



DeepL

Subscribe to DeepL Pro to translate larger documents.
Visit www.DeepL.com/pro for more information.

Vector units

Oguz Kaya

University' Paris-Saclay

Orsay, France

1 Principles

1 Principles

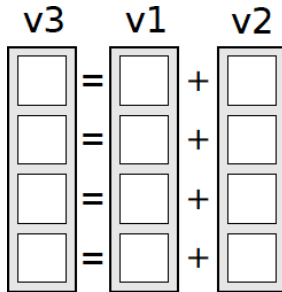
2 The AVX unit & the intrinseques

1 Principles

2 The AVX unit & the intrinseques

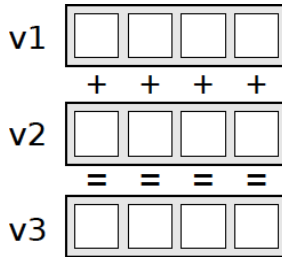
A scalar processor performs the operations ~~on~~ ^{by} each operation on scalar data (a scalar)

```
v3[0] = v1[0] + v2[0];  
v3[1] = v1[1] + v2[1];  
v3[2] = v1[2] + v2[2];  
v3[3] = v1[3] + v2[3];
```



A vector processor performs operations in parallel on a set of scalars, each operation being on vectors

```
v3 = v1 + v2;
```



Advantages/inconveniences of the vector

1 Advantage:

- Instructions n fois plus rapide, n la largeur de l'unité vectorielle si l'algorithme se prête au paradigme SIMD: une seule instruction simultanément sur plusieurs données
 - vector calculus scientific
 - simulations image
 - processing
 - automatic learning (neural networks)

2 Inconvénient:

- Gain only if the algorithm se prête au SIMD paradigm.
- Otherwise it is a loss of circuitry that could be used for other purposes (cores, cache, register, ...).
- Sometimes difficult (and dirty) to use by hand (automation possible).

- ~~At~~es 60: first Solomon project
- 1970s: ILLIAC IV, University of Illinois, then CDC,
- Cray, NEC, Hitachi, Fujitsu
- since 1996: scalar + vector processors, Pentium (MMX, SSE, AVX), PowerPC (AltiVec)
- now: in almost all architectures (x86, ARM, RISC-V)

General operation of SIMD units

Same as the operation of scalar registers, except that each operation is done on a set of data

- 1 Loading data from memory to SIMD registers
- 2 Operation on SIMD registers
- 3 Placement of the result in memory

- **Automatic vectorization:** The compiler tries to vectorize the code automatically.
- **Assembly programming:** Low-level programming, direct manipulation of registers and instructions.
- **Intrinsic:** Manipulation of variables by finite intrinsic functions in a `<mmmintrin.h>` file.
- **High-level libraries (OpenMP):** Self-vetted using the tips in the library.

1 Principles

2 The AVX unit & the intrinseques

Streaming SIMD Extensions (SSE)

- SSE - Pentium III, Athlon XP
- SSE2 - Pentium 4
- Willamette SSE3 - Pentium 4 Prescott
- SSSE3 - (Supplemental SSE3) Core Woodcrest
- SSE4.1 - Core 2 Penryn
- SSE4.2 - Core i7 Nehalem
- AVX - Sandy Bridge
- **AVX2** - Haswell
- AVX512 - Skylake and beyond

Streaming SIMD Extensions (SSE)

Each new version of SSE/AVX contains the previous SSE(N) = new instructions + SSE(N-1)

Intrinsics

The intrinsics are accessible through .h files that define: variable types

- functions operating on these variables

- SSE: `xmmintrin.h`
- SSE2: `emmintrin.h`
- SSE3: `pmmmintrin.h`
- SSSE3: `tmmintrin.h`
- SSE4.1: `smmintrin.h`
- SSE4.2: `nmmintrin.h`
- **AVX/AVX2: `immintrin.h`**

Run `lscpu` under Linux to display the ~~spec~~ version of e.

