

Parallel Programming Tutorial – SIMD

Bengisu Elis, M.Sc.

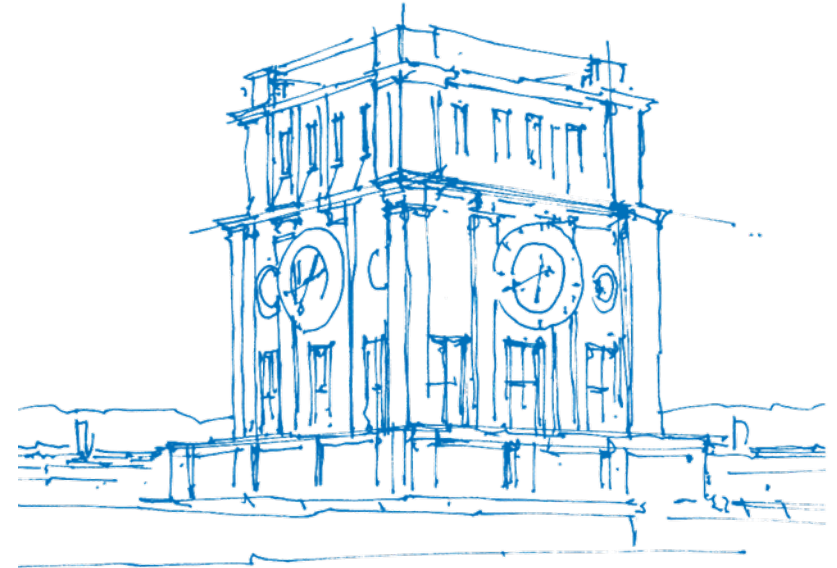
Philipp Czerner

Hasan Ashraf

Chair for Computer Architecture and Parallel Systems (Prof. Schulz)

Technical University of Munich

19 June 2019



TUM Uhrenturm

Organizational

- Speed up requirement for assignment is reduced to 120.
- `unit_test.c` added for assignment 6.
- Deadline for assignment 6 is on **26th June** (08:15 in the morning – before the tutorial!)
- Assignment 7 will be published today – deadline on **3rd July**
- Data dependency analysis and Loop transformation useful reading :
 - R. Allen and K Kennedy Optimizing compilers for modern architectures: A dependence based approach Morgan Kaufmann 2001.

