



Vectorization and AVX2 Intrinsics Programming



UNIVERSITÀ DEGLI STUDI
DI SALERNO

Biagio Cosenza
Department of Computer Science
University of Salerno, Italy

Outline

Introduction to Vector Instructions

Intel AVX2

SIMD Intrinsics Programming

Data types

Load and Store

Arithmetic

Slides and lecture notes by Biagio Cosenza under Creative Commons CC BY-SA 4.0 . All images report the source, author and the original license. You are free to copy and redistribute the material in any medium or format for any purpose, even commercially, but you must give appropriate credit. You are free to adapt, remix, transform, and build upon the material for any purpose, even commercially, but you must distribute your contributions under the same license as the original.



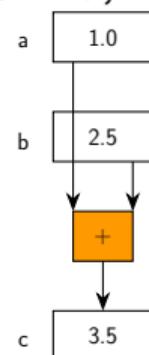
SIMD Instructions

Single Instruction, Multiple Data

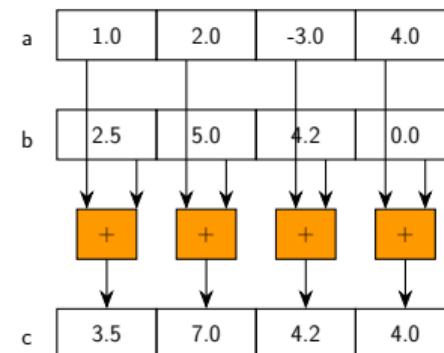
Vector instruction

- ▶ instruction that work on multiple data elements (pack) simultaneously
- ▶ typically: short, length fixed by the hardware, not necessarily fixed by the ISA

Example: $c = a + b;$



Scalar instruction



SIMD instruction

SIMD Extensions

ISA	SIMD Extension	Year
PA-RISC	Multimedia Acceleration eXtensions (MAX)	1994
SPARC	Visual Instruction Set (VIS)	1995
Alpha	Motion Video Instructions (MVI)	1996
MIPS	Multimedia Acceleration Instructions (MDMX)	1996
x86	MultiMedia eXtensions (MMX)	1996
PowerPC	AltiVec	1998
X86	Streaming SIMD Extensions (SSE, SSE2, SSE3, SSE4)	1999-2007
ARM	Advanced SIMD Extension (NEON)	2005
X86	Advanced Vector eXtensions (AVX, AVX2, AVX-512)	2011-2013-2017
ARM	Scalable Vector Extension (SVE, SVE2)	2020-2021
RISC-V	RVV: RISC-V "V" Vector Extension (0.7, 1.0)	2021-2023



Intel AVX, AVX2 and AVX-512

AVX (Advanced Vector Extensions)

- ▶ introduced with Intel's Sandy Bridge (2011)
- ▶ AVX extends SIMD operations to 256-bit registers
- ▶ improving floating-point performance for workloads such as multimedia and scientific computing

AVX2

- ▶ introduced with Haswell (2013)
- ▶ AVX2 enhances AVX by adding full 256-bit integer support
- ▶ new instructions such as gather for improved memory access patterns

AVX-512

- ▶ since Intel's Skylake-X (2017)
- ▶ AVX-512 doubles the register width to 512 bits
- ▶ introduces new mask registers for finer control over operations, and specialized subsets targeting HPC, cryptography, ...



