start;
5      int my_end = my_arg->my_end;
6      int my_increment = my_arg->my_increment;
7
8      step_arg_t *step_arg = (step_arg_t *)my_arg->func_arg;
9      int mj = step_arg->mj;
10     int mj2 = step_arg->mj2;
11     double *a = step_arg->a;
12     double *b = step_arg->b;
13     double *c = step_arg->c;
14     double *d = step_arg->d;
15     double *w = step_arg->w;
16     int sgn = step_arg->sgn;
17
18     double ambr;
19     double ambu;
20     int j;
21     int ja;
22     int jb;
23     int jc;
24     int jd;
25     int jw;
26     int k;
27     double wjw[2];
28
29     for (j = my_start; j < my_end; j += my_increment)
30     {
31         jw = j * mj;
32         ja = jw;
33         jb = ja;
34         jc = j * mj2;
35         jd = jc;
36
37         wjw[0] = w[jw * 2 + 0];
38         wjw[1] = w[jw * 2 + 1];
39
40         if (sgn < 0.0)
41         {
42             wjw[1] = -wjw[1];
43         }
44
45         for (k = 0; k < mj; k++)
46         {
47             c[(jc + k) * 2 + 0] = a[(ja + k) * 2 + 0] + b[(jb + k) * 2 + 0];
48             c[(jc + k) * 2 + 1] = a[(ja + k) * 2 + 1] + b[(jb + k) * 2 + 1];
49
50             ambr = a[(ja + k) * 2 + 0] - b[(jb + k) * 2 + 0];
51             ambu = a[(ja + k) * 2 + 1] - b[(jb + k) * 2 + 1];
52
53             d[(jd + k) * 2 + 0] = wjw[0] * ambr - wjw[1] * ambu;
54             d[(jd + k) * 2 + 1] = wjw[1] * ambr + wjw[0] * ambu;
55         }
56     }
57     return NULL;
58 }

```

`step` 函数在 `cfft2` 函数中被调用，我们将其中的 `step` 函数全部替换为 `parallel_for` 函数，开启线程执行 `step_pf`。

```

1 void cfft2(int n, double x[], double y[], double w[], double sgn)
2 /*
3  Purpose:
4   CFFT2 performs a complex Fast Fourier Transform.
5  Parameters:
6   Input, int N, the size of the array to be transformed.
7   Input/output, double X[2*N], the data to be transformed.
8   On output, the contents of X have been overwritten by work information.
9   Input, double W[N], a table of sines and cosines.
10  Input, double SGN, is +1 for a "forward" FFT and -1 for a "backward" FFT.
11  Output, double Y[2*N], the forward or backward FFT of X.
12 */
13 {
14     int j;
15     int m;
16     int mj;
17     int mj2;
18     int lj;
19     int tgle;
20     step_arg_t step_arg;
21
22     m = (int)(log((double)n) / log(1.99)); // m = log_2^n
23     mj = 1;
24     mj2 = 2 * mj;
25     lj = n / mj2;
26
27     tgle = 1; // Toggling switch for work array.
28     step_arg.mj = mj;
29     step_arg.mj2 = mj2;
30     step_arg.a = &x[0 * 2 + 0];
31     step_arg.b = &x[(n / 2) * 2 + 0];
32     step_arg.c = &y[0 * 2 + 0];
33     step_arg.d = &y[mj * 2 + 0];
34     step_arg.w = w;
35     step_arg.sgn = sgn;
36     parallel_for(0, lj, 1, step_pf, (void *)&step_arg, thread_count);
37     // step(n, mj, &x[0 * 2 + 0], &x[(n / 2) * 2 + 0], &y[0 * 2 + 0], &y[mj * 2 + 0], w, sgn);
38
39     if (n == 2)
40     {
41         return;
42     }
43
44     for (j = 0; j < m - 2; j++) // m = log_2^n, 这里的 step 要被重复执行很多次, 线程开销过大
45     {
46         mj = mj * 2;
47         mj2 = 2 * mj;
48         lj = n / mj2;
49         if (tgle)
50         {
51             step_arg.mj = mj;
52             step_arg.mj2 = mj2;
53             step_arg.a = &y[0 * 2 + 0];
54             step_arg.b = &y[(n / 2) * 2 + 0];
55             step_arg.c = &x[0 * 2 + 0];
56             step_arg.d = &x[mj * 2 + 0];
57             step_arg.w = w;
58             step_arg.sgn = sgn;
59
60             parallel_for(0, lj, 1, step_pf, (void *)&step_arg, thread_count);
61             // step(n, mj, &y[0 * 2 + 0], &y[(n / 2) * 2 + 0], &x[0 * 2 + 0], &x[mj * 2 + 0], w, sgn);
62             tgle = 0;
63         }
64         else
65         {
66             step_arg.mj = mj;
67             step_arg.mj2 = mj2;
68             step_arg.a = &x[0 * 2 + 0];
69             step_arg.b = &x[(n / 2) * 2 + 0];
70             step_arg.c = &y[0 * 2 + 0];
71             step_arg.d = &y[mj * 2 + 0];
72             step_arg.w = w;
73             step_arg.sgn = sgn;
74
75             parallel_for(0, lj, 1, step_pf, (void *)&step_arg, thread_count);
76             // step(n, mj, &x[0 * 2 + 0], &x[(n / 2) * 2 + 0], &y[0 * 2 + 0], &y[mj * 2 + 0], w, sgn);
77             tgle = 1;
78         }
79     }
80 }

```

