

IN4200 Partial Exam

Candidate number: 15350

April 2018

Computing PageRank

The program consist of a PE_function_15350.c file and an PE_main_15350.c file. The PE_function_15350.c consist of three functions: read_graph_from_file(filename), PageRank_iterations(d,eps) and top_n_webpages(xk,top_n).

The read_graph_from_file() function takes in a file name and reads the file, It is importanteant that the file have the same layout as the web_NotreDame.txt file!. The important variables that the function gets from the file is n_nodes and n_edges which comes from the line:

```
1 err = fscanf(file , "%*s %*s %d %*s %d" , &n_nodes , &n_edges );
```

Were it ignores the strings. The loop that saves the values to F_nodeId and T_nodeId, while excluding self-links is:

```
1 int s = 0;
2 int a,b;
3 while(s < n_edges) {
4     err = getline(&line , &len , file);
5     err = fscanf(file , "%d %d" , &a , &b);
6     if(a == b){
7         n_edges--;
8         continue;
9     }
10    F_nodeId[s] = a;
11    T_nodeId[s] = b;
12    s++;
13 }
```

n_nodes and n_edges gets used all the time in the program, while the other two is used to generate the A matrix as CRS format with arrays: row_ptr:

```
1 int cnt = 0;
2 for(int n=0; n < n_nodes; n++){
```

```

3     cnt += total.T[n]; //total.T is sum of equal values in T_nodeid
4     row_ptr[n+1] = cnt;
5 }

```

col_idx and values:

```

1  for (int j = 0; j < n_edges; j++) {
2      count = row_ptr[T_nodeId[j]] + temp[T_nodeId[j]];
3      col_idx[count] = F_nodeId[j];
4      values[count] = (double)1/total_F[F_nodeId[j]];
5      temp[T_nodeId[j]]++;
6  }

```

The function also finds the danglings point and location, which is used later in Page-Rank function.

The Page_Rank() function with input values d and ϵ , does the most of the calculation. The function iterates true a loop, were it calculates W and x and returns x , with the functions:

$$W^{k-1} = \sum_{m \in \mathcal{D}} x_m^{k-1}, \quad \mathcal{D} : \text{set of indices for all dangling webpages}$$

$$x^k = \frac{1 - d + d * W^{k-1}}{N} * \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} + d * Ax^{k-1}$$

and stops at the stopping criterion:

$$\max_{1 \leq j \leq N} |x_j^k - x_j^{k-1}| < \epsilon$$

The maximumly achievable (The big O) speed of the page iterations depends on the for and while loops: The first for loop gives $O(n)$, and from the while loop we get $O(k*w+n*n)$, all this gives: $O(n) + O(k*w+k*n*n) \approx O(n^2)$. The factors that prevent reaching the ideal speed is that the threads runs in different speed and, some threads need to wait for them to finish before they can do their job. The way to calculated the iterations is inside a while loop who runs up to the stopping criterion hits, runs parallel:

```

1  while (!array_ok) {
2      array_ok = true;
3      k++;
4      sum_wk = 0;
5      // calculating W
6      #pragma omp parallel for reduction(+:sum_wk)
7      for (int i=0; i<w_dang; i++){
8          sum_wk += xk_1[w_place[i]];
9      }
10
11     tamp = (1 - d + (d * sum_wk))/((double)n_nodes);

```

```

12     int j;
13     //calculating x
14     #pragma omp parallel for private(j)
15     for (int i = 0; i < n_nodes; i++) {
16         xk[i] = tmp; // calculating outside for loop since const
17         for (j = row_ptr[i]; j < row_ptr[i+1]; j++) {
18             xk[i] += d * (values[j] * xk_1[col_idx[j]]);
19         }
20         // Stopping crirerion
21         if (fabs(xk[i]-xk_1[i])>eps){
22             array_ok = false;
23         }
24     }
25     //updating values
26     tmp = xk_1;
27     xk_1 = xk;
28     xk = tmp;
29 }

```

The last function **Top_n_webpages()**, gets the x from Page_Rank() and a chosen value from the user that indicates how many top values they want to print out. It start with calculating the first top value and uses it to find the others, and it only goes true the loop top_n times n_nodes time. The function is build with 2 parallel run, the first one calculates one top value and saves it as max, it is important to have a another variable outside the parallel to hold on the idx value, and to make sure that all the threads give out the max value:

```

1  #pragma omp parallel
2  {
3      int idx=0;
4      #pragma omp for reduction(max:max)
5      for (int i = 0; i < n_nodes; i++) {
6          if(xk[i] > max) {
7              max = xk[i];
8              idx = i;
9          }
10     }
11     #pragma omp critical
12     {
13         if(xk[idx] == max){
14             idx2 = idx;
15         }
16     }
17 }
18
19 top_n_index[0] = idx2;
20 prevmax = max;
21 max = 0;

