

332

Advanced Computer Architecture

Chapter 7.1:

Vectors, vector instructions, vectorization and SIMD

March 2019

Paul H J Kelly

This section has contributions from Fabio Luporini (postdoc at Imperial) and Luigi Nardi (now a postdoc at Stanford).

Course materials online at

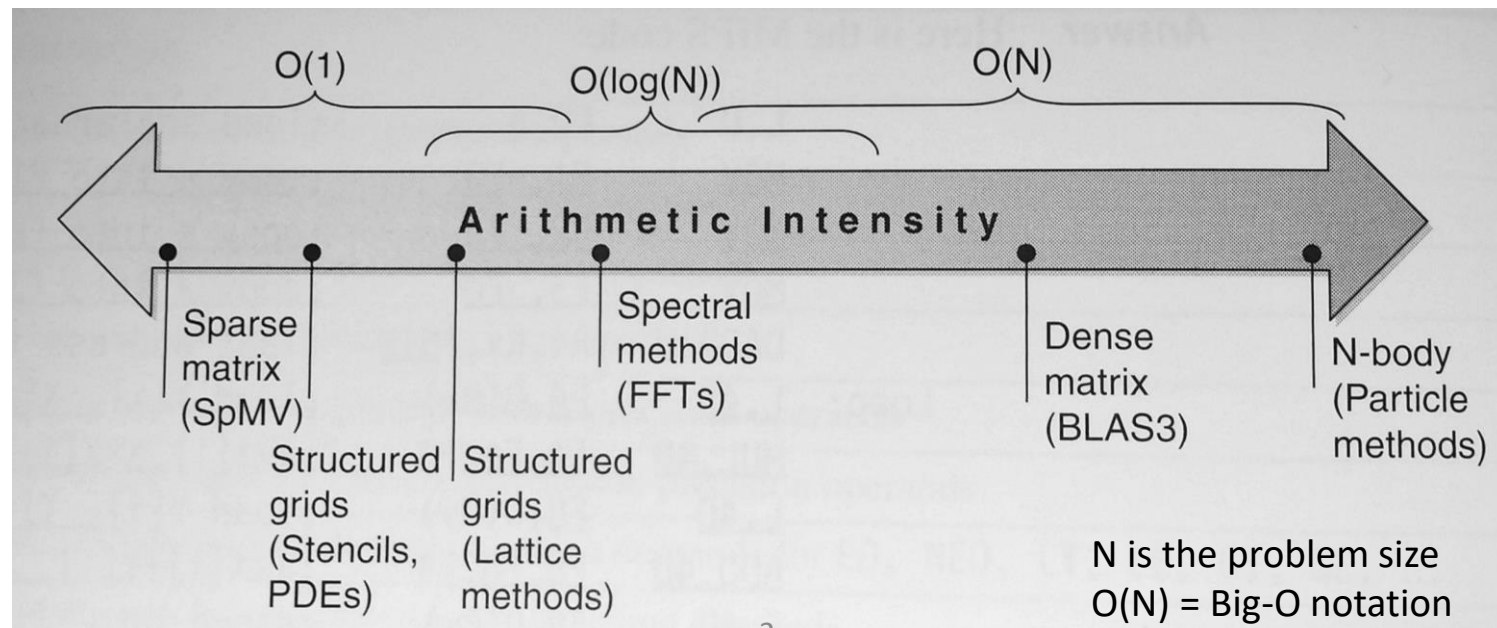
<http://www.doc.ic.ac.uk/~phjk/AdvancedCompArchitecture.html>

Arithmetic Intensity

E5-2690 v3 SP	CPU	416	68	~6	~24
E5-2690 v3 DP	CPU	208	68	~3	~24
K40 SP	GPU	4,290	288	~15	~60
K40 DP	GPU	1,430	288	~5	~40

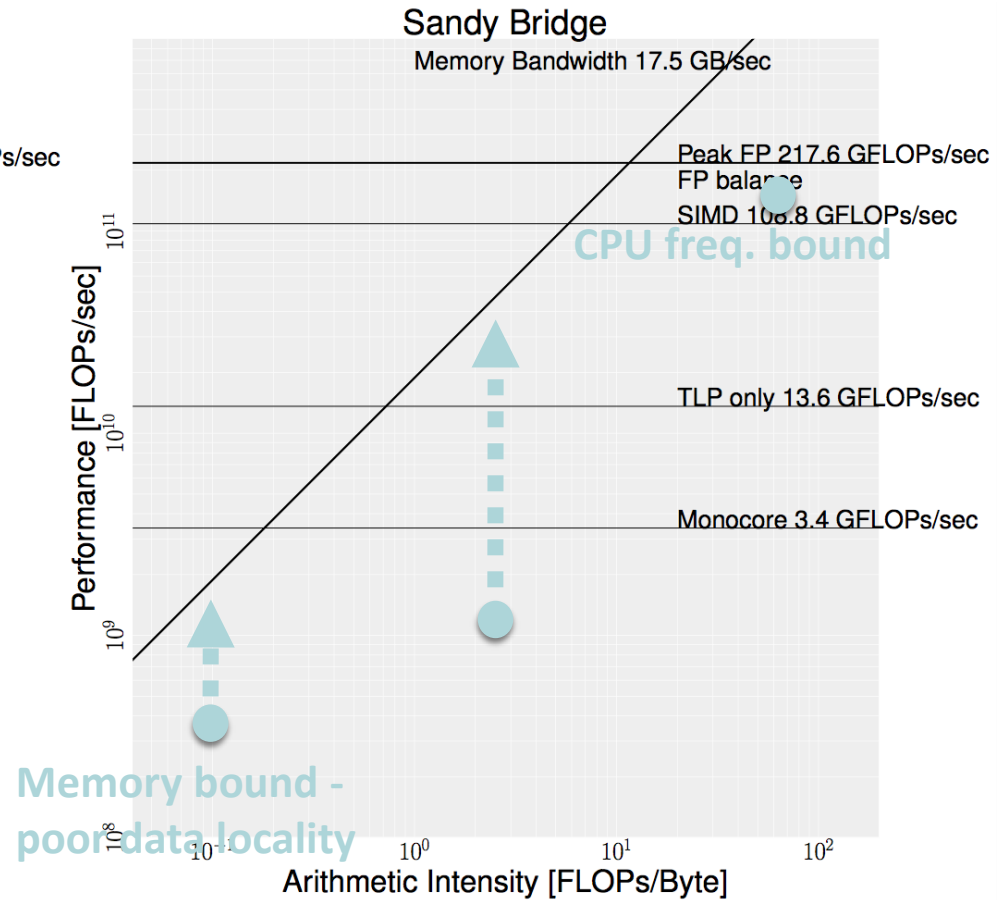
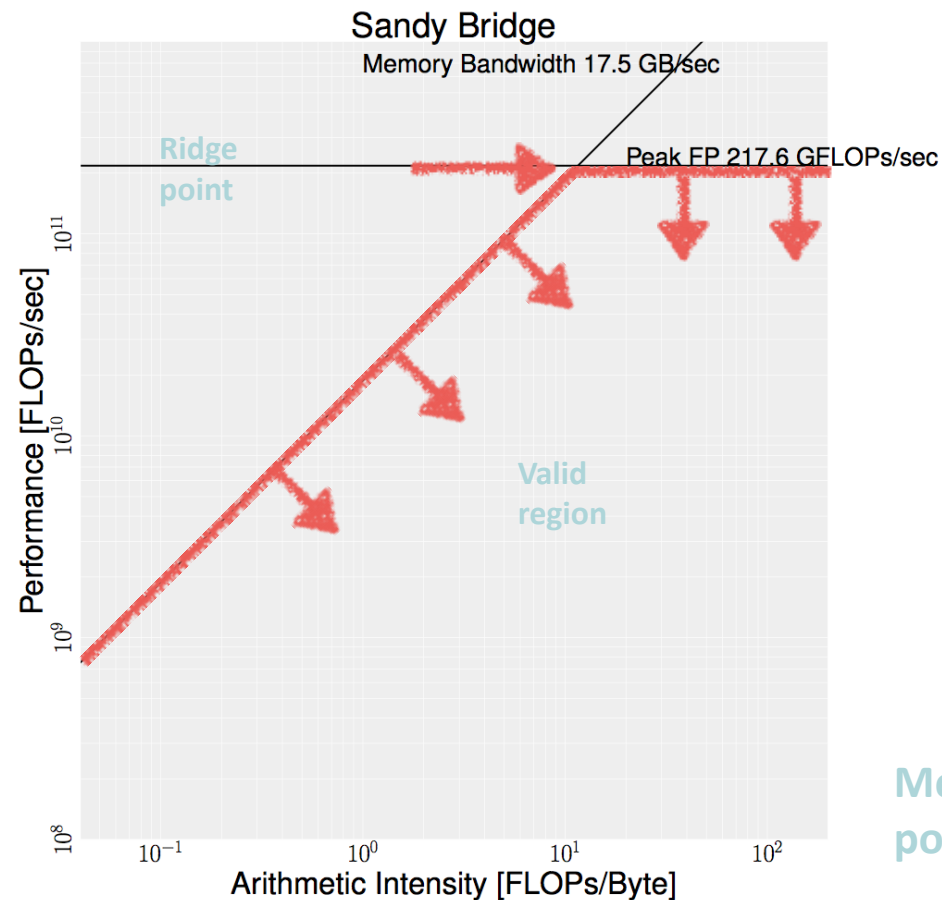
Without enough Ops/Word codes are likely to be bound by operand delivery
(SP: single-precision, 4B/word; DP: double-precision, 8B/word)

