

arm

SVE Deep Dive



arm

Arm C Language Extensions for SVE

Compiler Intrinsics for SVE

Arm C Language Extensions

Intrinsics and other features for supporting Arm features in C and C++

- ACLE extends C/C++ with Arm-specific features
 - Predefined macros: `__ARM_ARCH_ISA_A64`, `__ARM_BIG_ENDIAN`, etc.
 - Intrinsic functions: `__clz(uint32_t x)`, `__cls(uint32_t x)`, etc.
 - Data types: SVE, NEON and FP16 data types
- **ACLE for SVE** enables VLA programming with ACLE
 - *Nearly* one intrinsic per SVE instruction
 - Data types to represent the size-less vectors used for SVE intrinsics
- Intended for users that...
 - Want to hand-tune SVE code
 - Want to adapt or hand-optimize applications and libraries
 - Need low-level access to Arm targets

How to use ACLE

- Include the headers you need
 - `arm_acle.h` → for core ACLE
 - `arm_fp16.h` → to add scalar FP16 arithmetic
 - `arm_neon.h` → to add NEON intrinsics and data types
 - `arm_sve.h` → to add SVE intrinsics and data types
- Each of those require certain features at the compilation target
 - `arm_fp16.h` → Your target platform needs to support FP16 (-march=armv8-a+fp16)
 - `arm_neon.h` → Your target platform need to support NEON (-march=armv8-a+simd)
 - `arm_sve.h` → Your target platform need to support SVE (-march=armv8-a+sve)

SVE ACLE

SVE Arm C Language Extensions – aka *C intrinsics*

```
#include <arm_sve.h>
```

- VLA Data types:
 - `svfloat64_t`, `svfloat16_t`,
`svuint32_t`,...
- Predication:
 - Merging: `_m`
 - Zeroing: `_z`
 - Don't care: `_x`
 - Predicate type: `svbool_t`
- Use *C11 generics* for function overloading.
- Intrinsics are **not 1-1 with the ISA**.

Examples

```
svfloat32_t  
svadd[_n_f32]_z (svbool_t pg,  
                 svfloat32_t op1,  
                 float32_t op2);
```

```
svfloat16_t  
svsqrt_m (svfloat16_t inactive,  
          svbool_t pg,  
          svfloat16_t op)
```

Vectorizing a scalar loop with ACLE

$a[:] = 2.0 * a[:]$

Original Code

```
for (int i=0; i < N; ++i) {  
  
    a[i] = 2.0 * a[i];  
  
}
```

128-bit NEON vectorization with ACLE

```
int i;  
  
// vector loop  
for (i=0; (i<N-3) && (N&~3); i+=4) {  
    float32x4_t va = vld1q_f32(&a[i]);  
    va = vmulq_n_f32(va, 2.0);  
    vst1q_f32(&a[i], va)  
}  
  
// drain loop  
for (; i < N; ++i)  
    a[i] = 2.0 * a[i];
```

This is NEON,
not SVE!

Vectorizing a scalar loop with ACLE

```
a[:] = 2.0 * a[:]
```

```
for (int i=0; i < N; ++i) {  
    a[i] = 2.0 * a[i];  
}
```

SVE vectorization

```
for (int i = 0 ; i < N; i += ????????)          svcntw()  
{  
    svbool_t Pg = svwhilelt_b32(i, N);  
    svfloat32_t va = svld1(Pg, &a[i]);  
  
    va = svmul_x(Pg, va, 2.0);  
  
    svst1(Pg, &a[i], va);  
}
```

128-bit NEON vectorization

```
int i;  
  
// vector loop  
for (i=0; (i<N-3) && (N&~3); i+=4) {  
    float32x4_t va = vld1q_f32(&a[i]);  
    va = vmulq_n_f32(va, 2.0);  
    vst1q_f32(&a[i], va)  
}  
  
// drain loop  
for (; i < N; ++i)  
    a[i] = 2.0 * a[i];
```

Vectorizing a scalar loop with ACLE

`a[:] = 2.0 * a[:]`

```
for (int i=0; i < N; ++i) {  
    a[i] = 2.0 * a[i];  
}
```

SVE vectorization

```
for (int i = 0 ; i < N; i += ????????)          svcntw()  
{  
    svbool_t Pg = svwhilelt_b32(i, N);  
    svfloat32_t va = svld1(Pg, &a[i]);  
    va = svmul_x(Pg, va, 2.0);  
    svst1(Pg, &a[i], va);  
}
```

SVE vectorization with fewer branches

```
svbool_t all = svptrue_b32();  
svbool_t Pg;  
for (int i=0;  
    svptest_first(all,  
        Pg=svwhilelt_b32(i, N));  
    i += svcntw())  
{  
    svfloat32_t va = svld1(Pg, &a[i]);  
    va = svmul_x(Pg, va, 2.0);  
    svst1(Pg, &a[i], va);  
}
```


ACLE for SVE Cheat Sheet

- Vector types
 - `sv<datatype><datasize>_t`
 - `svfloat32_t`
 - `svint8_t`
 - `svuint16_t`
- Predicate types
 - `svbool_t`
- Functions
 - `svbase[disambiguator][type0][type1]...[predication]`
 - base is the lower-case name of an SVE instruction
 - *disambiguator* distinguishes between different forms of a function
 - *typeN* lists the types of vectors and predicates
 - *predication* describes the inactive elements in the result of a predicated operation

```
svfloat64_t svld1_f64(svbool_t pg, const float64_t *base)
svbool_t svwhilelt_b8(int64_t op1, int64_t op2)
svuint32_t svmla_u32_z(svbool_t pg, svuint32_t op1, svuint32_t op2, svuint32_t op3)
svuint32_t svmla_u32_m(svbool_t pg, svuint32_t op1, svuint32_t op2, svuint32_t op3)
```

02_ACLE/01_vecadd

See README.md for details

GCC 9.3

```
gcc --version
```

```
gcc (GCC) 9.3.0
```

```
gcc -fopt-info-all-vec -Ofast \  
-mcpu=native vec_add_acle.c
```

```
vec_add_acle.c:22:10: fatal error: arm_sve.h: No such file or  
directory
```

```
22 | #include <arm_sve.h>
```

```
| ^~~~~~
```

```
compilation terminated.
```

GCC 11

```
gcc --version
```

```
gcc (GCC) 11.0.0 20201025 (experimental)
```

```
gcc -fopt-info-all-vec -Ofast \  
-mcpu=native vec_add.c
```

```
# No errors
```

02_ACLE/01_vecadd

See README.md for details

vec_add_acle_arm.exe

```
whilelo p0.s, x8, x9
ld1w { z0.s }, p0/z, [x1, x8, lsl #2]
ld1w { z1.s }, p0/z, [x2, x8, lsl #2]
fadd z0.s, p0/m, z0.s, z1.s
st1w { z0.s }, p0, [x0, x8, lsl #2]
incw x8
cmp x8, #256, lsl #12
b.lo #-28 <vec_svadd_m+0x20>
ret
```

vec_add_arm.exe

```
whilelo p0.s, xzr, x9
b #24 <vec_add+0x60>
... (six nop instructions)
ld1w { z0.s }, p0/z, [x1, x8, lsl #2]
ld1w { z1.s }, p0/z, [x2, x8, lsl #2]
fadd z0.s, z1.s, z0.s
st1w { z0.s }, p0, [x0, x8, lsl #2]
incw x8
whilelo p0.s, x8, x9
b.mi #-24 <vec_add+0x60>
b #64 <vec_add+0xbc>
mov x8, xzr
b #28 <vec_add+0xa0>
```

arm

Working with SVE Instructions

SAXPY

```
subroutine saxpy(x,y,a,n)
real*4 x(n),y(n),a
do i = 1,n
    y(i) = a*x(i) + y(i)
enddo
```

Scalar [-march=armv8-a]

```
// x0 = &x[0], x1 = &y[0], x2 = &a, x3 = &n
saxpy_:
    ldrsw    x3, [x3]           // x3=*n
    mov     x4, #0             // x4=i=0
    ldr     s0, [x2]           // d0=*a
b .latch
.loop:
    ldr     s1, [x0,x4,ls1 2]   // s1=x[i]
    ldr     s2, [x1,x4,ls1 2]   // s2=y[i]
    fmaddd s2, s1, s0, s2       // s2+=x[i]*a
    str     s2, [x1,x4,ls1 2]   // y[i]=s2
add     x4, x4, #1           // i+=1
.latch:
cmp     x4, x3               // i < n
b.lt    .loop                // more to do?
ret
```

Key Operations

- whilelt constructs a **predicate** (p0) to dynamically map vector operations to vector data
- incw increments a scalar register (x4) by the number of float elements that fit in a vector register
- No drain loop! Predication handles the remainder

SVE [-march=armv8-a+sve]

```
// x0 = &x[0], x1 = &y[0], x2 = &a, x3 = &n
saxpy_:
    ldrsw    x3, [x3]           // x3=*n
    mov     x4, #0             // x4=i=0
whilelt p0.s, x4, x3         // p0=while(i++<n)
    ld1rw   z0.s, p0/z, [x2]   // p0:z0=bcast(*a)
.loop:
    ld1w    z1.s, p0/z, [x0,x4,ls1 2] // p0:z1=x[i]
    ld1w    z2.s, p0/z, [x1,x4,ls1 2] // p0:z2=y[i]
    fmla    z2.s, p0/m, z1.s, z0.s // p0?z2+=x[i]*a
    st1w    z2.s, p0, [x1,x4,ls1 2] // p0?y[i]=z2
incw     x4                   // i+=(VL/32)
.latch:
whilelt p0.s, x4, x3         // p0=while(i++<n)
b.first .loop                // more to do?
ret
```

How do you count by vector width?

No need for multi-versioning: one increment for all vector sizes

```
ld1w z1.s, p0/z, [x0,x4,ls1 2] // p0:z1=x[i]
ld1w z2.s, p0/z, [x1,x4,ls1 2] // p0:z2=y[i]
fm1a z2.s, p0/m, z1.s, z0.s // p0?z2+=x[i]*a
st1w z2.s, p0, [x1,x4,ls1 2] // p0?y[i]=z2
incw x4 // i+=(VL/32)
```

“Increment `x4` by the number of 32-bit lanes (`w`) that fit in a VL.”

