

High Performance Computing

Homework 1

Understanding the Programming Platform

Report

Peiyun Hu, 2010011297

Department of Computer Science & Technology

Sep. 22, 2012

1 Environment

1.1 Hardware

1.1.1 CPU

Model Name: Intel(R) Core(TM)2 Quad CPU @ 2.93GHz

CPU MHz: 1600.000

Cores: 4

But actually, only one core is involved in this task. What is more, the '1600 MHz' means that at which the processor is running right now, and '2.93 GHz' is the maximum CPU Speed.

1.1.2 Memory

Total Memory: 4055940 KB

1.1.3 Cache

Cache Size: 4096 KB

1.2 Software

1.2.1 OS

Linux version 2.6.38-8-generic

1.2.2 Compiler

gcc version 4.5.2 (Ubuntu/Linaro 4.5.2-8ubuntu4)

1.2.3 Compiler Options

Please refer to makefile.

1.2.4 Timing Method

As for the timing method, only meaningful calculation time are counted, which means that time spent in initialization and output are excluded from the final time.

Additionally, there are minute differences under different conditions. Firstly, in the vector-to-vector(\vec{a}, \vec{b}) computation, one float operations are involved in each pair of elements, which is a multiplication of $a_i \times b_i$. Secondly, in the matrix-to-vector calculation, two float operations need to be counted in each pair, which is a multiplication and an addition. Besides, matrix-to-matrix condition is similar with the m2v's condition.

Consequently, assuming that the problem number of rows and columns are both n, we could infer that:

$$\text{FLOPS}_{v2v} = \frac{n}{t} \quad (1)$$

$$\text{FLOPS}_{m2v} = \frac{2n^2}{t} \quad (2)$$

$$\text{FLOPS}_{m2m} = \frac{2n^3}{t} \quad (3)$$

2 Performance

2.1 Theoretical Peak Performance

As is shown in 1.1.1, it is known that A Core 2 Quad has 4 cores. HPC world uses the following formulae for node theoretical peak performance:

Node performance in GFLOPS = (CPU speed in GHz) x (number of CPU cores) x (CPU instruction per cycle) x (number of CPUs per node). [2]

In this sequential task, the node performance is as below:

$$\text{Performance} = 2.93\text{GHz} \times 1 \times 4 \times 1 = 11.72\text{GHz} \quad (4)$$

An important note about CPU instruction per cycle is 4, as is shown in [1].

Moreover, considering there are two indicators of cpu speed, we choose the higher one, because the cpu is smart enough to gear up when more calculation is needed. So the indicator '2.93 GHz' is more appropriate.

2.2 Experiment Results

2.2.1 Vector-Vector

The efficiency in varied size of dataset is shown as Figure 1. As we can see from the chart, in small dataset, the performance of float rivals that of double. Because the dataset is rather small that the time of timing is not short enough to be ignored, even timing possessed most of time consumed.

When dataset gets bigger, float is almost twice faster than double, owing to their length.

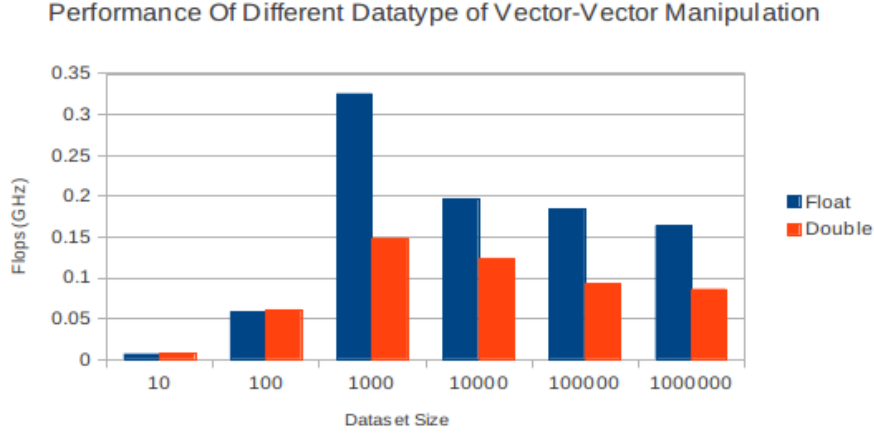


Figure 1: Performance of Vector-Vector Manipulation With Different Implements

2.2.2 Matrix-Vector

The efficiency in varied size of dataset is shown as Figure 2. As shown in the chart, the efficiency of float calculation is rather stable. And the double is increasing gradually, but stay lower than float operation.

2.2.3 Matrix-Matrix

The efficiency in varied size of dataset is shown as Figure 3. Two methods of multiplication are involved.