Advanced Vector Extensions (AVX)

16 256bit registers on x86-64: YMM0 - > YMM15

- memory transfers < - > AVX registers arithmetic
- operations
- operations logic tests
- permutations
-

Compilation with AVX

```
gcc program.c -o program -mavx2
```

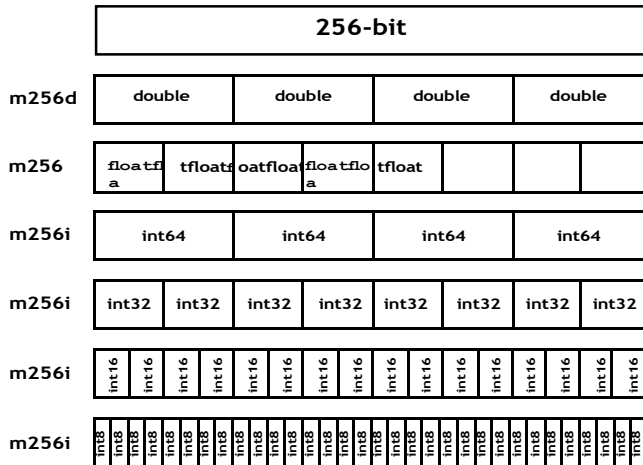
```
./program
```

Types of variables

- `__m256`: 8 floats
- `__m256d`: 4 doubles
- `__m256i`: 256/width integers (*width* = 8, 16, 32, 64)

There is a set of 'finite' instructions for each type.

AVX vectors



```
__m256 {operation}{non-alignment} {dataorganization}{datatype}...
```

```
__m256 mm256_add_ps( __m256 A, __m256 B)
```

`_mm256_loadu_ps256`

- s (single): single precision float (32bits)
- d (double): float double precision (64bits) i...
- (integer): integer
- p (packed): contiguous, operate on the whole
- vector s (scalar): operate only on one element
- (obsolete!) u (unaligned): data not aligned in
- memory l (low): low order bits
- h (high): most significant bits
- r (reversed): in reverse order

- 1 AVX variables declaration (registers): m256
__r1;
- 2 loading data from the memory to the AVX registers: float A[8]
__attribute__((aligned(16)));
...
r1 = mm256_load_ps(&A[0]);
- 3 operations on SIMD registers r1
= mm256_add_ps(r1, r1);
- 4 placing the content of the registers in memory:
_mm256_store_ps(&A[0], r1)

Transfers between the memory and the AVX

- **Floating:**

`__m256 mm256_load_ps(float *addr)`

`void mm256_store_ps(float *addr, __m256 a)`

- **Doubles:**

`__m256d mm256_load_pd(double *addr)`

`void mm256_store_pd(double*addr, __m256d a)`

- **Entity:**

`__m256i mm256_load_si256(m256i const * addr)`

`void mm256_store_si256(m256i const * addr, __m256i a)`

Transfers from memory to AVX registers

- aligned or not aligned:
_mm256_load. ... or mm256_loadu. ...
_mm256_blind. ... or mm256_blindu. . .
- per vector: $8 \times SP$, $4 \times DP$, ... mm256
_load ps, mm256 load pd, ... mm256
_store ps, mm256 store pd, ...

- The access to a variable is **aligned** if the address of this variable is a multiple of its size. This constraint is material.
- For AVX variables, the processor can perform efficient transfers of 32-byte (256-bit) AVX variables if it is aligned to 32 bytes.
- To get the alignment of a type: `__alignof_(type)`

A cache line is usually 32 to 64 bytes in size.

If the data ~~into~~ into an SSE register comes from a memory area that "straddles" 2 cache lines, then the 2 cache lines must be read in order to fill the SSE register, which implies a decrease in performance

Allocation of aligned donors

- 32 bytes static alignment: `int x __attribute`

`__((aligned(32))`

`float y[8] __attribute __ ((aligned (32)))`

- Dynamic alignment on 32 bytes:

`type* x = (type*) mm malloc(size*sizeof(type), 32);`

- Non-aligned acce's can be much slower (but not as much in the new generations).