```

80
81  if (tgle) // Last pass through data: move Y to X if needed.
82  {
83      ccopy(n, y, x);
84  }
85
86  mj = n / 2;
87  mj2 = mj * 2;
88  lj = n / mj2;
89  step_arg.mj = mj;
90  step_arg.mj2 = mj2;
91  step_arg.a = &x[0 * 2 + 0];
92  step_arg.b = &x[(n / 2) * 2 + 0];
93  step_arg.c = &y[0 * 2 + 0];
94  step_arg.d = &y[mj * 2 + 0];
95  step_arg.w = w;
96  step_arg.sgn = sgn;
97  parallel_for(0, lj, 1, step_pf, (void *)&step_arg, thread_count);
98  // step(n, mj, &x[0 * 2 + 0], &x[(n / 2) * 2 + 0], &y[0 * 2 + 0], &y[mj * 2 + 0], w, sgn);
99
100 return;
101 }

```

下面是运行结果。

```

fft_pf_test.txt M x
Lab > lab5 > fft > asset > fft_pf_test.txt
4  C/PARALLEL-FOR version
5
6  Demonstrate an implementation of the Fast Fourier Transform
7  of a complex data vector.
8
9  Number of threads =          12
10
11 Accuracy check:
12
13 FFT ( FFT ( X(1:N) ) ) == N * X(1:N)
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39 FFT_PARALLEL_FOR:
40 Normal end of execution.
41
42 03 December 2021 09:51:36 PM
43

```

	N	NITS	Error	Time	Time/Call	MFLOPS
	2	10000	7.859082e-17	2.421576e-01	1.210788e-05	0.825908
	4	10000	1.209837e-16	4.935291e-01	2.467646e-05	1.620978
	8	10000	6.820795e-17	7.275663e-01	3.637832e-05	3.298668
	16	10000	1.438671e-16	1.002143e+00	5.010717e-05	6.386311
	32	1000	1.331210e-16	3.922027e-01	1.961013e-04	4.079524
	64	1000	1.776545e-16	6.525185e-01	3.262592e-04	5.884891
	128	1000	1.929043e-16	9.471074e-01	4.735537e-04	9.460384
	256	1000	2.092319e-16	1.191531e+00	5.957656e-04	17.187968
	512	100	1.927488e-16	1.439251e-01	7.196256e-04	32.016650
	1024	100	2.312093e-16	1.904348e-01	9.521740e-04	53.771687
	2048	100	2.445006e-16	2.037586e-01	1.018793e-03	110.562214
	4096	100	2.476589e-16	2.468777e-01	1.234389e-03	199.094509
	8192	10	2.571250e-16	2.745175e-02	1.372587e-03	387.938866
	16384	10	2.736298e-16	3.632549e-02	1.816274e-03	631.446473
	32768	10	2.924127e-16	4.803912e-02	2.401956e-03	1023.166141
	65536	10	2.833553e-16	7.481372e-02	3.740686e-03	1401.582584
	131072	1	3.142312e-16	1.202956e-02	6.014778e-03	1852.290986
	262144	1	3.216005e-16	2.524719e-02	1.262360e-02	1868.957224
	524288	1	3.282664e-16	5.299924e-02	2.649962e-02	1879.549865
	1048576	1	3.284479e-16	1.229341e-01	6.146704e-02	1705.915969
	2097152	1	3.509548e-16	2.546743e-01	1.273371e-01	1729.275146

## (2) heated\_plate\_openmp 改造为基于 MPI 的进程并行应用

heated plate 是个数值计算问题，主要是通过迭代来计算 plate 上温度收敛状态的分布。

