
MATRIX MULTIPLIKATION MIT FPU

OPTIMIERUNG DER MATRIX MULTIPLIKATION MIT FPU UND
STREAMING-DATA BEFEHLEN (MMX, SSE, SSE2, AVX).

30. Juli 2021

Boris Foko Kouti (s0559792)

Inhaltsverzeichnis

Abbildungsverzeichnis	ii
Tabellenverzeichnis	iii
Code-Verzeichnis	iv
1 Einleitung	1
2 Grundlagen zu FPU und Intrinsics	2
2.1 Matrix Multiplikation	2
2.2 Floating Point Unit	3
2.3 SIMD-Operationen	4
3 Implementierung	6
3.1 Daten Generierung	6
3.2 Matrix Multiplikation Standard Version	7
3.3 Matrix Multiplikation mit FPU über Inline Assembler	7
3.4 Matrix Multiplikation mit MMX Intel Intrinsics	10
3.5 Matrix Multiplikation mit SSE Intel Intrinsics	11
3.6 Matrix Multiplikation mit SSE2 Intel Intrinsics	11
3.7 Matrix Multiplikation mit AVX Intel Intrinsics	12
4 Ergebnisse	13
4.1 Rechenzeiten der verschiedenen Varianten	13
4.2 VTune Messungen der Clocks und weiteres	14
4.3 VTune Messung des Energieverbrauchs	14
5 Zusammenfassung	15
Literatur	vi

Abbildungsverzeichnis

1	FPU Register Befüllung [4]	3
2	Evolution der Intel SIMD-Architekturen [7]	4
3	Darstellung der Rechenzeiten der verschiedenen Implementierungen (Standard, FPU, MMX, SSE, SSE2, AVX)	13

Tabellenverzeichnis

1	SIMD gleichzeitige Operanden nach Datentyp [8]	5
2	Clocks und weitere Messungen mit V-Tune in Debug Mode . .	14
3	Clocks und weitere Messungen mit V-Tune in Release Mode .	14
4	Messung des Energieverbrauchs	14

Code-Verzeichnis

1	Matrix Produkt Pseudo-Code	2
2	Load Add und Mul mit FPU	3
3	Generierung der Matrizen A und B	6
4	Init Funtionene für die Generierung der Matrizen A und B . .	6
5	Standard Version der Matrix Multiplikation	7
6	Matrix Multiplikation mit FPU	7
7	Matrix Multiplikation mit MMX	10
8	Matrix Multiplikation mit SSE	11
9	Matrix Multiplikation mit SSE2	11
10	Matrix Multiplikation mit AVX	12

1 Einleitung

Die Matrix Multiplikation in der linearen Algebra ist eine binäre Operation, die eine Matrix aus zwei Matrizen erzeugt. Bei der Matrixmultiplikation muss die Anzahl der Spalten der ersten Matrix gleich der Anzahl der Zeilen der zweiten Matrix sein. Die resultierende Matrix, das so genannte Matrixprodukt, hat die Anzahl der Zeilen der ersten und die Anzahl der Spalten der zweiten Matrix [1, 2].

Die meisten Algorithmen zur Berechnung dieses Produkts haben eine Komplexität von $O(n^3)$ bei $n * n$ Matrizen. Bei wachsenden n -Werte wird benötigt dieser Algorithmus immer mehr Zeit und Energie. Im Rahmen der Lehrveranstaltung MWP Wahlpflichtmodul/GIT Green IT sind wir beauftragt worden, dieser Algorithmus mit Hilfe von FPU (über die Floating-Point Execution Unit [3]) zu optimieren und gegebenenfalls Operationen über SIMD (Single instruction, multiple data) zu parallelisieren. Es geht also bei der vorliegenden Arbeit nicht an der ersten Linie darum die Komplexität dieses Algorithmus zu reduzieren (auch wenn über Unrolling mit SIMD oder FPU mache Vergleich-Operationen entfallen), sondern darum die einzelnen Operationen zu beschleunigen (FPU) oder zu parallelisieren (SIMD).

In den nächsten Abschnitten wird zunächst auf den Algorithmus für den Matrix-Produkt eingegangen, dann werden in dieser Reihenfolge FPU und SIMD (über Intel Intrinsics) kurz vorgestellt. Im Anschluss dazu werden sowohl die Umsetzung, als auch die erzielten Ergebnisse präsentiert. Zum Schluss wird noch einmal die gesamte Arbeit mit den wichtigsten Erkenntnissen zusammengefasst.

2 Grundlagen zu FPU und Intrinsics

Dieses Kapitel setzt die Grundlagen für weitere Bearbeitung dieser Thematik. Es wird sowohl auf dem Algorithmus für das Matrixprodukt, als auch auf FPU und SIMD mit Intel Intrinsics eingegangen.