Arithmetic intensity: Ops/Byte of DRAM traffic

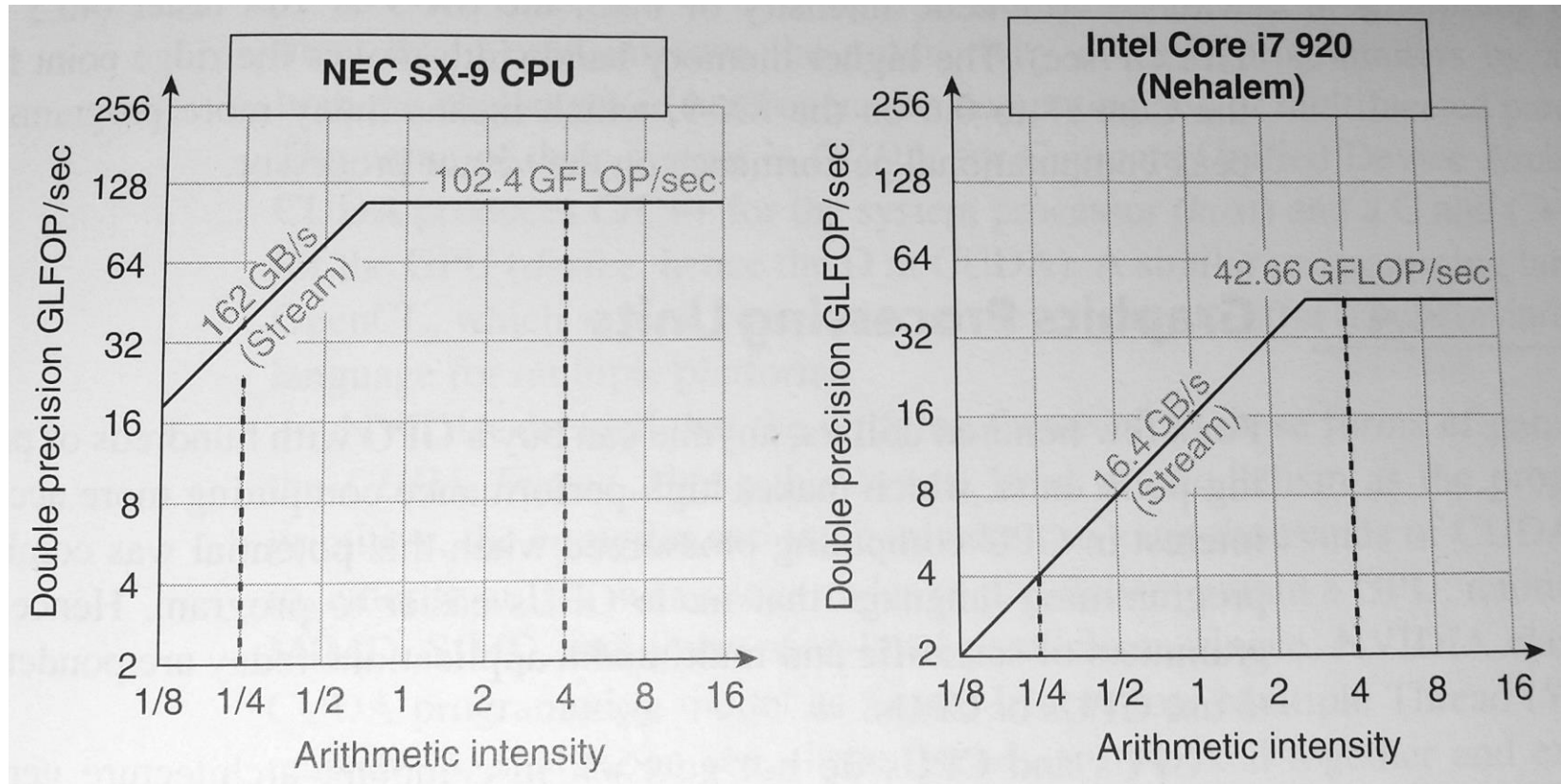


Roofline Model: Visual Performance Model

- Bound and bottleneck analysis (like Amdahl's law)
- Relates processor performance to off-chip memory traffic (bandwidth often the bottleneck)



Roofline Model: Visual Performance Model



- The ridge point offers insight into the computer's overall performance potential
- It tells you whether your application *should* be limited by memory bandwidth, or by arithmetic capability

Vector instruction set extensions

- Example: Intel's AVX512
- Extended registers ZMM0-ZMM31, 512 bits wide
 - Can be used to store 8 doubles, 16 floats, 32 shorts, 64 bytes
 - So instructions are executed in parallel in 64,32,16 or 8 “lanes”
- Predicate registers k0-k7 (k0 is always true)
 - Each register holds a predicate *per operand* (per “lane”)
 - So each k register holds (up to) 64 bits*
- Rich set of instructions operate on 512-bit operands

* k registers are 64 bits in the AVX512BW extension; the default is 16

AVX512: vector addition

- Assembler:

- `VADDPS zmm1 {k1}{z}, zmm2, zmm3`

- In C the compiler provides “vector intrinsics” that enable you to emit specific vector instructions, eg:

- `res = _mm512_maskz_add_ps(k, a, b);`

- Only lanes with their corresponding bit in `k1` are activated

- Two predication modes: *masking* and *zero-masking*

- With “zero masking” (shown above), inactive lanes produce zero
- With “masking” (omit “z” or “{z}”), inactive lanes do not overwrite their prior register contents

AVX512: vector addition

- Assembler:
 - `VADDPS zmm1 {k1}{z}, zmm2, zmm3`
- In C the compiler provides “vector intrinsics” that enable you to emit specific vector instructions, eg:
 - `res = _mm512_maskz_add_ps(k, a, b);`
- Only lanes with their corresponding bit in `k1` are activated
- Two predication modes: *masking* and *zero-masking*
 - With “zero masking” (shown above), inactive lanes produce zero
 - With “masking” (omit “z” or “{z}”), inactive lanes do not overwrite their prior register contents

More formally...

FOR $j \leftarrow 0$ TO $KL-1$

$i \leftarrow j * 32$

IF $k1[j]$ OR *no writemask*

THEN $DEST[i+31:i] \leftarrow SRC1[i+31:i] + SRC2[i+31:i]$

ELSE

IF *merging-masking* ; merging-masking

THEN * $DEST[i+31:i]$ remains unchanged*

ELSE ; zeroing-masking

$DEST[i+31:i] \leftarrow 0$

FI

FI;

ENDFOR;

Can we get the compiler to vectorise?

Compiler Explorer

C++ source #1

```
1 float c[1024];
2 float a[1024];
3 float b[1024];
4 void add ()
5 {
6     for (int i=0; i < 1024; i++)
7         c[i]=a[i]+b[i];
8 }
```

x86-64 gcc 5.4 (Editor #1, Compiler #1) C++

x86-64 gcc 5.4 -O3 -fopenmp

11010 .LX0: .text // \s+ Intel Demangle

Libraries + Add new...

```
1 _Z3addv:
2     xorl    %eax, %eax
3 .L2:
4     movaps  a(%rax), %xmm0
5     addq    $16, %rax
6     addps   b-16(%rax), %xmm0
7     movaps  %xmm0, c-16(%rax)
8     cmpq    $4096, %rax
9     jne     .L2
10    rep ret
11 b:
12     .zero   4096
13 a:
14     .zero   4096
15 c:
16     .zero   4096
```

Output (0/1) g++ (GCC-Explorer-Build) 5.4.0 - cached (4432)

In sufficiently simple cases, no problem:
Gcc reports:
test.c:6:3: note: loop vectorized

Secure | https://godbolt.org

Compiler Explorer

Editor

Diff View

More

Share

Other

C++ source #1

A

Save/Load

Add new

C++

```
1 float c[1024];
2 float a[1024];
3 float b[1024];
4 void add (int N)
5 {
6     for (int i=0; i < N; i++)
7         c[i]=a[i]+b[i];
8 }
```

11010

.LX0:

.text

//

\s+

Intel

Demangle

Libraries

Add new

```
1 .L3addi:
2     testl %edi, %edi
3     jle .L1
4     leal -4(%rdi), %edx
5     leal -1(%rdi), %ecx
6     shrl $2, %edx
7     addl $1, %edx
8     cmpl $2, %ecx
9     leal 0(%rdx,4), %eax
10    jbe .L9
11    xorl %ecx, %ecx
12    xorl %esi, %esi
13
14    .L5:
15        movaps a(%rcx), %xmm0
16        addl $1, %esi
17        addq $16, %rcx
18        addps b-16(%rcx), %xmm0
19        movaps %xmm0, c-16(%rcx)
20        cmpl %esi, %edx
21        ja .L5
22        cmpl %edi, %eax
23        je .L12
24
25    .L3:
26        movslq %eax, %rdx
27        movss b(%rdx,4), %xmm0
28        addss a(%rdx,4), %xmm0
29        movss %xmm0, c(%rdx,4)
30        leal 1(%rax), %edx
31        cmpl %edx, %edi
32        jle .L1
33        movslq %edx, %rdx
34        addl $2, %eax
35        movss a(%rdx,4), %xmm0
36        cmpl %eax, %edi
37        addss b(%rdx,4), %xmm0
38        movss %xmm0, c(%rdx,4)
39        jle .L1
40        cltq
41        movss a(%rax,4), %xmm0
42        addss b(%rax,4), %xmm0
43        movss %xmm0, c(%rax,4)
44        ret
45
46    .L1:
47        rep ret
48
49    .L12:
50        rep ret
51
52    .L9:
53        xorl %eax, %eax
54        jmp .L3
55
56 b:
57     .zero 4096
58 a:
59     .zero 4096
60 c:
61     .zero 4096
```

Output (0/3)

g++ (GCC-Explorer-Build) 5.4.0 - 377ms (5760B)

If the trip count is not known to be divisible by 4:

gcc reports:
test.c:6:3: note: loop vectorized
test.c:6:3: note: loop turned into non-loop; it never loops.
test.c:6:3: note: loop with 3 iterations completely unrolled

Basically the same vectorised code as before

Three copies of the non-vectorised loop body to mop up the additional iterations in case N is not divisible by 4

```

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

If the alignment of the operand pointers is not known:

gcc reports:

test.c:6:3: note: loop vectorized

test.c:6:3: note: loop peeled for vectorization to enhance alignment

test.c:6:3: note: loop turned into non-loop; it never loops.

test.c:6:3: note: loop with 3 iterations completely unrolled

test.c:1:6: note: loop turned into non-loop; it never loops.

test.c:1:6: note: loop with 4 iterations completely unrolled