The first one is simple and intuitive. After reading in the matrix, do the multiplication directly like below:

```

1  for(i=0; i<a_rows; i++)
2      for(j=0; j<b_cols; j++)
3          for(k=0; k<a_cols; k++)

```

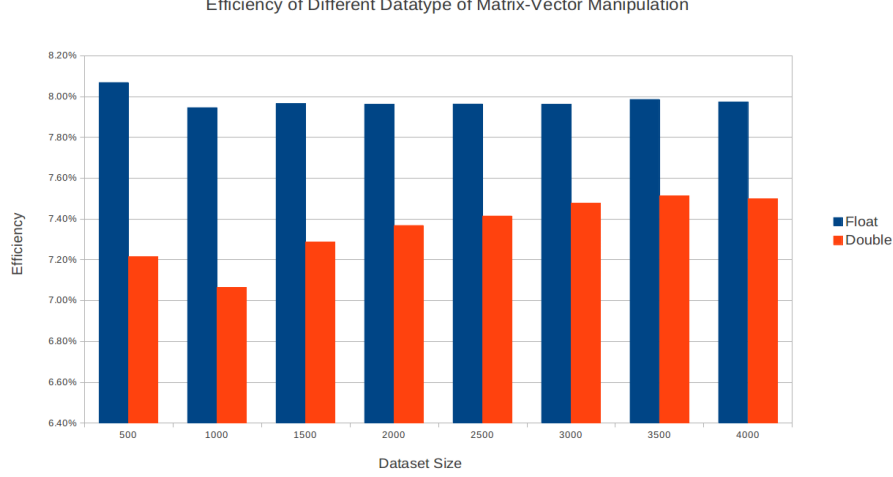


Figure 2: Efficiency of Matrix-Vector Manipulation With Different Implements

```
4 | res[i*a_cols+j] += v1[i*a_cols+k] * v2[k*b_cols+j];
```

But when considering taking advantage of cache, we can make it better by preprocessing the second matrix v_2 . By transposing v_2 , we can achieve a better hit rate when accessing v_2 , owing to the linear arrangement of the array in the physical storage. And the code is a little different from above:

```
1 | for(i=0; i<a_rows; i++)
2 |     for(j=0; j<b_cols; j++)
3 |         for(k=0; k<a_cols; k++)
4 |             res[i][j] += v1[i*a_cols+k] * v2[j*a_cols + k];
```

We can notice that the index of each array now increment by 1. And, we got almost 4 times faster than that of former method as is shown in 3.

What is more **interesting** is, in Figure 3, the float implement of first method is nearly on the same level with the second method. To interpret this, we should review 1.1.3. The cache can hold:

$$\frac{4096 * 1024}{500^2 * 2 * 4} = 2.097 \quad (5)$$

32-bit single-precision floating numbers. That is to say, without the optimization of memory access, CPU could still reach a high hit rate.

But cache is not big enough to store double-precision floating matrixs of this size at one time. That is why we see a sharp contrast between performance of float and performance of double, in the first method.

3 Further Discussion

Besides, noticing that SSE, which was firstly introduced in Pentium III, added eight new-bit registers known as XMM0 through XMM7, the cpu could handle 4 32-bit float at one time. Owing to SSE2 introduced later since Pentium Series, the usage of XMM is expanded to 64-bit double-precision mode. Additionally, the AMD64(originally called x86-64) added further eight registers from XMM8 to XMM15.

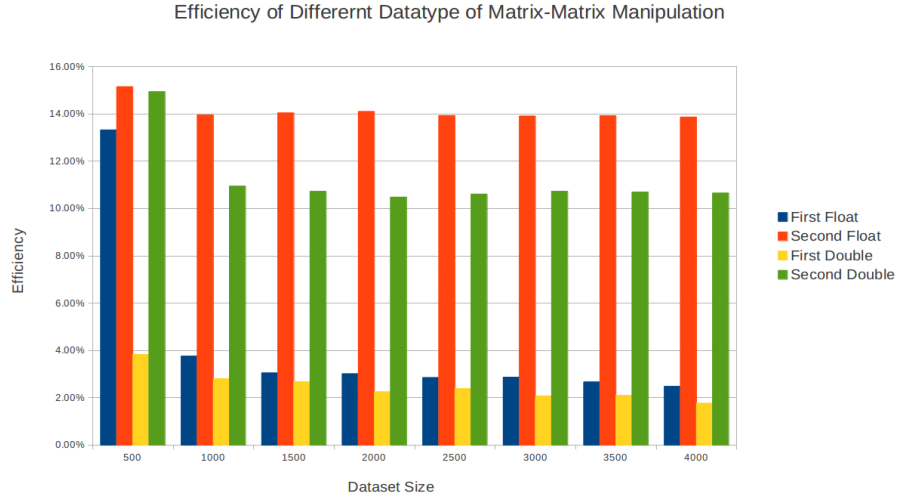


Figure 3: Efficiency of Matrix-Matrix Manipulation With Different Implements

So with SSE programmed inline, we could further enhance the performance. However, it is a little embarrassing in this task, for I did not find a easy way to mutiply using SSE.

