

Parallel Programming Tutorial - SIMD

Vincent Bode, M.Sc.

Chair for Computer Architecture and Parallel Systems (CAPS)

Technical University Munich

May 26, 2020



TUM Uhrenturm

Public Service Announcements

- The current assignment is extended to Friday night. Here is some supporting material:
 - The exercise support wiki on [GitLab](#).
 - Help from the Q&A session (link to recording [here](#) or on moodle).
- New RocketChat channel for discussing high-performance solutions: [#parprog-need-for-speed](#).
- There was some feedback asking us to include some more examples in the exercises.
- Additional (unofficial) resource: parprog comprehensive question bank [here](#) (source: [GitHub](#)).

Introduction to SIMD

Execution Models

What mental model do you have for the execution of this code on a single core?

```
1 // a and b are float arrays
2 for( int i = 0; i < 128; i++ ) {
3     c[i] = a[i] * b[i];
4 }
```

Sequential Execution

Iteration	0
Operand 1	
Operation	
Operand 2	
Result	

Sequential Execution

Iteration	0
Operand 1	
Operation	*
Operand 2	
Result	

Sequential Execution

Iteration	0
Operand 1	0
Operation	*
Operand 2	
Result	

Sequential Execution

Iteration	0
Operand 1	0
Operation	*
Operand 2	0
Result	

Sequential Execution

Iteration	0
Operand 1	0
Operation	*
Operand 2	0
Result	0

The diagram illustrates the sequential execution of a multiplication operation. It consists of a table with five rows: Iteration, Operand 1, Operation, Operand 2, and Result. The values in the table are 0, 0, *, 0, and 0 respectively. A vertical line separates the labels from the values. A downward arrow points from the '0' in the 'Operand 2' row to the '0' in the 'Result' row, indicating the flow of data from the operand to the result.

Sequential Execution

Iteration	0	1
Operand 1	<div>0</div>	
Operation	*	
Operand 2	<div>0</div>	
Result	<div>0</div>	

The diagram illustrates the sequential execution of a multiplication operation. It consists of a table with five rows: Iteration, Operand 1, Operation, Operand 2, and Result. The first column (Iteration) has values 0 and 1. The second column (Operand 1) has a box containing 0 for iteration 0. The third column (Operation) has a box containing * for iteration 0. The fourth column (Operand 2) has a box containing 0 for iteration 0. The fifth column (Result) has a box containing 0 for iteration 0. A vertical arrow points from the Operand 2 box of iteration 0 to the Result box of iteration 0.

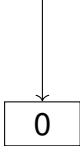
Sequential Execution

Iteration	0	1
Operand 1	<div>0</div>	
Operation	*	*
Operand 2	<div>0</div>	
Result	<div>0</div>	

The diagram illustrates the sequential execution of a multiplication operation. It consists of a table with two columns representing iterations 0 and 1. In iteration 0, Operand 1 is 0, the Operation is multiplication (*), and Operand 2 is 0. The Result of this operation is 0. In iteration 1, the Operation is again multiplication (*), but Operand 1 and Operand 2 are not specified, and the Result is also not shown. A vertical arrow points from the Operand 2 box in iteration 0 to the Result box in iteration 0, indicating the flow of data from the operand to the result.

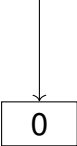
Sequential Execution

Iteration	0	1
Operand 1	<div>0</div>	<div>1</div>
Operation	*	*
Operand 2	<div>0</div>	
Result	<div>0</div>	



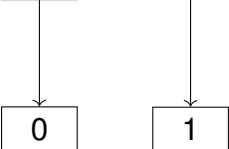
Sequential Execution

Iteration	0	1
Operand 1	0	1
Operation	*	*
Operand 2	0	1
Result	0	



Sequential Execution

Iteration	0	1
Operand 1	<div>0</div>	<div>1</div>
Operation	*	*
Operand 2	<div>0</div>	<div>1</div>
Result	<div>0</div>	<div>1</div>



```
graph TD; subgraph Iteration_0 [Iteration 0]; O1_0[0]; O2_0[0]; R_0[0]; end; subgraph Iteration_1 [Iteration 1]; O1_1[1]; O2_1[1]; R_1[1]; end; O2_0 --> R_0; O2_1 --> R_1;
```

