

# 1. Mandatory Assignment

Aligned SSE: yes

Unaligned SSE: no

Mark Jervelund

(Mjerv15)

Parallel Computing (DM818)



UNIVERSITY OF  
SOUTHERN DENMARK

IMADA

October 20, 2017 at 10

# Contents

Project description . . . . .	1
Introduction . . . . .	2
Work load . . . . .	2
Issues . . . . .	2
Design . . . . .	2
Advice from lecturer . . . . .	2
Own design choices . . . . .	2
Implementation . . . . .	4
Square-dgemm . . . . .	4
Blocking functions . . . . .	4
do-block . . . . .	4
Core-dyn . . . . .	5
Core-4x4 . . . . .	5
Testing . . . . .	5
Comparison . . . . .	5
Conclusion . . . . .	5
appendix . . . . .	6
Goto paper . . . . .	6
Error, Valgrind result and stack overflow post . . . . .	6
Code . . . . .	6

## Project description

To describe the services an operating system provides to users, processes, and other systems

To discuss the various ways of structuring an operating system

To explain how operating systems are installed and customized and how they boot

## Introduction

This report contains documentation for the first mandatory assignment for DM818 Parallel Computing, in the first assignment we were tasked with designing a high performance single threaded matrix multiplication program.

## Work load

For the following project i ended up working alone as i didn't really know anyone in the course as I'm a bachelor student, the people i knew from other courses were already in groups and the remainder that i talked to were phds as were not allowed to form groups.

## Issues

There is currently a bug with the freeing the memory allocation, this means that the program will work one iteration. I've tried debugging with gdb and Valgrind, but i haven't been able to figure out where the bug is yet. I assume it's an overflow, or an out of place memory happening in the B Array.

## Design

### Advice from lecturer

For designing the algorithm we were given 2 examples of how its naively implemented, a naive 3 for loop implementation, and a naive blocked version.

Furthermore we were pointed heavy towards the a paper by Goto on high performance matrix multiplication, which describes how you should design the blocking for the best possible performance.

and as a third helping hand we were given a lecture by Jacob who explained how he implemented it, and what we should do to get the highest performance.

### Own design choices

I made the design choice to follow Jacobs implementation. this meant writing functions that pack the larger matrix into smaller sub-matrices that are then used for the matrix matrix multiplication.

### Block size and Cache

For figuring out the block size, the size of the largest possible matrix that can fit in level 2 cache was calculated.

$$\text{sqr}((256 * 1024\text{byte})/64\text{bit}) = 181.019335984 \quad (1)$$

181 isn't a good size when handling matrix blocking, The number is therefore rounded down to nearest  $2^n$  which is 128.

$$128 * 128 * 64\text{bit} = 128\text{Kbytes} \quad (2)$$

An equation to solve this was also shown doing the lecture. The this is used the 128 by 128 for the main matrix, and 4 by 128 for the remaining.

$$8 * (m_c * k_c + m_c * n_c + N_c * k_c) = 8 * (128 * 128 + 128 * 4 + 4 * 128) = 139264\text{bytes} \quad 136\text{Kbytes} \quad (3)$$

After doing the calculations we can see that the program will need at least 136 KByte of cache to store the data needed. level2 cache in most modern cpu's is 256 KB, so this will suffice. A reason for not pushing it a using maybe a 256 by 128 block will mean that all 256 KB of the level2 cache would be used by the main matrix, this could cause some severed trotting issues due to it being pushed to level3 cache or main memory..

## Intrinsic

I only decided on only implementing intrinsic for a core of 4x4, and just using a naive for loop implementation for the smaller blocks, as programming all the possible different blocks alone would take too much time.

This project was designed using avx2 instruction set, this set contains 6 data types

1. `__m128` contains 4 floats
2. `__m128d` contains 2 doubles
3. `__m128i` contains 8 integers
4. `__m256` contains 8 floats
5. `__m256d` contains 4 doubles
6. `__m256i` contains 16 integers

A generalization for the naming scheme of data types is `__m<size><type>` where size is 128, 256 or 512 when you count the newest avx512 in.

The naming convention on the functions for the instructions follow the same scheme.

`__<size>_<name>_<type>` Its almost the same except a few more data types can be used here.

The data type the project will use is `pd` which is packed doubles.

However a large array of data types are supplied to use with the intrinsic instructions.