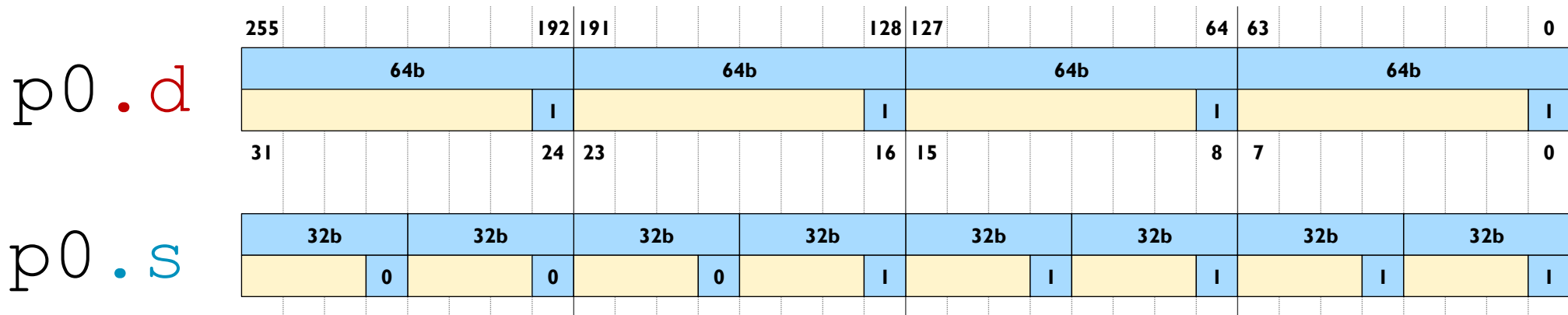
VLA Increment and Count

<code>incb x0, mul3, mul #2</code>	<code>x0 += 2 x (largest mul3 <= VL.b)</code>
<code>incd z0.d, pow2</code>	<code>each lane += largest pow2 <= VL.d</code>
<code>incp x0, p0.s</code>	<code>x0 += # active lanes</code>
<code>incp z0.h, p0</code>	<code>each vector lane += # active lanes of p</code>
<code>cntw x0</code>	<code>x0 = VL.s</code>
<code>cntp x0, p0, p1.s.</code>	<code>x0 = # active lanes of (p0 && p1.s)</code>

Predicates: Active Lanes vs Inactive Lanes

Predicate registers track lane activity

- 16 predicate registers (P0-P15)
- 1 predicate bit per 8 vector bits (lowest predicate bit per lane is significant)
- On **load**, active elements update the destination
- On **store**, inactive lanes leave destination unchanged (p0/**m**) or set to 0's (p0/**z**)



Predicate Condition Flags

SVE is a *predicate-centric* architecture

- Predicates support complex nested conditions and loops.
- Predicate generation also sets condition flags.
- Reduces vector loop management overhead.

Overloading the A64 NZCV condition flags

Condition Test	A64 Name	SVE Alias	SVE Interpretation
Z=1	EQ	NONE	No active elements are true
Z=0	NE	ANY	Any active element is true
C=1	CS	NLAST	Last active element is not true
C=0	CC	LAST	Last active element is true
N=1	MI	FIRST	First active element is true
N=0	PL	NFRST	First active element is not true
C=1 & Z=0	HI	PMORE	More partitions: some active elements are true but not the last one
C=0 Z=1	LS	PLAST	Last partition: last active element is true or none are true
N=V	GE	TCONT	Continue scalar loop
N!=V	LT	TSTOP	Stop scalar loop

Initialization when vector length is unknown

- Vectors cannot be initialized from compile-time constant, so...
 - INDEX `Zd.S, #1, #4` : `Zd = [1, 5, 9, 13, 17, 21, 25, 29]`
- Predicates cannot be initialized from memory, so...
 - PTRUE `Pd.S, MUL3` : `Pd = [(T, T, T), (T, T, T), F, F]`
- Vector loop increment and trip count are unknown at compile-time, so...
 - INCD `Xi` : increment scalar `Xi` by # of 64b dwords in vector
 - WHILELT `Pd.D, Xi, Xe` : next iteration predicate `Pd = [while i++ < e]`
- Vectors stores to stack must be dynamically allocated and indexed, so...
 - ADDVL `SP, SP, #-4` : decrement stack pointer by (`4*VL`)
 - STR `Zi, [SP, #3, MUL VL]` : store vector `Z1` to address (`SP+3*VL`)

arm

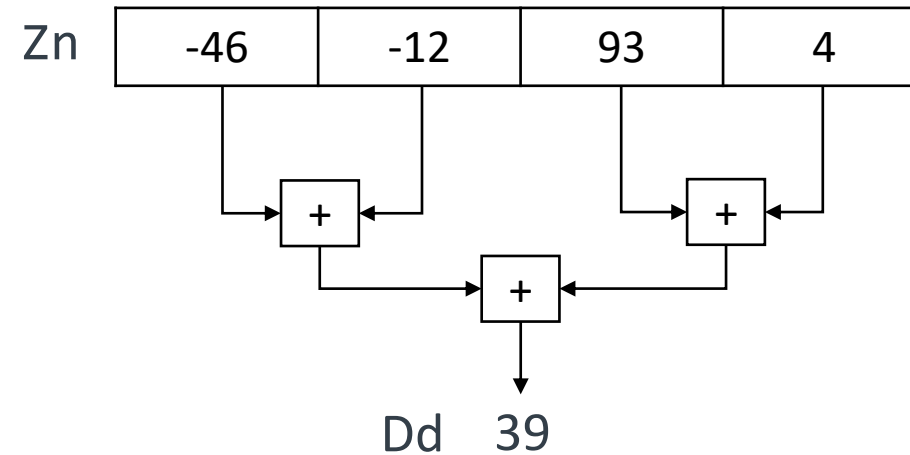
Horizontal
Reductions

SVE includes a rich set of horizontal operations

The operation happens across the lanes of a vector

- Addition, maximum, minimum, bitwise AND, OR, XOR ...

SADDV <Dd>, <Pg>, <Zn.T>



C code without intrinsics

Compile with `-march=armv8-a+sve` or `-mcpu=native` if your CPU has SVE

```
double ddot (double *a, double *b, int n)
{
    double sum = 0.0;
    for ( int i = 0; i < n; i++ ) {
        sum += a[i] * b[i];
    }
    return sum;
}
```

- Accumulates in a scalar
- No horizontal reductions

```
cmp        w2, #1
b.lt      #52 <ddot+0x38>
mov       w9, w2
mov       x8, xzr
whilelo   p0.d, xzr, x9
fmov     d0, xzr
ld1d     { z1.d }, p0/z, [x0, x8, lsl #3]
ld1d     { z2.d }, p0/z, [x1, x8, lsl #3]
incd     x8
fmul     z1.d, z1.d, z2.d
fadda    d0, p0, d0, z1.d
whilelo   p0.d, x8, x9
b.mi     #-24 <ddot+0x18>
ret
fmov     d0, xzr
ret
```

C code with SVE intrinsics

Vector accumulator; horizontal reduction

```
#include <arm_sve.h>

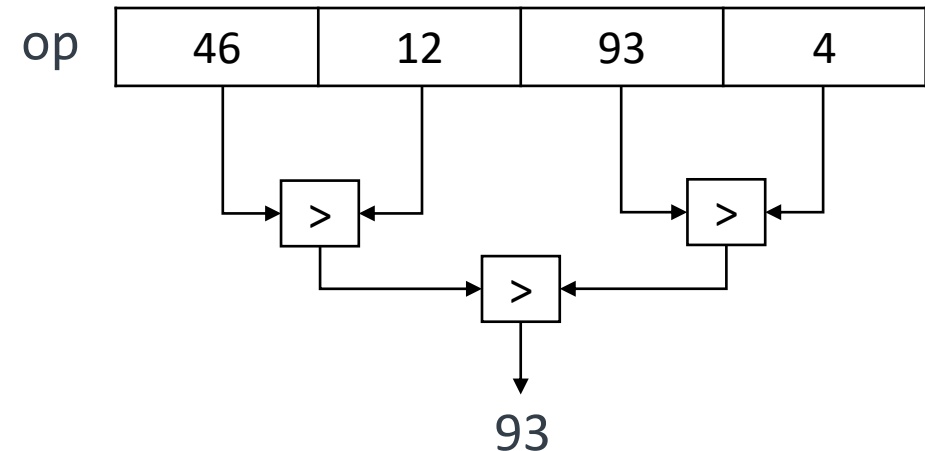
double ddot (double *a, double *b, int n) {
    svfloat64_t svsum = svdup_f64(0.0);
    svbool_t pg;
    svfloat64_t sva, svb, svsum;
    for (int i = 0; i < n; i += svcntd()) {
        pg = svwhilelt_b64(i, n);
        sva = svld1_f64(pg, &a[i]);
        svb = svld1_f64(pg, &b[i]);
        svsum = svmla_f64_m(pg, svsum, sva, svb);
    }
    return svaddv_f64(svptrue_b64(), svsum);
}
```

```
cmp        w2, #1
b.lt      #60 <ddot+0x40>
mov       w8, wzr
cntd     x9
mov      z0.d, #0
whilelt  p0.d, w8, w2
sxtw     x10, w8
ld1d    { z1.d }, p0/z, [x0, x10, lsl #3]
ld1d    { z2.d }, p0/z, [x1, x10, lsl #3]
add      w8, w8, w9
cmp      w8, w2
fmla    z0.d, p0/m, z1.d, z2.d
b.lt    #-28 <ddot+0x14>
ptrue   p0.d
faddv   d0, p0, z0.d
ret
mov     z0.d, #0
ptrue  p0.d
faddv  d0, p0, z0.d
ret
```

02_ACLE/02_hreduce

- Most vector instructions operate on a lane-by-lane basis
- SVE includes a rich set of horizontal operations where the operation happens across the lanes of a vector
- Examples of such operations are addition, maximum, minimum, bitwise AND, OR, XOR

```
suint32_t svmaxv(svbool_t pg, svuin32_t op)
```



02_ACLE/02_hreduce

See README.md for details

Reduce to scalar on every iteration

```
uint32_t max_scalar(uint32_t array[SIZE])
{
    uint32_t max = 0.0;

    uint32_t vl = svcntw();
    svbool_t p32_all = svptrue_b32();

    for (int i=0; i<SIZE; i+=vl) {
        // load array elements from memory
        svuint32_t vec = svld1(p32_all, &array[i]);
        // get max within vector
        uint32_t vecmax = svmaxv(p32_all, vec);
        if (vecmax > max) {
            max = vecmax;
        }
    }

    return max;
}
```