Sequential Execution

Iteration	0	1	2
Operand 1	0	1	
Operation	*	*	
Operand 2	0	1	
Result	0	1	

Sequential Execution

Iteration	0	1	2
Operand 1	0	1	
Operation	*	*	*
Operand 2	0	1	
Result	0	1	

Sequential Execution

Iteration	0	1	2
Operand 1	0	1	2
Operation	*	*	*
Operand 2	0	1	
Result	0	1	

Sequential Execution

Iteration	0	1	2
Operand 1	0	1	2
Operation	*	*	*
Operand 2	0	1	2
Result	0	1	

Sequential Execution

Iteration	0	1	2
Operand 1	0	1	2
Operation	*	*	*
Operand 2	0	1	2
Result	0	1	4

Sequential Execution

Iteration	0	1	2	3
Operand 1	0	1	2	
Operation	*	*	*	
Operand 2	0	1	2	
Result	0	1	4	

Sequential Execution

Iteration	0	1	2	3
Operand 1	0	1	2	
Operation	*	*	*	*
Operand 2	0	1	2	
Result	0	1	4	

Sequential Execution

Iteration	0	1	2	3
Operand 1	0	1	2	3
Operation	*	*	*	*
Operand 2	0	1	2	
Result	0	1	4	

Sequential Execution

Iteration	0	1	2	3
Operand 1	0	1	2	3
Operation	*	*	*	*
Operand 2	0	1	2	3
Result	0	1	4	

The diagram illustrates the sequential execution of a multiplication operation. It consists of a table with five rows: Iteration, Operand 1, Operation, Operand 2, and Result. The columns represent iterations 0, 1, 2, and 3. In each iteration, Operand 1 and Operand 2 are multiplied. The result of the multiplication is shown in the Result row. Arrows point from the Operand 2 boxes to the Result boxes for iterations 0, 1, and 2, indicating the flow of data.

Iteration	0	1	2	3
Operand 1	0	1	2	3
Operation	*	*	*	*
Operand 2	0	1	2	3
Result	0	1	4	

Sequential Execution

Iteration	0	1	2	3
Operand 1	0	1	2	3
Operation	*	*	*	*
Operand 2	0	1	2	3
Result	0	1	4	9

Sequential Execution

Iteration	0	1	2	3
Operand 1	0	1	2	3
Operation	*	*	*	*
Operand 2	0	1	2	3
Result	0	1	4	9

Single Core Vectorized

Iteration	0-4
Operand 1	
Operation	
Operand 2	
Result	

Single Core Vectorized

Iteration	0-4
Operand 1	
Operation	*
Operand 2	
Result	

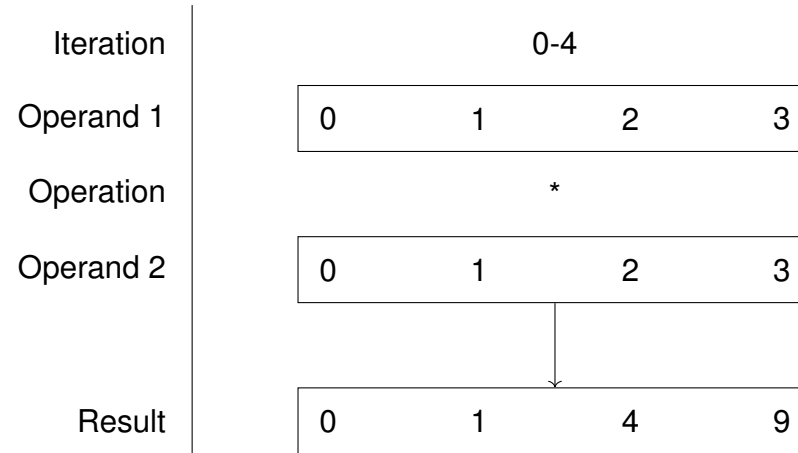
Single Core Vectorized

Iteration	0-4			
Operand 1	0	1	2	3
Operation	*			
Operand 2				
Result				

Single Core Vectorized

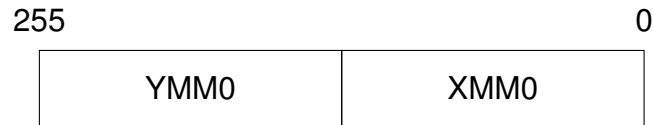
Iteration	0-4			
Operand 1	0	1	2	3
Operation	*			
Operand 2	0	1	2	3
Result				

