

Vector SIMD Instructions

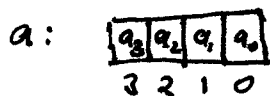
①

Contents

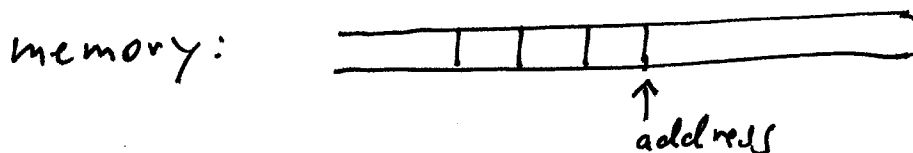
- basic instructions 2-5
- useful building blocks
 - load/store 6-7
(gather/scatter)
 - reordering 8-9
 - matrix-vector product 10

Information

vector indexing:



in most intrinsics, the order of operands matter



← address increases

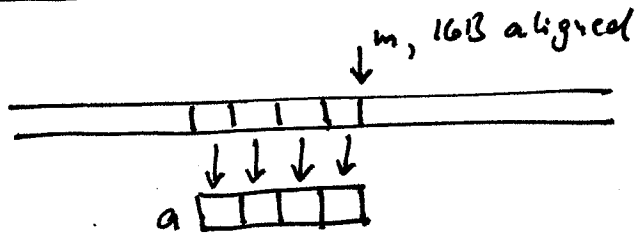
SIMD extensions timeline: SSE, SSE2, SSE3, SSSE3, SSE4

We focus on single precision float, 4-way

1 vector = 128 bit = 16 B, data type -- m128

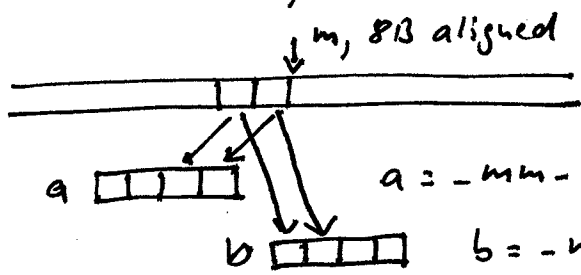
Unless stated otherwise, instructions are SSE or later

Load and Store

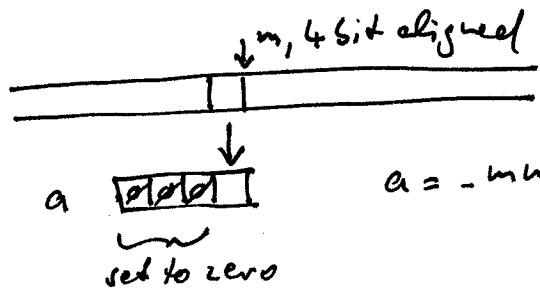


$a = -mm_load_ps(m);$ aligned
 $a = -mm_loadu_ps(m);$ unaligned (avoid)

$a = p[i];$ if p is: $--m128 * p$



$a = -mm_loadl_pi(a, m);$ (keeps upper half)
 $b = -mm_loadh_pi(b, m);$ (keeps lower half)



$a = -mm_load_ss(m)$

stores are analogous

Constants

c:

4.0	3.0	2.0	1.0
-----	-----	-----	-----

$c = -mm_set_ps(4.0, 3.0, 2.0, 1.0);$

d:

1.0	1.0	1.0	1.0
-----	-----	-----	-----

$d = -mm_set1_ps(1.0);$

e:

∅	∅	∅	1.0
---	---	---	-----

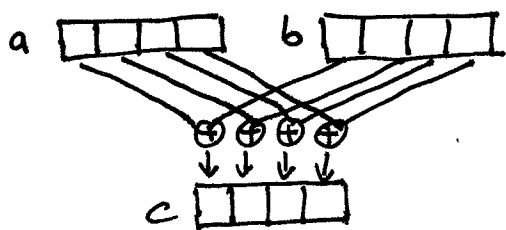
$e = -mm_set_ss(1.0);$

f:

∅	∅	∅	∅
---	---	---	---

$f = -mm_setzero_ps();$

Vector arithmetic



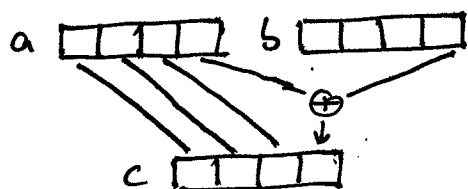
$$c = \text{-mm-add-ps}(a, b); \quad "a+b"$$

analogous:

$$c = \text{-mm-sub-ps}(a, b); \quad "a-b"$$

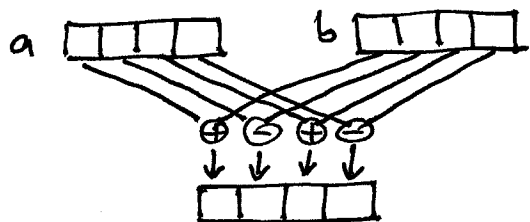
$$c = \text{-mm-mul-ps}(a, b); \quad "a \cdot b"$$

Scalar arithmetic

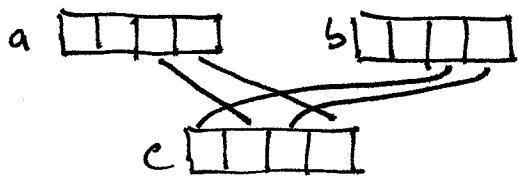


$$c = \text{-mm-addss}(a, b);$$

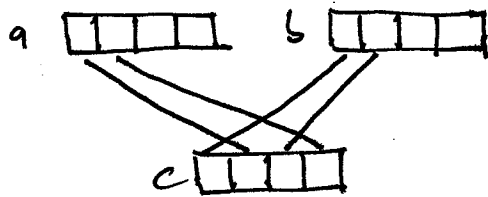
AddSub (SSE3 and later)



Reorder Instructions

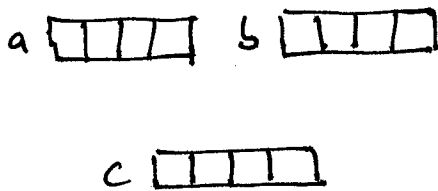


$$c = \text{_mm_unpacklo_ps}(a, b)$$



$$c = \text{_mm_unpackhi_ps}(a, b)$$

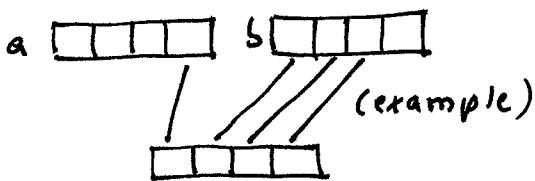
shuffle:



$$c = \text{_mm_shuffle_ps}(a, b, \text{_MM_SHUFFLE}(2, 4, 1, 3))$$

$$\begin{aligned} c_0 &= a_i & i, j, k, l \in \{0, \dots, 3\} \\ c_1 &= a_j \\ c_2 &= b_k \\ c_3 &= b_l \end{aligned}$$

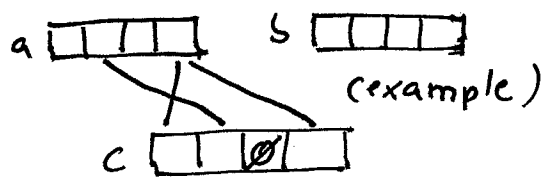
align: (SSE3 and later)



"any 4 consecutive elements of the concatenation of a and b go into c"

`_mm_alignr_epi8` use with `_mm_castsi128_ps`

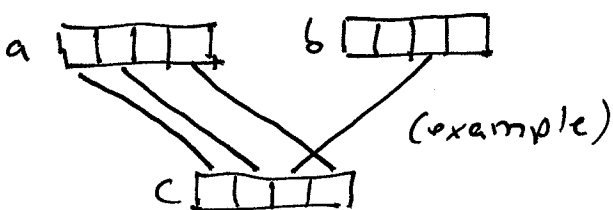
shuffle: (SSE3 and later)



"c is filled in each position with any element from a or 0, as specified by b"

`_mm_shuffle_epi8`

blend: (SSE4 and later)