Reduce to vector; final reduction to scalar

```
uint64_t max(uint64_t array[SIZE])
{
    // initialize max vector with zeros
    svuint64_t max = svdup_u64(0);

    // get number of 64-bit elements in vector
    uint64_t vl = svcntd();

    // all true mask - assuming no partial vector mask needed for simplicity
    svbool_t p64_all = svptrue_b64();

    for (int i=0; i<SIZE; i+=vl) {
        // load array elements from memory
        svuint64_t vec = svld1(p64_all, &array[i]);
        // get max between loaded vector to max
        max = svmax_m(p64_all, max, vec);
    }

    // return max across values within the vector
    return svmaxv(p64_all, max);
}
```


arm

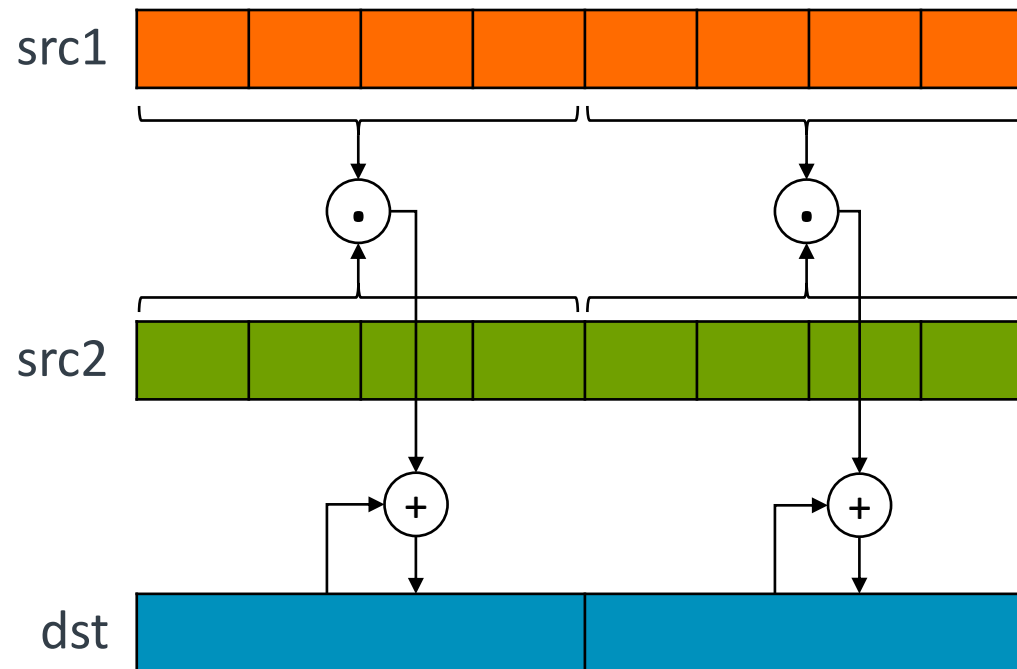
Low-precision
Dot Product
with Widening

02_ACLE/03_dotprod

- Dot product for low-precision value with widening

```
svint32_t svdot(svint8_t src1, svint8_t src2)
```

```
svint64_t svdot(svint16_t src1, svint16_t src2)
```



02_ACLE/03_dotprod

dotprod_acle_arm.exe

```
000000000400720 dot_product:
400720: e8 03 1f aa      mov     x8, xzr
400724: 00 c0 b8 25      mov     z0.s, #0
400728: e0 e3 18 25      ptrue  p0.b
40072c: 05 00 00 14      b      #20 <dot_product+0x20>
400730: 1f 20 03 d5      nop
400734: 1f 20 03 d5      nop
400738: 1f 20 03 d5      nop
40073c: 1f 20 03 d5      nop
400740: 01 40 08 a4      ld1b   { z1.b }, p0/z, [x0, x8]
400744: 22 40 08 a4      ld1b   { z2.b }, p0/z, [x1, x8]
400748: 28 50 28 04      addvl  x8, x8, #1
40074c: 1f 01 10 71      cmp    w8, #1024
400750: 20 04 82 44      udot   z0.s, z1.b, z2.b
400754: 63 ff ff 54      b.lo  #-20 <dot_product+0x20>
400758: e0 e3 98 25      ptrue  p0.s
40075c: 00 20 81 04      uaddv  d0, p0, z0.s
400760: 00 00 26 1e      fmov   w0, s0
400764: c0 03 5f d6      ret
400768: 06 00 00 14      b      #24 <main>
40076c: 1f 20 03 d5      nop
400770: 1f 20 03 d5      nop
400774: 1f 20 03 d5      nop
400778: 1f 20 03 d5      nop
40077c: 1f 20 03 d5      nop
```

dotprod_arm.exe

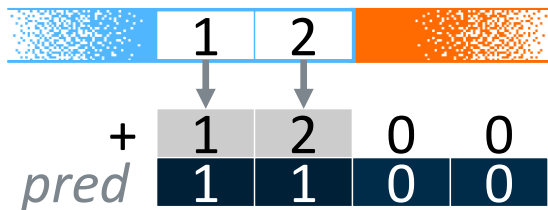
```
0000000004006c0 dot_product:
4006c0: 09 80 80 52      mov     w9, #1024
4006c4: e8 03 1f aa      mov     x8, xzr
4006c8: e0 1f 29 25      whilelo p0.b, xzr, x9
4006cc: 00 c0 b8 25      mov     z0.s, #0
4006d0: 04 00 00 14      b      #16 <dot_product+0x20>
4006d4: 1f 20 03 d5      nop
4006d8: 1f 20 03 d5      nop
4006dc: 1f 20 03 d5      nop
4006e0: 01 40 08 a4      ld1b   { z1.b }, p0/z, [x0, x8]
4006e4: 22 40 08 a4      ld1b   { z2.b }, p0/z, [x1, x8]
4006e8: 28 50 28 04      addvl  x8, x8, #1
4006ec: 00 1d 29 25      whilelo p0.b, x8, x9
4006f0: 40 04 81 44      udot   z0.s, z2.b, z1.b
4006f4: 64 ff ff 54      b.mi  #-20 <dot_product+0x20>
4006f8: e0 e3 98 25      ptrue  p0.s
4006fc: 00 20 81 04      uaddv  d0, p0, z0.s
400700: 00 00 66 9e      fmov   x0, d0
400704: c0 03 5f d6      ret
400708: 06 00 00 14      b      #24 <main>
40070c: 1f 20 03 d5      nop
400710: 1f 20 03 d5      nop
400714: 1f 20 03 d5      nop
400718: 1f 20 03 d5      nop
40071c: 1f 20 03 d5      nop
```

arm

Vector Partition with the First-faulting Register (FFR)

Vector Partitioning: when a vector spans a protected region

With VLA, we don't always know what data we may touch



- **Software-managed speculative vectorisation**
 - Create sub-vectors (partitions) in response to data and dynamic faults
- **First-fault load allows access to safely cross a page boundary**
 - First element is mandatory but others are a “speculative prefetch”
 - Dedicated FFR predicate register indicates successfully loaded elements
- **Allows uncounted loops with break conditions**
 - Load data using first-fault load
 - Create a *before-fault* partition from FFR
 - Test for break condition
 - Create a *before-break* partition from condition predicate
 - Process data within partition
 - Exit loop if break condition was found.

Vectorizing strlen

A vector load could result in a segfault if the vector spans protected memory.

Source Code

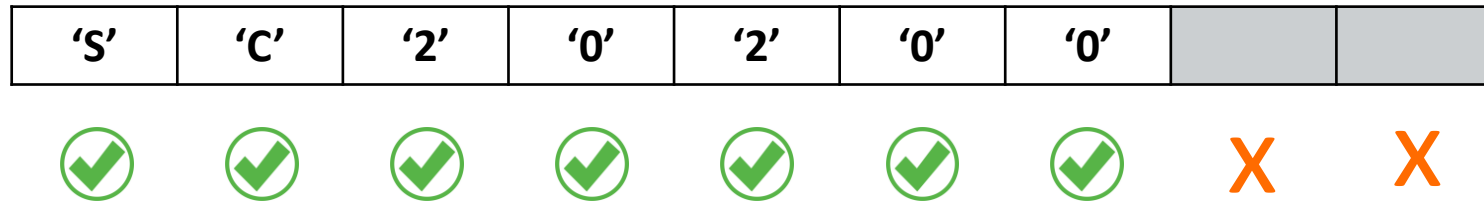
```
int strlen(const char *s) {  
    const char *e = s;  
    while (*e) e++;  
    return e - s;  
}
```

Scalar [-march=armv8-a]

```
// x0 = s  
strlen:  
    mov     x1, x0                // e=s  
.loop:  
    ldrb   x2, [x1],#1           // x2=*e++  
    cbnz  x2, .loop             // while(*e)  
.done:  
    sub   x0, x1, x0 // e-s  
    sub   x0, x0, #1 // return e-s-1  
    ret
```

Fault-tolerant Speculative Vectorization

- Some loops have dynamic exit conditions that prevent vectorization
 - E.g. the loop breaks on a particular value of the traversed array



- The access to unallocated space **does not trap** if it is not the first element
 - Faulting elements are stored in the first-fault register (FFR)
 - Subsequent instructions are predicated using the FFR information to operate only on successful element accesses

Vectorizing strlen

Partitioning off protected memory

SVE [-march=armv8-a+sve]

```
strlen:
    mov     x1, x0                // e=s
    ptrue  p0.b                  // p0=true
.loop:
    setffr                                // ffr=true
    ldff1b z0.b, p0/z, [x1]         // p0:z0=ldff(e)
    rdffr  p1.b, p0/z              // p0:p1=ffr
    cmpeq  p2.b, p1/z, z0.b, #0    // p1:p2>(*e==0)
    brkbs  p2.b, p1/z, p2.b        // p1:p2=until(*e==0)
    incp  x1, p2.b                 // e+=popcnt(p2)
    b.last .loop                  // last active=>!break
    sub    x0, x1, x0              // return e-s
    ret
```

```
int strlen(const char *s) {
    const char *e = s;
    while (*e) e++;
    return e - s;
}
```

Scalar [-march=armv8-a]

```
// x0 = s
strlen:
    mov     x1, x0
.loop:
    ldrb   x2, [x1],#1
    cbnz  x2, .loop
.done:
    sub    x0, x1, x0
    sub    x0, x0, #1
    ret
```


Optimized strlen (SVE)

```
strlen:
    mov     x1, x0
    ptrue  p0.b
.loop:
    setffr
    ldff1b z0.b, p0/z, [x1]
    rdffr  p1.b, p0/z
    cmpeq  p2.b, p1/z, z0.b, #0
    brkbs  p2.b, p1/z, p2.b
    incp   x1, p2.b
    b.last .loop
    sub    x0, x1, x0
    ret
```