1. `ps` packed floats
2. `pd` packed doubles
3. `epi8/16/32/64` vectors containing signed integers of size 8 to 84
4. `epu8/16/32/64` vectors containing unsigned integers of size 8 to 84
5. `si128/si256/si512` unspecified vector of size 128, 256, or 512
6. `m<size><type>` when using input vectors that are different than return type.

There are too many different functions to list so I'll only list the ones i have used in this project

1. `mul` - Multiplication
2. `add` - Adds two vectors
3. `store` - Stores a vector from registers to memory
4. `setzero` - sets all entries in a vector to 0.0
5. `load` - loads data from memory into registers.
6. `set1` - sets all values in a vector to specified value.

## Implementation

### Square-dgemm

The initial function `square_dgemm` gets the 3 matrices and `M` in as the 4 arguments. `M` is directly saved as a global variable.

3 sub matrices are then reserved, `Ablock`, `Bblock`, and `C`, block. `Ablock` and `Cblock` are allocated in a way that is aligned in memory by 32-bytes. this is done so the intrinsic instructions will run faster when storing or loading information.

The first for loop blocks `B` into blocks of size  $KC * M(128 * \text{size of matrix})$ . The second for loop blocks `A` into blocks of 128 by 128 so we keep as much of `A` in cache as possible.

`Prepare_block` is then called with the blocks we want to multiply.

```

1 void square_dgemm(int M, double *A, double *B, double *C) {
2     lda = M;
3     Ablock = (double *) _mm_malloc(MC * KC * sizeof(double), 32); //128*128
4     Bblock = (double *) malloc(lda * KC * sizeof(double)); // M * 128
5     Cblock = (double *) _mm_malloc(MC * NR * sizeof(double), 32); //128*4
6
7     for (unsigned int k = 0; k < lda; k += KC) {
8         packBBlock(B + k, std::min(KC, lda - k));
9         for (unsigned int i = 0; i < lda; i += MC) {
10             packAblock(A + i + k * lda, std::min(MC, lda - i), std::min(KC, lda - k));
11             Prepare_block(C + i, std::min(MC, lda - i), lda, std::min(KC, lda - k));
12         }
13     }
14     _mm_free(Ablock);
15     free(Bblock);
16     _mm_free(Cblock);
17 }

```

### Blocking functions

The program has 2 block packing functions. and 1 unpacking, the 2 packing functions are quite close in usecase and what they do.

`PackA` packs up to 128x128 of matrix `A` into `Ablock`, and `Packb` packs up to 128\*lda into `BlockB`.

The last packing function is for unblocking is for unpacking `Cblock` into `C`

### do-block

The `do_block` functions contains the 2 inner for loops that split `A` into smaller blocks, and the logic that decides if its uses the high performance multiplication function, or the naive dynamic implementation.

```

1 void Do_block(double *C, unsigned int M, unsigned int N, unsigned int K) {
2     double *B = Bblock;
3     for (unsigned int n = 0; n < N; n += NR) {
4         for (unsigned int m = 0; m < M; m += MR) {
5             unsigned int Max_M = std::min(NR, M - m);
6             unsigned int Max_N = std::min(MR, N - n);
7             if (Max_M == MR && Max_N == NR) {
8                 core_4_4(Ablock + m * K, B, Cblock + m, K);
9             } else {
10                 core_dyn(Ablock + m, B, Cblock + m, M, N, K);
11             }
12         }
13         B += NR * K;
14         unpackCBlock(C + n * lda, M, std::min(NR, N - n));
15     }
16 }

```

## Core-dyn

This is simply a standard matrix matrix multiplication function using 3 loops, and no intrinsics.

## Core-4x4

The core 4x4 function is the high performance implementation of matrix matrix multiplication. It performs 32 double precision floating point operations per cycle using fewer operations.

The first optimization that the intrinsic functions does for us is that we load 4 doubles into registers from memory at the time.

after this we load a single double from B into 4 positions in a vector again only using one instruction.

this pattern keeps going with all our operations. we can load blocks into registers a lot faster, and we can compute 4 numbers at the time using this function.

```

1  __m256d c0 = _mm256_setzero_pd();
2  __m256d a1 = _mm256_load_pd(A + k * MR);
3  __m256d b00 = _mm256_set1_pd(B[k]);
4  c0 = _mm256_add_pd(c0, _mm256_mul_pd(b00, a1));
5  _mm256_store_pd(C, c0);