一个 plate 用一个  $M \times N$  的矩形区域来表示，程序中设置  $M = N = 500$ ，开始时每个节点都赋予一个初始的温度，程序中设置左边界、右边界和下边界的温度都为 100，上边界的温度为 0，中间位置的温度是边界温度的均值（所有边界节点温度的和除以边界节点数）。这个过程是初始化的过程，不好并行，我因此用 master 进程来初始化，在 master 进程中使用 openmp 并行。



迭代的过程非常简单，对面中间位置的节点（不在边界的节点），每轮更新它的温度为其相邻节点温度的均值（上下左右四个节点），即：

$$w[center] = \frac{w[north] + w[south] + w[east] + w[west]}{4}, \text{ for each iteration}$$

一直进行下去，直到相对误差（上次迭代结果与本次迭代结果相比）小于一个阈值，误差的计算是采取所有节点中误差最大值作为本次迭代的误差。

首先是 master 进程（rank 0）初始化 plate。内部用 openmp 并行加快速度。比较简单，这里不再赘述。

```
1  if (my_rank == 0) // I'm the master
2  {
3      printf("\n");
4      printf("HEATED_PLATE_MPI\n");
5      printf("  C/MPI version\n");
6      printf("  A program to solve for the steady state temperature distribution\n");
7      printf("  over a rectangular plate.\n");
8      printf("\n");
9      printf("  Spatial grid of %d by %d points.\n", M, N);
10     printf("  The iteration will be repeated until the change is <= %e\n", epsilon);
11     printf("  Number of processes = %d\n", comm_sz);
12
13     int i, j;
14     #pragma omp parallel shared(w) private(i, j)
15     {
16         // 初始化 w 边界的值
17         #pragma omp for
18         for (i = 1; i < M - 1; i++)
19         {
20             w[i][0] = 100.0; // left border
21             w[i][N - 1] = 100.0; // right border
22         }
23         #pragma omp for
24         for (j = 0; j < N; j++)
25         {
26             w[M - 1][j] = 100.0; // bottom border
27             w[0][j] = 0.0; // top border
28         }
29
30         // 计算边界 w 值之和，加到 mean 上
31         #pragma omp for reduction(+ : mean)
32         for (i = 1; i < M - 1; i++)
33         {
34             mean += w[i][0] + w[i][N - 1];
35         }
36         #pragma omp for reduction(+ : mean)
37         for (j = 0; j < N; j++)
38         {
39             mean += w[M - 1][j] + w[0][j];
40         }
41     }
42
43     // 算均值，除以边界节点的个数
44     mean = mean / (double)(2 * M + 2 * N - 4);
45     printf("\n");
46     printf("  MEAN = %f\n", mean);
47
48     // 初始化 w 内部为 mean
49     #pragma omp parallel shared(mean, w) private(i, j)
50     {
51         #pragma omp for
52         for (i = 1; i < M - 1; i++)
53             for (j = 1; j < N - 1; j++)
54                 w[i][j] = mean;
55     }
56 }
```

然后，每个进程计算自己需要计算的范围，这里是按行划分的。注意到，第 0 行和第 M-1 行是不需要计算的，因为它们属于边界节点，温度是固定的。

`my_bound` 函数会根据 `start`、`end`、进程数 `comm_sz` 和进程号 `my_rank` 来计算该进程需要计算的范围。可以处理不能被整除的情况。