Suboptimal implementation

```
    setffr          /* initialize FFR */
    ptrue p2.b      /* all ones; loop invariant */
    mov x1, 0       /* initialize length */

    /* Read a vector's worth of bytes, stopping on first fault. */
0:   ldff1b z0.b, p2/z, [x0, x1]
    rdffrs p0.b, p2/z
    b.nlast 2f

    /* First fault did not fail: the whole vector is valid. Avoid
    depending on the contents of FFR beyond the branch. */
    incb x1, all    /* speculate increment */
    cmpeq p1.b, p2/z, z0.b, 0 /* loop if no zeros */
    b.none 0b
    decb x1, all    /* undo speculate */

    /* Zero found. Select the bytes before the first and count them. */
1:   brkbs p0.b, p2/z, p1.b
    incp x1, p0.b
    mov x0, x1
    ret

    /* First fault failed: only some of the vector is valid. Perform the
    comparison only on the valid bytes. */
2:   cmpeq p1.b, p0/z, z0.b, 0
    b.any 1b

    /* No zero found. Re-init FFR, increment, and loop. */
    setffr
    incp x1, p0.b
    b 0b
```

03_SVE/D3.1_sve_strcmp

See 03_SVE/SVE-SVE2-programming-examples-REL-01.pdf

```
ptrue    p5.b
setffr
mov      x5, #0
.L_loop:
ldff1b  z0.b, p5/z, [str1, x5]
ldff1b  z1.b, p5/z, [str2, x5]
rdffrs  p7.b, p5/z
b.nlast .L_fault
incb    x5
cmpeq   p0.b, p5/z, z0.b, #0
cmpne   p1.b, p5/z, z0.b, z1.b
.L_test:
orrs    p4.b, p5/z, p0.b, p1.b
b.none  .L_loop
.L_fault:
incp    x5, p7.b
setffr
cmpeq   p0.b, p7/z, z0.b, #0
cmpne   p1.b, p7/z, z0.b, z1.b
```

```
Initialize predicate
Initialize FFR to all true
Initialize loop counter to 0
- BEGIN LOOP BODY -
Load str1 with first-fault
Load str2 with first-fault
Read FFR and place active elements in predicate p7
If first active element is not true handle string remainder
Increment loop counter by number of byte lanes in SVE register
Compare full vector of str1 chars to 0 (look for end of string)
Compare full vector of str1 to str2
- BEGIN LOOP CONDITION -
Bitwise OR of predicates to check for differences and str end
Predicate condition: no active elements are true (keep looping)
- BEGIN reminder strcmp -
Increment loop counter by number of true predicate elements
Reset FFR
Compare partial vector of str1 chars to 0
Compare partial vector of str1 to str2
```

arm

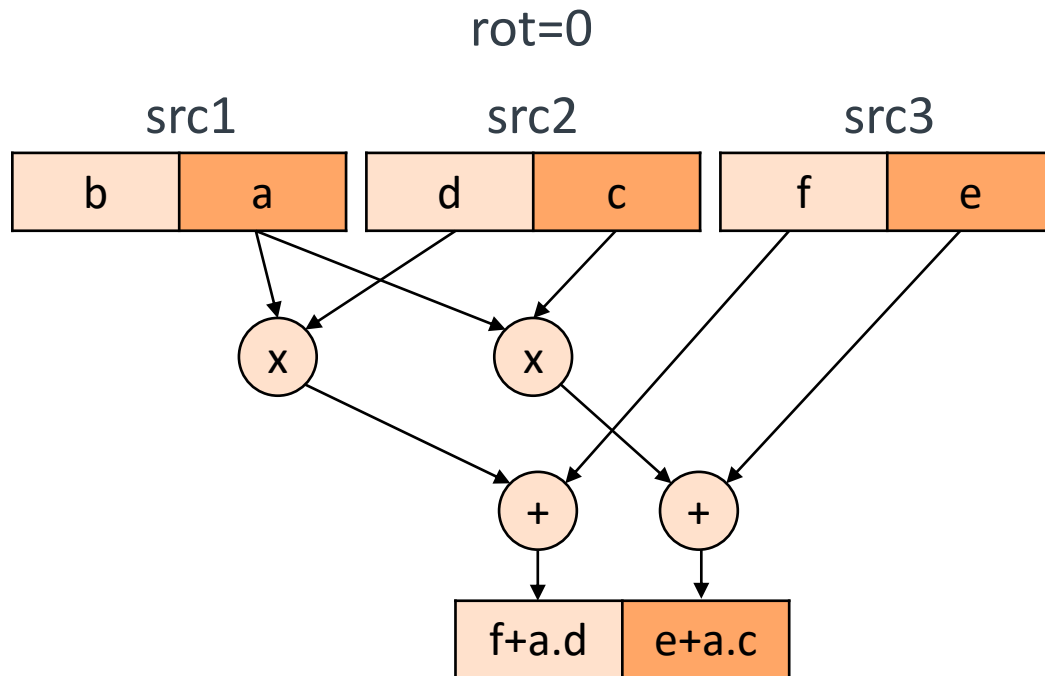
Complex
Arithmetic

02_ACLE/04_dotprod_complex

Complex arithmetic with FMLA or FCMLA

Complex multiply:

$$(a+ib).(c+id) = (ac-bd)+i(ad+bc)$$

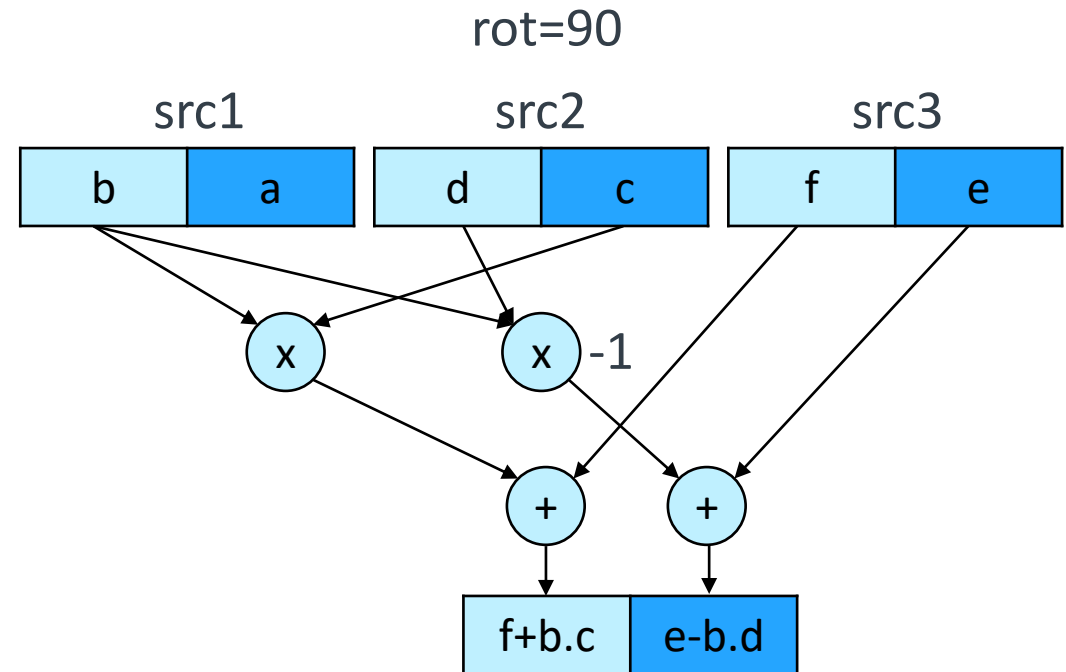


FCMLA in SVE works for 4 rotations:

- 0, 90, 180 and 270

Complex multiply-add needs a pair of instructions

```
svcmmla_z(pg, src3, src1, src2, 0);  
svcmmla_z(pg, src3, src1, src2, 90);
```



02_ACLE/04_dotprod_complex

See README.md for details

ACLE code using FMLA

```
for (int i=0; i<SIZE; i+=v1) {
    svfloat32x2_t va = svld2(p32_all, (float32_t*)&a[i]);
    svfloat32x2_t vb = svld2(p32_all, (float32_t*)&b[i]);
    svfloat32x2_t vc = svld2(p32_all, (float32_t*)&c[i]);

    vc.v0 = svmla_m(p32_all, vc.v0, va.v0, vb.v0); //c.re += a.re * b.re
    vc.v1 = svmla_m(p32_all, vc.v1, va.v1, vb.v0); //c.im += a.im * b.re
    vc.v0 = svmls_m(p32_all, vc.v0, va.v1, vb.v1); //c.re -= a.im * b.im
    vc.v1 = svmla_m(p32_all, vc.v1, va.v0, vb.v1); //c.im += a.re * b.im

    svst2(p32_all, (float32_t*)&c[i], vc);
}
```

ACLE code using FCMLA

```
for (int i=0; i<SIZE; i+=v1/2) {
    svfloat32_t va = svld1(p32_all, (float32_t*)&a[i]);
    svfloat32_t vb = svld1(p32_all, (float32_t*)&b[i]);
    svfloat32_t vc = svld1(p32_all, (float32_t*)&c[i]);

    vc = svcmla_m(p32_all, vc, va, vb, 0); //c += a * b
    vc = svcmla_m(p32_all, vc, va, vb, 90);

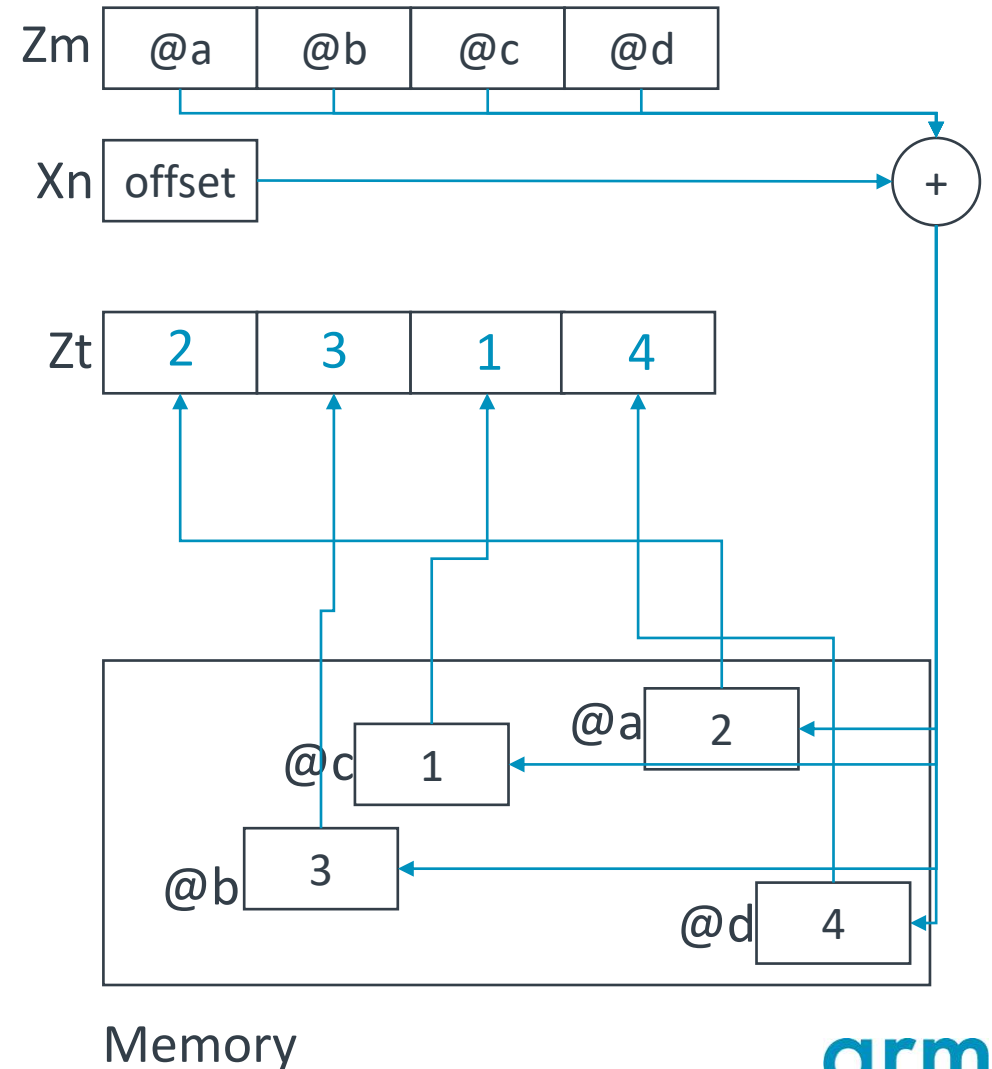
    svst1(p32_all, (float32_t*)&c[i], vc);
}
```