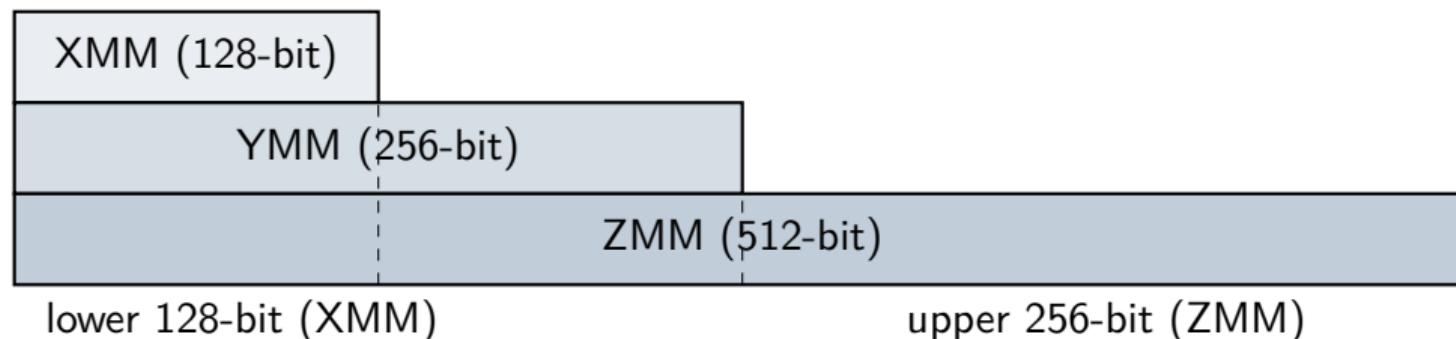
AVX Registers

Vector registers of different size, from SSE to AVX-512

SSE XMM0 — XMM15, 128-bit registers

AVX & AVX2 YMM0 — YMM15, 256-bit registers

AVX-512 ZMM0 — ZMM31, 512-bit registers



AVX/AVX2 Registers and Data Type

Vector registers can be mapped on different data types

YMM register can contain up to 256 bits of data, e.g.

- ▶ four doubles (4×64 bits = 256)
- ▶ eight floats (8×32 bits = 256)
- ▶ eight ints (8×32 bits = 256)

255	YMM0								0
	float	float	float	float	float	float	float	float	
	double		double		double		double		
	int32	int32	int32	int32	int32	int32	int32	int32	
	16	16	16	16	16	16	16	16	
	8	8	8	8	8	8	8	8	

YMM registers are versatile and can be used for both floating-point and integer operations depending on the instruction operating on them.



Programming Models for SIMD

From high to low level

Automatic vectorization

- ▶ SIMD instruction automatically generated by the compiler

SIMD-optimized libraries

- ▶ e.g., FFTW, oneMKL, oneDNN

Data parallel programming models with SIMD support

- ▶ e.g., OpenCL, SYCL, ispc, TBB, OpenMP SIMD pragma
- ▶ implicit use of vector

SIMD intrinsics

- ▶ explicit use of vector, manually write code with SIMD instruction
- ▶ leave register allocation and instruction scheduling to the compiler

Assembly

- ▶ instruction selection, scheduling and register allocation by the developer



SIMD Intrinsics Programming

Low-level C functions for SIMD programming

- ▶ an intrinsic function is mapped on one (sometime more than one) vectorial instruction

Example:

- ▶ the intrinsic function `_mm256_add_ps()` maps directly to `vaddps`

Combining the performance of assembly with the convenience of a high-level function

- ▶ register allocation and instruction scheduling performed by the compiler
- ▶ instruction selection enforced by intrinsics



Intel Intrinsics Guide

Overview: <https://software.intel.com/sites/landingpage/IntrinsicsGuide>

The Intel Intrinsics Guide is an interactive reference tool for Intel intrinsic instructions, which are C style functions that provide access to many Intel instructions - including Intel® SSE, AVX, AVX-512, and more - without the need to write assembly code. ×

MMX
 SSE
 SSE2
 SSE3
 SSSE3
 SSE4.1
 SSE4.2
 AVX
 AVX2
 FMA
 AVX-512
 KNC
 AMX
 SVML
 Other

Application-Targeted
 Arithmetic
 Bit Manipulation
 Cast
 Compare
 Convert
 Cryptography
 Elementary Math

Functions
 General Support
 Load
 Logical
 Mask

mm_search ?