```

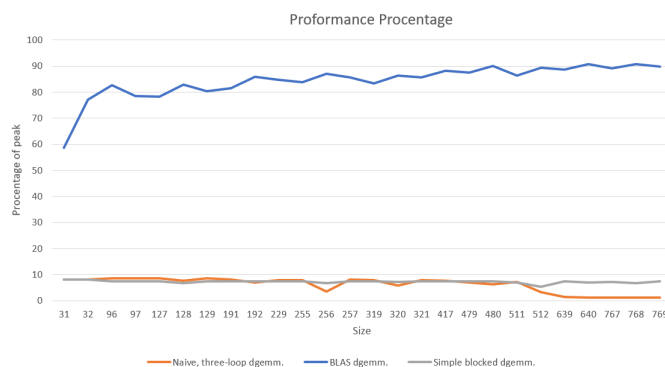
## Testing

For testing the program i was unable to get the benchmark to work due to a memory allocation/free error in the code. i had it run once but there were some problems in other parts of the code, but the performance was around 5000 mflops. but I'm currently unsure if that result could be trusted so I'm only mentioning it here.

The only testing I've done is to check if it returns the correct values, this was done using a testing program i made by removing some parts of the Benchmark program and changing it so it compared to the other implementations we were given as well as printing the output to the terminal for manual verification.

## Comparison

The comparison of the 3 implementations we were given.



## Conclusion

The Matrix matrix multiplication assignment proved to be a bit harder than anticipated. I thought I had it working, but found out with 2 days left that there was an error i was unable to find. I did however get to learn a lot of new things, and i got to get better at using the tools for debugging and finding out where the issue in the code is. furthermore I've learned how use intrinsic and how it can boost a programs run time by a large factor.

## appendix

### Goto paper

<https://dl.acm.org/citation.cfm?id=1356053>

### Error, Valgrind result and stack overflow post

```

1 *** Error in './benchmark-dgemm': munmap_chunk(): invalid pointer: 0x000000000eca9b0 ***
2 benchmark-dgemm: malloc.c:2405: sysmalloc: Assertion '(old_top == initial_top (av) && old_size
   == 0) || ((unsigned long) (old_size) >= MINSIZE && prev_inuse (old_top) && ((unsigned long)
   old_end & (pagesize - 1)) == 0)' failed.
3
4
5 ==13219== Your program just tried to execute an instruction that Valgrind
6 ==13219== did not recognise. There are two possible reasons for this.
7 ==13219== 1. Your program has a bug and erroneously jumped to a non-code
8 ==13219== location. If you are running Memcheck and you just saw a
9 ==13219== warning about a bad jump, it's probably your program's fault.
10 ==13219== 2. The instruction is legitimate but Valgrind doesn't handle it,
11 ==13219== i.e. it's Valgrind's fault. If you think this is the case or
12 ==13219== you are not sure, please let us know and we'll try to fix it.
13 ==13219== Either way, Valgrind will now raise a SIGILL signal which will
14 ==13219== probably kill your program.
15
16 https://stackoverflow.com/questions/2987207/why-do-i-get-a-c-malloc-assertion-failure

```

## Code

```

1 //
2 // Created by jervelund on 9/21/17.
3 //
4
5
6 #include <immintrin.h>
7 #include <stdio.h>
8 #include <algorithm>
9 #include <malloc.h>
10 #include <iostream>
11
12 /*
13  In case you're wondering, dgemm stands for Double-precision, GEneral Matrix-Matrix
14  multiplication.
15 */
16
17 const char *dgemm_desc = "mjerv15 blocked dgemm.";
18
19 // #define min(a,b) (((a)<(b))?(a):(b))
20
21 unsigned int NR = 4;
22 unsigned int MR = 4;
23
24 unsigned int KC = 128;
25 unsigned int MC = 128;
26 unsigned int lda; //size lda*lda of matrix
27
28 double *Ablock;
29 double *Bblock;
30 double *Cblock;