arm

Data Movement

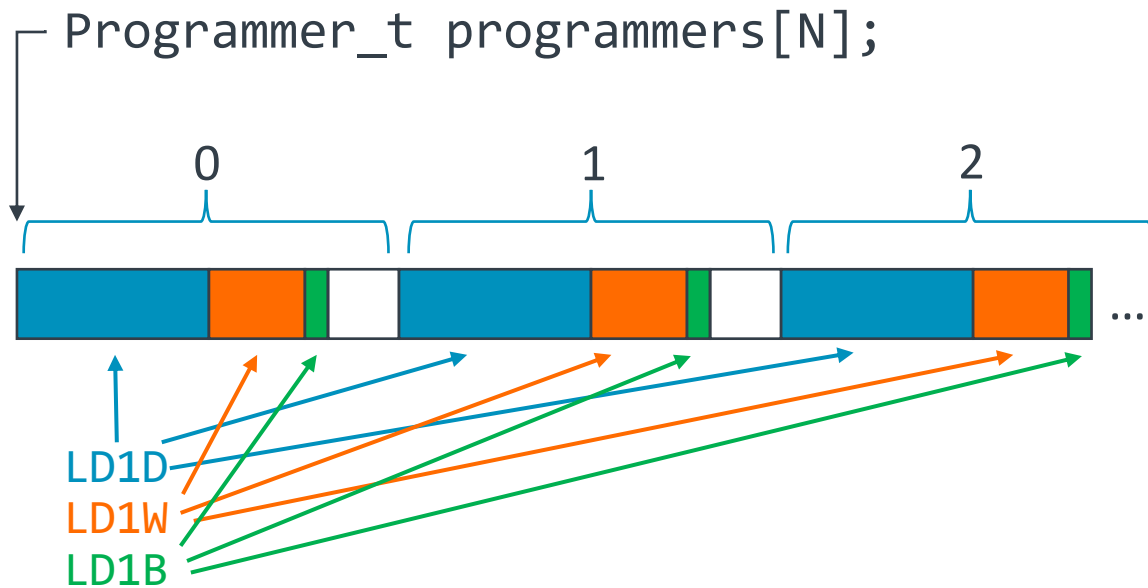
Gather/Scatter Operations

- Enable vectorization of codes with non-adjacent accesses on adjacent lanes
- Examples:
 - Outer loop vectorization
 - Strided accesses (larger than +1)
 - Random accesses
- Performance implementation dependent
 - Worst case one separate access per element
- LD1D <Zt>.D, Ps/Z [<Xn>, <Zm>.D]

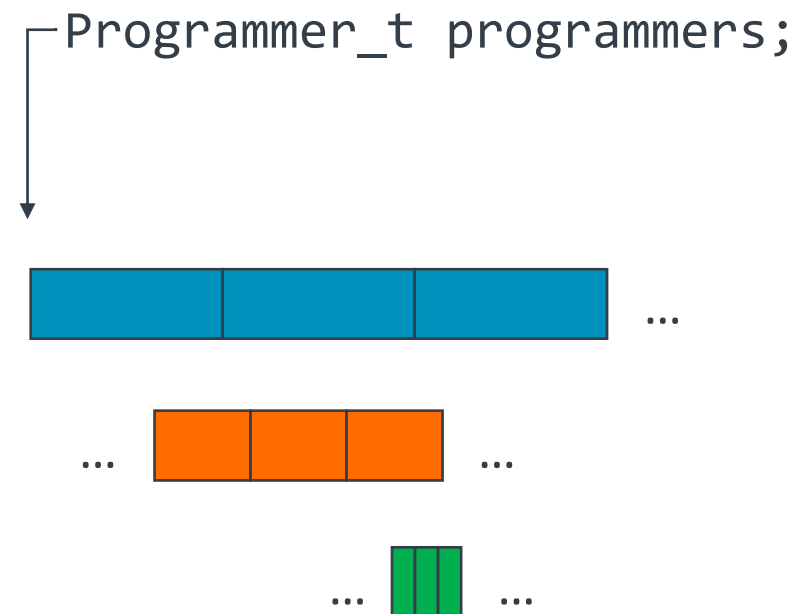


Array of Structures vs. Structure of Arrays

```
typedef struct {  
    uint64_t num_projects;  
    float caffeine;  
    bool vim_nemacs;  
} Programmer_t;
```



```
typedef struct {  
    uint64_t num_projects[N];  
    float caffeine[N];  
    bool vim_nemacs[N];  
} Programmer_t;
```



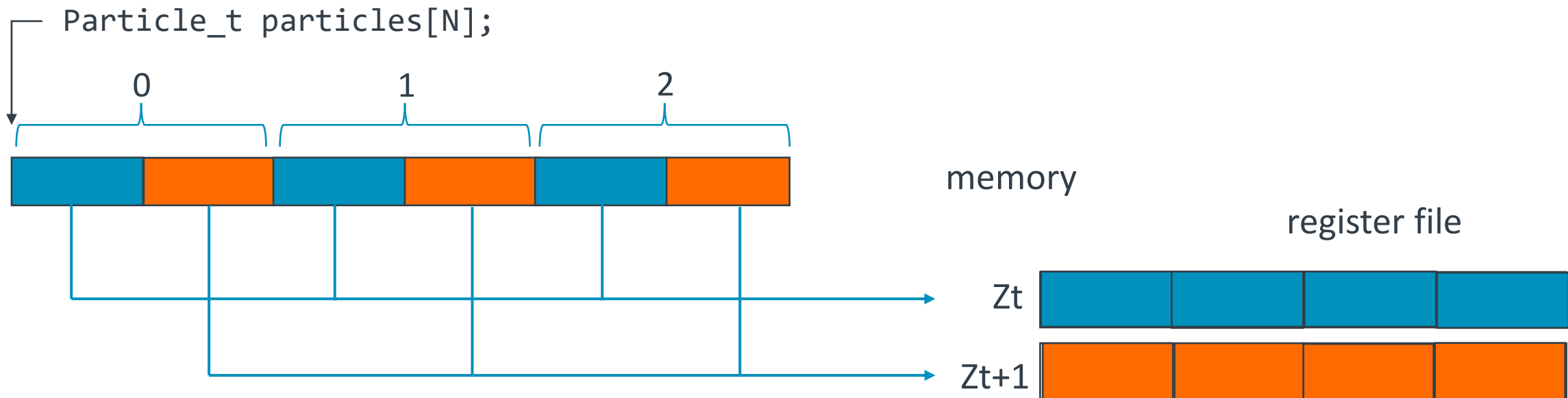
02_ACLE/05_gather

- Use SVE vector tuples to access structure
- Performance depends on u-arch and memory system

```
typedef struct {  
    uint32_t x;  
    uint32_t y;  
} Particle_t;
```

Structure loads:

```
LD2W {<Zt>.H, <Zt+1>.H}, P0/Z, [X0]
```



02_ACLE/05_gather

See README.md for details

Original code

```
typedef struct {
    int32_t x;
    int32_t y;
} Particle_t;

void move(Particle_t p[SIZE], int32_t x, int32_t y)
{
    for (int i=0; i<SIZE; i++) {
        p[i].x += x;
        p[i].y += y;
    }
}
```

Use Vector Tuples

```
void move(Particle_t p[SIZE], int32_t x, int32_t y)
{
    uint32_t v1 = svcntw();

    svbool_t p32_all = svptrue_b32();

    for (int i=0; i<SIZE; i+=v1) {
        svint32x2_t vp = svld2(p32_all, (int32_t*)&p[i]);

        vp.v0 = svadd_m(p32_all, vp.v0, x);
        vp.v1 = svadd_m(p32_all, vp.v1, y);

        svst2(p32_all, (int32_t*)&p[i], vp);
    }
}
```