__m256i_mm256_abs_ep116 (__m256i a)	vpabsw
__m256i_mm256_abs_ep132 (__m256i a)	vpabsd
__m256i_mm256_abs_ep18 (__m256i a)	vpabsb
__m256i_mm256_add_ep116 (__m256i a, __m256i b)	vpadde
__m256i_mm256_add_ep132 (__m256i a, __m256i b)	vpaddd
__m256i_mm256_add_ep164 (__m256i a, __m256i b)	vpaddq
__m256i_mm256_add_ep18 (__m256i a, __m256i b)	vpaddb
__m256d_mm256_add_pd (__m256d a, __m256d b)	vaddpd
__m256_mm256_add_ps (__m256 a, __m256 b)	vaddps
__m256i_mm256_adds_ep116 (__m256i a, __m256i b)	vpaddsw
__m256i_mm256_adds_ep18 (__m256i a, __m256i b)	vpaddsb
__m256i_mm256_adds_epu16 (__m256i a, __m256i b)	vpadusw
__m256i_mm256_adds_epu8 (__m256i a, __m256i b)	vpadusb
__m256d_mm256_addsub_pd (__m256d a, __m256d b)	vaddsubpd
__m256_mm256_addsub_ps (__m256 a, __m256 b)	vaddsubps
__m256i_mm256_alignr_ep18 (__m256i a, __m256i b, const int imm8)	vpalignr
__m256d_mm256_and_pd (__m256d a, __m256d b)	vandpd
__m256_mm256_and_ps (__m256 a, __m256 b)	vandps
__m256i_mm256_and_s1256 (__m256i a, __m256i b)	vpand
__m256d_mm256_andnot_pd (__m256d a, __m256d b)	vandnpd
__m256_mm256_andnot_ps (__m256 a, __m256 b)	vandnps
__m256i_mm256_andnot_s1256 (__m256i a, __m256i b)	vpandn
__m256i_mm256_avg_epu16 (__m256i a, __m256i b)	vpavgw
__m256i_mm256_avg_epu8 (__m256i a, __m256i b)	vpavgb



Intel Intrinsics Guide

vaddps

```
__m256 _mm256_add_ps (__m256 a, __m256 b)
```

vaddps

Synopsis

```
__m256 _mm256_add_ps (__m256 a, __m256 b)
#include <immintrin.h>
Instruction: vaddps ymm, ymm, ymm
CPUID Flags: AVX
```

Description

Add packed single-precision (32-bit) floating-point elements in **a** and **b**, and store the results in **dst**.

Operation

```
FOR j := 0 to 7
    i := j*32
    dst[i+31:i] := a[i+31:i] + b[i+31:i]
ENDFOR
dst[MAX:256] := 0
```

Performance

Architecture	Latency	Throughput (CPI)
Icelake	4	0.5
Skylake	4	0.5
Broadwell	3	1
Haswell	3	1
Ivy Bridge	3	1



Intrinsic to Instruction Mapping

Typically one-to-one mapping, but there are exceptions: the intrinsic `_mm256_fmadd_ps` maps to three instructions

- ▶ `vfmadd132ps`
- ▶ `vfmadd213ps`
- ▶ `vfmadd231ps`

```
_m256 _mm256_fmadd_ps (_m256 a, _m256 b, _m256 c)           vfmadd132ps,...
```

Synopsis

```
_m256 _mm256_fmadd_ps (_m256 a, _m256 b, _m256 c)
#include <immintrin.h>
Instruction: vfmadd132ps ymm, ymm, ymm
              vfmadd213ps ymm, ymm, ymm
              vfmadd231ps ymm, ymm, ymm
CPUID Flags: FMA
```

Description

Multiply packed single-precision (32-bit) floating-point elements in `a` and `b`, add the intermediate result to packed elements in `c`, and store the results in `dst`.

Operation

```
FOR j := 0 to 7
    i := j*32
    dst[i+31:i] := (a[i+31:i] * b[i+31:i]) + c[i+31:i]
ENDFOR
dst[MAX:256] := 0
```

Performance

Architecture	Latency	Throughput (CPI)
Icelake	4	0.5
Skylake	4	0.5
Knights Landing	6	0.5
Broadwell	5	0.5
Haswell	5	0.5



AVX/AVX2 Data Types

Each type starts with two underscores (`__`), an `m`, and the width of the vector in bits

- ▶ AVX/AVX2 registers are 256 bits, AVX512 registers are 512 bits

A vector type ending in `d` contains doubles; if it doesn't have a suffix, it contains floats

