

Umeå University
Department of Computing Science

Parallel Programming 7.5 p
5DV152

Exercises, Chapter/Topic 2

Submitted 2017-01-26
Author: Lorenz Gerber (dv15lgr@cs.umu.se lozger03@student.umu.se)
Instructor: Lars Karlsson / Mikael Ränner

Contents

1	Introduction	1
2	2.4 - Counting pages	1
3	2.8 - Hardware multithreading and caches	1
4	2.10 - Communication overhead	1
5	2.16 - Speedup efficiency	1
6	2.19 - Scalability	3
7	2.20 - Linear speedup and strong scalability	3
8	2.23 - Alternative algorithm for computing histogram	5
	References	5

1 Introduction

This report is part of the mandatory coursework. It describes the solutions for several chosen exercises from the course book [1].

2 2.4 - Counting pages

The 12 lower bits are used for within page addressing and the 20 upper bits of the 32 bit address are used for pages. With 20 bits, 2^{20} pages (1'046'576) can be addressed. 12 bits corresponds to a page size of 4Kb. Hence the total virtual memory is 4Gb.

3 2.8 - Hardware multithreading and caches

Caching operates on whole cache lines. Hence, the chance that another process/thread changes something in a specific cache line increases with cache size and number of processes/threads. The specific situation that can happen is called 'false sharing': When one process changes data in a cache line, there is no possibility to check or know for another process which data exactly was changed. It can very well be the case that to the current process unrelated data was changed and a cache reload would not be needed. However, cache has to be reloaded.

4 2.10 - Communication overhead

The calculations can be split in two parts: the instructions and the communication. For both a) and b) the instructions will take the same time:

$\text{instructions} / \text{cores} / \text{instructions_per_second}$

Hence $10^{12} \div 1000 \div 10^6 = 1000\text{sec}$. For communications, the time is calculated as

$\text{messages_to_send} * \text{time_per_message}$

while messages to send is described as $10^9(p - 1)$.

a) if sending one message takes 10^{-9} seconds, the communication with 1000 cores will take: $10^9(1000 - 1) \times 10^{-9} = 999\text{sec}$, together with the actual calculation, $1'999\text{sec} = 33.32\text{min}$.

b) if sending one message takes 10^{-3} seconds, the communication with 1000 cores will take: $10^9(1000 - 1) \times 10^{-3} = 9.99 \times 10^{14}\text{s}$, or several trillion years...

5 2.16 - Speedup efficiency

To calculate the hypothetical speed up between a serial and parallel program with the run time formulas $T_{\text{serial}} = n^2$ and $T_{\text{parallel}} = \frac{n^2}{p} + \log_2(p)$ where $T = [\text{ms}]$, a short program was written that calculates the speedup $S = \frac{T_{\text{serial}}}{T_{\text{parallel}}}$. The obtained data was then plotted (figure 1) using the R statistical environment [2].

```
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
```

2(5)

```
#define START_N 10
#define START_P 1
#define N_LENGTH 6
#define P_LENGTH 8

int main (void) {

    FILE *f = NULL;
    if((f = fopen("speedup.txt", "w"))<0){
        perror("fopen");
        exit(EXIT_FAILURE);
    }

    FILE *g = NULL;

    if((g = fopen("efficiency.txt", "w"))<0){
        perror("fopen");
        exit(EXIT_FAILURE);
    }

    double n = START_N;
    double p = START_P;

    double speedUp = 0;
    double efficiency = 0;

    for(int i = 0 ; i < N_LENGTH; i++){
        for(int j = 0; j < P_LENGTH; j++){

            speedUp = n*n / ((n * n) / p + (log(p)/log(2)));
            efficiency = n*n / (p * (n * n) / p + (log(p)/log(2)));
            fprintf(f, "%lf\t", speedUp);
            fprintf(g, "%lf\t", efficiency);
            printf("n = %lf, p = %lf, speedup = %lf, efficiency = %lf\n", n, p, speedUp, efficiency);
            p*=2;
        }

        fprintf(f, "\n");
        fprintf(g, "\n");
        n*=2;
        p = START_P;
    }

    fclose(f);
    fclose(g);
    return 0;
}
```

a) Figure 1 shows graphs for constant n (upper graph) and constant p (lower graph) respectively. From the upper graph, it can be seen that for constant p , the growth of speedup for larger and larger problem instances n decreases and eventually closes in to a horizontal asymptote. The lower graph shows, that for large problem sizes n , increasing the number of cores p leads to almost linear speedups, while when we increase the number of cores for small problem instances, the speed up levels off very soon.

b) Efficiency is defined as $E = \frac{T_{serial}}{pT_{parallel}}$. In the current question, $T_{parallel}$ is defined as $\frac{T_{serial}}{p} + T_{overhead}$. Hence efficiency E translates here to $\frac{T_{serial}}{p(\frac{T_{serial}}{p} + T_{overhead})} = \frac{T_{serial}}{T_{serial} + pT_{overhead}}$. From the last expression, it can be seen, that if p is held constant and T_{serial} grows faster than $T_{overhead}$, the ratio will increase, hence the efficiency increases. The opposite follows, that if $T_{overhead}$ grows monotonically faster than T_{serial} , it will sooner or later dominate the denominator, the ratio will decrease, hence the efficiency decrease.