Single Core Vectorized



SIMD

- This is SIMD: Single Instruction Multiple Data.
- To support this, we need vector registers and Instruction Set support.
- Intel has had multiple extensions to its SIMD instructions.
- For example, SSE operates with 128 bit registers. AVX with 256 bit registers:



Intrinsics

Compilers can automatically vectorize your code in some cases. However, we will focus on Intel Intrinsics, a set of C style functions that can help you vectorize your code without writing assembly.

Intrinsics define datatypes and operations on these datatypes. For AVX, in GCC, these are defined in the header `<immintrin.h>`.

Intrinsics

AVX Datatypes:

- `__m256` can hold eight 32 bit floating point numbers (float)
- `__m256d` can hold four 64 bit floating point numbers (double)
- `__m256i` can hold eight, 32 bit integer (int) OR four, 64 bit integers

AVX function examples:

- Functions for loading / storing data e.g. `_mm256_loadu_ps`
- Arithmetic functions e.g. `_mm_add_ps` add packed floating point numbers
- Byte manipulation functions like `_mm_movelh_ps`.

Intrinsics

Some terminology in the Intrinsics Guide:

- Packed (`_ps`) operations operate on the entire vector.
- Scalar operations (`_sd`) operate on the least significant data element (bits 0-31 for floats).
- Latency: Number of clock cycles from beginning of execution to end of execution.
 - Unit: cycles (lower is better)
- Throughput: Rate of execution of instructions.
 - Unit: Cycles per Instruction (CPI, lower is better) or Instructions per Cycle (IPC, higher is better)

These functions and datatypes are much closer to assembly than to normal C. You necessarily have to think about how to use these vector instructions efficiently to make full use of the hardware, which you also need to consider. It is recommended that you consult the [Intel Intrinsics Guide](#).

Simple Example

Simple Example

```

1  #include <stdlib.h>
2  #include <stdint.h>
3  #include <stdio.h>
4
5  int main() {
6      int size = 20000;
7      int iterations = 100000;
8
9      uint32_t* a = (uint32_t*)malloc(size * sizeof(uint32_t));
10     for (int i = 0; i < size; ++i) a[i] = i;
11
12     for (int iter = 0; iter < iterations; ++iter) {
13         for (int i = 0; i < size; ++i) {
14             a[i] = ((a[i] * a[i]) >> 1) ^ a[i];
15         }
16     }
17
18     uint32_t sum = 0;
19     for (int i = 0; i < size; ++i) sum += a[i];
20     printf("%x\n", sum);
21 }

```

Simple Example

Zoom in on the interesting code:

```
1          a[i] = ((a[i] * a[i]) >> 1) ^ a[i];
```

Simple Example

C is much too hard to read, so look at the assembly

```
1      a[i] = ((a[i] * a[i]) >> 1) ^ a[i];
```

```
1      mov edx, DWORD PTR [rcx]
2      mov esi, edx
3      imul esi, edx
4      shr esi
5      xor edx, esi
6      mov DWORD PTR [rcx-4], edx
```

You can explore the assembly yourself at https://godbolt.org/z/nh-KJ_

Simple Example

Rewrite the code to clearly show the individual operations.

```

1      uint32_t* a_ptr = a + i;
2      uint32_t a_i    = *a_ptr;
3      uint32_t mul     = a_i * a_i;
4      uint32_t srl     = mul >> 1;
5      uint32_t xor     = srl ^ a_i;
6      *a_ptr           = xor;

```

```

1      mov edx, DWORD PTR [rcx]
2      mov esi, edx
3      imul esi, edx
4      shr esi
5      xor edx, esi
6      mov DWORD PTR [rcx-4], edx