- ▶ `_m128i` and `_m256i` vectors contain integers
- ▶ `_m256i` may contain 32 `chars`, 16 `shorts`, 8 32-bit `ints`, or 4 64-bit `longs`, either signed or unsigned

<code>__m128</code>	128-bit vector containing 4 floats
<code>__m128d</code>	128-bit vector containing 2 doubles
<code>__m128i</code>	128-bit vector containing integers
<code>__m256</code>	256-bit vector containing 8 floats
<code>__m256d</code>	256-bit vector containing 4 doubles
<code>__m256i</code>	256-bit vector containing integers



AVX/AVX2 Function Naming Conventions

Function name: `_mm<bit_width>_<name>_<data_type>`

- ▶ `<bit_width>` the size of the vector returned by the function; for 128-bit vectors, this is empty; for 256-bit vectors, this is set to 256
- ▶ `<name>` describes the operation performed by the intrinsic
- ▶ `<data_type>` is the data type of the function's primary arguments

`<data_type>` identifies the content of the input values:

- ▶ `ps` packed single-precision (floats)
- ▶ `pd` packed double-precision (doubles)
- ▶ `epi8/epi16/epi32/epi64` contain 8-bit/16-bit/32-bit/64-bit signed integers
- ▶ `epu8/epu16/epu32/epu64` contain 8-bit/16-bit/32-bit/64-bit unsigned integers
- ▶ `si128/si256` unspecified 128-bit vector or 256-bit vector
- ▶ `m128/m128i/m128d/m256/m256i/m256d` identifies input vector types when they're different than the type of the returned vector



AVX Example: Vector Difference

```
1 #include <immintrin.h>
2 #include <stdio.h>
3
4 int main() {
5     /* Initialize the two argument vectors */
6     __m256 evens = _mm256_set_ps(2.0, 4.0, 6.0, 8.0, 10.0, 12.0, 14.0, 16.0);
7     __m256 odds  = _mm256_set_ps(1.0, 3.0, 5.0, 7.0, 9.0, 11.0, 13.0, 15.0);
8
9     /* Compute the difference between the two vectors */
10    __m256 result = _mm256_sub_ps(evens, odds);
11
12    /* Display the elements of the result vector */
13    float* f = (float*)&result;
14    printf("%f %f %f %f %f %f %f %f\n", f[0], f[1], f[2], f[3], f[4], f[5], f[6], f[7]);
15    return 0;
16 }
```

Compile with specific AVX flag:

```
1 > gcc -mavx -o hello_avx hello_avx.c
2 1.000000 1.000000 1.000000 1.000000 1.000000 1.000000 1.000000
```

 <https://godbolt.org/z/nMjrvbcG4>

⚠️ In `_mm256_set_ps`, the ordering of the input parameters is right-to-left



Initialization to Scalar Values

256-bit register can be arranged in different ways

- ▶ by 32- or 64-bit float
- ▶ by 8-, 16-, 32- or 64-bit integer (either signed or unsigned)

<code>_mm256_setzero_ps/pd</code>	Returns a floating-point vector filled with zero
<code>_mm256_setzero_si256</code>	Returns an integer vector whose bytes are set to zero
<code>_mm256_set1_ps/pd</code>	Fill a vector with a floating-point value
<code>_mm256_set1_epi8/epi16/epi32/epi64</code>	Fill a vector with an integer
<code>_mm256_set_ps/pd</code>	Initialize a vector with 8 floats (<code>ps</code>) or 4 doubles (<code>pd</code>)
<code>_mm256_set_m128/m128d/m128i</code>	Initialize a 256-bit vector with two 128-bit vectors
<code>_mm256_setr_ps/pd</code>	Initialize a vector with 8 floats (<code>ps</code>) or 4 doubles (<code>pd</code>) in reverse order
<code>_mm256_setr_epi8/epi16/epi32/epi64</code>	Initialize a vector with integers in reverse order