```
1 /* Each process compute my_first_m and my_last_m */
2 int my_first_m, my_last_m, my_m;
3 my_bound(1, M - 1, comm_sz, my_rank, &my_first_m, &my_last_m, &my_m);
```

```
1 void my_bound(int start, int end, int comm_sz, int my_rank, int *my_first_m, int *my_last_m, int *my_m)
2 {
3     if (my_rank == comm_sz - 1) // if I'm the last process
4     {
5         *my_first_m = start + (end - start) / comm_sz * my_rank;
6         *my_last_m = end; // set my_last_m as M-1
7     }
8     else
9     {
10        *my_first_m = start + (end - start) / comm_sz * my_rank;
11        *my_last_m = *my_first_m + (end - start) / comm_sz;
12    }
13    *my_m = *my_last_m - *my_first_m;
14 }
```

接下来，我们先分析 openmp 版本的 heated plate 迭代部分的代码。

```
1 #pragma omp parallel shared(u, w) private(i, j)
2 {
3     // Save the old solution in U.
4     #pragma omp for
5     for (i = 0; i < M; i++)
6         for (j = 0; j < N; j++)
7             u[i][j] = w[i][j];
8
9     // Determine the new estimate of the solution at the interior points.
10    // The new solution W is the average of north, south, east and west neighbors.
11    #pragma omp for
12    for (i = 1; i < M - 1; i++)
13        for (j = 1; j < N - 1; j++)
14            w[i][j] = (u[i - 1][j] + u[i + 1][j] + u[i][j - 1] + u[i][j + 1]) / 4.0;
15 }
16
17 // C and C++ cannot compute a maximum as a reduction operation.
18
19 // Therefore, we define a private variable MY_DIFF for each thread.
20 // Once they have all computed their values, we use a CRITICAL section
21 // to update DIFF.
22 diff = 0.0;
23 #pragma omp parallel shared(diff, u, w) private(i, j, my_diff)
24 {
25     my_diff = 0.0;
26     #pragma omp for
27     for (i = 1; i < M - 1; i++)
28         for (j = 1; j < N - 1; j++)
29             if (my_diff < fabs(w[i][j] - u[i][j]))
30                 my_diff = fabs(w[i][j] - u[i][j]);
31
32     #pragma omp critical
33     if (diff < my_diff)
34         diff = my_diff;
35 }
```

首先，迭代开始时，用一个  $M \times N$  的二维数组 `u` 暂存 `w`，这一步是非常耗时的，因为每次迭代都要做一次复制，大量的内存存取造成巨大的开销。考虑到要使用 MPI 并行，如果将整个 `w` 广播给所有的进程，每个进程还要将 `w` 复制到 `u`，开销是非常巨大的，所以我们可以针对这一部分进行优化。我的想法是，master 进程保留 `w` 的值，由于是按行分配计算任务的，所以可以将每个进程需要的 `w` 的部分发送给它，而不是广播整个 `w`。

于是，每个进程可以用一个 `w_buf` 来暂存需要的 `w`，`w_buf` 可以定义为 `double w_buf[my_m + 2][N]`，注意行数为 `my_m + 2` 的原因是计算 `my_m` 行还需要边界的上下两行温度。

但是，这么做还不够，由于在每次迭代更新温度时，是按照上次迭代的温度来更新的，所以如果只用 `w_buf` 一个缓冲区进行收发数据的话，还需要一个 `u_buf` 来暂存接收到的 `w_buf`，再根据 `u_buf` 暂存的值将计算结果保存到 `w_buf`，这又需要将 `w_buf` 复制到 `u_buf`，造成巨大开销。

进一步的，每个进程还可以用一个 `my_w` 来保存自己本次迭代计算的结果，于是就可以直接利用 `w_buf` 来计算，将结果保存到 `my_w` 即可，不需要保存旧值的额外操作。

最后，还可以发现，每个进程计算误差 `my_diff` 和 `my_w[i][j]` 可以合并到一起，这样会将两次双重循环的开销减少到一次。

下面是每个进程都会定义的 `my_w` 和 `w_buf`，分别作为发送缓冲区和接收缓冲区。