```

1 123adP5i_5_1
2 testl %eax, %eax
3 pushq %r13
4 pushq %r12
5 pushq %r10
6 pushq %r9
7 .L1:
8 movs %r13, %rax
9 leal 1(%r13), %r10
10 andl $15, %eax
11 shrq $2, %rax
12 radd %rax
13 andl $3, %eax
14 cmpl %r10, %eax
15 cmovs %r10, %rax
16 cmpl $4, %r10
17 .L2:
18 movl %r10, %eax
19 .L3:
20 movss (%r13), %xmm0
21 cmpl $1, %eax
22 movl $1, %r10
23 addss (%r10), %xmm0
24 movss %xmm0, (%r13)
25 .L5:
26 movss 4(%r13), %xmm0
27 cmpl $2, %eax
28 movl $2, %r10
29 addss 4(%r10), %xmm0
30 movss %xmm0, 4(%r13)
31 .L5:
32 movss 8(%r13), %xmm0
33 cmpl $3, %eax
34 movl $3, %r10
35 addss 8(%r10), %xmm0
36 movss %xmm0, 8(%r13)
37 .L5:
38 movss 12(%r13), %xmm0
39 cmpl $4, %r10
40 addss 12(%r10), %xmm0
41 movss %xmm0, 12(%r13)
42 .L5:
43 cmpl %eax, %r10
44 .L1:
45 .L4:
46 subl %eax, %r10
47 movl %eax, %r10
48 movl %eax, %r11
49 leal -4(%r10), %r10
50 subl %eax, %r10
51 shrq $2, %r10
52 addl $1, %r10
53 cmpl $2, %r10
54 leal 0(%r10,4), %r10
55 .L7:
56 leaq 0(%r10,4), %rax
57 vorl %r10, %rax
58 leaq (%r13,%rax), %r13
59 leaq (%r10,%rax), %r12
60 leaq (%r13,%rax), %r11
61 vorl %r10, %rax
62 .L9:
63 movups (%r13,%rax), %xmm0
64 addl $1, %r10
65 addps 0(%r13,%rax), %xmm0
66 movups %xmm0, (%r13,%rax)
67 addl $15, %rax
68 cmpl %r10, %r10
69 .L9:
70 addl %r10, %r10
71 cmpl %r10, %r10
72 .L1:
73 .L7:
74 movsld %r10, %rax
75 movss (%r13,%rax,4), %xmm0
76 addss (%r10,%rax,4), %xmm0
77 movss %xmm0, (%r13,%rax,4)
78 leal 1(%r13), %r10
79 cmpl %r10, %r10
80 .L1:
81 addl $2, %r10
82 movss (%r13,%rax,4), %xmm0
83 cmpl %r10, %r10
84 addss (%r10,%rax,4), %xmm0
85 movss %xmm0, (%r13,%rax,4)
86 .L1:
87 movsld %r10, %r10
88 movss (%r13,%r10,4), %xmm0
89 addss (%r10,%r10,4), %xmm0
90 movss %xmm0, (%r13,%r10,4)
91 .L1:
92 .L1:
93 popq %r10
94 popq %r12
95 popq %r13
96 ret
97 .L12:
98 testl %eax, %eax
99 .L3:
100 vorl %r10, %r10
101 .L4:
    
```

Three copies of the non-vectorised loop body to align the start address of the vectorised code on a 32-byte boundary

Basically the same vectorised code as before

Three copies of the non-vectorised loop body to mop up the additional iterations in case N is not divisible by 4

Output (0/6) g++ (GCC-Explorer-Build) 5.4.0 - 578ms (6934B)

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

1 22addps5,5,1
2 testl %eax, %eax
3 jg .L27
4 leaq 16(%rax), %rax
5 leaq 16(%rax), %rax
6 leaq 16(%rax), %rax
7 cmov %rax, %rax
8 scasd %rax
9 cmov %rax, %rax
10 scasd %rax
11 ori %eax, %rax
12 leaq 16(%rax), %rax
13 cmov %rax, %rax
14 scasd %rax
15 cmov %rax, %rax
16 scasd %rax
17 ori %rax, %eax
18 testb %al, %al
19 je .L1
20 cmov %eax, %eax
21 jbe .L3
22 movl %rax, %rax
23 pushl %rax
24 pushl %rax
25 andl %eax, %eax
26 pushl %rax
27 pushl %rax
28 shrl %eax, %eax
29 negl %eax
30 andl %eax, %eax
31 cmov %rax, %eax
32 cmov %eax, %eax
33 xorl %rax, %rax
34 testl %eax, %eax
35 je .L4
36 movss (%rax), %xmm0
37 cmov %eax, %eax
38 movl %eax, %eax
39 addss (%rax), %xmm0
40 movss %xmm0, (%rax)
41 je .L4
42 movss 4(%rax), %xmm0
43 cmov %eax, %eax
44 movl %eax, %eax
45 addss 4(%rax), %xmm0
46 movss %xmm0, 4(%rax)
47 je .L4
48 movss 8(%rax), %xmm0
49 cmov %eax, %eax
50 addss 8(%rax), %xmm0
51 movss %xmm0, 8(%rax)
52 .L4:
53 subl %eax, %rax
54 salq %eax, %rax
55 xorl %rax, %rax
56 leaq -4(%rax), %rax
57 leaq (%rax,%rax), %rax
58 leaq (%rax,%rax), %rax
59 xorl %eax, %eax
60 addl %rax, %rax
61 shrl %eax, %rax
62 addl %eax, %rax
63 leaq 8(%rax,%rax), %rax
64 .L7:
65 movq (%rax,%rax), %xmm0
66 addl %eax, %eax
67 addq 8(%rax,%rax), %xmm0
68 movq %xmm0, (%rax,%rax)
69 addl %eax, %rax
70 cmov %rax, %eax
71 je .L7
72 addl %eax, %rax
73 cmov %eax, %rax
74 je .L1
75 movsld %rax, %rax
76 movss (%rax,%rax), %xmm0
77 addss (%rax,%rax), %xmm0
78 movss %xmm0, (%rax,%rax)
79 leal 16(%rax), %eax
80 cmov %eax, %eax
81 je .L1
82 cmsg
83 addl %eax, %rax
84 movss (%rax,%rax), %xmm0
85 cmov %eax, %eax
86 addss (%rax,%rax), %xmm0
87 movss %xmm0, (%rax,%rax)
88 .L1:
89 movsld %rax, %rax
90 movss (%rax,%rax), %xmm0
91 addss (%rax,%rax), %xmm0
92 movss %xmm0, (%rax,%rax)
93 .L1:
94 popq %rax
95 popq %rax
96 popq %rax
97 popq %rax
98 .L27:
99 rep ret
100 .L3:
101 xorl %eax, %eax
102 .L12:
103 movss (%rax,%rax), %xmm0
104 addss (%rax,%rax), %xmm0
105 movss %xmm0, (%rax,%rax)
106 addl %eax, %eax
107 cmov %eax, %eax
108 je .L12
109 rep ret

If the pointers might be aliases:

gcc reports:

- test.c:6:3: note: loop vectorized
- test.c:6:3: note: loop versioned for vectorization because of possible aliasing
- test.c:6:3: note: loop peeled for vectorization to enhance alignment
- test.c:6:3: note: loop turned into non-loop; it never loops.
- test.c:6:3: note: loop with 3 iterations completely unrolled
- test.c:1:6: note: loop turned into non-loop; it never loops.
- test.c:1:6: note: loop with 3 iterations completely unrolled

```
gcc reports:
test.c:6:3: note: loop vectorized
test.c:6:3: note: loop versioned for vectorization because of
possible aliasing
test.c:6:3: note: loop peeled for vectorization to enhance alignment
test.c:6:3: note: loop turned into non-loop; it never loops.
test.c:6:3: note: loop with 3 iterations completely unrolled
test.c:1:6: note: loop turned into non-loop; it never loops.
test.c:1:6: note: loop with 3 iterations completely unrolled
```

Non-vector version of the loop for the case when c might overlap with a or b

What to do if the compiler just won't vectorise your loop? Option #1: ivdep pragma

```
void add (float *c, float *a, float *b)
{
    #pragma ivdep
    for (int i=0; i <= N; i++)
        c[i]=a[i]+b[i];
}
```

IVDEP (Ignore Vector DEpendencies) compiler hint.

Tells compiler “Assume there are no loop-carried dependencies”

This tells the compiler vectorisation is *safe*: it might still not vectorise

What to do if the compiler just won't vectorise your loop? Option #2: **OpenMP 4.0 pragmas**

loopwise:

```
void add (float *c, float *a, float *b)
{
    #pragma omp simd
    for (int i=0; i <= N; i++)
        c[i]=a[i]+b[i];
}
```

Indicates that the loop can be transformed into a SIMD loop
(i.e. the loop can be executed concurrently using SIMD instructions)

functionwise:

```
#pragma omp declare simd
void add (float *c, float *a, float *b)
{
    *c=*a+*b;
}
```

"declare simd" can be applied to a function to enable
SIMD instructions at the function level from a SIMD loop

Tells compiler “vectorise this code”. It might still not do it...

What to do if the compiler just won't vectorise your loop? Option #2: SIMD intrinsics:

```
void add (float *c, float *a, float *b)
{
    __m128* pSrc1 = (__m128*) a;
    __m128* pSrc2 = (__m128*) b;
    __m128* pDest = (__m128*) c;
    for (int i=0; i <= N/4; i++)
        *pDest++ = _mm_add_ps(*pSrc1++, *pSrc2++);
}
```

Vector instruction lengths are hardcoded in the data types and intrinsics

This tells the compiler which specific vector instructions to generate. This time it really will vectorise!

What to do if the compiler just won't vectorise your loop? Option #3: SIMT

Basically... think of each lane as a thread

Or: vectorise an *outer* loop:

```
#pragma omp simd
For (int i=0; i<N; ++i) {
    if(){...} else {...}
    for (int j=....) {...}
    while(...) {...}
    f(...)
}
```

In the body of the vectorised loop, each lane executes a different iteration of the loop – *whatever* the loop body code does

More later – when we look at GPUs

Summary Vectorisation Solutions

1. Indirectly through **high-level libraries**/code generators
2. **Auto-vectorisation** (eg use “-O3 -mavx2 -fopt-info” and hope it vectorises):
 - code complexity, sequential languages and practices get in the way
 - Give your **compiler hints** and hope it vectorises:
 - C99 "restrict" (implied in FORTRAN since 1956)
 - #pragma ivdep
3. **Code explicitly**:
 - In assembly language
 - SIMD instruction intrinsics
 - OpenMP 4.0 #pragma omp simd
 - Kernel functions:
 - OpenMP 4.0: #pragma omp declare simd
 - OpenCL or CUDA: more later

History:

Intel x86 ISA extended with SIMD

		width	Int.	SP	DP
1997	MMX	64	✓		
1999	SSE	128	✓	✓(x4)	
2001	SSE2	128	✓	✓	✓(x2)
2004	SSE3	128	✓	✓	✓
2006	SSSE 3	128	✓	✓	✓
2006	SSE 4.1	128	✓	✓	✓
2008	SSE 4.2	128	✓	✓	✓
2011	AVX	256	✓	✓(x8)	✓(x4)
2013	AVX2	256	✓	✓	✓
<i>future</i>	AVX-512	512	✓	✓(x16)	✓(x8)

ATPESC 2014, James Reinders: <http://extremecomputingtraining.anl.gov/files/2014/08/20140804-1030-1115-ATPESC-Argonne-Reinders.2.pdf>

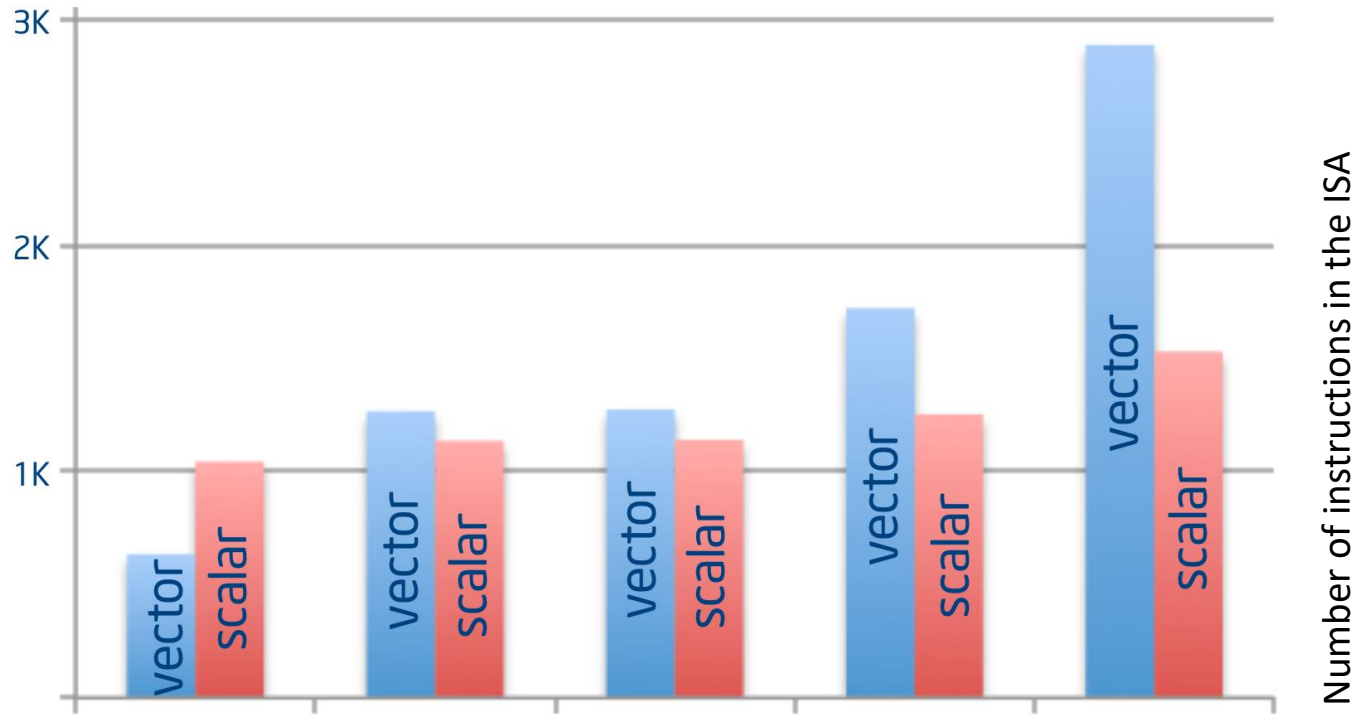
History:

Intel x86 ISA ex

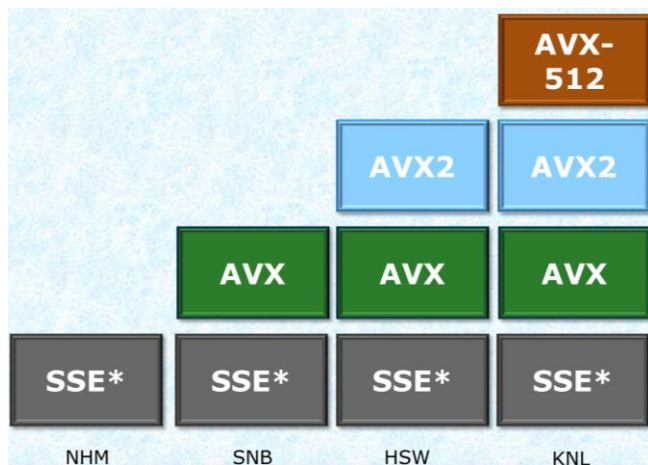
		width
1997	MMX	64
1999	SSE	128
2001	SSE2	128
2004	SSE3	128
2006	SSSE 3	128
2006	SSE 4.1	128
2008	SSE 4.2	128
2011	AVX	256
2013	AVX2	256
future	AVX-512	512

- Wider registers
(from 32 to 512 bits)
- More registers
- Richer instruction set
(predication, FMAs, gather, scatter, ...)
- Easier exploitation
(better compiler support, high-level functions, libraries...)

Growth in vector instructions on Intel



Backwards compatibility accumulation



ATPESC 2014, James Reinders:

<http://extremecomputingtraining.anl.gov/files/2014/08/20140804-1030-1115-ATPESC-Argonne-Reinders.2.pdf>

Elena Demikhovskiy (Intel): <http://lvm.org/devmtg/2013-11/slides/Demikhovskiy-Poster.pdf>

Issues inherent in the vector model

Example 1

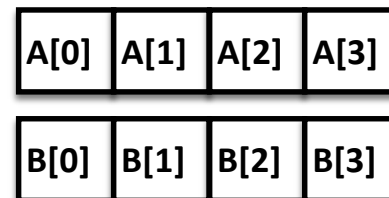
```
double A[N], B[N], C[N]
for i = 0 to N, i++
    C[i] = sqrt(A[i] + B[i])
```

SIMD version