2.1 Matrix Multiplikation

Zur mathematischen Definition einer Matrixmultiplikation werden hier zwei Matrizen A ($m * n$) und B ($n * p$) betrachtet:

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \cdot & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} \text{ und } B = \begin{pmatrix} b_{11} & b_{12} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & b_{2p} \\ \dots & \dots & \cdot & \dots \\ b_{n1} & b_{n2} & \dots & b_{np} \end{pmatrix}$$

Das Produkt dieser beiden Matrizen C = AB von der Größe ($m * p$) ist definiert als:

$$C = \begin{pmatrix} c_{11} & c_{12} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2p} \\ \dots & \dots & \cdot & \dots \\ c_{m1} & c_{m2} & \dots & c_{mp} \end{pmatrix}, \text{ sodass } c_{ij} = a_{i1}b_{1j} + a_{i2}b_{2j} + \dots + a_{in}b_{nj} = \sum_{k=1}^n a_{ik}b_{kj},$$

mit $i = 1, \dots, m$ und $j = 1, \dots, p$.

Dieser Vorgang wird algorithmisch in dem Pseudocode 1 wiedergegeben.

```
1 Input: matrices A and B
2 Let C be a new matrix of the appropriate size
3 For i from 1 to m:
4     For j from 1 to p:
5         Let sum = 0
6         For k from 1 to n:
7             Set sum = sum + A[i][k] x B[k][j]
8         Set C[i][j] = sum
9 Return C
```

Code 1: Matrix Produkt Pseudo-Code

Die Standard-Implementierung (zum Beispiel in c++ oder c) der Operation $sum = sum + A[i][k] * B[k][j]$ bei realen Zahlen oder Gleitkommazahlen stellt sich eher als ineffizient und rechenaufwändiger heraus. Dies kann durch Einsatz der Floating-Point Execution Unit (FPU) verbessert werden. Im nächsten Abschnitt wird auf FPU genau eingegangen.

2.2 Floating Point Unit

Die arithmetisch logische Einheit (ALU) eines normalen Prozessors kann nur mit ganzzahlige Werte arbeiten. Dies mag bei manchen Operationen ausreichen, ist aber oft notwendig mit Dezimal- oder Fließkomma-Zahlen zu arbeiten. Dies kann zwar noch mit AssemblerROUTINEN basierend auf die Standard ALU-Operationen umgesetzt werden, allerdings sehr ineffizient. Zur effizienteren Lösung dieses Problem wurde ein hochspezialisierter Co-Prozessor der "Floating Point Unit" entwickelt. FPU ins Deutsch Fließkomma-Recheneinheit ermöglicht über seine acht 80-bits Registers (ST0 bis ST7) die Verarbeitung von Zahlen mit Nachkommastellen. Operanden sowie Ergebnisse werden in den acht verfügbaren Registern geladen und gespeichert. Das Laden eines Wertes in dem Register ST0 erfolgt immer über den Befehl **fld**. Dieser lädt den übergebenen Wert in ST0. Sollte aber davor schon ein Wert in ST0 gewesen sein, so wird dieser zunächst nach ST1 verschoben und erst danach wird der neue Wert in ST0 geladen. Dieser Prozess wird in der Abbildung 1 veranschaulicht. Im Debug Modus lassen sich diese Register auch in Microsoft Visual Studio über *Menü* → *Debuggen* → *Fenster* → *Register* → *Kontext – Menü*(rechteMaustaste) → *FloatingPoint*.

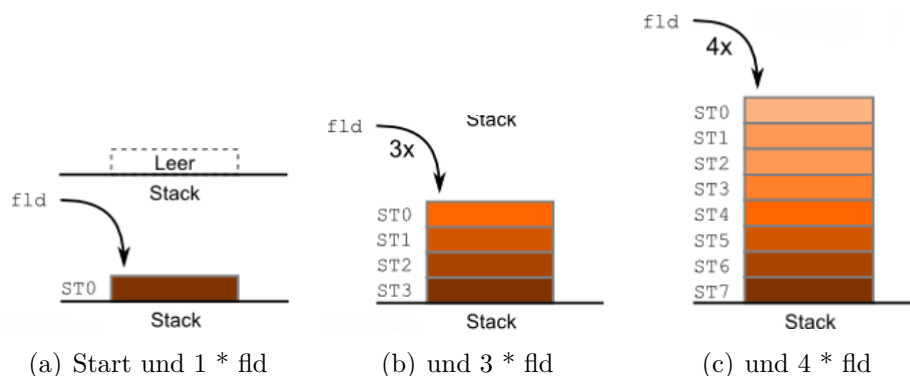


Abbildung 1: FPU Register Befüllung [4]