Introduction to SIMD

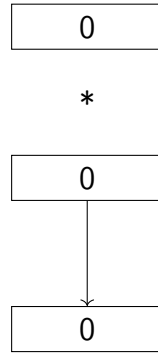
Execution Models

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

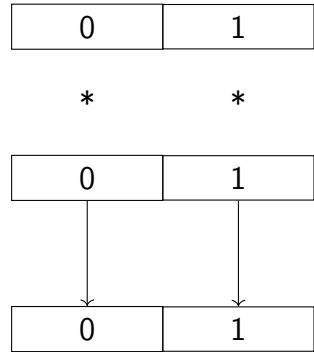
```

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



Sequential Execution

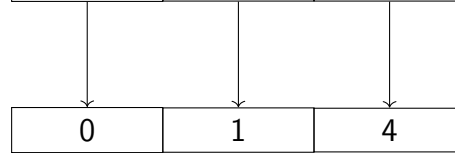


Sequential Execution

0	1	2
---	---	---

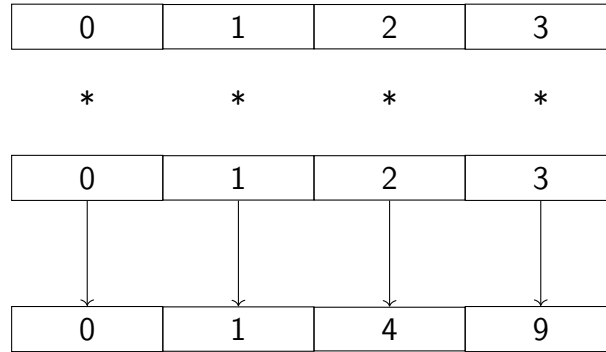
* *

0	1	2
---	---	---

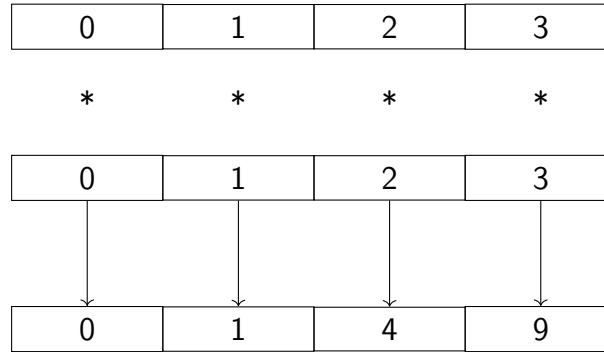


0	1	4
---	---	---

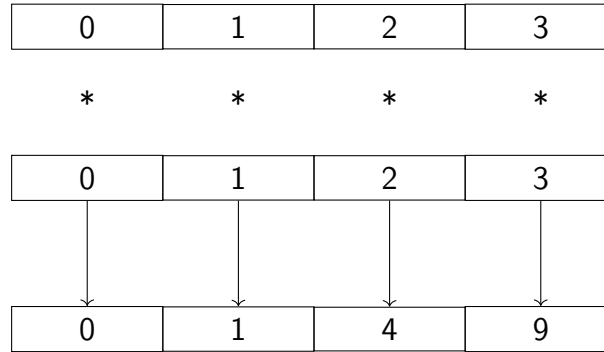
Sequential Execution



Sequential Execution

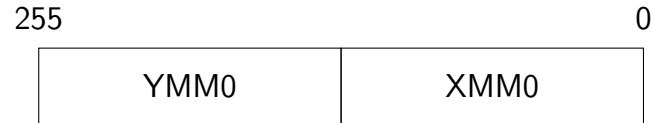


Single Core Vectorized



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 8, 32 bit integer (int) OR 4, 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 to perform an instruction.
- Throughput: Number of clock cycles you need to wait to start independent instructions.

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, c_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, c_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, c_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, c_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?