```
1 // set my_w as send buffer, the computation results will be saved in my_w
2 double my_w[my_m][N];
3 for (int i = 0; i < my_m; i++)
4 {
5     my_w[i][0] = 100;
6     my_w[i][N-1] = 100;
7 }
8
9 // set w_buf as recv buffer
10 double w_buf[my_m + 2][N];
```

然后进入迭代循环，每次计算开始前，master 进程先把各个进程需要的 `w` 发送给它们，slave 进程用 `w_buf` 接收。当然，`master` 也要计算，将第一个部分 copy 到自己的 `w_buf` 中。

```
1 // Do not send the whole w
2 if (my_rank == 0)
3 {
4     for (int slave = 1; slave < comm_sz; slave++)
5     {
6         int slave_first_m, slave_last_m, slave_m;
7         my_bound(1, M - 1, comm_sz, slave, &slave_first_m, &slave_last_m, &slave_m);
8         MPI_Send(&w[slave_first_m - 1][0], (slave_m + 2) * N, MPI_DOUBLE, slave, 0, MPI_COMM_WORLD);
9     }
10    memcpy(&w_buf[0][0], &w[my_first_m - 1][0], sizeof(double) * (my_m + 2) * N);
11 }
12 else
13 {
14     MPI_Recv(&w_buf[0][0], (my_m + 2) * N, MPI_DOUBLE, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
15 }
```

然后是计算部分，每个进程将计算的结果保存到 `my_w` 中，同时，某个位置计算完成后，还会更新误差 `my_diff`。

```

1 // Compute my_w for each process and my_diff at the same time
2 // Note that my_w[i][j] <=> w_buf[i+1][j]
3 double my_diff = 0.0;
4 for (int i = 0; i < my_m; i++)
5 {
6     for (int j = 1; j < N - 1; j++)
7     {
8         my_w[i][j] = (w_buf[i][j] + w_buf[i+2][j] + w_buf[i+1][j-1] + w_buf[i+1][j+1]) / 4.0;
9         my_diff = MAX(my_diff, fabs(my_w[i][j] - w_buf[i+1][j]));
10    }
11 }

```

最后是 slave 进程向 master 发送自己的计算结果，这里数组下标位置要细心一点，master 将结果汇总到 `w` 上，这样，下次迭代开始时，又将 `w` 分配给每个进程的 `w_buf`。 `w` 更新完成之后，对 `my_diff` 进行 reduce 操作，选取最大的 `my_diff` 到 `diff`。这里使用 `MPI_Allreduce`，因为每个进程都要进行循环条件的判断，误差是否小于阈值。

```

1 if (my_rank != 0) // if I'm slave, send my_w to master
2 {
3     // printf("I'm process %d, my begin row = %d\n", my_rank, my_first_m);
4     MPI_Send(&my_w[0][0], my_m * N, MPI_DOUBLE, 0, 0, MPI_COMM_WORLD);
5 }
6 else // if I'm the master, recv my_w from slaves
7 {
8     for (int slave = 1; slave < comm_sz; slave++)
9     {
10        int slave_first_m, slave_last_m, slave_m;
11        my_bound(1, M - 1, comm_sz, slave, &slave_first_m, &slave_last_m, &slave_m);
12        MPI_Recv(&w[slave_first_m][0], slave_m * N, MPI_DOUBLE, slave, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
13    }
14    // master should also copy its own my_w to w
15    memcpy(&w[1][0], &my_w[0][0], sizeof(double) * my_m * N);
16 }
17
18 // Reduce the maximum my_diff to diff, every process will get a copy
19 MPI_Allreduce(&my_diff, &diff, 1, MPI_DOUBLE, MPI_MAX, MPI_COMM_WORLD);

```

在 `heated-plate-mpi` 目录下，`make` 后再执行 `make test`，会依次运行 OpenMP 版本和 MPI 版本的 heated plate。默认开启最大的线程数和进程数，我的 CPU 是 6 核 12 线程的，OpenMP 的线程数为 12，MPI 进程数为 6。

下面是运行的结果。两者迭代的次数和记录的误差完全相同，最后的 Wallclock time 相差只有 0.5 秒，而 OpenMP 的线程数是 MPI 进程数的两倍（可能与 Intel 的超线程技术的实现有关），MPI 版本的性能可以说非常好。