Sollte für eine Operation einen Wert benötigt sein, der in ST_N (mit $0 < N \leq 7$) liegt, kann der Werte in ST0 mit dem in ST_N über den Befehl **fst Ziel** vertauscht werden (Ziel ist hier ST_N). Die Operation $Setsum = sum + A[i][k]xB[k][j]$ könnte dann wie folgt aussehen:

```

1 fld    dword ptr[esi + ecx * 4] // Load A[i][k]
2 fmul   dword ptr[edi + edx * 4] // Mult B[k][j]
3 mov    ecx, dword ptr[MATRIX_C] // Load C[i][j]
4 fadd   dword ptr[ecx + eax * 4] // C[i][j] + A[i][k]*B[k][j]

```

Code 2: Load Add und Mul mit FPU

2.3 SIMD-Operationen

Single instruction, multiple data (SIMD) ist eine Art der Parallelverarbeitung aus der Flynn Taxonomie [5]. SIMD beschreibt Computer mit mehreren Verarbeitungseinheiten, die dieselbe Operation für mehreren Datensätzen gleichzeitig durchführen. Dies basiert auf dem Prinzip der Parallelität auf Datenebene [6]. Diese Operationen haben sich über die Jahren von MMX in 1997 bis Core und weitere heute. Ein kleiner Auszug dieser Evolution für den Intel Prozessor in in der Abbildung 2 zu sehen.

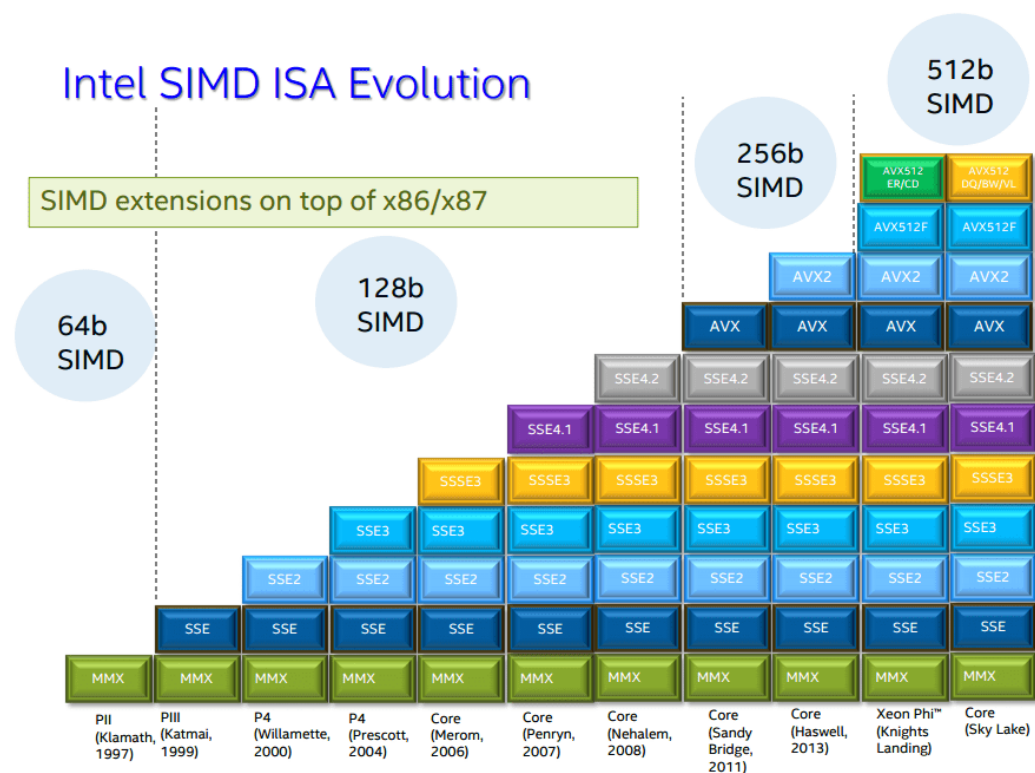


Abbildung 2: Evolution der Intel SIMD-Architekturen [7]

Die erste Variante MMX nutzt die acht 80-bits Register von FPU. Diese werden hier allerdings MM0 bis MM7 genannt. Damit lässt sich eine Operation gleichzeitig auf 8 *char* (8-Bit), 4 *short* (16-Bit), zwei *int* (32-Bit) und ein *int64* (64-Bit) durchführen. Diese Informationen über MMX und weitere sind in der Tabelle 1 zusammengefasst.

Tabelle 1: SIMD gleichzeitige Operanden nach Datentyp [8]

	char	short	int	int64	float	double
MMX	8	4	2	1	0	0
SSE	0	0	0	0	4	0
SSE2	16	8	4	2	4	2
AVX	16	8	4	2	8	4
AVX2	32	16	8	4	8	4
AVX512F	32	16	16	8	16	8
AVX512BW	64	32	16	8	16	8

Im Debug Modus lassen sich die Register für MMX und SSE auch in Microsoft Visual Studio über *Menü → Debuggen → Fenste → Register → Kontext – Menü(rechteMaustaste) → MMX(oderSSE)*.