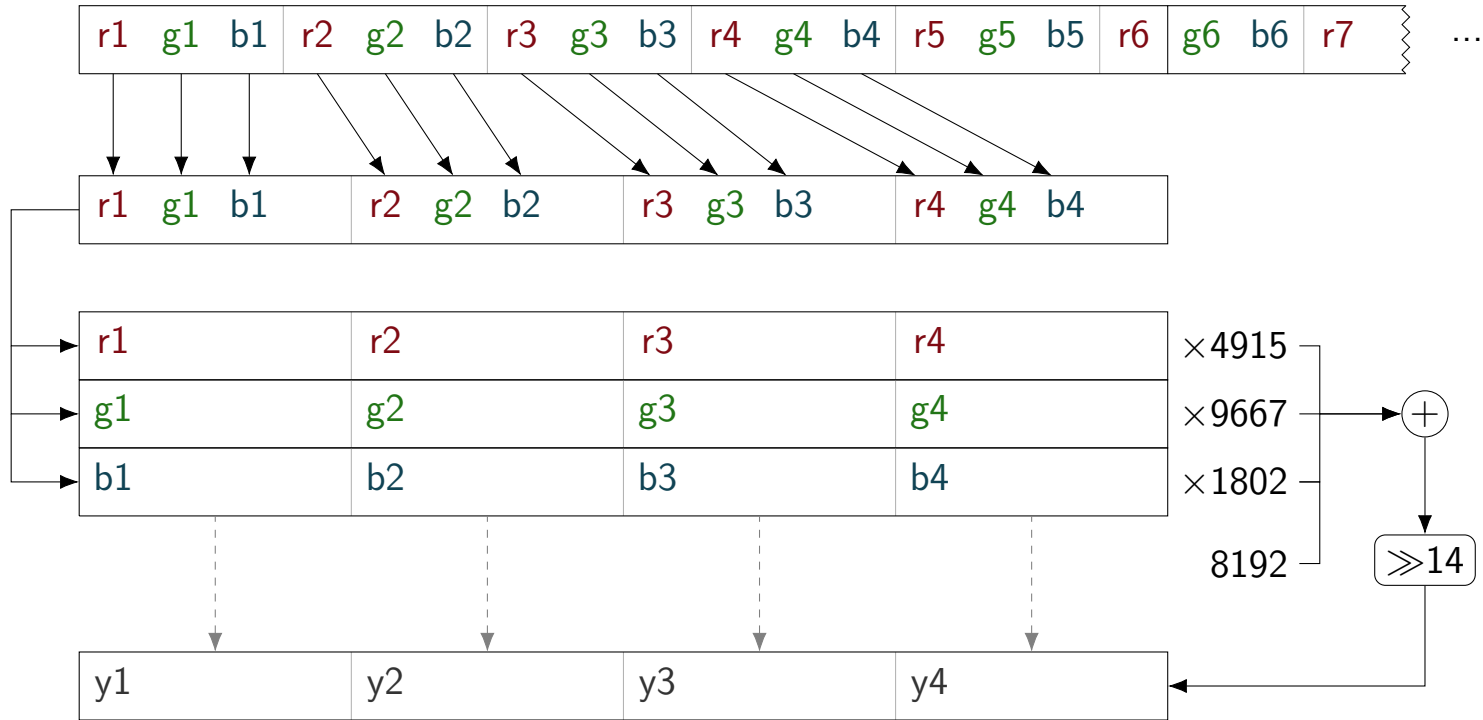
Speedup?

~ 3.6

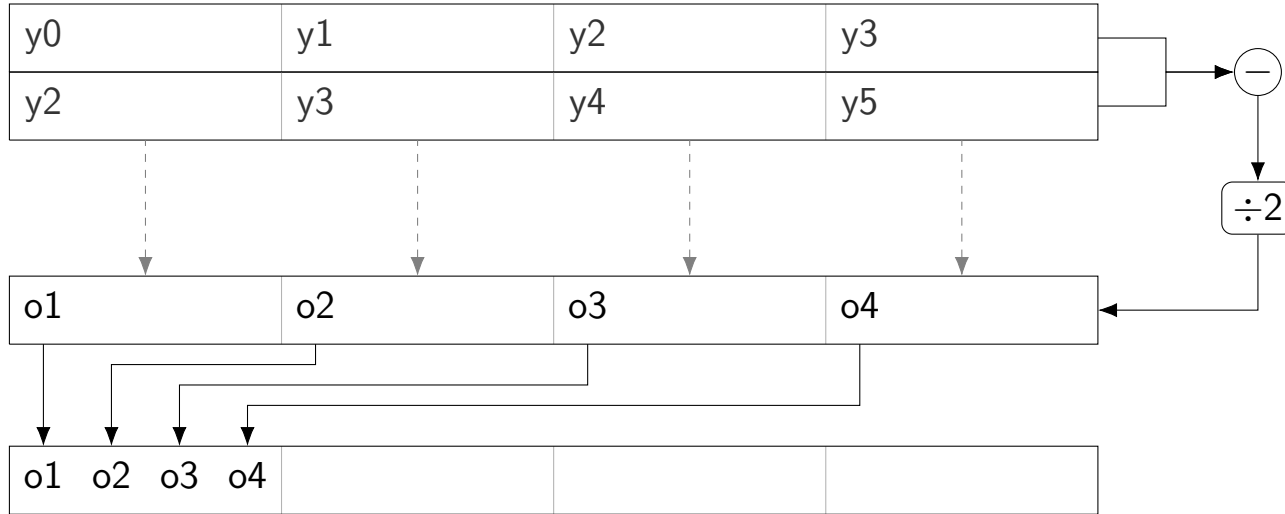
Complicated Example

```
1  for (int y = 0; y < height; ++y) {
2      for (int x = 1; x+1 < width; ++x) {
3          uint8_t* p = source + 3 * (width*y + x);
4          int8_t* q = target +      width*y + x ;
5          int u_l = (p[-3]*4915 + p[-2]*9667 + p[-1]*1802 + 8192) >> 14;
6          int u_r = (p[ 3]*4915 + p[ 4]*9667 + p[ 5]*1802 + 8192) >> 14;
7          *q = (int8_t)((u_r - u_l) / 2);
8      }
9  }
```

This code converts an image from RGB into greyscale, then computes the difference between adjacent pixels.



Combine values of multiple iterations:



Then do it three times more, for o_5 to o_{16} . In total, 3×16 bytes read and 16 bytes written.

Assignment 7 – SIMD

Assignment 7 – SIMD

- Transposed matrix multiplication $C := AB^T$
- That means entry c_{ij} is 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
- No parallelism this time!
- Required speedup is 2.3
- Deadline on 3rd July

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 7, 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, Fr 08-10, 01.06.020