- `_mm256_add_pd(_m256d A, _m256d B)` - (Add-Packed-Double)
 - Inputs: [A0, A1, A2, A3], [B0, B1, B2, B3]
 - Output: [A0 + B0, A1 + B1, A2 + B2, A3 + B3]
- `_mm256_add_ps(_m256 A, _m256 B)` - (Add-Packed-Single)
 - Inputs: [A0, ... , A7], [B0, ... , B7]
 - Output: [A0 + B0, ... , A7 + B7]
- `_mm256_add_epi64(_m256i A, _m256i B)` - (Add-Packed-Int64)
 - Inputs: [A0, A1, A2, A3], [B0, B1, B2, B3]
 - Output: [A0 + B0, A1 + B1, A2 + B2, A3 + B3]
- `_mm256_add_epi32(_m256i A, _m256i B)` - (Add-Packed-Int32)
- `_mm256_add_epi16(_m256i A, _m256i B)` - (Add-Packed-Int16)
- - - - - -

Arithmetic operations

mm256 add_epi8(m256i A, m256i B) - (Add-Packed-Int8)

Same for mm256 sub . .

- `_mm256_mul_pd(_m256d A, m256d B)` - (Multiply-Packed-Double)
 - Inputs: `[A0, A1, A2, A3]`, `[B0, B1, B2, B3]`
 - Output: `[A0 * B0, A1 * B1, A2 * B2, A3 * B3]`
- `_mm256_mul_ps(_m256 A, m256 B)` - (Multiply-Packed-Single)
 - Inputs: `[A0, ... , A7]`, `[B0, ... , B7]`
 - Output: `[A0 * B0, ... , A7 * B7]`
- `_mm256_mul_epi64(_m256i A, m256i B)` - (Multiply-Packed-Int64)
 - Inputs: `[A0, A1, A2, A3]`, `[B0, B1, B2, B3]`
 - Output: `[A0 * B0, A1 * B1, A2 * B2, A3 * B3]`
- `_mm256_mul_epi32(_m256i A, m256i B)` - (Multiply-Packed-Int32)
- `_mm256_mul_epi16(_m256i A, m256i B)` - (Multiply-Packed-Int16)
- - - - - -

Arithmetic operations

mm256 mul epi8(m256i A, m256i B) - (Multiply-Packed-Int8)

Ditto for mm256 div ... ; beware, the division is **slow!**

- `_mm256_fmadd_pd(_m256 A, __m256 B, __m256 C) -`
(Fused-Multiply-Add-Packed-Double)
 - `[A0, A1, A2, A3], [B0, B1, B2, B3], [C0, C1, C2, C3]` Output: [
■ `A0 * B0 + C0, A1 * B1 + C1, A2 * B2 + C2, A3 * B3 + C3]`
- `_mm256_fmadd_ps(_m256 A, __m256 B, __m256 C)`
- (Fused-Multiply-Add-Packed-Single)
 - `[A0, ... , A7], [B0, ... , B7], [C0, ... , C7]` Output:
■ `[A0 * B0 + C0, ... , A7 * B7 + C7]`

Same for `mm256 fmsub . .`

- `_mm256_hadd_pd(A, B)` - (Horizontal-Add-Packed-Double)
 - `[A0, A1, A2, A3], [B0, B1, B2, B3]`
 - Output: `[A0 + A1, B0 + B1, A2 + A3, B2 + B3]`
- `_mm256_hadd_ps(A, B)` (Horizontal-Add-Packed-Single)
 - `[A0, ..., A7], [B0, ..., B7]`
 - Output: `[A0 + A1, A2 + A3, B0 + B1, B2 + B3, A4 + A5, ..., B6 + B7]`
- `_mm256_hadd_epi32(_mm256i A, _mm256i B, _mm256i C)` - (Horizontal-Add-Packed-Int32)
- `_mm256_hadd_epi16(_mm256i A, _mm256i B, _mm256i C)` - (Horizontal-Add-Packed-Int16)
- Same as `_mm256_hsub` . .