3 Implementierung

In diesem Kapitel werden die wichtigsten Implementierungsschritte vom Lösungsaufbau, über die Generierung der Matrizen (A und B) bis hin zu den Optimierungen mit FPU, MMX, SSE, SSE2 und AVX besprochen. Zur Umsetzung des hier vorgestellten Programms wurde die Programmiersprache C verwendet. Der Code wurde mit Hilfe der IDE MS Visual Studio 2019 entwickelt, kompiliert, debuggt und getestet.

3.1 Daten Generierung

```
1 // Declaration of matrix for double and float operations
2 double *MATRIX_D_A, *MATRIX_D_B, *MATRIX_D_C;
3 float *MATRIX_F_A, *MATRIX_F_B, *MATRIX_F_C;
4 // Memory allocation with malloc
5 MATRIX_D_A = malloc(N * N * sizeof(double));
6 MATRIX_D_B = malloc(N * N * sizeof(double));
7 MATRIX_D_C = malloc(N * N * sizeof(double));
8 MATRIX_F_A = malloc(N * N * sizeof(float));
9 MATRIX_F_B = malloc(N * N * sizeof(float));
10 MATRIX_F_C = malloc(N * N * sizeof(float));
11 // Init Matrix with random value on false and 0 on true
12 initMatrixD(MATRIX_D_A, false, N);
13 initMatrixD(MATRIX_D_B, false, N);
14 initMatrixD(MATRIX_D_C, true, N);
15 // Copy Double values to Float values
16 copyMatrixDF(MATRIX_D_A, MATRIX_F_A, N);
17 copyMatrixDF(MATRIX_D_B, MATRIX_F_B, N);
18 copyMatrixDF(MATRIX_D_C, MATRIX_F_C, N);
```

Code 3: Generierung der Matrizen A und B

Die im Codestück 3 verwendete Funktion *initMatrixD* wird hier implementiert:

```
1 void initMatrixD(double* MATRIX, bool initToZero, int N) {
2     if (initToZero) {
3         // Init MATRIX Value to 0 (see project)
4     } else {
5         for (int i = 0; i < N; i++) {
6             for (int j = 0; j < N; j++) {
7                 MATRIX[i * N + j] = random_double(MIN, MAX);
8             }
9         }
10    }
11 }
```

Code 4: Init Funktionene für die Generierung der Matrizen A und B

3.2 Matrix Multiplikation Standard Version

Die Standard Version der Matrix Multiplikation ist analog zum Algorithmus 1 implementiert. Zur Messung der Rechenzeit wird ein *clock_t* aus der Klasse *time.h* verwendet.

```
1 int multiply(const double* MATRIX_A, const double* MATRIX_B,
2             double* MATRIX_C, int N) {
3     int msec = 0;
4     clock_t start, finish;
5     start = clock();
6
7     for (int i = 0; i < N; i++) {
8         for (int j = 0; j < N; j++) {
9             // Multiply the row of MATRIX_A by the column of
10            // MATRIX_B to get the row of MATRIC_C.
11            for (int k = 0; k < N; k++) {
12                MATRIX_C[i * N + j] += MATRIX_A[i * N + k] * MATRIX_B
13                [k * N + j];
14            }
15        }
16    }
17    finish = clock();
18    msec = 1000.0 * (finish - start) / CLOCKS_PER_SEC;
19    return msec;
20 }
```