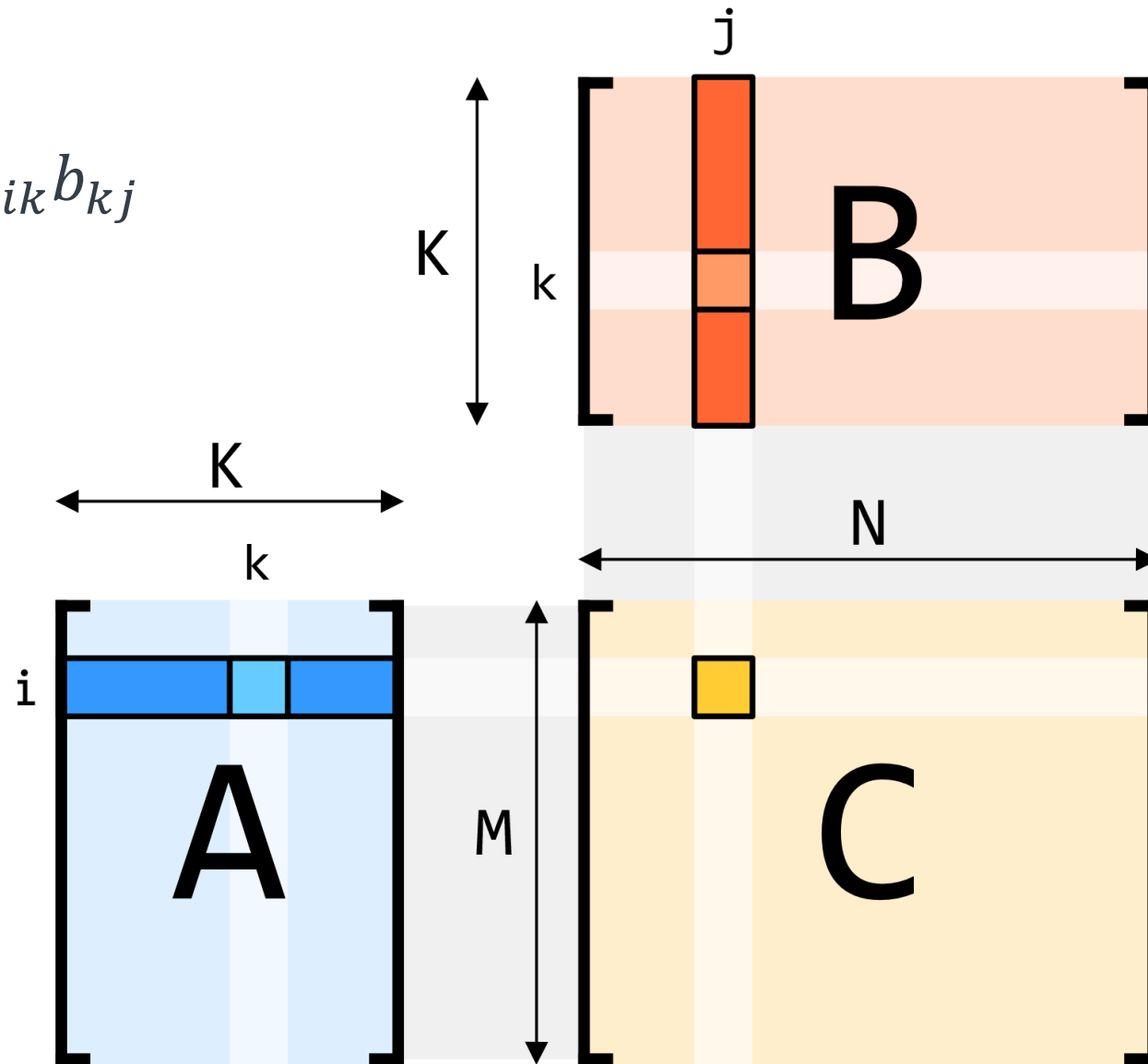
arm

SVE Load Replicate and GEMM

Neat theory; YMMV

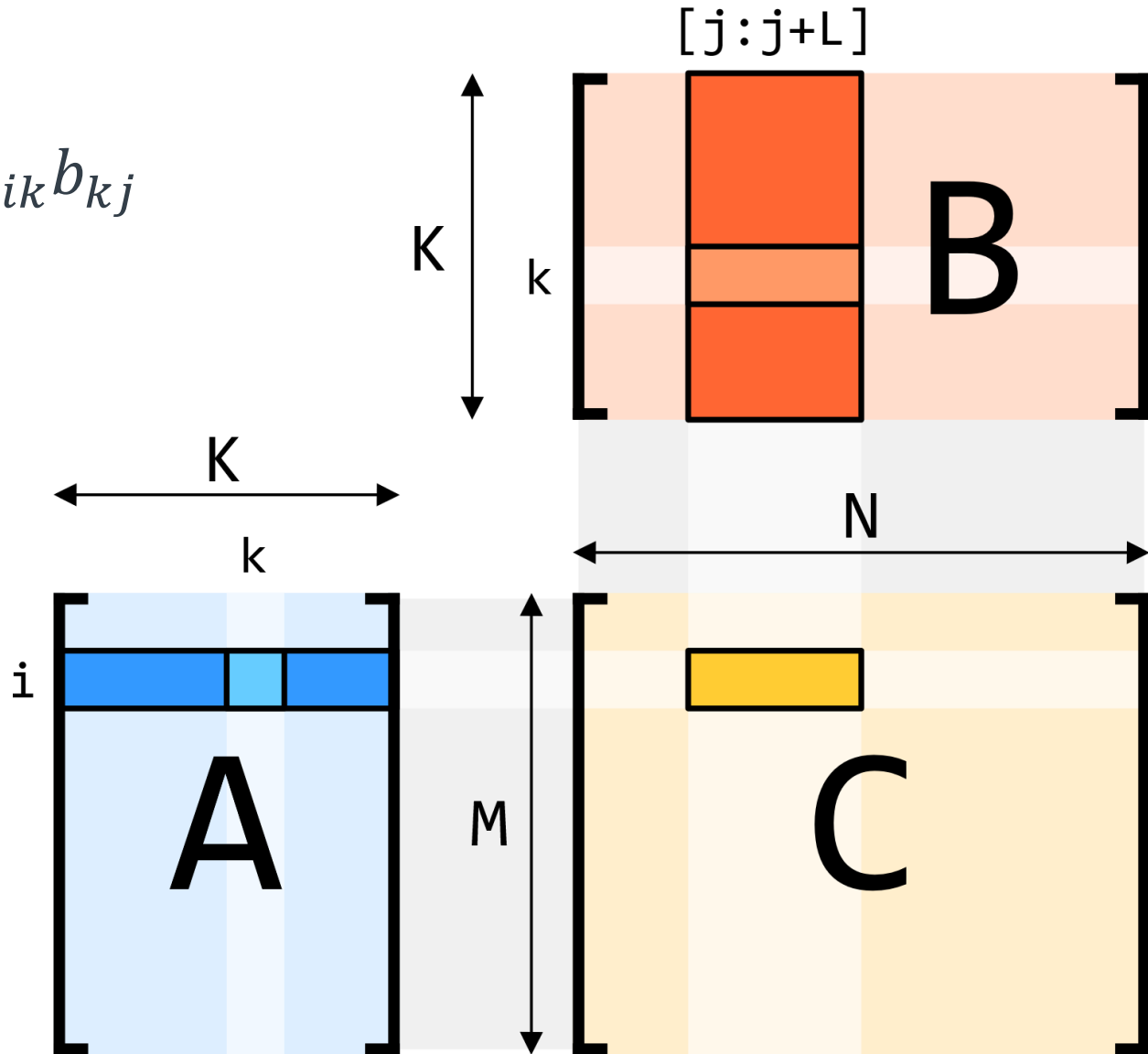
GEMM

$$c_{ij} += \sum_{k=0}^{K-1} a_{ik} b_{kj}$$



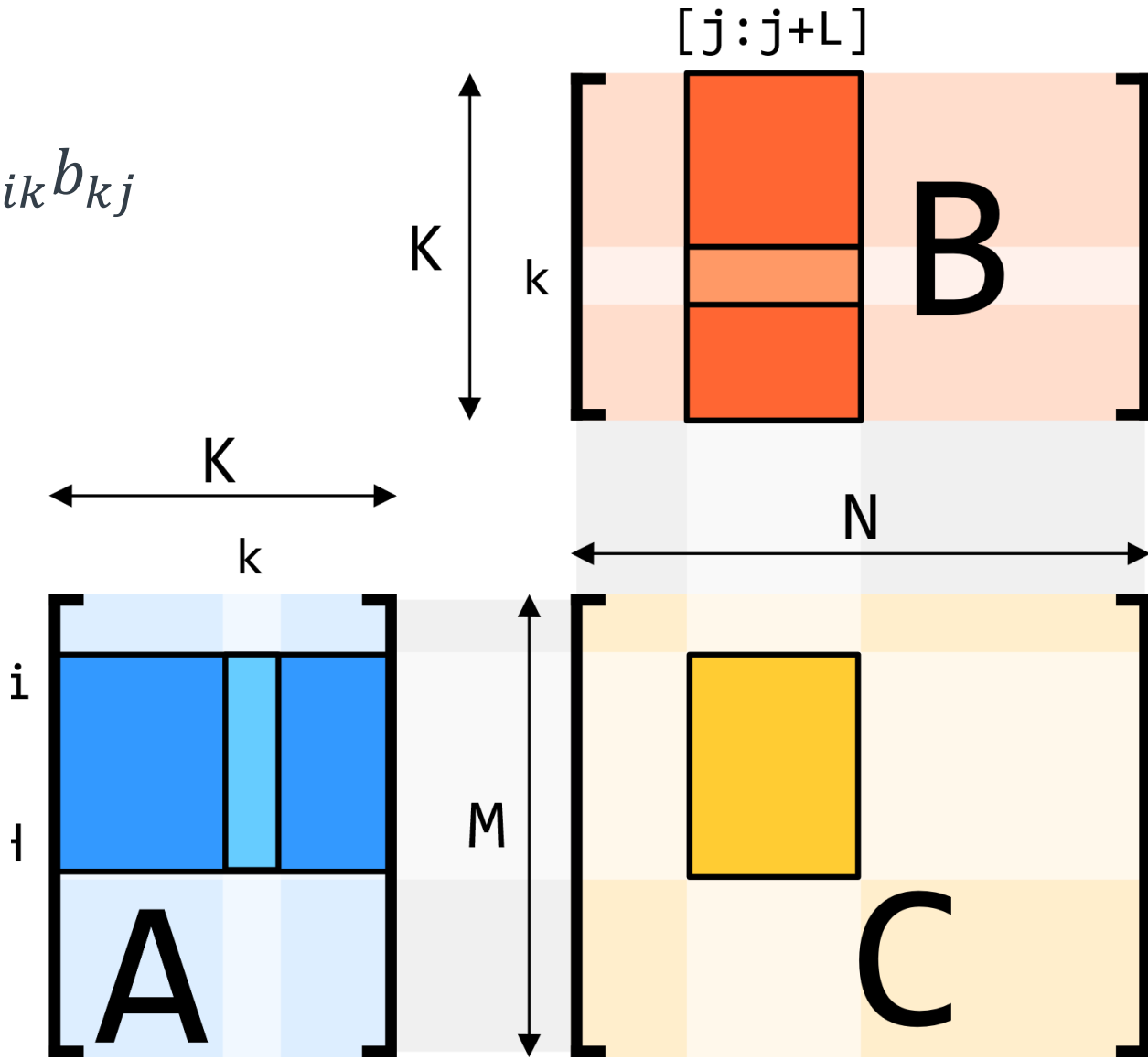
GEMM – step 1 – vectorize along the columns