```

Simple Example

Do the actual vectorisation.

```

1      __m128i* a_ptr = a + i
2      __m128i a_i    = *a_ptr
3      __m128i mul     = a_i * a_i
4      __m128i srl     = mul >> 1
5      __m128i xor     = srl ^ a_i
6      *a_ptr = xor
    
```

```

1      mov edx, DWORD PTR [rcx]
2      mov esi, edx
3      imul esi, edx
4      shr esi
5      xor edx, esi
6      mov DWORD PTR [rcx-4], edx
    
```


Simple Example

Do the actual vectorisation.

```

1      __m128i* a_ptr = a + i ???
2      __m128i a_i    = *a_ptr
3      __m128i mul     = a_i * a_i
4      __m128i srl     = mul >> 1
5      __m128i xor     = srl ^ a_i
6      *a_ptr = xor
    
```

```

1      mov edx, DWORD PTR [rcx]
2      mov esi, edx
3      imul esi, edx
4      shr esi
5      xor edx, esi
6      mov DWORD PTR [rcx-4], edx
    
```

Simple Example

Do the actual vectorisation.

```

1      __m128i* a_ptr = (__m128i*)(a + i);
2      __m128i a_i    = *a_ptr ???
3      __m128i mul     = a_i * a_i
4      __m128i srl     = mul >> 1
5      __m128i xor     = srl ^ a_i
6      *a_ptr = xor
    
```

```

1      mov edx, DWORD PTR [rcx]
2      mov esi, edx
3      imul esi, edx
4      shr esi
5      xor edx, esi
6      mov DWORD PTR [rcx-4], edx
    
```

Simple Example

Do the actual vectorisation.

```

1      __m128i* a_ptr = (__m128i*)(a + i);
2      __m128i a_i    = _mm_load_si128(a_ptr);
3      __m128i mul     = a_i * a_i ???
4      __m128i srl     = mul >> 1
5      __m128i xor     = srl ^ a_i
6      *a_ptr = xor
    
```

```

1      mov edx, DWORD PTR [rcx]
2      mov esi, edx
3      imul esi, edx
4      shr esi
5      xor edx, esi
6      mov DWORD PTR [rcx-4], edx
    
```

Simple Example

Do the actual vectorisation.

```

1      __m128i* a_ptr = (__m128i*)(a + i);
2      __m128i a_i    = _mm_load_si128(a_ptr);
3      __m128i mul     = _mm_mullo_epi32(a_i, a_i);
4      __m128i srl     = mul >> 1 ???
5      __m128i xor     = srl ^ a_i
6      *a_ptr = xor
    
```

```

1      mov edx, DWORD PTR [rcx]
2      mov esi, edx
3      imul esi, edx
4      shr esi
5      xor edx, esi
6      mov DWORD PTR [rcx-4], edx
    
```

Simple Example

Do the actual vectorisation.

```

1      __m128i* a_ptr = (__m128i*)(a + i);
2      __m128i a_i    = _mm_load_si128(a_ptr);
3      __m128i mul     = _mm_mullo_epi32(a_i, a_i);
4      __m128i srl     = _mm_srli_epi32(mul, 1);
5      __m128i xor     = srl ^ a_i ???
6      *a_ptr = xor
    
```

```

1      mov edx, DWORD PTR [rcx]
2      mov esi, edx
3      imul esi, edx
4      shr esi
5      xor edx, esi
6      mov DWORD PTR [rcx-4], edx
    
```

Simple Example

Do the actual vectorisation.

```

1      __m128i* a_ptr = (__m128i*)(a + i);
2      __m128i a_i    = _mm_load_si128(a_ptr);
3      __m128i mul     = _mm_mullo_epi32(a_i, a_i);
4      __m128i srl     = _mm_srli_epi32(mul, 1);
5      __m128i xor     = _mm_xor_si128(srl, a_i);
6      *a_ptr = xor ???