```

heated_plate_openmp_test.txt M x
Lab > lab5 > heated-plate-mpi > heated_plate_openmp_test.txt
1
2 HEATED_PLATE_OPENMP
3 C/OpenMP version
4 A program to solve for the steady state temperature distribution
5 over a rectangular plate.
6
7 Spatial grid of 500 by 500 points.
8 The iteration will be repeated until the change is <= 1.000000e-03
9 Number of processors available = 12
10 Number of threads = 12
11
12 MEAN = 74.949900
13
14 Iteration Change
15
16 1 18.737475
17 2 9.368737
18 4 4.098823
19 8 2.289577
20 16 1.136684
21 32 0.568281
22 64 0.282805
23 128 0.141777
24 256 0.070808
25 512 0.035427
26 1024 0.017707
27 2048 0.008856
28 4096 0.004428
29 8192 0.002210
30 16384 0.001043
31
32 16955 0.001000
33
34 Error tolerance achieved.
35 Wallclock time = 8.774925
36
37 HEATED_PLATE_OPENMP:
38 Normal end of execution.
39

```

```

heated_plate_mpi_test.txt M x
Lab > lab5 > heated-plate-mpi > heated_plate_mpi_test.txt
1
2 HEATED_PLATE_MPI
3 C/MPI version
4 A program to solve for the steady state temperature distribution
5 over a rectangular plate.
6
7 Spatial grid of 500 by 500 points.
8 The iteration will be repeated until the change is <= 1.000000e-03
9 Number of processes = 6
10
11 MEAN = 74.949900
12
13 Iteration Change
14
15 1 18.737475
16 2 9.368737
17 4 4.098823
18 8 2.289577
19 16 1.136684
20 32 0.568281
21 64 0.282805
22 128 0.141777
23 256 0.070808
24 512 0.035427
25 1024 0.017707
26 2048 0.008856
27 4096 0.004428
28 8192 0.002210
29 16384 0.001043
30
31 16955 0.001000
32
33 Error tolerance achieved.
34 Wallclock time = 9.201528
35
36 HEATED_PLATE_OPENMP:
37 Normal end of execution.
38

```

另外，我还将优化前的 MPI 版本的运行结果 backup 了一下，下面是 backup 版本的结果。（优化前指的是，没有用 `w_buf` 作为接收 `w` 必要部分的缓冲区，和没有用 `my_w` 作为计算结果的缓冲区、避免 `w_buf` 的复制，这两部分的优化）

可以看到，没有用缓冲区的墙上时间是 42.5 秒，比我优化后的 9.2 秒慢了接近 5 倍。实际上，我在写代码的时候，先考虑到的是 `w_buf`，再考虑到的是 `my_w`，而前者的优化为程序带来了 2 倍的性能提升，后者为程序带来了 2.5 倍的提升。

由于我发送和接收只用到了 `w`，所以没必要用 MPI pack/unpack 来打包和解包，这对于我优化后的程序并不能带来性能提升，所以没有实现。