Code 5: Standard Version der Matrix Multiplikation

3.3 Matrix Multiplikation mit FPU über Inline Assembler

Die Matrix Multiplikation mit FPU verläuft analog zu der Standard Version mit dem Unterschied, dass diese zum Einen als Inline Assembler [9, 10] geschrieben ist und zum Anderen bei Multiplikation und Addition Operationen die FPU Register und Befehle einsetzt.

```
1 /*
2 By default start first loop
3 LN1 First for loop increment counter
4 LN2 First for loop compare i with N
5 After LN2 start the second loop
6 LN3 Second for loop increment counter
7 LN4 Second for loop compare i with N
8 After LN4 start the third loop
9 LN5 Third for loop increment counter
10 LN6 Multiplication and Addition
```

```

11 LN7 Close third for loop
12 LN8 Close second for loop
13 LN9 Close first for loop
14 LN10, LN11 Clean and return
15 */
16 int multiply_fpu(float* MATRIX_A, float* MATRIX_B, float*
    MATRIX_C, int N) {
17     int msec = 0;
18     clock_t start, finish;
19     start = clock();
20
21     __asm {
22     //
23         push        ebp
24         sub         esp, 0E4h
25         push        ebx
26         push        esi
27         push        edi
28         lea         edi, [ebp - 0E4h]
29         mov         ecx, 39h
30         mov         eax, 0CCCCCCCCh
31         rep stos    dword ptr es : [edi]
32         mov         dword ptr[ebp - 8], 0
33         jmp LN2
34     LN1 :
35         mov         eax, dword ptr[ebp - 8]
36         add         eax, 1
37         mov         dword ptr[ebp - 8], eax
38     LN2 :
39         mov         eax, dword ptr[ebp - 8]
40         cmp         eax, dword ptr[N]
41         jge LN10
42         mov         dword ptr[ebp - 14h], 0
43         jmp LN4
44     LN3 :
45         mov         eax, dword ptr[ebp - 14h]
46         add         eax, 1
47         mov         dword ptr[ebp - 14h], eax
48     LN4 :
49         mov         eax, dword ptr[ebp - 14h]
50         cmp         eax, dword ptr[N]
51         jge LN9
52         mov         dword ptr[ebp - 20h], 0
53         jmp LN6
54     LN5 :
55         mov         eax, dword ptr[ebp - 20h]
56         add         eax, 1
57         mov         dword ptr[ebp - 20h], eax
58     LN6 :

```

```

59     mov     eax, dword ptr[ebp - 20h]
60     cmp     eax, dword ptr[N]
61     jge     LN8
62     mov     eax, dword ptr[ebp - 8]
63     imul    eax, dword ptr[N]
64     add     eax, dword ptr[ebp - 14h]
65     mov     ecx, dword ptr[ebp - 8]
66     imul    ecx, dword ptr[N]
67     add     ecx, dword ptr[ebp - 20h]
68     mov     edx, dword ptr[ebp - 20h]
69     imul    edx, dword ptr[N]
70     add     edx, dword ptr[ebp - 14h]
71     mov     esi, dword ptr[MATRIX_A]
72     mov     edi, dword ptr[MATRIX_B]
73     fld     dword ptr[esi + ecx * 4]
74     fmul    dword ptr[edi + edx * 4]
75     mov     ecx, dword ptr[MATRIX_C]
76     fadd    dword ptr[ecx + eax * 4]
77     mov     edx, dword ptr[ebp - 8]
78     imul    edx, dword ptr[N]
79     add     edx, dword ptr[ebp - 14h]
80     mov     eax, dword ptr[MATRIX_C]
81     fstp    dword ptr[eax + edx * 4]
82     LN7 :
83     jmp     LN5
84     LN8 :
85     jmp     LN3
86     LN9 :
87     jmp     LN1
88     LN10 :
89     pop     edi
90     LN11 :
91     pop     esi
92     pop     ebx
93     add     esp, 0E4h
94     cmp     ebp, esp
95     mov     esp, ebp
96     pop     ebp
97     ret
98 }
99
100 finish = clock();
101 msec = 1000.0 * (finish - start) / CLOCKS_PER_SEC;
102
103 return msec;
104 }

```

Code 6: Matrix Multiplikation mit FPU

3.4 Matrix Multiplikation mit MMX Intel Intrinsics

Die MMX kann keine Fließkommazahlen verarbeiten, daher werden wir short(int16) als Input und int32 als Ergebnis verwenden. MMX Verarbeitet 4 short gleichzeitig. MMX ist aus diesem Grund nicht für die Optimierung relevant. Es wird hier trotzdem beispielhaft aufgeführt.

```
1 int multiply_mmx(const short* MATRIX_A, const short* MATRIX_B,
2                 int* MATRIX_C, int N)
3 {
4     int msec = 0;
5     clock_t start, finish;
6     start = clock();
7     __m64 a_line, b_line;
8
9     for (int i = 0; i < N; i++) {
10         for (int j = 0; j < N; j++) {
11             __m64 sum = _mm_setzero_si64(); // init sum to zero
12             int sum_down = 0;
13             // add 4 short at the same time using the MMX
14             // equivalent to asm pmaddwd
15             _m_paddwd function
16             for (int k = 0; k < N; k+=4) {
17                 a_line = _mm_set_pi16(MATRIX_A[i * N + k],
18                                     MATRIX_A[i * N + k + 1], MATRIX_A[i * N + k + 2],
19                                     MATRIX_A[i * N + k + 3]);
20                 b_line = _mm_set_pi16(MATRIX_B[k * N + j],
21                                     MATRIX_B[(k + 1) * N + j], MATRIX_B[(k + 2) * N + j],
22                                     MATRIX_B[(k + 3) * N + j]);
23                 sum = _m_paddw(sum, _m_pmaddwd(a_line, b_line));
24             }
25             sum_down = _mm_cvtsi64_si32(sum); // save low 32 bits
26             sum = _m_psrlqi(sum, 32); // shift right on 32 bits
27             sum_down += _mm_cvtsi64_si32(sum); // save low 32 bits
28             MATRIX_C[i * N + j] = sum_down;
29         }
30     }
31     // Clear the MMX registers and MMX state
32     _m_empty();
33     _mm_empty();
34     finish = clock();
35     msec = 1000.0 * (finish - start) / CLOCKS_PER_SEC;
36
37     return msec;
38 }
```