```



```

31 //          From Lecture by Jacob
32 //
33 //          +-----+
34 //          |               |
35 //          |               |
36 //          |               |
37 //          |               |
38 //          |               |
39 //          |      K      |      B      |
40 //          |               |               |
41 //          |               |               |
42 //          |               |               |
43 //          |               |               |
44 //          |               |               |
45 //          |               |               |
46 //          |               |               |
47 //          |               |               |
48 //          |               |               |
49 //          |               |               |
50 //          |               |               |
51 //          |               |               |
52 //          |               |               |
53 //          |               |               |
54 //          |      A      |               |
55 // M|      |               |      C      |
56 //          |               |               |
57 //          |               |               |
58 //          |               |               |
59 //          |               |               |
60 //          |               |               |
61 //          |               |               |
62 //          +-----+ +-----+
63 //debugging
64 void printMatrix2(double *matrix, int MatrixSize) {
65     std::cout << ("Starting Matrix") << '\n';
66     for (int i = 0; i < MatrixSize * MatrixSize; i++) {
67         std::cout << matrix[i] << ' ';
68         if ((i + 1) % MatrixSize == 0) {
69             printf("\n");
70         }
71     }
72     std::cout << ("Ending matrix") << '\n';
73 }
74
75 void printdouble(char *text, double *X, int arraySize) {
76     std::cout << text << '\n';
77     for (int i = 0; i < arraySize; i++) {
78         double a = X[i];
79
80         std::cout << a << ' ';
81     }
82     std::cout << '\n';
83 }
84
85
86 void packAblock(double *A, unsigned int M, unsigned int K) {
87     unsigned int a = 0;
88     for (unsigned int m = 0; m < M; m += MR) {
89         unsigned int MMax = std::min(MR, M - m);
90         for (unsigned int k = 0; k < K; k++) {
91             for (unsigned int i = 0; i < MMax; i++) {
92                 Ablock[a++] = A[m + i + k * lda];
93             }
94         }
95     }
96 }

```

```

94     }
95 }
96 }
97 }
98
99 void packBBlock(double *B, unsigned int K) {
100     unsigned int b = 0;
101     for (unsigned int n = 0; n < lda; n++) {
102         for (unsigned int k = 0; k < K; k++) {
103             Bblock[b++] = B[k + n * lda];
104         }
105     }
106 }
107
108
109 //TODO currently broken, currently fixed
110 void unpackCBlock(double *C, unsigned int M, unsigned int N) {
111     // printMatrix2(Cblock,4 );
112     for (unsigned int n = 0; n < NR; n++) {
113         for (unsigned int i = 0; i < M; i++) {
114             // int lookingat = i + n * MC; ///Debugging value used when debugging with gdb
115             // int storeingat = i + n * lda; ///Debugging value used when debugging with gdb
116             // double x = Cblock[lookingat]; ///Debugging value used when debugging with gdb
117             C[i + n * lda] = Cblock[i + n * MC];
118         }
119     }
120 }
121
122
123 //This is pretty what much Jacob did, except i wrote it with AVX2
124 //He said we should let the compiler take care of how to handle this. eg.
125 // __m256d c01x0-3 could be written as __m256d c23 = _mm_set_pd(0.0,0.0,0.0,0.0); and same goes
126 // for most of this function.
127 void core_4_4(double *A, double *B, double *C, unsigned int K) {
128     // std::cout << "K is = " << K << '\n';
129     __m256d c0 = _mm256_setzero_pd();
130     __m256d c1 = _mm256_setzero_pd();
131     __m256d c2 = _mm256_setzero_pd();
132     __m256d c3 = _mm256_setzero_pd();
133
134     for (unsigned int k = 0; k < K - 1; k += 2) {
135         // printdouble(A+k*MR,4);
136         // printdouble(A+(2+k)*MR,4);
137
138         __m256d a1 = _mm256_load_pd(A + k * MR);
139         // printdouble("A0123k0",A + (k) * MR, 4);
140
141         __m256d a2 = _mm256_load_pd(A + (k + 1) * MR);
142
143         // printdouble("A0123k1",A+ (k+1) * MR, 4);
144
145         __m256d b00 = _mm256_set1_pd(B[k]);
146         __m256d b01 = _mm256_set1_pd(B[k + K]);
147         __m256d b02 = _mm256_set1_pd(B[k + K * 2]);
148         __m256d b03 = _mm256_set1_pd(B[k + K * 3]);
149
150         __m256d b10 = _mm256_set1_pd(B[k + 1]);
151         __m256d b11 = _mm256_set1_pd(B[k + K + 1]);
152         __m256d b12 = _mm256_set1_pd(B[k + K * 2 + 1]);
153         __m256d b13 = _mm256_set1_pd(B[k + K * 3 + 1]);
154
155         //do left
156         c0 = _mm256_add_pd(c0, _mm256_mul_pd(b00, a1));