References

- [1] Intel. *Intel 64 and IA-32 Architectures Optimization Reference Manual*, November 2009. 2.1 Intel Core Microarchitecture and Enhanced Intel Core Microarchitecture, Intel Wide Dynamic Execution, Page 30.
- [2] NOVATTE. How to calculate theoretical peak performance of a cpu-based hpc system. <http://novatte.com/blog/2012/03/how-to-calculate-theoretical-peak-performance-of-a-cpu-based-hpc-system/>.

Appendix

A Implement of Vector-Vector

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <time.h>
4
5
6  #define start_time clock_gettime(CLOCK_MONOTONIC, &start);
7  #define end_time clock_gettime(CLOCK_MONOTONIC, &finish);
8
9  struct timespec start, finish;
10
11 int main(int argc, char *argv[]){
12     int output_switch;
13     char *input_file_name;
14     char *output_file_name;

```

```

15     FILE* pfile;
16
17     /* resolve the input arguments */
18     if (argc < 2){
19         printf("Error! the number of arguments is less(%d)!\n",
20             argc);
21         system("pause");
22         return 0;
23     }
24     input_file_name = argv[1]; // get input file name
25
26     output_switch = 0;        // output results or not
27     if (argc > 2) {
28         output_switch = 1;
29         output_file_name = argv[2];
30     }
31
32     if ((pfile = fopen(input_file_name, "rb")) == NULL){
33         printf("file: %s can not be opened \n", input_file_name);
34         system("pause");
35         return 0;
36     }
37
38     int size;
39     float *v1, *v2, *res;
40
41     // set size
42     fread((void*)&size, sizeof(int), 1, pfile);
43
44     // alloc
45     v1 = (float*)malloc(size * sizeof(float));
46     v2 = (float*)malloc(size * sizeof(float));
47     res = (float*)malloc(size * sizeof(float));
48
49     // read in data
50     fread((void*)v1, sizeof(float), size, pfile);
51     fread((void*)v2, sizeof(float), size, pfile);
52
53     // set start
54     start_time
55
56     // calculate
57     int i;
58     for(i=0; i<size; i++)
59         res[i] = v1[i] * v2[i];
60
61     // set end
62     end_time
63     printf("%d %.16lf\n", size, finish.tv_sec-start.tv_sec + (
64         double)(finish.tv_nsec - start.tv_nsec) / 1000000000.0);
65
66     #ifndef SILENT
67     for (i = 0; i < a_rows; i++){
68         printf("%f ", v1[i]);
69     }
70     printf("\n");

```

```

69 #endif
70
71 /* output results */
72 if (output_switch){
73     if ((pfile = fopen(output_file_name, "wb")) == NULL){
74         printf("file: %s can not be opened \n",
75             output_file_name);
76         system("pause");
77         return 0;
78     }
79     fwrite(res, sizeof(float), size, pfile);
80     fclose(pfile);
81 }

```

B Implement of Matrix-Vector

```

82
83 #include <stdio.h>
84 #include <stdlib.h>
85 #include <time.h>
86 #include <memory.h>
87
88 #define start_time clock_gettime(CLOCK_MONOTONIC, &start);
89 #define end_time clock_gettime(CLOCK_MONOTONIC, &finish);
90
91 struct timespec start, finish;
92
93 int main(int argc, char *argv[]){
94     int output_switch;
95     char *input_file_name;
96     char *output_file_name;
97     FILE* pfile;
98
99     /* resolve the input arguments */
100     if (argc < 2){
101         printf("Error! the number of arguments is less(%d)!\n",
102             argc);
103         system("pause");
104         return 0;
105     }
106     input_file_name = argv[1]; // get input file name
107
108     output_switch = 0; // output results or not
109     if (argc > 2) {
110         output_switch = 1;
111         output_file_name = argv[2];
112     }
113
114     if ((pfile = fopen(input_file_name, "rb")) == NULL){
115         printf("file: %s can not be opened \n", input_file_name);
116         system("pause");
117         return 0;

```