Logical and comparative operations

- `_mm256_{and, or, xor, ...} {ps, pd, si256}(A,`
- `_B) mm256 cmp {ps, pd}(A, B, op)`
 - `op = CMP_LT OS: Comparison $A < B$ op =`
 - `= CMP_LE OS: Comparison $A \leq B$ op =`
 - `CMP_EQ OS: Comparison $A == B$ op =`
 - `CMP_GT OS: Comparison $A > B$ op =`
 - `CMP_GE OS: Comparison $A \geq B$`
- `_mm256_cmpeq_epi{64, 32, 16, 8}(m256i A, __m256i B)`
- `_mm256_cmpgt_epi{64, 32, 16, 8}(m256i A, __m256i B)`

For `C = mm256 cmp. . (A, B)`, if `A[i] op B[i]` is correct, `C[i] = 0b11. 1`. If it is false, `C[i] = 0b00. 0`.

The instruction allows to obtain any permutation of data in a register AVX. The permutation is specified by an immediate entry of which bits are used as indices.

- `__m256d B = mm256_permute4x64_pd(m256d A, const int imm8)`
 $B[0] = A[\text{imm8}[1:0]]$
 $B[1] = A[\text{imm8}[3:2]]$
 $B[2] = A[\text{imm8}[5:4]]$
 $B[3] = A[\text{imm8}[7:6]]$
- `__m256i B = mm256_permute4x64_epi64(m256i A, const int imm8)`
 $B[0] = A[\text{imm8}[1:0]]$; $B[1] = A[\text{imm8}[3:2]]$
 $B[2] = A[\text{imm8}[5:4]]$; $B[3] = A[\text{imm8}[7:6]]$
- `__m256 B = mm256_permutevar8x32_ps(m256 A, __m256i idx)` $B[i] = A[\text{idx}[(32 * i) + 2 : (32 * i)]]$
- `__m256i B = mm256_permutevar8x32_epi32(m256i A, __m256i idx)` $B[i] = A[\text{idx}[(32 * i) + 2 : (32 * i)]]$

Assignment of constants

- `_mm256_set_ps(float f0, . . . , float f7)`
- `_mm256_set_pd(double d0, double d1, double d2, double d3)`
- `_mm256_set_epi64x(int64_t i0, . . . , int64_t i3)`
- `_mm256_set_epi32(int i0, . . . , int i7)`
- `_mm256_set_epi16(short i0, . . . , short i15)`
- `_mm256_set_epi8(char i0, . . . , char i31)`
- There is also the variant `_mm256_set1_ps(float f)` which assigns the same value to all the elements of the vector (same for double, `int{64, 32, 16, 8}`)

Tips for good performance

- Organize data correctly to avoid non-contiguous access.
- Stay in the cache as much as possible.
- Align the data (less important in the ~~new~~ versions).
- Stay as long as possible in the registers (i.e., avoid reading/writing on the memory)

The initial structure could not ~~lead~~ to the vectorization. In this case it would be necessary to rethink it.

- Array of structures (AoS): xyzwxyzwxyzw
- Structure of Arrays (SoA): xxxxxxxxxx . . . yyyyyyyy . . .
zzzzzzzzzz . . . wwwwwwwwwwww . . .
- Hybrid structure:
xxxxyyyzzzzwwwxxxxyyyzzzzwww .

- Vector units allow SIMD operation mode on a data vector. The state of the AVX
- instruction set is rather chaotic (just like x86-64) Instructions not compatible with all
- data types
- Substantial gain potential (16-32x with AVX2)