$$c_{ij} += \sum_{k=0}^{K-1} a_{ik} b_{kj}$$

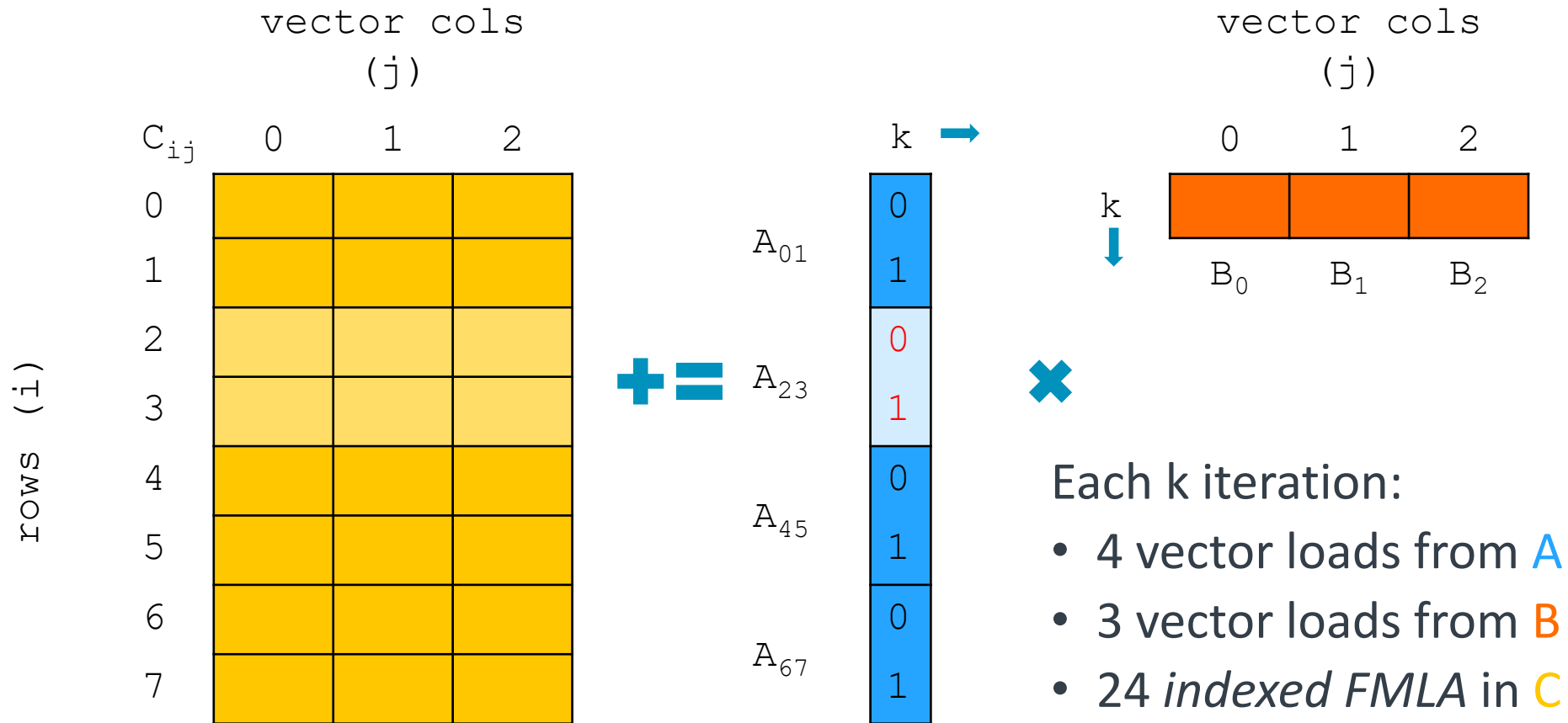


GEMM – step 2 – unroll along rows

$$c_{ij} += \sum_{k=0}^{K-1} a_{ik} b_{kj}$$



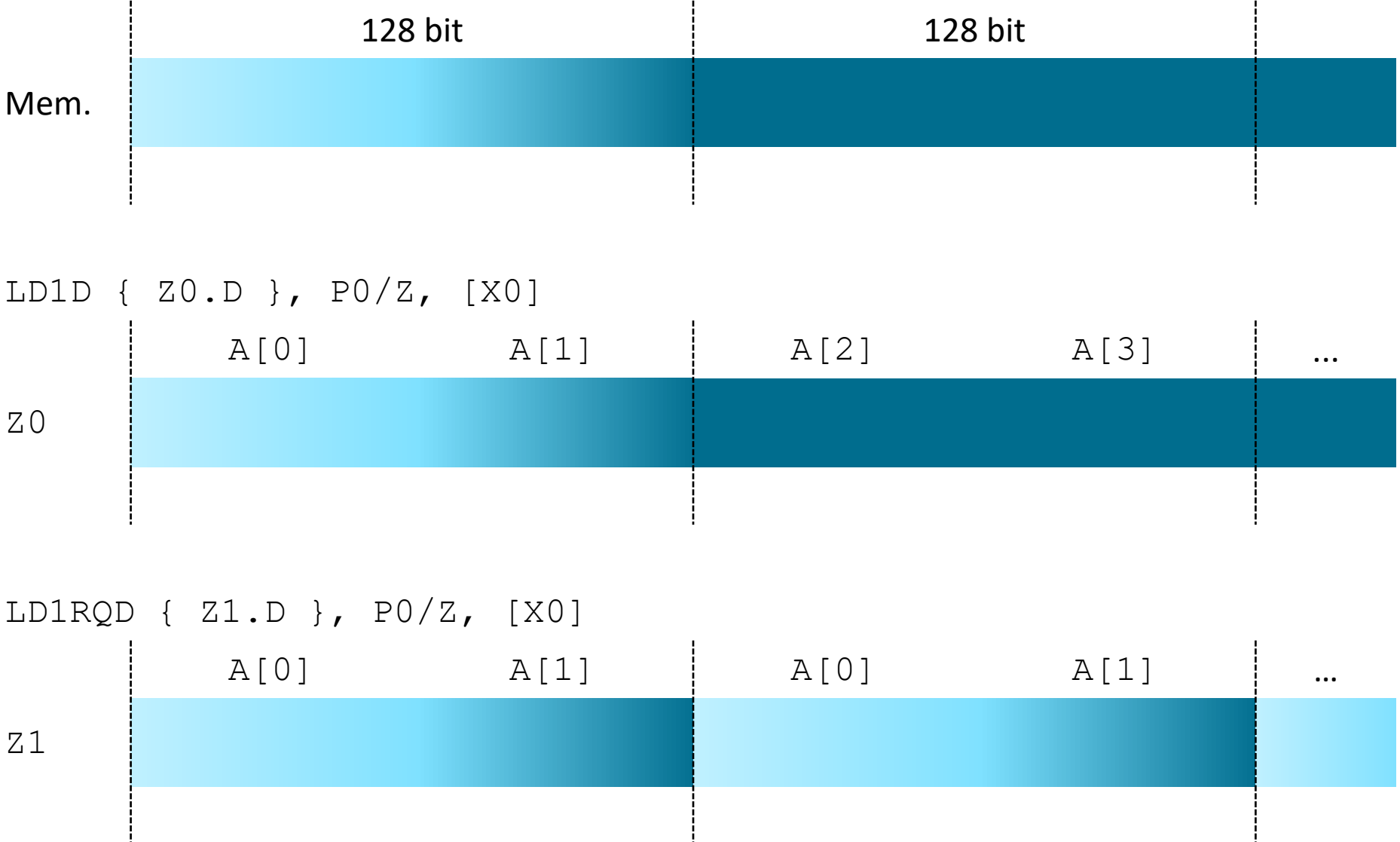
DGEMM kernel (NEON, 128-bit, 24 accumulators + 4 + 3)



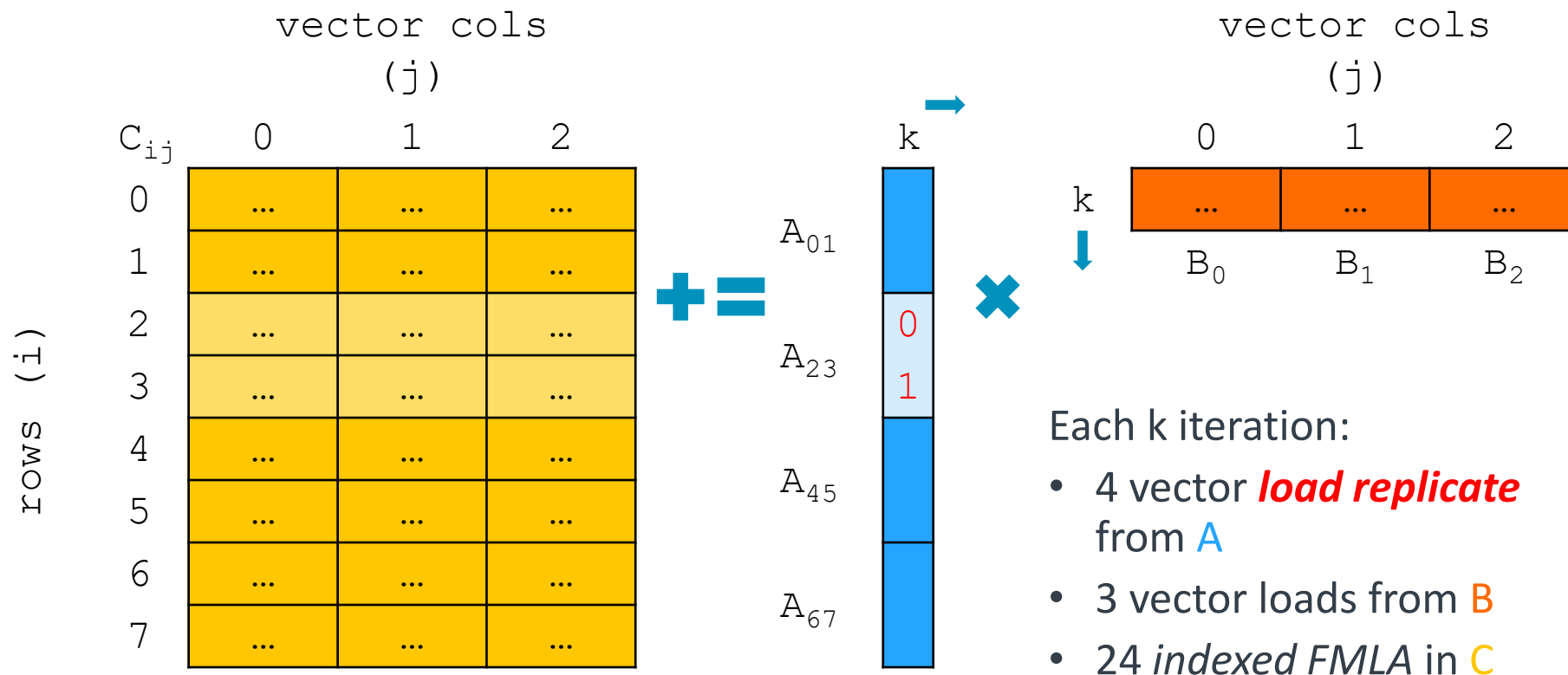
```
fmla C20.2d, B0.2d, A23.d[0] fmla C30.2d, B0.2d, A23.d[1]
fmla C21.2d, B1.2d, A23.d[0] fmla C31.2d, B1.2d, A23.d[1]
fmla C22.2d, B2.2d, A23.d[0] fmla C32.2d, B2.2d, A23.d[1]
```