Load from Memory

Instruction for aligned, unaligned and masked load

Different kind of load instructions

- ▶ aligned load
- ▶ unaligned load
- ▶ masked load

<code>_mm256_load_ps/pd</code>	Loads a floating-point vector from an aligned memory location.
<code>_mm256_load_si256</code>	Loads an integer vector from an aligned memory location.
<code>_mm256_loadu_ps/pd</code>	Loads a floating-point vector from an unaligned memory location.
<code>_mm256_loadu_si256</code>	Loads an integer vector from an unaligned memory location.
<code>_mm_maskload_ps/pd</code> <code>_mm256_maskload_ps/pd</code>	Load portions of a 128-bit/256-bit floating-point vector.
<code>_mm_maskload_epi32/64</code> <code>_mm256_maskload_epi32/64</code>	Load portions of a 128-bit/256-bit integer vector.



Memory Alignment

Data must be aligned at an N-byte boundary:

MMX 64-bit (8 bytes)

SSE, Altivec 128-bit (16 bytes)

AVX 256-bit (32 bytes)

Handling alignment

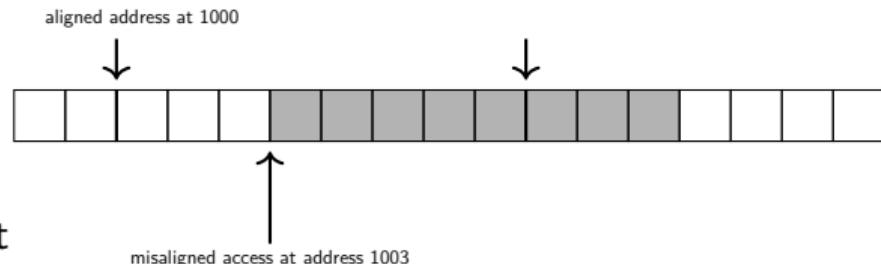
- ▶ software: the programmer has to do it
- ▶ OS exception: very slow
- ▶ hardware: transparent to the programmer

Aligned loads are faster because they can be fetched in a single memory access

- ▶ data is fetched directly into the SIMD register without needing additional work
- ▶ unaligned accesses could involve multiple memory accesses

The compiler aligns `__mm256` (local and global) to 32-byte boundaries

- ▶ **A** to align integer, float, or double arrays, you must use `__declspec(align(32))`



Memory Alignment

When loading data into vectors, memory alignment is important

- ▶ each `_mm256_load_*` intrinsic accepts a memory address that must be aligned on a 32-byte boundary, i.e., the address must be divisible by 32
- ▶ any attempt to load unaligned data with `_m256_load_*` produces a segmentation fault

Example with C11 aligned allocation `aligned_alloc` (alternative are POSIX

`posix_memalign`, Windows `_aligned_malloc` and Intel `__mm_malloc` / `__mm_free`):

```
1 float* aligned_floats = (float*)aligned_alloc(32, 64 * sizeof(float));
2 //... initialize data ...
3 __m256 vec = _mm256_load_ps(aligned_floats);
```

If the data isn't aligned at a 32-bit boundary, the `_m256_loadu_*` functions should be used instead, for example:

```
1 float* unaligned_floats = (float*)malloc(64 * sizeof(float));
2 //... initialize data ...
3 __m256 vec = _mm256_loadu_ps(unaligned_floats);
```



Masked Load

Suppose you want to process a float array of 11 float in AVX

- ▶ 11 isn't divisible by 8, the last five floats of the second `_m256` vector need to be set to zero so they don't affect the computation

This selective loading can be accomplished with the `_maskload_` functions, which accept two arguments:

1. a memory address
2. an integer vector with the same number of elements as the returned vector (called **mask**). For each element in the integer vector whose **highest** bit is one, the corresponding element in the returned vector is read from memory. If the **highest bit** in the integer vector is zero, the corresponding element in the returned vector is set to zero



Masked Load Example

```
1 int array[8]    = { 1, 2, 3, 4, 5, 6, 7, 8 };
2 __m256i mask   = _mm256_setr_epi32(-2, -42, 10, 0, -100, 0, 4, 8);
3 __m256i result = _mm256_maskload_epi32(array, mask);
4
5 int* res = (int*)\&result;
6 printf("%d %d %d %d %d %d %d\n",
7     res[0], res[1], res[2], res[3], res[4], res[5], res[6], res[7]);
```