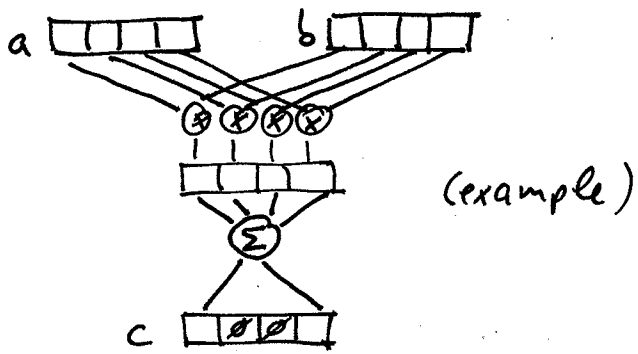


"c is filled in each position with any element from a or b from the same position"

`_mm_blend_ps`

Dot product (SSE4 and later)

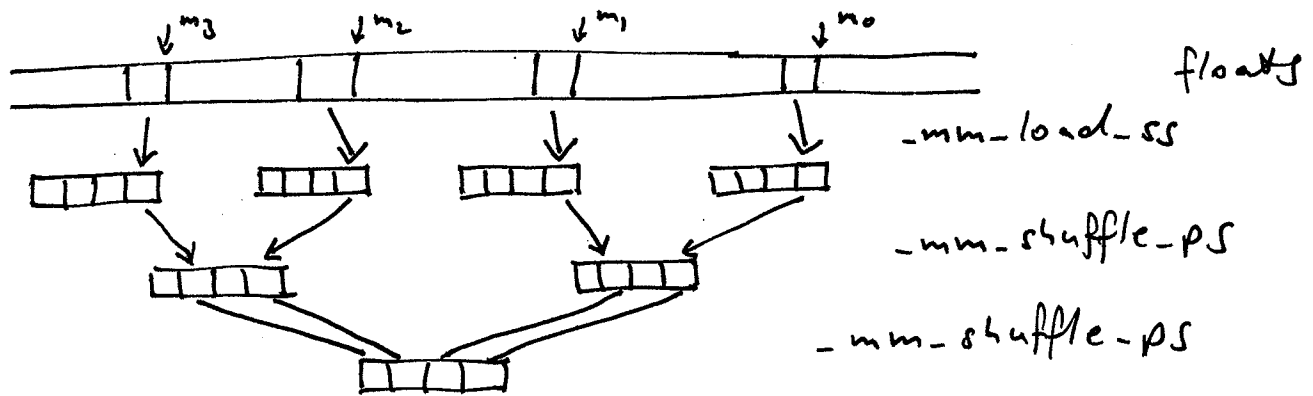
5



"computes the pointwise product of a and b and writes an arbitrary sum of the resulting numbers into selected elements of c — the others are set to zero"

-mm_dp-ps(a, b, mask)

Load 4 real numbers from arbitrary memory locations



7 instructions, this is the right way

Note:

- whenever possible avoid this by restructuring the algorithm or data to have aligned vector loads `-mm-load-ps`
- the above should be equivalent to the following but a.) the above is safer; b.) be aware that the below are 7 instructions

```
float f[20] = {...};
--m128 vf = -mm-set-ps(f[3], f[5], f[17], f[13]);
```

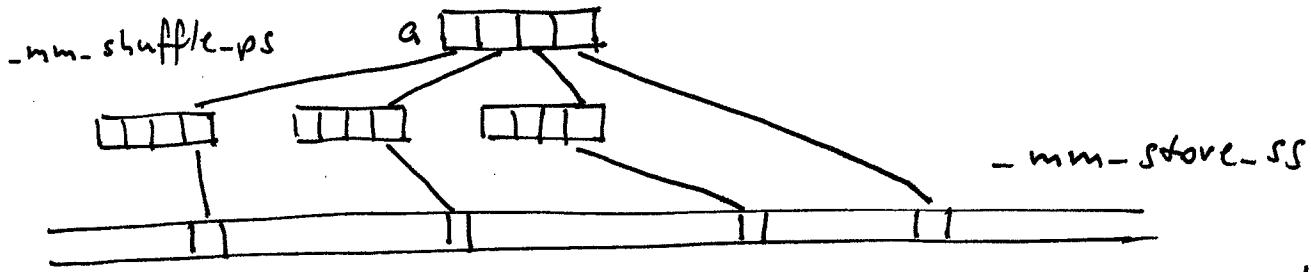
Don't do this:

```
float f[20] = {...};
--declspec(aligned(16)) g[4];
--m128 vf;
g[0] = f[3];
g[1] = f[17];
g[2] = f[5];
g[3] = f[13];
vf = -mm-load-ps(g);
```