```
loop: VLOAD av, A[i:v]
      VLOAD bv, B[i:v]
      VADD cv, bv, av
      VSQRT cv, cv
      VSTORE C[i:v], cv
      INCR i
      IF i < N/v: loop
```

Notation:

- $:v$ indicates that the assembly operation is over v elements
- subscript v indicates that the register is actually a vector register, hosting v elements



E.g. $v=4$

Simple issues: bad array size

```
loop: VLOAD av, A[i:v]  
      VLOAD bv, B[i:v]  
      VADD cv, bv, av  
      VSQRT cv, cv  
      VSTORE C[i:v], cv  
      INCR i  
      IF i<N/v: loop
```

Issue 1: N might not be a multiple of the vector length v

or

N is known only at runtime

Simple issues: bad array size

```
loop: VLOAD av, A[i:v]  
      VLOAD bv, B[i:v]  
      VADD cv, bv, av  
      VSQRT cv, cv  
      VSTORE C[i:v], cv  
      INCR i  
      IF i < N/v: loop  
      IF N%v == 0: exit  
peel: LOAD a, A[v*i + 0]  
      ...  
      ...  
exit: ...
```

Issue 1: N might not be a multiple of the vector length v

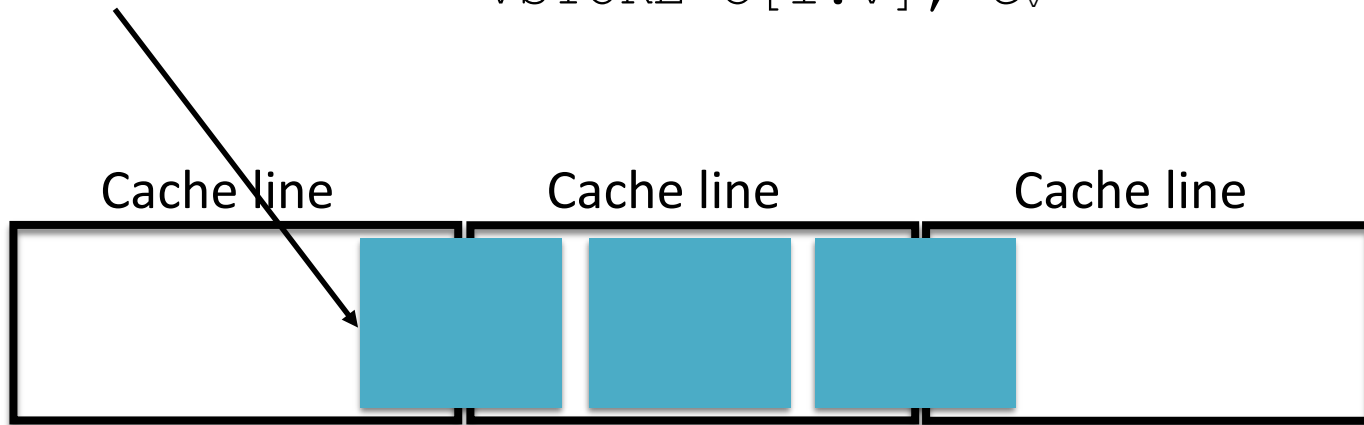
or

- Or not known at compile-time

Medium issues: data alignment

