

Lab 2: Parallel Longest Common Subsequence

姜海天 19307110022

0.1 Problem formulation

This report is to record the performance optimization for longest common subsequence. Let $A[0..M-1]$ be a string of length M and $B[0..N-1]$ be a string of length N . Then the state-transition equation is:

$$dp(i, j) = \begin{cases} 0, & i < 0 \text{ or } j < 0 \\ dp(i-1, j-1) + 1, & A[i] = B[j] \\ \max\{dp(i, j-1), dp(i-1, j)\}, & A[i] \neq B[j] \end{cases}$$

The naive implementation is

```
1 int naive_lcs_2d(const std::string A, const std::string B) {
2     const int M = A.length();
3     const int N = B.length();
4     int * dp = (int *)calloc(M * N, sizeof(int));
5     for (int i = 0; i < M; i++)
6     {
7         for (int j = 0; j < N; j++)
8         {
9             int up = (i > 0) ? dp[(i - 1) * N + j] : 0;
10            int left = (j > 0) ? dp[i * N + j - 1] : 0;
11            int upleft = (i > 0 && j > 0) ? dp[(i - 1) * N + j - 1] : 0;
12            if (A[i] == B[j])
13                dp[i * N + j] = upleft + 1;
14            else
15                dp[i * N + j] = max(up, left);
16        }
17    }
18    int res = dp[M * N - 1];
19    free(dp);
20    return res;
21 }
```

0.2 Test environment

All the code are tested on the testing machine of this course. It has a 8-core, 8-thread [Intel Xeon E5-2650](#) CPU.
L1d cache: 32K, L1i cache: 32K, L2 cache: 256K, L3 cache: 20MiB

1 1-D Space Usage for LCS

One important observation in the naive 2-D implementation is, in each iteration of the outer loop we only need values from all columns of the previous row. So there is no need to store all rows in our dp matrix, we can just store two rows at a time and use them. In that way, used space will be reduced from `dp[m+1][n+1]` to `dp[2][n+1]`. Below is the implementation of the above idea. Here I use the transpose of the dp matrix to have a better cache hit rate. Notice that few codes need to be revised compared with the 2-D version.

```
1 int lcs_basic(const char * _A, const char * _B, int M, int N) {
2     int (*dp)[2] = (int(*)[2])calloc(N * 2, sizeof(int));
3     bool bi;
4     for (int i = 0; i < M; ++i) {
5         bi = i & 1;
6         for (int j = 0; j < N; ++j) {
7             int up = (i > 0) ? dp[j][1 - bi] : 0;
8             int left = (j > 0) ? dp[j - 1][bi] : 0;
9             int upleft = (i > 0 && j > 0) ? dp[j - 1][1 - bi] : 0;
10            if (_A[i] == _B[j]) {
11                dp[j][bi] = upleft + 1;
12            } else {
13                dp[j][bi] = max(up, left);
14            }
15        }
16    }
17    int res = dp[N - 1][bi];
18    free(dp);
19    return res;
20 }
```

2 1-D by Anti-Diagonal

In this case, we can fill the dp matrix by looping over $s=i+j$ from 0 to $M+N-2$. The inner loop can be $t=j-i$ from $-M+1$ to $N-1$. To make the index bigger than 0, we can let $t=j-i+M-1$ ranging from 0 to $M+N-2$. So now we can have the sequential version of the anti-diagonal method, and there is no data dependency in the inner loop.

```
1 int lcs_ad(const char * _A, const char * _B, int M, int N) {
2     int* dp = (int*)calloc(M + N - 1, sizeof(int));
3     for (int s = 0; s < M + N - 1; ++s) { // s = i + j
4         int start = (s < M) ? (M - s - 1) : (s - M + 1);
5         int end = (s < N) ? (M + s + 1) : (N * 2 + M - s - 1);
6         for (int t = start; t < end; t += 2) { // t = j - i + M - 1
7             int i = (s + M - t - 1) / 2;
8             int j = (s + t + 1 - M) / 2;
```

```

9         if (_A[i] == _B[j]) {
10             ++dp[t];
11         } else {
12             dp[t] = max(dp[t-1], dp[t+1]);
13         }
14     }
15 }
16 int res = dp[N-1];
17 free(dp);
18 return res;
19 }

```

3 Parallel Version

To parallelize the above code, we can just add the OpenMP primitive for the inner loop.

