

ROYAL INSTITUTE OF TECHNOLOGY



INTRODUCTION TO HIGH-PERFORMANCE  
COMPUTING, DN2258

---

# Final project, parallel search Draft

---

Sindri Magnússon, [sindrim@kth.se](mailto:sindrim@kth.se), 871209-7156  
Brynjar Smári Bjarnasson, [bsbj@kth.se](mailto:bsbj@kth.se), 840824-4690

September 28, 2012

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	The data . . . . .	2
<b>2</b>	<b>Theoretical performance estimation</b>	<b>2</b>
<b>3</b>	<b>Results</b>	<b>2</b>
3.1	OpenMP . . . . .	2
3.2	MPI . . . . .	4
<b>A</b>	<b>Code</b>	<b>4</b>
A.1	OpenMP . . . . .	4

# 1 Introduction

In today's scientific environment the amount of data that is generated is increasing like never before. One consequence of these revolutionary time is the ability to process and work with the data in acceptable time. We plan to look further into one of the most important problem in the data universe which is searching.

embarrasly parallel

## 1.1 The data

To isolate the main problem and to be able to analyse the results in as direct way as possible we decided to use as simple data as we could think of. Another benefit of a simple data set is the ability of creating a data with a large number of rows. The data set is a file containing in every line an id and a sequence of six uniform random integers on the interval  $[0, 9]$ . All of the id's and all the sequences of random numbers have the same number of characters so the number of characters in each line is fixed.

Here below we see an example of the data with 1000 rows.

```
1 0000000 124523
2 0000001 561313
3 0000002 921836
4 ...
5 9999999 196142
```

example\_data.txt

The data we used for testing contains  $10^9$  lines and each line has 17 characters. Since a character is one byte the file is  $17 \times 10^9$  bytes (17 Gbyte).

## 2 Theoretical performance estimation

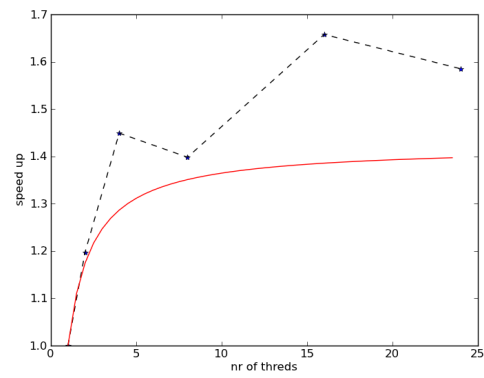
Since each comparison is independent we see that the problem of searching is actually embarrassingly parallel. However we need to read the data file from disk which is a sequential operation. So if we only look at the searching part of the algorithm then we would expect a linear speedup. But if we look at the whole thing and let  $P$  be the time percent of the search in the parallel part then from Amdahl's law we expect the speedup on  $N$  threads to be:

$$S_N = \frac{1}{(1 - P) + \frac{P}{N}}$$

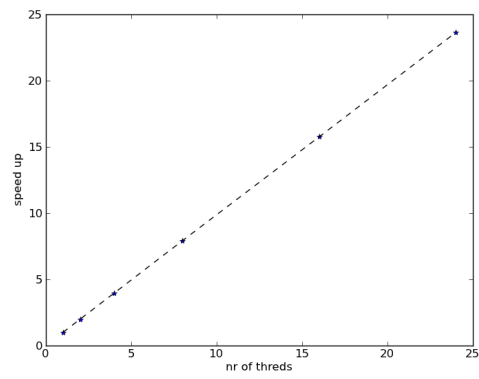
## 3 Results

### 3.1 OpenMP

Some graphics of the result we got. We have still not been able to run it on Lindgren but these results are from our personal computers.



(a) The red-line is the theoretical speed-up from Ahmads-law, the stars are show the acctual speed-up.