Code 7: Matrix Multiplikation mit MMX

3.5 Matrix Multiplikation mit SSE Intel Intrinsics

Über SSE Register lassen sich 4 Floats gleichzeitig verarbeiten. SSE ist dadurch mindestens 4 mal schneller als der Standard.

```
1
2 int multiply_sse(const float* MATRIX_A, const float* MATRIX_B,
3               float* MATRIX_C, int N)
4 {
5     int msec = 0;
6     clock_t start, finish;
7     start = clock();
8     for (int i = 0; i < N; i++) {
9         for (int j = 0; j < N; j += 4) {
10             __m128 sum = _mm_load_ps(MATRIX_C + i * N + j);
11             for (int k = 0; k < N; k += 4) {
12                 sum = _mm_add_ps(_mm_mul_ps(
13                     _mm_set1_ps(MATRIX_A[i * N + k]),
14                     _mm_load_ps(MATRIX_B + k * N + j)), sum);
15             }
16             _mm_store_ps(MATRIX_C + i * N + j, sum);
17         }
18     }
19     finish = clock();
20     msec = 1000.0 * (finish - start) / CLOCKS_PER_SEC;
21     return msec;
22 }
```

Code 8: Matrix Multiplikation mit SSE

3.6 Matrix Multiplikation mit SSE2 Intel Intrinsics

Da SSE2 dieselbe Anzahl von Float (nämlich 4) wie SSE gleichzeitig verarbeiten kann, ist hier keine Verbesserung zu erwarten. Statt dessen wurde hier Testweise die Operation auf Double (2 Double gleichzeitig) durchgeführt.

```
1 int multiply_sse2(const double* MATRIX_A, const double*
2               MATRIX_B, double* MATRIX_C, int N)
3 {
4     int msec = 0;
5     clock_t start, finish;
6     start = clock();
7     for (int i = 0; i < N; i++) {
8         for (int j = 0; j < N; j += 2) {
9             __m128d sum = _mm_load_pd(MATRIX_C + i * N + j);
10            for (int k = 0; k < N; k++) {
11                sum = _mm_add_pd(_mm_mul_pd(
12                    _mm_set_pd1(MATRIX_A[i * N + k]),
13                    _mm_load_pd(MATRIX_B + k * N + j)), sum);
14            }
15        }
16    }
17    finish = clock();
18    msec = 1000.0 * (finish - start) / CLOCKS_PER_SEC;
19    return msec;
20 }
```



```

13     }
14     _mm_store_pd(MATRIX_C + i * N + j, sum);
15 }
16 }
17 finish = clock();
18 msec = 1000.0 * (finish - start) / CLOCKS_PER_SEC;
19 return msec;
20 }

```

Code 9: Matrix Multiplikation mit SSE2

3.7 Matrix Multiplikation mit AVX Intel Intrinsics

AVX ermöglicht die Verarbeitung von 8 Double und 16 Float. AVX ist also mindestens 16 mal schneller als die Standard Multiplikation mit Float.

```

1 int multiply_avx(const float* MATRIX_A, const float* MATRIX_B,
2                 float* MATRIX_C, int N)
3 {
4     int msec = 0;
5     clock_t start, finish;
6     start = clock();
7     const int block_width = N >= 256 ? 512 : 256;
8     const int block_height = N >= 512 ? 8 : N >= 256 ? 16 : 32;
9     for (int row_offset = 0; row_offset < N; row_offset +=
10         block_height) {
11         for (int column_offset = 0; column_offset < N;
12             column_offset += block_width) {
13             for (int i = 0; i < N; ++i) {
14                 for (int j = column_offset; j < column_offset +
15                     block_width && j < N; j += 8) {
16                     __m256 sum = _mm256_load_ps(MATRIX_C + i * N + j);
17                     for (int k = row_offset; k < row_offset +
18                         block_height && k < N; ++k) {
19                         sum = _mm256_fmadd_ps(_mm256_set1_ps(MATRIX_A[i *
20                             N + k]), _mm256_load_ps(MATRIX_B + k * N + j), sum);
21                     }
22                     _mm256_store_ps(MATRIX_C + i * N + j, sum);
23                 }
24             }
25         }
26     }
27     finish = clock();
28     msec = 1000.0 * (finish - start) / CLOCKS_PER_SEC;
29     return msec;
30 }

```

Code 10: Matrix Multiplikation mit AVX

4 Ergebnisse

Dieser Kapitel befasst sich mit den Ergebnissen, die im Laufe der mit der Anwendung durchgeführten Tests und Messungen gesammelt wurden. Diese Ergebnisse bestehen aus den gemessenen Rechenzeiten der verschiedenen Implementierungen, die mit VTune [11] gemessenen Clocks und den entsprechenden Energieverbrauch. Die in diesem Kapitel präsentierten Messungen sind auf einem Intel(R) Core(TM) i7-10510U CPU @ 1.80GHz, 2.30 GHz mit 8 logische Kern durchgeführt worden.