```

117     }
118
119     int a_rows, a_cols, b_cols;
120     float *v1, *v2, *res;
121
122     // set size
123     fread((void*)&a_rows, sizeof(int), 1, pfile);
124     fread((void*)&a_cols, sizeof(int), 1, pfile);
125
126     // alloc
127     v1 = (float*)malloc(a_rows*a_cols * sizeof(float));
128
129     v2 = (float*)malloc(a_cols * sizeof(float));
130     res = (float*)malloc(a_cols * sizeof(float));
131
132     // read in data
133     fread((void*)v1, sizeof(float), a_rows*a_cols, pfile);
134     fread((void*)v2, sizeof(float), a_cols, pfile);
135     // memset
136     memset(res, 0.0, a_rows * sizeof(float));
137
138     // set start
139     start_time
140
141     // calculate
142     int i, j;
143     for(i=0; i<a_rows; i++)
144         for(j=0; j<a_cols; j++)
145             res[i] += v1[i*a_cols+j] * v2[j];
146
147     // set end
148     end_time
149     printf("%d %d %.16lf\n", a_rows, a_cols, finish.tv_sec-start.
        tv_sec + (float)(finish.tv_nsec - start.tv_nsec) /
        1000000000.0);
150
151     printf("a: %d %d \n", a_rows, a_cols);
152     #ifndef SILENT
153     for (i = 0; i < a_rows; i++){
154         for (j = 0; j < a_cols; j++){
155             printf("%f ", v1[i*a_cols+j]);
156             printf("\n");
157         }
158     }
159     #endif
160
161     /* output results */
162     if (output_switch){
163         if ((pfile = fopen(output_file_name, "wb")) == NULL){
164             printf("file: %s can not be opened \n",
165                 output_file_name);
166             system("pause");
167             return 0;
168         }
169         fwrite(res, sizeof(float), a_cols*b_cols, pfile);
170         fclose(pfile);
171     }

```



```
170
171 }
```

C Implement of Matrix-Matrix

```
172
173 #include <stdio.h>
174 #include <stdlib.h>
175 #include <time.h>
176 #include <memory.h>
177
178 #define start_time clock_gettime(CLOCK_MONOTONIC, &start);
179 #define end_time clock_gettime(CLOCK_MONOTONIC, &finish);
180
181 struct timespec start, finish;
182
183 int main(int argc, char *argv[]){
184     int output_switch;
185     char *input_file_name;
186     char *output_file_name;
187     FILE* pfile;
188
189     /* resolve the input arguments */
190     if (argc < 2){
191         printf("Error! the number of arguments is less(%d)!\n",
192             argc);
193         system("pause");
194         return 0;
195     }
196     input_file_name = argv[1]; // get input file name
197
198     output_switch = 0; // output results or not
199     if (argc > 2) {
200         output_switch = 1;
201         output_file_name = argv[2];
202     }
203
204     if ((pfile = fopen(input_file_name, "rb")) == NULL){
205         printf("file: %s can not be opened \n", input_file_name);
206         system("pause");
207         return 0;
208     }
209
210     float *v1, *v2, *res;
211     int a_rows, a_cols, b_cols;
212
213     // set size
214     fread((void*)&a_rows, sizeof(int), 1, pfile);
215     fread((void*)&a_cols, sizeof(int), 1, pfile);
216     fread((void*)&b_cols, sizeof(int), 1, pfile);
217
218     // alloc
219     v1 = (float*)malloc(a_rows* a_cols * sizeof(float));
220     v2 = (float*)malloc(a_cols* b_cols* sizeof(float));
```

```

219     res = (float*)malloc(a_rows*b_cols* sizeof(float));
220
221 // read in data
222     fread((void*)v1, sizeof(float), a_rows*a_cols, pfile);
223     fread((void*)v2, sizeof(float), a_cols*b_cols, pfile);
224
225     fclose(pfile);
226
227 // memset
228     memset(res, 0.0, a_rows * b_cols * sizeof(float));
229 // set start
230     start_time
231
232 // calculate
233     int i, j, k;
234     for(i=0; i<a_rows; i++)
235         for(j=0; j<b_cols; j++)
236             for(k=0; k<a_cols; k++)
237                 res[i*a_cols+j] += v1[i*a_cols+k] * v2[k*b_cols+j];
238
239 // set end
240     end_time
241     printf("%d %d %d %.16lf\n", a_rows, a_cols, b_cols, finish.
        tv_sec-start.tv_sec + (double)(finish.tv_nsec - start.
        tv_nsec) / 1000000000.0);
242
243     printf("a: %d %d \n", a_rows, a_cols);
244     printf("b: %d %d \n", a_cols, b_cols);
245     printf("c: %d %d \n", a_rows, b_cols);
246 #ifndef SILENT
247     for (i = 0; i < a_rows; i++){
248         for (j = 0; j < a_cols; j++){
249             printf("%f ", v1[i*a_cols+j]);
250             printf("\n");
251         }
252     }
253 #endif
254
255 /* output results */
256 if (output_switch){
257     if ((pfile = fopen(output_file_name, "wb")) == NULL){
258         printf("file: %s can not be opened \n",
259             output_file_name);
260         system("pause");
261         return 0;
262     }
263     fwrite(res, sizeof(float), a_cols*b_cols, pfile);
264     fclose(pfile);
265 }

```