```

```

1      mov edx, DWORD PTR [rcx]
2      mov esi, edx
3      imul esi, edx
4      shr esi
5      xor edx, esi
6      mov DWORD PTR [rcx-4], edx

```

Simple Example

Do the actual vectorisation.

```

1      __m128i* a_ptr = (__m128i*)(a + i);
2      __m128i a_i    = _mm_load_si128(a_ptr);
3      __m128i mul     = _mm_mullo_epi32(a_i, a_i);
4      __m128i srl     = _mm_srli_epi32(mul, 1);
5      __m128i xor     = _mm_xor_si128(srl, a_i);
6      _mm_store_si128(a_ptr, xor);

```

```

1      mov edx, DWORD PTR [rcx]
2      mov esi, edx
3      imul esi, edx
4      shr esi
5      xor edx, esi
6      mov DWORD PTR [rcx-4], edx

```

Simple Example

Do the actual vectorisation.

```

1      __m128i* a_ptr = (__m128i*)(a + i);
2      __m128i a_i    = _mm_load_si128(a_ptr);
3      __m128i mul     = _mm_mullo_epi32(a_i, a_i);
4      __m128i srl     = _mm_srli_epi32(mul, 1);
5      __m128i xor     = _mm_xor_si128(srl, a_i);
6      _mm_store_si128(a_ptr, xor);

```

```

1
2      vmovdqa xmm0, XMMWORD PTR [rdx]
3      vpmulld xmm1, xmm0, xmm0
4      vpsrld xmm1, xmm1, 1
5      vpxor xmm0, xmm0, xmm1
6      vmovaps XMMWORD PTR [rdx-16], xmm0

```



```

1  #include <stdlib.h>
2  #include <stdint.h>
3  #include <stdio.h>
4
5
6  int main() {
7      int size = 20000;
8      int iterations = 100000;
9
10     uint32_t* a = (uint32_t*)malloc(size * sizeof(uint32_t))
11     for (int i = 0; i < size; ++i) a[i] = i;
12
13     for (int iter = 0; iter < iterations; ++iter) {
14         for (int i = 0; i < size; ++i) {
15             __m128i* a_ptr = (__m128i*)(a + i);
16             __m128i a_i = _mm_load_si128(a_ptr);
17             __m128i mul = _mm_mullo_epi32(a_i, a_i);
18             __m128i srl = _mm_srli_epi32(mul, 1);
19             __m128i xor = _mm_xor_si128(srl, a_i);
20             _mm_store_si128(a_ptr, xor);
21         }
22     }
23
24     uint32_t sum = 0;
25     for (int i = 0; i < size; ++i) sum += a[i];
26     printf("%x\n", sum);
27 }

```

Does it work?

```

1  #include <stdlib.h>
2  #include <stdint.h>
3  #include <stdio.h>
4
5
6  int main() {
7      int size = 20000;
8      int iterations = 100000;
9
10     uint32_t* a = (uint32_t*)malloc(size * sizeof(uint32_t))
11     for (int i = 0; i < size; ++i) a[i] = i;
12
13     for (int iter = 0; iter < iterations; ++iter) {
14         for (int i = 0; i < size; i += 4) {
15             __m128i* a_ptr = (__m128i*)(a + i);
16             __m128i a_i = _mm_load_si128(a_ptr);
17             __m128i mul = _mm_mullo_epi32(a_i, a_i);
18             __m128i srl = _mm_srli_epi32(mul, 1);
19             __m128i xor = _mm_xor_si128(srl, a_i);
20             _mm_store_si128(a_ptr, xor);
21         }
22     }
23
24     uint32_t sum = 0;
25     for (int i = 0; i < size; ++i) sum += a[i];
26     printf("%x\n", sum);
27 }

```

Does it work?

```

1  #include <stdlib.h>
2  #include <stdint.h>
3  #include <stdio.h>
4  #include <immintrin.h>
5
6  int main() {
7      int size = 20000;
8      int iterations = 100000;
9
10     uint32_t* a = (uint32_t*)malloc(size * sizeof(uint32_t))
11     for (int i = 0; i < size; ++i) a[i] = i;
12
13     for (int iter = 0; iter < iterations; ++iter) {
14         for (int i = 0; i < size; i += 4) {
15             __m128i* a_ptr = (__m128i*)(a + i);
16             __m128i a_i = _mm_load_si128(a_ptr);
17             __m128i mul = _mm_mullo_epi32(a_i, a_i);
18             __m128i srl = _mm_srli_epi32(mul, 1);
19             __m128i xor = _mm_xor_si128(srl, a_i);
20             _mm_store_si128(a_ptr, xor);
21         }
22     }
23
24     uint32_t sum = 0;
25     for (int i = 0; i < size; ++i) sum += a[i];
26     printf("%x\n", sum);
27 }