4.1 Rechenzeiten der verschiedenen Varianten

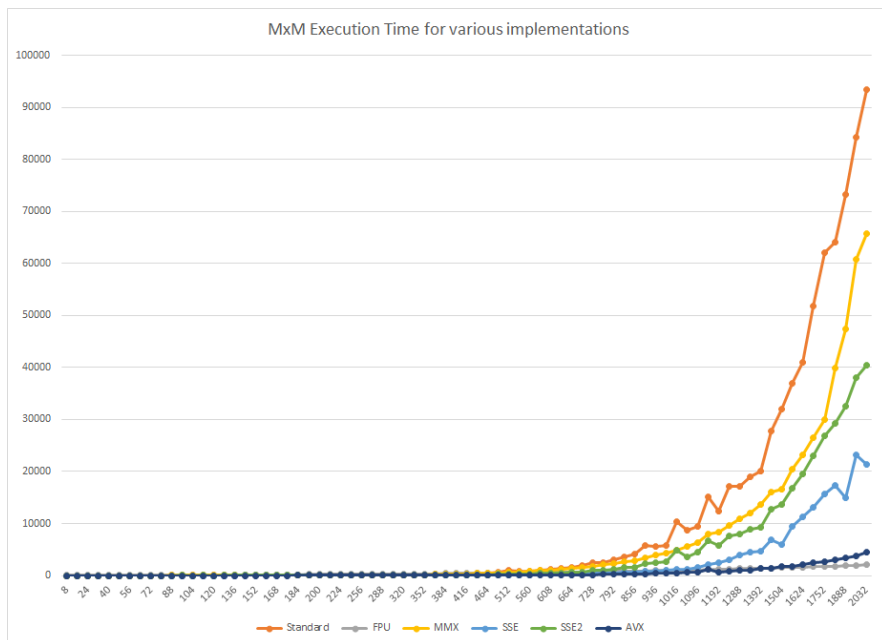


Abbildung 3: Darstellung der Rechenzeiten der verschiedenen Implementierungen (Standard, FPU, MMX, SSE, SSE2, AVX)

Die Kurven in Abb 3 sind aus den Daten aus dem Projekt MxM_Streaming_Data abgeleitet. Die “statistic.xlsx” Datei dazu wird mit dieser Dokumentation mit geliefert. Zu den Ergebnissen selbst muss unterzeichnet werden, dass MMX und SSE2 hier nicht von großer Relevanz sind, da bei dem Ersten keine Floating-Point Operationen möglich sind (wurde ausschließlich für Testzwecke implementiert), und bei dem Zweiten hier SSE2 Double verwendet wurden (da SSE2 genauso wie SSE 4 Float verarbeiten kann, besteht da kein Mehrwert). Auf der Grafik sind zwei klare Gewinner zu erkennen,

nämlich FPU und AVX, während die Standard Version exponentiell steigt, sind FPU und AVX eher linear mit einer sehr kleinen Steigung.

4.2 VTune Messungen der Clocks und weiteres

Hier wird ein Auszug aus den Messungen vorgestellt. Die vollständige Messungen sind als Kommentar in dem Code abgegeben worden.

Tabelle 2: Clocks und weitere Messungen mit V-Tune in Debug Mode

Function	CPU Time	Clockticks	Instructions Retired	CPI Rate
multiply	14.600s	48,727,800,000	29,050,200,000	1.677
multiply_fpu	6.874s	27,883,800,000	29,030,400,000	0.961
multiply_mmxx	5.000s	18,208,800,000	18,304,200,000	0.995
multiply_sse2	5.267s	17,299,800,000	12,385,800,000	1.397
multiply_sse	1.249s	4,352,400,000	6,179,400,000	0.704
multiply_avx	0.468s	1,742,400,000	4, 163,400,000	0.419

Tabelle 3: Clocks und weitere Messungen mit V-Tune in Release Mode

Function	CPU Time	Clockticks	Instructions Retired	CPI Rate
multiply	3.166s	13,559,400,000	8,616,600,000	1.574
multiply_mmxx	3.101s	12,585,600,000	6,741,000,000	1.867
multiply_sse2	1.590s	5,414,400,000	4,842,000,000	1.118
multiply_sse	0.495s	1,733,400,000	2,421,000,000	0.716
multiply_avx	0.151s	552,600,000	2,188,800,000	0.252
multiply_fpu	keine Messung da es Inline ASM			

4.3 VTune Messung des Energieverbrauchs

Eine Messung des Energieverbrauchs ist ausschließlich für die Standard-Umsetzung und die optimierte FPU Variante durchgeführt worden.