6 2.19 - Scalability

The efficiency for the given T_{serial} and $T_{parallel}$ can be expressed as:

$$E = \frac{n}{p(\frac{n}{p} + \log_2(p))} = \frac{n}{n + p \times \log_2(p)}$$

To understand the problem, we substitute $p \times \log_2(p)$ with U . Then the expression for efficiency E reads $\frac{n}{n+U}$. It becomes obvious that if U increases by p , then n has to increase by $\frac{p \times U}{U}$ to maintain the same ratio E :

$$E_{constant} = \frac{n \frac{p \times U}{U}}{n \frac{p \times U}{U} + p \times U}$$

From this follows:

$$\frac{n \frac{k \times \log(k \times p)}{\log(p)}}{n \frac{k \times \log(k \times p)}{\log(p)} + k \times p \times \log_2(k \times p)}$$

- Hence, the sought factor for n is $\frac{k \times \log(k \times p)}{\log(p)}$
- If $p = 8$ then $16 = p \times 2$, hence $k = 2, p = 8$. Filling these values into the expression from a), the numeric result is $\frac{8}{3}$.

7 2.20 - Linear speedup and strong scalability

A linear speed-up is when a parallel program using p cores runs p times faster than the serial one. So yes, by definition, if a parallel program has a linear speed-up, it is strongly scalable. However, in practice it will hardly be obtainable to have a perfect linear speed-up. In this case, amdahl's law shows that if not 100% of the code can be made parallel, it can be at most weakly scalable.

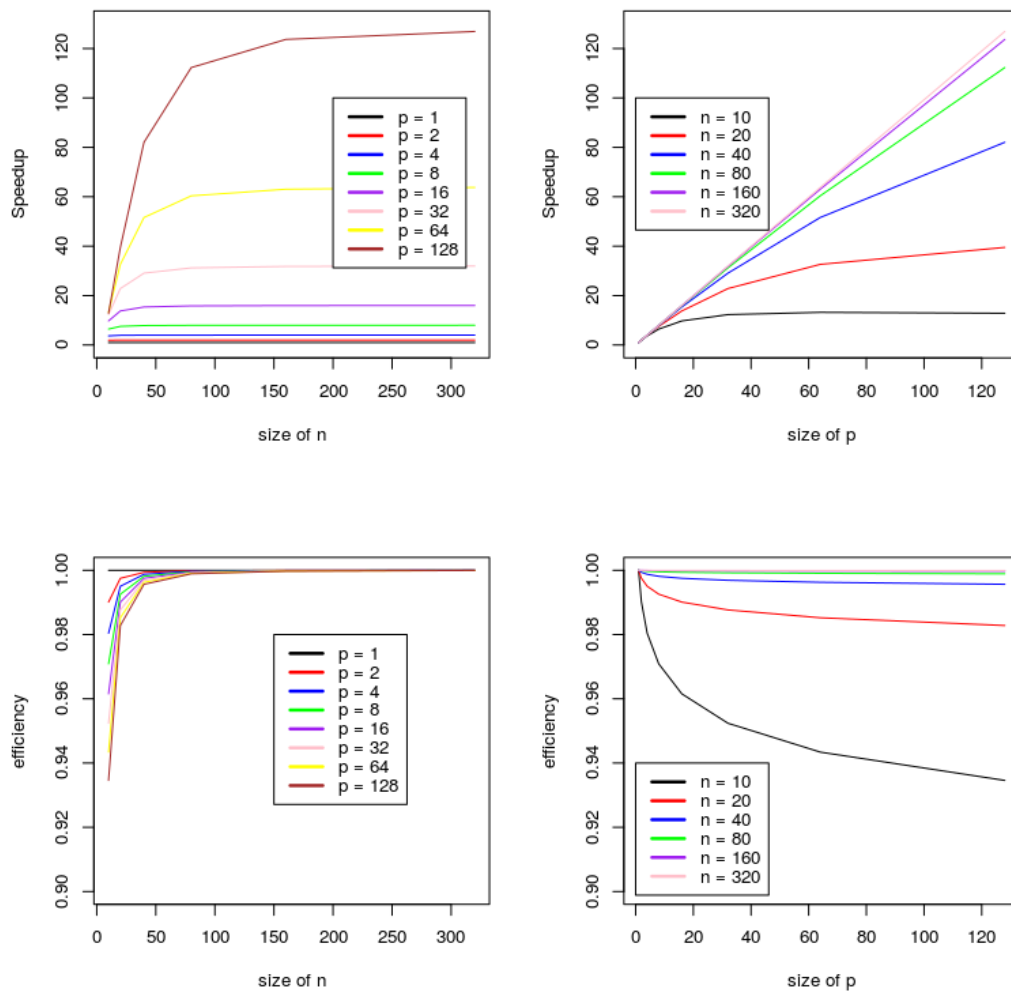


Figure 1: This figure shows the speedup for constant p , respectively constant n .

8 2.23 - Alternative algorithm for computing histogram

Assuming `data_count` c and `bin_counts` b to be $c \gg b$, agglomerating over bins should result in fewer but longer/larger tasks. Many small tasks could be an advantage from a practical point of view when we want to scale up, as we can easily distribute the tasks on many cores. Further, when we agglomerate over data, we can sometimes save time if a low index bin is a hit (in case of using linear search from low to high). When we aggregate from bins, we always need to iterate over full length c . Hence in most cases, it will be slower. While being slower, each run should be quite well predictable in run-time. This again could be an advantage when we run parallel on just a few cores.

References

- [1] P.S. Pacheco. *An Introduction to Parallel Programming*. Morgan Kaufman, 2011.
- [2] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2016.