Output:

```
1 1 2 0 0 5 0 0 0
```

Elements are zeroed out when the *highest bit* is not set in the corresponding mask element `setr` loads the mask in reversed order, e.g.:

- ▶ -100 is **11111111 11111111 11111111 10011100**
- ▶ 8 is **00000000 00000000 00000000 00000100**
- ▶ the most significant bit is 0 for non-negative number



Reversed Order in Set

_mm256_set_epi32 vs _mm256_setr_epi32

```
__m256 _mm256_set_ps (float e7, float e6, ...
float e5, float e4, float e3, float e2, float
e1, float e0)
```

Synopsis

```
__m256 _mm256_set_ps (float e7, float e6, float e5,
float e4, float e3, float e2, float e1, float e0)
#include <immintrin.h>
```

Instruction: Sequence

CPUID Flags: AVX

Description

Set packed single-precision (32-bit) floating-point elements in `dst` with the supplied values.

Operation

```
dst[31:0] := e0
dst[63:32] := e1
dst[95:64] := e2
dst[127:96] := e3
dst[159:128] := e4
dst[191:160] := e5
dst[223:192] := e6
dst[255:224] := e7
dst[MAX:256] := 0
```

```
__m256 _mm256_setr_ps (float e7, float e6, ...
float e5, float e4, float e3, float e2, float
e1, float e0)
```

Synopsis

```
__m256 _mm256_setr_ps (float e7, float e6, float
e5, float e4, float e3, float e2, float e1, float
e0)
#include <immintrin.h>
```

Instruction: Sequence

CPUID Flags: AVX

Description

Set packed single-precision (32-bit) floating-point elements in `dst` with the supplied values in reverse order.

Operation

```
dst[31:0] := e7
dst[63:32] := e6
dst[95:64] := e5
dst[127:96] := e4
dst[159:128] := e3
dst[191:160] := e2
dst[223:192] := e1
dst[255:224] := e0
dst[MAX:256] := 0
```



Store into Memory

Instruction for aligned, unaligned and masked store

Different kind of store instructions

- ▶ aligned store
- ▶ unaligned store
- ▶ masked store

`_mm256_store_ps/pd`

Store a floating-point vector into an aligned memory

`_mm256_store_si256`

Store an integer vector into an aligned memory

`_mm256_storeu_ps/pd`

Store a floating-point vector into an unaligned memory

`_mm256_storeu_si256`

Store an integer vector into an unaligned memory

`_mm_maskstore_ps/pd` `_mm256_maskstore_ps/pd`

Store portions of a 128-bit/256-bit floating-point vector

`_mm_maskstore_epi32/64` `_mm256_maskstore_epi32/64`

Store portions of a 128-bit/256-bit integer vector



Arithmetic Intrinsics

Addition, subtraction and saturated operations

Packed add and subtract, with saturated version for integers.

`_mm256_add_ps/pd`

Add two floating-point vectors

`_mm256_sub_ps/pd`

Subtract two floating-point vectors

`_mm256_add_epi8/16/32/64`

Add two integer vectors

`_mm256_sub_epi8/16/32/64`

Subtract two integer vectors

`_mm256_adds_epi8/16` `_mm256_adds_epu8/16` Add two integer vectors with saturation
`_mm256_subs_epi8/16` `_mm256_subs_epu8/16` Subtract two integer vectors with saturation

The `s` stands for saturated arithmetic: `_adds_` / `_subs_`

- ▶ limited to a fixed range between a minimum and maximum value
- ▶ clamp the result to the minimum/maximun value that can be stored