(b) Speed up of the search part

## 3.2 MPI

# A Code

## A.1 OpenMP

```
1 #include <stdio.h>
2 #include <string.h>
3 #include <time.h>
4 #include <math.h>
5 #include <omp.h>
6 #include <stdlib.h>
7 long long ae_load_file_to_memory(const char *filename, char
    **result)
8 {
9     // filename: path to the file to read
10    // result: pointer to character array that contains the
        content of the file.
11    long long size = 0;
12    FILE *f = fopen(filename, "rb");
13    if (f == NULL)
14    {
15        *result = NULL;
16        return -1; // -1 means file opening fail
17    }
18    fseek(f, 0, SEEK_END);
19    size = ftell(f);
20    fseek(f, 0, SEEK_SET);
21    *result = (char *)malloc(size+1);
22    if (size != fread(*result, sizeof(char), size, f))
23    {
24        free(*result);
25        return -2; // -2 means file reading fail
26    }
27    fclose(f);
28    (*result)[size] = 0;
29    return size;
30 }
31
32 int read_file(char* input_file, char* key, int result_size_block,
    long long nr_lines, int line_size){
33     // input_file: is direction to a data file.
34     // key: the search string.
35     // result_size_block: is the block size of the array that
36     // contains the results which contains the id
37     // of the matching data.
38     // nr_lines: nr of lines in file.
39     // line_size: size of line in char
40     //
41     // It is assumed that all lines in input_file have the same
        length and are on the form: "0000001 123123\n" where first
        is an index and then the search values
42
43     char line[line_size]; // line in file
44     long long i;
45     int result_size = result_size_block; // initialized result size
46     int *result;
47     int result_counter;
48     int key_len = strlen(key);
49
50     result = malloc(sizeof(*result)*result_size);
```

```

51 #pragma omp parallel for private(i,line)
    shared(result,result_counter,result_size,key,result_size_block)
52 for ( i = 0; i < nr_lines; i++ ) {
53     strncpy(line,&input_file[i*line_size],line_size);
54
55     line[line_size - 1] = '\0';
56     // line contains a single line from the file. We compare the
        last part of the line to the search key.
57     if ( strcmp(&line[line_size-(key.len+1)], key) == 0){
58         // if we have used all spaces in result_counter we need to
            reallocate and increase the size.
59         if (result_size == result_counter ){
60             result_size = result_size + result_size_block;
61             result = realloc(result , result_size*sizeof(*result) );
62             if (result == NULL){
63                 printf("Error reallocating memory\n");
64                 exit(1);
65             }
66         }
67         // we keep the line number where we found the key
        result[result_counter] = atoi(strtok(line,"\t\n"));
68         result_counter++;
69     }
70 }
71 }
72 return result_counter;
73 }
74
75
76 int main( int argc , const char* argv[] )
77 {
78     double start,end;
79     double dif;
80
81     char* file_name = "../data/file.txt";
82     char* result;
83     char* search_key = "123123";
84
85     long long nr_bytes;
86     long long i;
87     long long nr_lines;
88
89     int string_size = strlen(search_key);
90     int block_size;
91     int read_count;
92     int line_size;
93
94     start = omp_get_wtime();
95     nr_bytes = ae_load_file_to_memory(file_name,&result);
96     end = omp_get_wtime();
97     dif = end-start;
98     printf("LoadFile: %f\n", dif);
99
100    // assume each line in file is equally long. here we get the
        line size.
101    for ( i=0 ; i<nr_bytes ; i++ ){
102        if ( result[i]=='\n' ){
103            line_size = i+1;
104            break;
105        }
106    }
107
108    //

```

```

109     nr_lines = nr_bytes / line_size;
110
111     block_size = 1.5*(nr_lines/pow(20,string_size));
112     read_count;
113     start = omp_get_wtime();
114     read_count =
        read_file(result,search_key,block_size,nr_lines,line_size);
115     end = omp_get_wtime();
116     dif = end - start;
117     printf("Search: %f\n",dif);
118     printf("result found: %i\n", read_count);
119
120 }

```

../src/search\_openmp.c