```

The next loop goes true the program again and compares the value with new max and so on.

```

1 for (int i = 1; i < top_n; i++) {
2     #pragma omp parallel
3     {
4         int idx = 0;
5         #pragma omp for reduction(max:max)
6         for (int j = 0; j < n_nodes; j++) {
7             if(xk[j] > max && xk[j] < prevmax) {
8                 max = xk[j];
9                 idx = j;
10            }
11        }
12        #pragma omp critical
13        {
14            if(xk[idx] == max){
15                idx2 = idx;
16            }
17        }
18    }
19    top_n_index[i] = idx2;
20    prevmax = max;
21    max = 0;
22 }

```

Results and OPenmp

As shown in the README.md file it is many ways to run the program, but i used the same optimization option, "-O2 -fopenmp" in every run with the same damping constant and only top 10 values were printed. I have also tested the program without any OpenMP, which takes a bit more time. When running the program with $\epsilon = 1e-10$ we get 89 iterations before the stopping criterie hits, and when using bigger ϵ it will be less iterations and more iteration when using smaller ϵ . The only changes was the OMP_NUM_THREADS. My pc have 4 cpu cores, and each cpu core can use two threads, that means that i should not run more then 8 threads. A run of the program with OMP_NUM_THREADS=8 (which gives the best running time):

```

15350 > export OMP_NUM_THREADS=8
15350 > gcc-8 -O2 -fopenmp PE_main_15350.c
15350 > ./a.out web-NotreDame.txt 0.85 1e-10 10
Reading file...
File read and closed...
Finding the CRS arrays ...
Finding dangling points ...
Total dangling points are: 188795
Iterating true x_k...
Calculating xk ...

Total iterations : 89
Finding top 10 values...
  PAGE      | MAX VALUE
  1964      | 0.005627
    1       | 0.005405
 124803     | 0.003326

```

```

12130 | 0.002857
191268 | 0.002749
32831 | 0.002732
3452 | 0.002590
83607 | 0.002460
1974 | 0.002376
142733 | 0.002341
Reading graph from file: took 0.487371 seconds
PageRank iterations: took 0.234013 seconds
Top n webpages: took 0.003070 seconds
ALL DONE!!, TOTAL TIME = 0.724474 seconds

```

When running the program with 1,2 and 4 Threads, we see that the Page iteration function uses less and less time:

```

15350 > export OMP_NUM_THREADS=1
15350 > ./a.out web-NotreDame.txt 0.85 1e-10 10
Reading graph from file: took 0.481889 seconds
PageRank iterations: took 0.363762 seconds
Top n webpages: took 0.003147 seconds
ALL DONE!!, TOTAL TIME = 0.848815 seconds

```

```

15350 > export OMP_NUM_THREADS=2
15350 > ./a.out web-NotreDame.txt 0.85 1e-10 10
Reading graph from file: took 0.485228 seconds
PageRank iterations: took 0.278124 seconds
Top n webpages: took 0.003123 seconds
ALL DONE!!, TOTAL TIME = 0.766499 seconds

```

```

15350 > export OMP_NUM_THREADS=4
15350 > ./a.out web-NotreDame.txt 0.85 1e-10 10
Reading graph from file: took 0.489218 seconds
PageRank iterations: took 0.238720 seconds
Top n webpages: took 0.002538 seconds
ALL DONE!!, TOTAL TIME = 0.730501 seconds

```

There is no point to try more threads, it will only take more times to load the threads then running the program. We see that it almost takes same time for 8 and 4 number of threads. Here is a run of the program with 10000 top n values with and without OpenMP:

```

15350 > gcc-8 PE_main_15350.c -fopenmp
15350 > export OMP_NUM_THREADS=4
15350 > ./a.out web-NotreDame.txt 0.85 1e-10 10000
Reading graph from file: took 0.515665 seconds
PageRank iterations: took 0.607381 seconds
Top n webpages: took 5.016241 seconds
ALL DONE!!, TOTAL TIME = 6.139315 seconds

```

```

15350 > gcc-8 test2.c
15350 > ./a.out web-NotreDame.txt 0.85 1e-10 10000
WITHOUT OPENMP
Reading graph from file: took 0.529530 seconds
PageRank iterations: took 1.130368 seconds
Top n webpages: took 9.358309 seconds
ALL DONE!!, TOTAL TIME = 11.018255 seconds

```

The time using OpenMP is much better than without openmp!.