```

```

156     c1 = _mm256_add_pd(c1, _mm256_mul_pd(b01, a1));
157     c2 = _mm256_add_pd(c2, _mm256_mul_pd(b02, a1));
158     c3 = _mm256_add_pd(c3, _mm256_mul_pd(b03, a1));
159
160     c0 = _mm256_add_pd(c0, _mm256_mul_pd(b10, a2));
161     c1 = _mm256_add_pd(c1, _mm256_mul_pd(b11, a2));
162     c2 = _mm256_add_pd(c2, _mm256_mul_pd(b12, a2));
163     c3 = _mm256_add_pd(c3, _mm256_mul_pd(b13, a2));
164
165 }
166 //edge case.
167 if (K % 2 == 1) {
168     unsigned int k = K - 1;
169     __m256d a0 = _mm256_load_pd(A + k * MR);
170
171     __m256d b00 = _mm256_set1_pd(B[k]);
172     __m256d b01 = _mm256_set1_pd(B[k + K]);
173     __m256d b02 = _mm256_set1_pd(B[k + K * 2]);
174     __m256d b03 = _mm256_set1_pd(B[k + K * 3]);
175
176     c0 = _mm256_add_pd(c0, _mm256_mul_pd(b00, a0));
177     c1 = _mm256_add_pd(c1, _mm256_mul_pd(b01, a0));
178     c2 = _mm256_add_pd(c2, _mm256_mul_pd(b02, a0));
179     c3 = _mm256_add_pd(c3, _mm256_mul_pd(b03, a0));
180 }
181
182 //TODO something is broken here. its been fixed
183 //Debugging value used when debugging with gdb
184 // double *tmp = (double *) _mm_malloc(NR * sizeof(double), 32);
185 //
186 // _mm256_store_pd(tmp, c0123x0);
187 //
188 // _mm256_store_pd(tmp, c0123x1);
189 //
190 // _mm256_store_pd(tmp, c0123x2);
191 //
192 // _mm256_store_pd(tmp, c0123x3);
193
194 // _mm_free(tmp);
195
196 _mm256_store_pd(C, c0);
197 _mm256_store_pd(C + MC, c1);
198 _mm256_store_pd(C + 2 * MC, c2);
199 _mm256_store_pd(C + 3 * MC, c3);
200
201
202 }
203
204 //works.
205 void core_dyn(double *A, double *B, double *C, unsigned int M, unsigned int N, unsigned int K)
206 {
207     for (unsigned int j = 0; j < N; j++) {
208         for (unsigned int i = 0; i < M; i++) {
209             double cij = C[j * lda + i];
210             for (int k = 0; k < K; ++k) {
211                 cij += A[k * lda + i] * B[j * lda + k];
212             }
213             C[j * MC + i] += cij;
214         }
215     }
216 }
217 void Do_block(double *C, unsigned int M, unsigned int N, unsigned int K) {

```

```

218     double *B = Bblock;
219     for (unsigned int n = 0; n < N; n += NR) {
220         for (unsigned int m = 0; m < M; m += MR) {
221             unsigned int Max_M = std::min(NR, M - m);
222             unsigned int Max_N = std::min(MR, N - n);
223             if (Max_M == MR && Max_N == NR) {
224                 core_4_4(Ablock + m * K, B, Cblock + m, K);
225             } else {
226                 core_dyn(Ablock + m, B, Cblock + m, M, N, K);
227             }
228         }
229         B += NR * K;
230         unpackCBlock(C + n * lda, M, std::min(NR, n - n));
231     }
232 }
233
234 void square_dgemm(int M, double *A, double *B, double *C) {
235     lda = M;
236     Ablock = (double *) _mm_malloc(MC * KC * sizeof(double), 32); //128*128
237     Bblock = (double *) malloc(lda * KC * sizeof(double)); // M * 128
238     Cblock = (double *) _mm_malloc(MC * NR * sizeof(double), 32); //128*4
239
240     for (unsigned int k = 0; k < lda; k += KC) {
241         packBBlock(B + k, std::min(KC, lda - k));
242         for (unsigned int i = 0; i < lda; i += MC) {
243             packAblock(A + i + k * lda, std::min(MC, lda - i), std::min(KC, lda - k));
244             Do_block(C + i, std::min(MC, lda - i), lda, std::min(KC, lda - k));
245         }
246     }
247     _mm_free(Ablock);
248     free(Bblock);
249     _mm_free(Cblock);
250 }

```