```
heated_plate_mpi_test.text.backup x
Lab > lab5 > heated-plated-mpi > heated_plate_mpi_test.text.backup
1
2 HEATED_PLATE_MPI
3 C/MPI version
4 A program to solve for the steady state temperature distribution
5 over a rectangular plate.
6
7 Spatial grid of 500 by 500 points.
8 The iteration will be repeated until the change is <= 1.000000e-03
9 Number of processes = 6
10
11 MEAN = 74.949900
12
13 Iteration  Change
14
15      1 18.737475
16      2  9.368737
17      4  4.098823
18      8  2.289577
19     16  1.136604
20     32  0.568201
21     64  0.282805
22    128  0.141777
23    256  0.070808
24    512  0.035427
25   1024  0.017707
26   2048  0.008856
27   4096  0.004428
28   8192  0.002210
29  16384  0.001043
30
31  16955  0.001000
32
33 Error tolerance achieved.
34 Wallclock time = 42.484814
35
36 HEATED_PLATE_OPENMP:
37 Normal end of execution.
38
```

### (3) fft\_pf 和 fft\_openmp 的性能对比

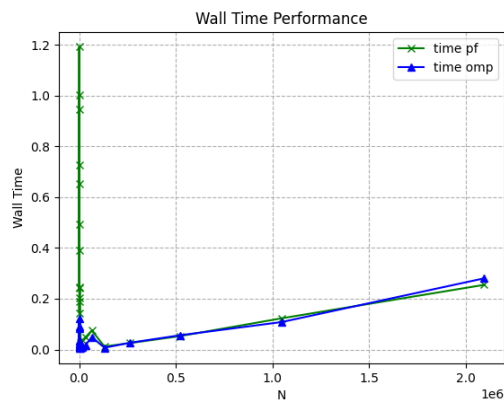
在 `fft` 下执行 `make test`，会在 `asset` 生成误差、运行时间、MFLOPS 的文件。任何并行优化都是在保证程序计算正确的前提下进行的，可以查看生成的文件中的 `error` 一列，我们的 `parallel_for` 每次迭代的误差与 `openmp` 完全相同，说明计算结果没有问题。

下面是可视化后的结果。

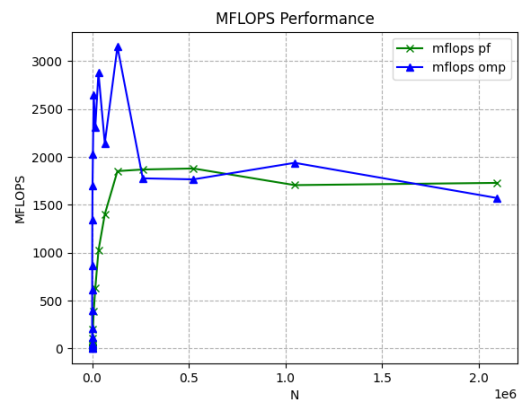
首先是墙上时间，可以看到在 `N` 较小时（大约是小于 50），我们的 `parallel_for` 运行时间远大于 `OpenMP` 的时间，但是在 `N > 1000` 时，我们的 `parallel_for` 有着与 `OpenMP` 几乎相同甚至略优的运行时间。这说明我们的程序在可扩展性上还是很好的，但是在 `N` 较小时，没有处理好线程数的问题，导致线程开销过大。

然后是 MFLOPS，对于每秒的浮点运算次数，大约在 `N < 2000` 时，`OpenMP` 的 MFLOPS 比 `parallel_for` 要快大约 1.3 ~ 1.5 倍。在 `N > 2000` 时，`parallel_for` 和 `OpenMP` 相差较小，大约相同。

## Wall Time



## MFLOPS



在 `fft` 下，执行 `make memory-test`，会使用 `valgrind` 的 `massif` 工具对程序进行内存测试。然后执行 `make print`，会将保存到 `asset` 下的文件打印出来。

parallel\_for 的结果:

[illegible]

openmp 的结果:



- <https://stackoverflow.com/questions/6481005/how-to-obtain-the-number-of-cpus-cores-in-linux-from-the-command-line>
- <https://valgrind.org/docs/manual/ms-manual.html>
- <https://stackoverflow.com/questions/52063507/how-to-use-valgrinds-massif-out-file-option-correctly>