} mem → register → mem round trip
⇒ expensive

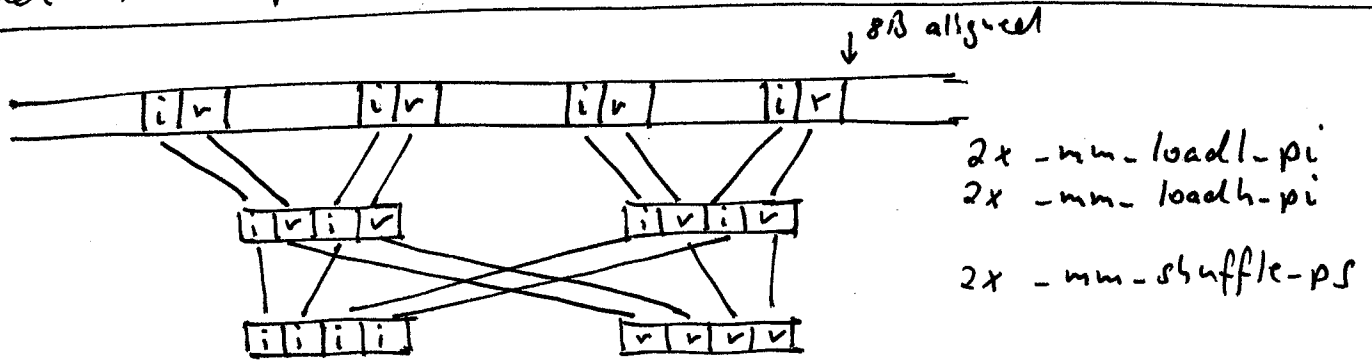
(same problem with unions and pointers)

Store 4 real numbers to arbitrary memory locations (7)



7 instructions, shorter critical path than load

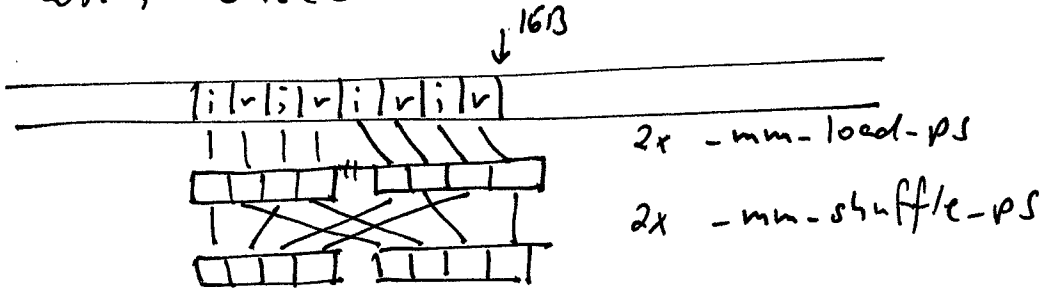
Load 4 complex numbers (= 4 pairs of real numbers)



6 instructions

store analogous

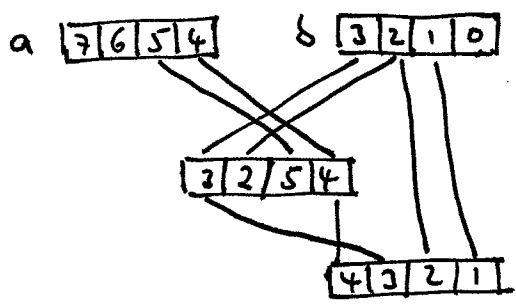
same with consecutive data:



4 instructions

store analogous

Shift by 1

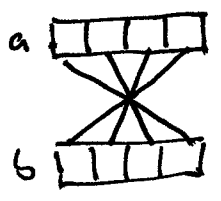


`_mm_shuffle_ps`
`_mm_shuffle_ps`

2 instructions

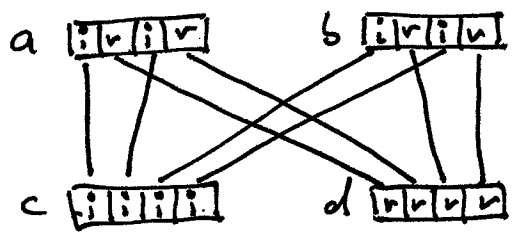
SSE3 and later: `_mm_alignr_epi8` + casts 1 instruction