Tabelle 4: Messung des Energieverbrauchs

Executable	Debug Energy Consumption (mJ)	RELEASE Energy Consumption (mJ)
MxM_Standard	140,469.360	59,315.186
Executable	130,589.722	2,218.933

5 Zusammenfassung

Ziel dieser Arbeit ist es gewesen, mit FPU und SIMD-Operatoren die Matrix Multiplikation zu optimieren. Es wurde im Kapitel 2 zunächst den Versuch unternommen, die Matrix Multiplikation selbst sowohl mathematisch als auch algorithmisch genau zu definieren. Dann wurden abwechselnd FPU und SIMD mit Intel Intrinsics genau erörtert. So konnte der Grundstein für die Implementierung der verschiedenen Varianten (FPU und SIMD mit MMX, SSE, SSE2 und AVX) gelegt werden. Die Umsetzung erfolgte im C Code mit der IDE MS Visual Studio 2019 (siehe Kapitel 3). Mit den aufgenommenen Messungen konnte gezeigt werden, wie effizienter und schneller die FPU Variante war im Vergleich zu der Standard-Version, sowohl was die Laufzeit angeht, die Anzahl der Clocks, als auch den Energieverbrauch. Es konnte durch die Kurven in Abb. 3 eine gewisse Tendenz, was das Verhältnis der Varianten zu der Größe der Input Matrizen angeht. Die Standard Version verhält sich exponentiell, während die FPU und AVX Varianten eher linear mit einer sehr kleinen Steigung verlaufen.

Literatur

- [1] *Math Insight: Multiplying matrices and vectors*. URL: https://mathinsight.org/matrix_vector_multiplication. (visited on 25.07.2021).
- [2] *Algebra Symbols: A comprehensive collection of 225+ symbols used in algebra, categorized by subject and type into tables along with each symbol's name, usage and example*. URL: <https://mathvault.ca/hub/higher-math/math-symbols/algebra-symbols/>. (visited on 25.07.2021).
- [3] S. F. Anderson u. a. "The IBM System/360 Model 91: Floating-Point Execution Unit". In: *IBM Journal of Research and Development* 11.1 (1967), S. 34–53. DOI: 10.1147/rd.111.0034.
- [4] Thorsten Thormählen. *Technische Informatik I FPU, MMX, SSE, x86-64*. URL: https://www.mathematik.uni-marburg.de/~thormae/lectures/ti1/ti_10_3_ger_web.html#1. (visited on 25.07.2021).
- [5] Michael J. Flynn. "Some Computer Organizations and Their Effectiveness". In: *IEEE Transactions on Computers* C-21.9 (1972), S. 948–960. DOI: 10.1109/TC.1972.5009071.
- [6] W. Daniel Hillis und Guy L. Steele. "Data Parallel Algorithms". In: *Commun. ACM* 29.12 (Dez. 1986), S. 1170–1183. ISSN: 0001-0782. DOI: 10.1145/7902.7903. URL: <https://doi.org/10.1145/7902.7903>.
- [7] Mostafa Soliman. "State-of-the-Art in Processor Architecture". In: (Jan. 2016), S. 13–14. DOI: 10.13140/RG.2.1.3105.8000. URL: https://www.researchgate.net/profile/Mostafa_Soliman3/publication/288993171_State-of-the-%20Art_in_Processor_Architecture/links/56881ca308ae19758398eb08/State-of-the-Art-in-Processor-Architecture.pdf.
- [8] *Intel® Intrinsics Guide: The Intel® Intrinsics Guide is a reference tool for Intel intrinsics*. URL: <https://software.intel.com/sites/landingpage/IntrinsicsGuide/#>. (visited on 25.07.2021).
- [9] *MASM for x64 (ml64.exe): Visual Studio includes both 32-bit and 64-bit hosted versions of Microsoft Assembler (MASM) to target x64 code*. URL: <https://docs.microsoft.com/en-us/cpp/assembler/masm/masm-for-x64-ml64-exe?view=msvc-160>. (visited on 25.07.2021).
- [10] *asm declaration: asm-declaration gives the ability to embed assembly language source code within a C++ program*. URL: <https://en.cppreference.com/w/cpp/language/asm>. (visited on 25.07.2021).

- [11] *Intel® VTune™ Profiler: Quickly Find and Fix Performance Bottlenecks and Realize All the Value of Your Hardware*. URL: <https://software.intel.com/content/www/us/en/develop/tools/oneapi/components/vtune-profiler.html#gs.7n99l7>. (visited on 25.07.2021).