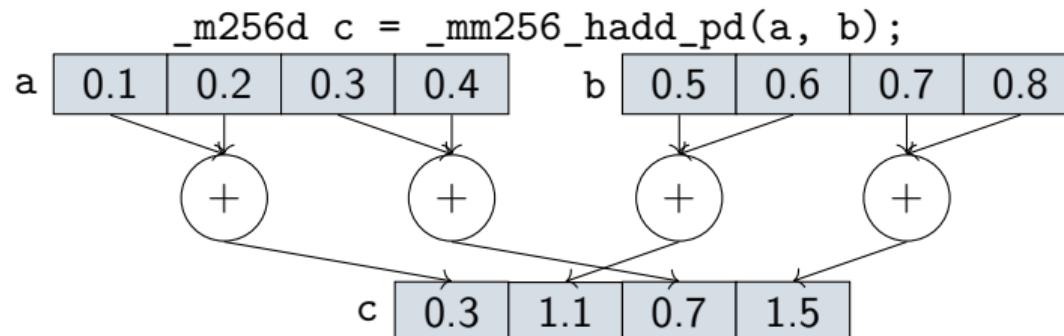
Horizontal Addition and Subtraction

Horizontal functions `_hadd_` / `_hsub_` perform addition and subtraction horizontally

- ▶ add or subtract adjacent elements within each vector
- ▶ results are stored in an interleaved fashion

Example: `_mm256_hadd_pd` horizontally adds double vectors a and b

- ▶ horizontally add adjacent pairs of double-precision (64-bit) floating-point elements in a and b, and pack the results in dst



AVX Example: Horizontal Addition

```
1 #include <immintrin.h>
2 #include <stdio.h>
3 int main() {
4     __m256 vec_a = _mm256_setr_pd(0.1, 0.2, 0.3, 0.4);
5     __m256 vec_b = _mm256_setr_pd(0.5, 0.6, 0.7, 0.8);
6
7     /* Compute the horizontal add between the two vectors */
8     __m256 result = _mm256_hadd_pd(vec_a, vec_b);
9
10    /* Display the elements of the result vector */
11    double* d = (double*)&result;
12    printf("%f %f %f %f \n", d[0], d[1], d[2], d[3]);
13    return 0;
14 }
```

Compile with specific AVX flag:

 <https://godbolt.org/z/edTs13je9>

```
1 > gcc -mavx horizontal_add.c
2 0.300000 1.100000 0.700000 1.500000
```



Multiplication and Division

For floating point and integer

Floating point multiplication (`ps` and `pd`) works like other arithmetic instructions.

However, integer multiplication is subject to overflow

<code>_mm256_mul_ps/pd</code>	Multiply two floating-point vectors
<code>_mm256_mul_ep32/epu32</code>	Multiply the lowest four elements of vectors containing 32-bit integers
<code>_mm256_mullo_epi16/32</code>	Multiply integers and store low halves
<code>_mm256_mulhi_epi16/epu16</code>	Multiply integers and store high halves
<code>_mm256_mulhrs_epi16</code>	Multiply 16-bit elements to form 32-bit elements
<code>_mm256_div_ps/pd</code>	Divide two floating-point vectors

⚠ If you multiply 2 N-bit numbers, the results can take up to $2N$ bits

Low and high bits:

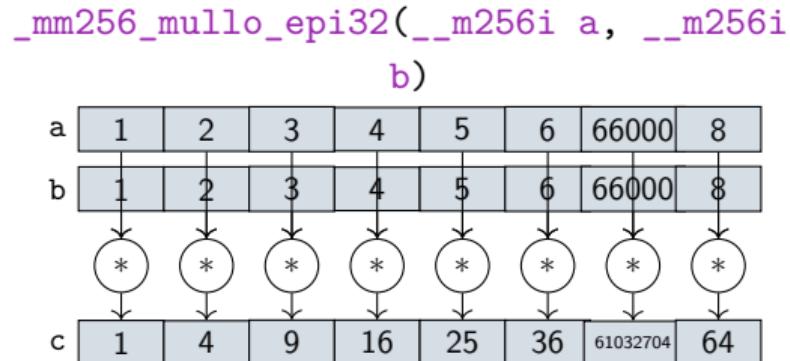
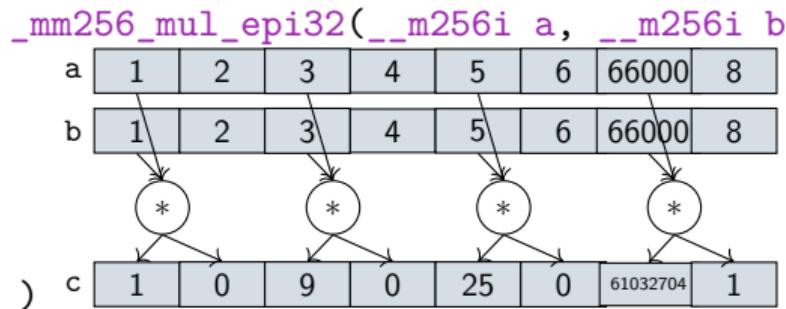
- ▶ `_mullo_` multiply every element but only store the low half of each product
- ▶ `_mulhi_` multiply every element but only store the high half of each product
- ▶ ⚠ there is no `mulhi_epi32`