```
loop: VLOAD av, A[i:v]  
      VLOAD bv, B[i:v]  
      ...  
      VSTORE C[i:v], cv
```

Base address
of array A



E.g.: AVX on Sandy Bridge: Cache line: 64B, vector length: 32B, double: 8B

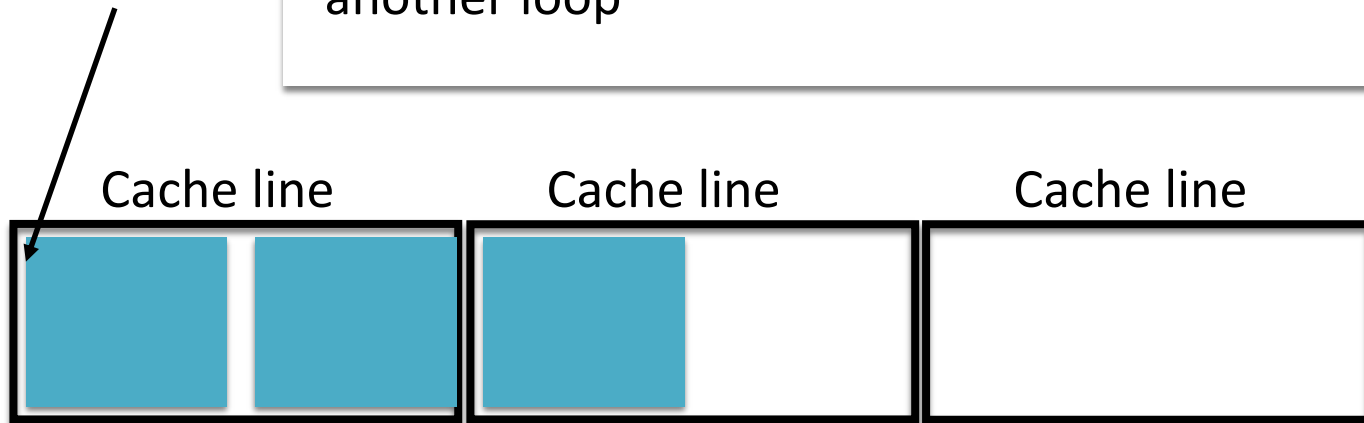
Issue 2: Memory accesses should be aligned to page and cache boundaries

Medium issues: data alignment

Solution: change the allocation point of A

- Use of special mallocs or special array qualifiers
- **Global transformation:** might affect alignment in another loop

Base address
of array A



E.g.: AVX on Sandy Bridge: Cache line: 64B, vector length: 32B, double: 8B

Issue 2: Memory accesses should be aligned to page
and cache boundaries
(tricky with stencils, for $i \in \{B[i]=A[i-1]+A[i]+A[i+1]\}$)

Advanced issues: bad access patterns

Example 2

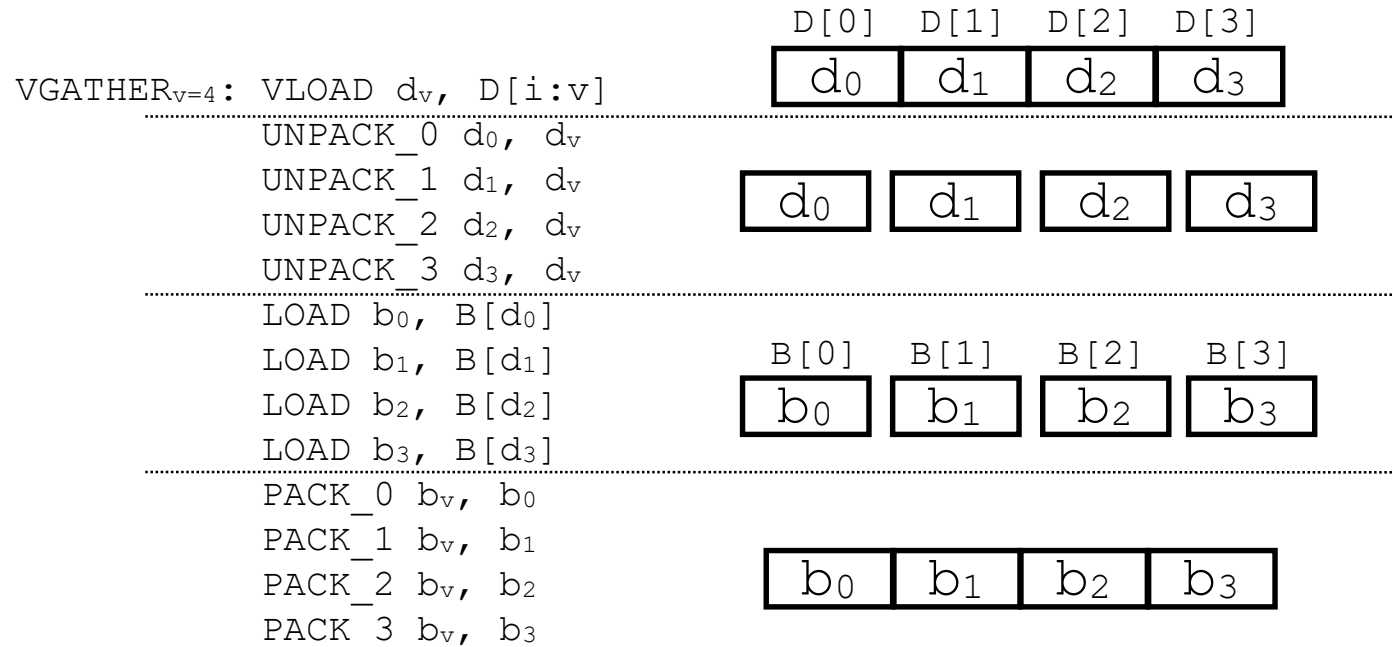
```
double A[N], B[N], C[N], D[N]
for i = 0 to N, i++
    C[i] = A[2*i] + B[D[i]]
```

SIMD version

```
loop: VLOAD av, A[i], stride=2
      VGATHER bv, B, D[i:v]
      VADD cv, bv, av
      VSTORE C[i:v], cv
incr: INCR i
      IF i<N/v: loop
```

Advanced issues: bad access patterns

$B[D[i]] \implies \text{VGATHER } b_v, B, D[i:v]$



Issue 3: regardless of the ISA the (micro-)interpretation of these instructions is expensive

Each $B[d_i]$ might fall on a different cache line, perhaps even a different page – vgather execution time depends on whether “coalescing” occurs

Advanced issues: branch divergence

```
for i = 0 to 63, i++  
    if A[i] > 0  
        B[i] = A[i] * 4
```

Solution: **Predication through masking**

.....

loop:

```
VLOAD      av, A[i:v]  
VCMP_P     Rmask, av, R0  
VMUL_P     bv{Rmask}, av, R4  
VSTORE_P   B[i:v]{Rmask}, bv  
VRESET_P   Rmask  
INCR       Ri  
CMP        Ri < 64/v: loop
```

Add a new boolean vector register (the vector mask register)

- Operates on elements whose corresponding bit in the mask is 1
- Requires ISA extension to set the mask register

Vector execution alternatives

Implementation may execute n-wide vector operation with an n-wide ALU
– or maybe in smaller, m-wide blocks

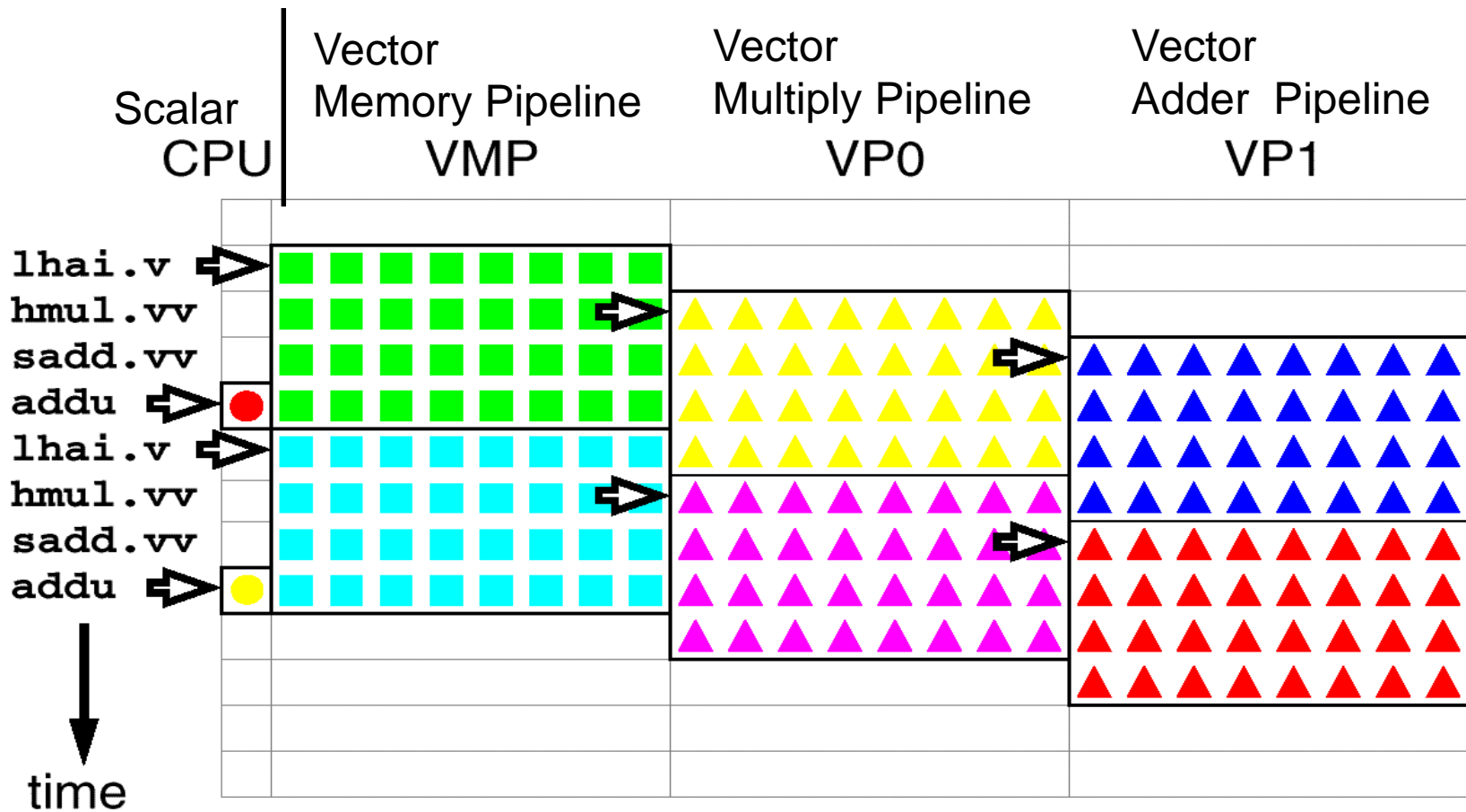
vector pipelining:

- Consider a simple static pipeline
- Vector instructions are executed serially, element-by-element, using a pipelined FU
- We have several pipelined FUs
- “vector chaining” – each word is forwarded to the next instruction as soon as it is available
- FUs form a long pipelined chain

uop decomposition:

- Consider a dynamically-scheduled o-o-o machine
- Each n-wide vector instruction is split into m-wide uops at decode time
- The dynamic scheduling execution engine schedules their execution, possibly across multiple FUs
- They are committed together

Vector pipelining – “chaining”



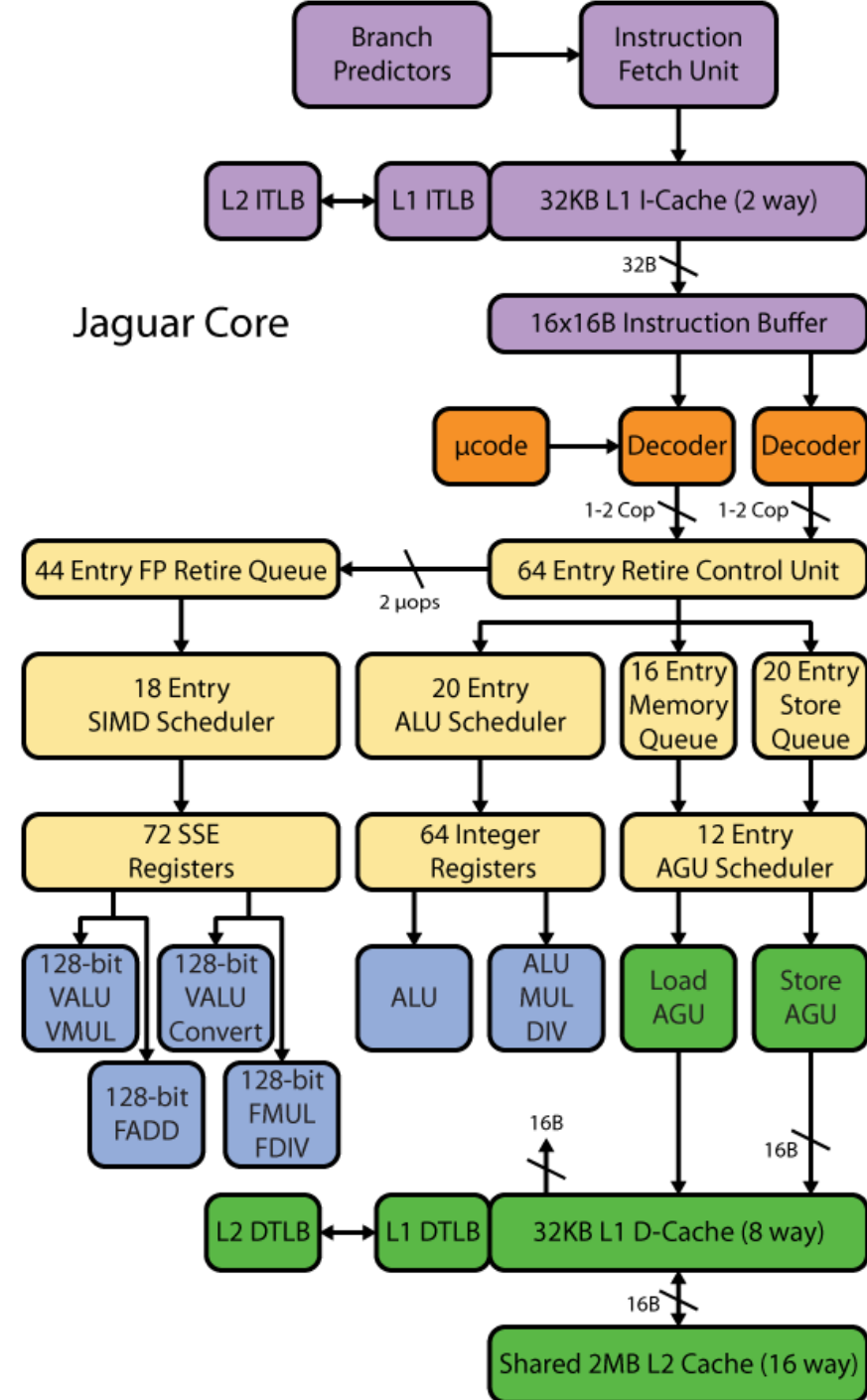
● ■ ▲ Operations
➡ Instruction issue

- Vector FUs are 8-wide - each 32-wide vector instruction is executed in 4 blocks
- Forwarding is implemented block-by-block
- So memory, mul, add and store are chained together into one continuously-active pipeline

Uop decomposition - example

AMD Jaguar

- Low-power 2-issue dynamically-scheduled processor core
- Supports AVX-256 ISA
- Has two 128-bit vector ALUs
- 256-bit AVX instructions are split into two 128-bit uops, which are scheduled independently
- Until retirement
- A “zero-bit” in the rename table marks a register which is known to be zero
- So no physical register is allocated and no redundant computation is done



SIMD Architectures: discussion

- **Reduced Turing Tax: more work, fewer instructions**
- **Relies on compiler or programmer**
- **Simple loops are fine, but many issues can make it hard**
- **“lane-by-lane” predication allows conditionals to be vectorised, but branch divergence may lead to poor utilisation**
- **Indirections can be vectorised on some machines (vgather, vscatter) but remain hard to implement efficiently unless accesses happen to fall on a small number of distinct cache lines**
- **Vector ISA allows broad spectrum of microarchitectural implementation choices**
- **Intel’s vector ISA has grown enormous as vector length has been successively increased**
- **ARM’s “scalable vector extension” (SVE) is an ISA design that hides the vector length (by using a special loop branch)**

Topics we have not had time to cover

ARM's SVE:

- ➡ a vector ISA that achieves binary compatibility across machines with different vector width uop decomposition

Matrix registers and matrix instructions

- ➡ Eg Nvidia's "tensor cores"

Pipelined vector architectures:

- ➡ The classical vector supercomputer

Whole-function vectorisation, ISPC, SIMT

- ➡ Vectorising nested conditionals
- ➡ Vectorising non-innermost loops
- ➡ Vectorising loops containing while loops

SIMT and the relationship/similarities with GPUs

- ➡ Coming!