```

Does it work?

```

1  #include <stdlib.h>
2  #include <stdint.h>
3  #include <stdio.h>
4  #include <immintrin.h>
5
6  int main() {
7      int size = 20000;
8      int iterations = 100000;
9
10     uint32_t* a = (uint32_t*)aligned_alloc(16, size * sizeof(uint32_t))
11     for (int i = 0; i < size; ++i) a[i] = i;
12
13     for (int iter = 0; iter < iterations; ++iter) {
14         for (int i = 0; i < size; i += 4) {
15             __m128i* a_ptr = (__m128i*)(a + i);
16             __m128i a_i = _mm_load_si128(a_ptr);
17             __m128i mul = _mm_mullo_epi32(a_i, a_i);
18             __m128i srl = _mm_srli_epi32(mul, 1);
19             __m128i xor = _mm_xor_si128(srl, a_i);
20             _mm_store_si128(a_ptr, xor);
21         }
22     }
23
24     uint32_t sum = 0;
25     for (int i = 0; i < size; ++i) sum += a[i];
26     printf("%x\n", sum);
27 }

```

Does it work?

```

1  #include <stdlib.h>
2  #include <stdint.h>
3  #include <stdio.h>
4  #include <immintrin.h>
5
6  int main() {
7      int size = 20000;
8      int iterations = 100000;
9
10     uint32_t* a = (uint32_t*)aligned_alloc(16, size * sizeof(uint32_t))
11     for (int i = 0; i < size; ++i) a[i] = i;
12
13     for (int iter = 0; iter < iterations; ++iter) {
14         for (int i = 0; i < size; i += 4) {
15             __m128i* a_ptr = (__m128i*)(a + i);
16             __m128i a_i = _mm_load_si128(a_ptr);
17             __m128i mul = _mm_mullo_epi32(a_i, a_i);
18             __m128i srl = _mm_srli_epi32(mul, 1);
19             __m128i xor = _mm_xor_si128(srl, a_i);
20             _mm_store_si128(a_ptr, xor);
21         }
22     }
23
24     uint32_t sum = 0;
25     for (int i = 0; i < size; ++i) sum += a[i];
26     printf("%x\n", sum);
27 }

```

Yes!

Speedup?

Speedup?

~ 3.6

Assignment 5 – SIMD

Assignment 5 – SIMD

- Transposed matrix multiplication $C \Leftarrow \alpha AB^T + \beta C$
- That means we always calculate the dot product of the i -th row of A and the j -th row of B
 - Why not normal multiplication?
- Your task is to *manually* vectorise the code using Intel intrinsics
- Required speedup is 2

The server only supports AVX!
(No AVX2, no AVX512.)

The server only supports AVX!
(No AVX2, no AVX512.)

The server only supports AVX!
(No AVX2, no AVX512.)

Assignment 5, Hints

- Only use instructions up until AVX (this includes all of the SSE extensions)
 - If you try to use anything that is not supported, the compiler should generate an extremely unhelpful error message
 - In case you manage to convince the compiler to compile your code anyway, the executable will crash on the server (but probably run fine locally)
 - The slides of Micheal Klemm's talk include some AVX512 instructions, take care
- Use the [Intel Intrinsics Guide](#) to find out which instructions to use
- Note that the input is not a even multiple of the vector size, so you have to process the remainder
- Inspect the assembly to find out what the compiler is actually doing, either the old-fashioned way (`gcc -S ... > out.s` and inspect the file) or via [Compiler Explorer](#)
- The server's CPU is an [Intel Xeon E5-2680 v0 2.70 GHz](#) (Sandy Bridge)
- Also remember the Q&A sessions, Mo. 14:00.

Content Questions

This is a placeholder.

<insert your question here>

What we covered today

- Using intrinsics for SIMD
- Assignment 05