Low and High with Integer Multiplication

mul vs mullo

⚠ `_mm256_mul_ep32` only multiply the low signed 32-bit integers from each packed 64-bit element, and store the signed 64-bit in the result.



Note that in C/C++, integer overflow is

- ▶ unsigned integer: modulo power of two
- ▶ signed integer: undefined behavior

Example: Low and High with Integer Multiplication

```
1 #include <immintrin.h>
2 #include <stdio.h>
3 int main() {
4     // Initialize the two integer vectors
5     __m256i a = _mm256_setr_epi32(1, 2, 3, 4, 5, 6, 66000, 8);
6     __m256i b = _mm256_setr_epi32(1, 2, 3, 4, 5, 6, 66000, 8);
7
8     // Compute multiply with mullow and mul
9     __m256i lo = _mm256_mullo_epi32(a, b);
10    __m256i hi = _mm256_mul_epi32(a, b);
11
12    // Print results, note overflow on the 7th element
13    int32_t* i = (int32_t*)&hi;
14    printf("%d %d %d %d %d %d %d %d\n", i[0], i[1], i[2], i[3], i[4], i[5], i[6], i[7]);
15    i = (int32_t*)&lo;
16    printf("%d %d %d %d %d %d %d %d\n", i[0], i[1], i[2], i[3], i[4], i[5], i[6], i[7]);
17    return 0;
18 }
```

Compile with specific AVX flag:

```
1 > gcc -fno-tree-vectorize -mavx2 mullo_add.c
2 1 0 9 0 25 0 61032704 1
3 1 4 9 16 25 36 61032704 64
```

 <https://godbolt.org/z/MEYPb8YaY>



FMA: Fused Multiply Add Instruction

FMA provides greater speed and accuracy than performing multiplication and addition separately. In particular, for accuracy

- ▶ when you multiply two floating-point values, a and b, the result is really `round(a*b)`, where `round(x)` returns the floating-point value closest to x
- ▶ with FMA: instead of `round(round(a*b)+c)`, they return `round(a*b+c)`

`_mm_fmadd_ps/pd` Multiply two vectors and add the product to a third vector
`_mm256_fmadd_ps/pd` (`res = a * b + c`)

`_mm_fmsub_ps/pd` Multiply two vectors and subtract a vector from the product
`_mm256_fmsub_ps/pd` (`res = a * b - c`)



FMA Example

Example without FMA support

```
1 void multiply_and_add(const float*a, const float*b, const float*c, float*d) {  
2     for(int i=0; i<8; i++) {  
3         d[i] = a[i] * b[i];  
4         d[i] = d[i] + c[i];  
5     }  
6 }
```

Example with AVX FMA

```
1 __m256 multiply_and_add(__m256 a, __m256 b, __m256 c) {  
2     return _mm256_fmadd_ps(a, b, c);  
3 }
```



Exercise: AVX2 Intrinsics Programming

Write a C/C++ program that

1. implements a function `vec_add(float *a, float *b, float *c, int size)` that add the elements of two arrays and writes the results in a third array of the same size. You can assume that the input is not aligned and the size a multiple of 8.
2. implements a function `vec_add_aligned`, which is the same as of (1), but assume the arrays are aligned
3. implements a function `vec_add_anysize`, which is the same of (2), but support input/output arrays of any size.

For all three functions, evaluate the performance against a scalar code compiler in `-O0` and a scalar code compiled in `-O3`.