By just adding the `#pragma omp parallel for`, I can get the following result:

<pre> 19307110022 @ ubuntu-ty-server in ~/pcpo_lab-1 make test for i in 1 2 3 4 5; do \ ./main ./input/\$i.txt ; \ done ----- M: 624 N: 946 lcs_1d(correct): 4.445 (ms) lcs_ad(correct): 3.136 (ms) lcs_ad_parallel(correct): 5.379 (ms) ----- M: 660 N: 489 lcs_1d(correct): 1.455 (ms) lcs_ad(correct): 1.539 (ms) lcs_ad_parallel(correct): 4.353 (ms) ----- M: 483 N: 841 lcs_1d(correct): 1.833 (ms) lcs_ad(correct): 1.948 (ms) lcs_ad_parallel(correct): 4.402 (ms) ----- M: 435 N: 762 lcs_1d(correct): 1.494 (ms) lcs_ad(correct): 1.584 (ms) lcs_ad_parallel(correct): 4.319 (ms) ----- M: 570 N: 995 lcs_1d(correct): 3.473 (ms) lcs_ad(correct): 3.109 (ms) lcs_ad_parallel(correct): 5.188 (ms) </pre>	<pre> 19307110022 @ ubuntu-ty-server in ~/pcpo_lab-1 make test for i in 1 2 3 4 5; do \ ./main ./input/\$i.txt ; \ done ----- M: 7680 N: 3409 lcs_1d(correct): 118.450 (ms) lcs_ad(correct): 124.063 (ms) lcs_ad_parallel(correct): 46.785 (ms) ----- M: 9458 N: 3648 lcs_1d(correct): 155.624 (ms) lcs_ad(correct): 163.290 (ms) lcs_ad_parallel(correct): 59.344 (ms) ----- M: 6915 N: 4006 lcs_1d(correct): 125.369 (ms) lcs_ad(correct): 131.233 (ms) lcs_ad_parallel(correct): 47.776 (ms) ----- M: 8088 N: 9716 lcs_1d(correct): 355.098 (ms) lcs_ad(correct): 372.226 (ms) lcs_ad_parallel(correct): 109.163 (ms) ----- M: 6041 N: 9933 lcs_1d(correct): 442.253 (ms) lcs_ad(correct): 284.016 (ms) lcs_ad_parallel(correct): 87.684 (ms) </pre>	<pre> 19307110022 @ ubuntu-ty-server in ~/pcpo_lab-1 make test for i in 1 2 3 4 5; do \ ./main ./input/\$i.txt ; \ done ----- M: 63214 N: 66079 lcs_1d(correct): 18916.234 (ms) lcs_ad(correct): 19792.404 (ms) lcs_ad_parallel(correct): 4387.698 (ms) ----- M: 42247 N: 54947 lcs_1d(correct): 10518.976 (ms) lcs_ad(correct): 10996.296 (ms) lcs_ad_parallel(correct): 2595.207 (ms) ----- M: 82455 N: 62726 lcs_1d(correct): 23423.327 (ms) lcs_ad(correct): 24512.797 (ms) lcs_ad_parallel(correct): 5499.096 (ms) ----- M: 57969 N: 67519 lcs_1d(correct): 17727.095 (ms) lcs_ad(correct): 18548.083 (ms) lcs_ad_parallel(correct): 4280.226 (ms) ----- M: 95102 N: 97740 lcs_1d(correct): 42154.234 (ms) lcs_ad(correct): 44110.214 (ms) lcs_ad_parallel(correct): 9525.971 (ms) </pre>
---	---	--

We can see that for small M and N below 1K, the overhead of creating threads takes the majority of time so that the parallel version takes longer to run. But for M and N below 10K, the parallel version takes dramatic less time to run, and the speedup is about **3**. When the length of string goes larger to below 100K, the speedup comes to about **4.5**.

Considering different threads will write to nearby positions of the memory simultaneously in the parallel version, the written value in the cache line of one processor may cause the cache line in other processor to be invalidated. And this kind of situation may happen back and forth so that the invalidation storm causes lots of communication in the hardware. One thought is to unroll the loop and make it parallel at the same time. So that one processor will be in charge of the whole cache line, and do not interfere with other processors. In OpenMP, this is simply a longer primitive: `#pragma omp parallel for schedule(static 8)`. But it actually makes the computation time longer.