Reverse vector



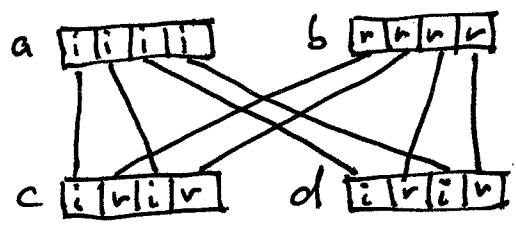
`b = _mm_shuffle_ps(a, a, _MM_SHUFFLE(0, 3, 2, 1));`

Interleaved complex → split complex



`c = _mm_shuffle_ps(b, a, _MM_SHUFFLE(3, 1, 2, 1));`
`d = _MM_SHUFFLE(2, 0, 3, 0);`

Split complex → interleaved complex



`c = _mm_unpackhi_ps(b, a);`
`d = _mm_unpacklo_ps(b, a);`

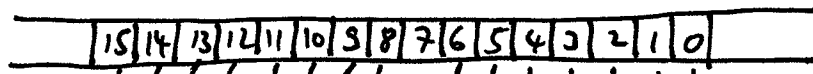
Transposition: 4x4 matrix

4x4 matrix:

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

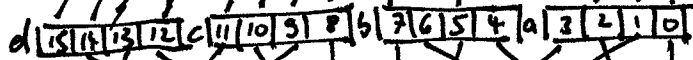
= A

in memory:



16B aligned
↓

4 aligned loads



4 shuffles

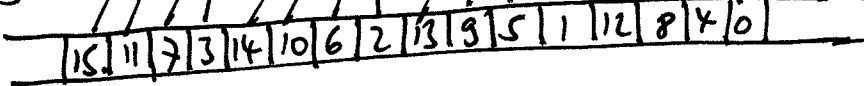


4 shuffles



4 aligned stores

in memory:



as matrix:

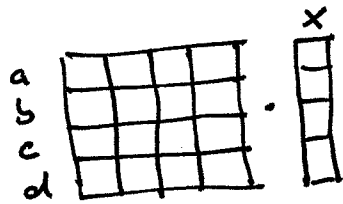
0	4	8	12
1	5	9	13
2	6	10	14
3	7	11	15

= A^T

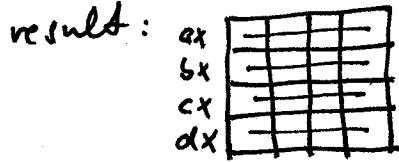
⊗ done by the macro `_MM_TRANSPOSE4_PS(a, b, c, d);`

8 instructions

Matrix-vector product

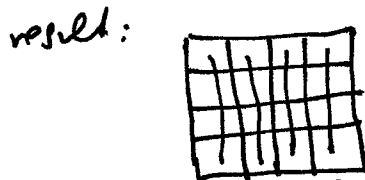


1. step: 4 vector products ax, bx, cx, dx (4 instructions)



SSE:

2. step: transpose (8 instructions)

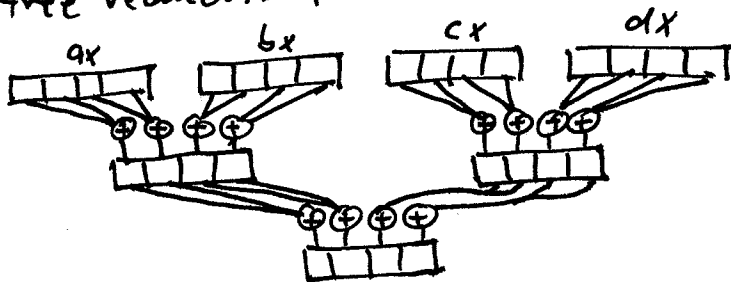


3. step: sum rows (3 instructions)

total: 15 instructions

SSE2:

2. step: tree reduction



3 instructions
(`_mm_hadd_ps`)

total: 7 instructions

SSE4: has dot product instruction but still 7 instructions are needed (exercise)