SVE load **Replicate** **Quadword** instructions: LD1RQ [BHWD]

```
double *A;
```



DGEMM kernel (SVE, *LEN* x 128-bit, 24 accumulators + 4 + 3)



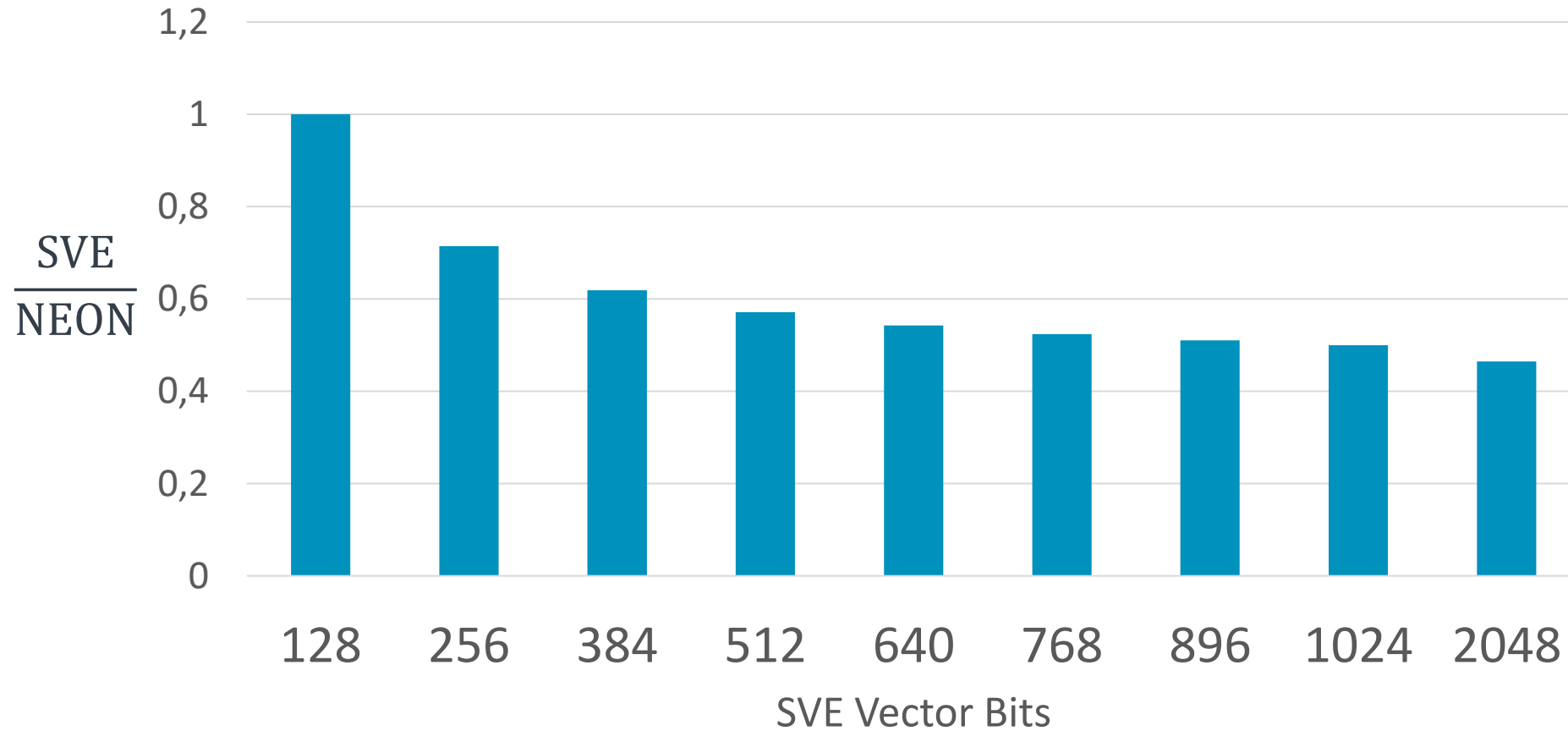
```
fmla C20.d, B0.d, A23.d[0]
fmla C21.d, B1.d, A23.d[0]
fmla C22.d, B2.d, A23.d[0]
```

```
fmla C30.d, B0.d, A23.d[1]
fmla C31.d, B1.d, A23.d[1]
fmla C32.d, B2.d, A23.d[1]
```

DGEMM: SVE vs NEON

	Each k iteration	Bytes in one k iteration	Total C area computed	SVE/NEON A and B data reads for same C area
SVE	<ul style="list-style-type: none"> 4 vector load replicate from A 3 vector loads from B 	$4 \times 128b$ + $3 \times \text{LEN} \times 128b$	$24 \times 128b \times \text{LEN}$	$\frac{\text{SVE}}{\text{NEON}} = \frac{4 + 3 \times \text{LEN}}{7}$
NEON	<ul style="list-style-type: none"> 4 vector loads from A 3 vector loads from B 	$7 \times 128b$	$24 \times 128b$	

DGEMM: SVE more memory efficient than LEN times NEON

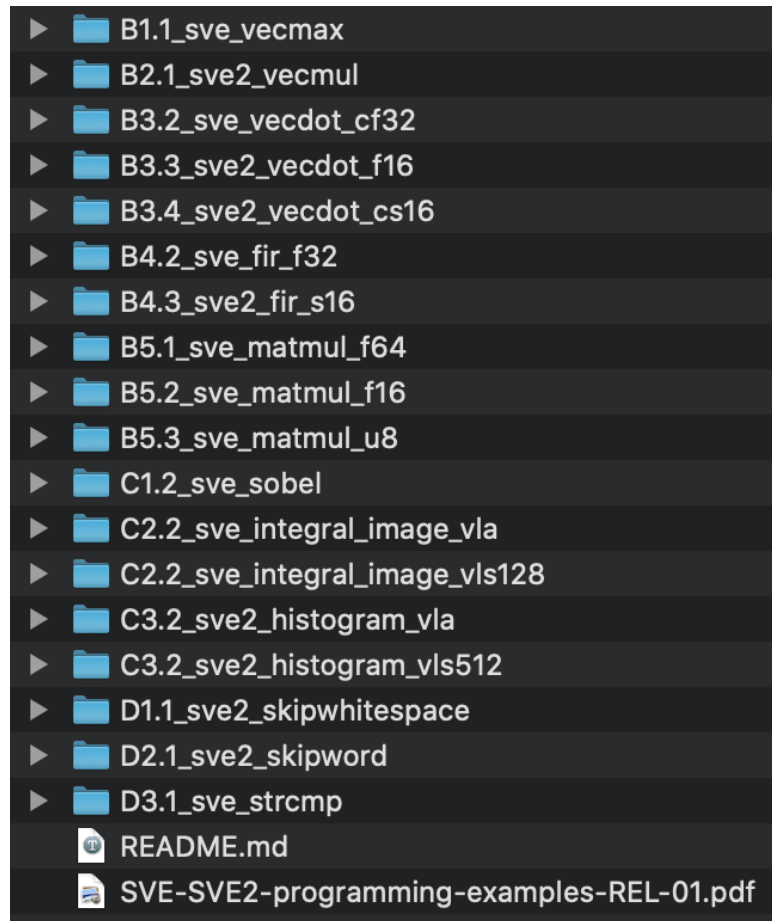


arm

Micro-examples

03_SVE: Micro-examples for individual SVE instructions

Based on <https://developer.arm.com/documentation/dai0548/latest>



Using the Micro-examples

- PDF document describes each example
- Directory name indicates section in the PDF
- Includes both SVE and SVE2 examples
- May need ArmIE to run some examples, e.g. SVE2 examples require ArmIE to run on A64FX

03_SVE: Micro-examples for individual SVE instructions

Based on <https://developer.arm.com/documentation/dai0548/latest>

- ▶ B1.1_sve_vecmax
- ▶ B2.1_sve2_vecmul
- ▶ B3.2_sve_vecdot_cf32
- ▶ B3.3_sve2_vecdot_f16
- ▶ B3.4_sve2_vecdot_cs16
- ▶ B4.2_sve_fir_f32
- ▶ B4.3_sve2_fir_s16
- ▶ B5.1_sve_matmul_f64
- ▶ B5.2_sve_matmul_f16
- ▶ B5.3_sve_matmul_u8
- ▶ C1.2_sve_sobel
- ▶ C2.2_sve_integral_image_vla
- ▶ C2.2_sve_integral_image_vls128
- ▶ C3.2_sve2_histogram_vla
- ▶ C3.2_sve2_histogram_vls512
- ▶ D1.1_sve2_skipwhitespace
- ▶ D2.1_sve2_skipword
- ▶ D3.1_sve_strcmp
- ▶ README.md
- ▶ SVE-SVE2-programming-examples-REL-01.pdf

Contents		
SVE Programming Examples		
Release information	ii	
Non-Confidential Proprietary Notice	iii	
Preface		
About this book	ix	
Using this book	x	
Conventions	xi	
Typographical conventions	xi	
Numbers	xi	
Pseudocode descriptions	xi	
Assembler syntax descriptions	xi	
Rules-based writing	xii	
Identifiers	xii	
Examples	xiii	
Additional reading	xiii	
Feedback	xiv	
Feedback on this book	xiv	
Part A Arm Scalable Vector Extension (SVE) overview		
Chapter A1 Introduction		
A1.1 Overview	17	
A1.2 Register file	18	
A1.2.1 Vector register file	18	
A1.2.2 Predicate register file	18	
A1.2.3 First Fault Register	19	
A1.3 Predicate Condition Flags	20	
A1.3.1 Overview	20	
A1.3.2 AArch64 Condition Codes and Flags	20	
A1.3.3 SVE Condition Codes and Flags	20	
Chapter A2 Compiler support		
A2.1 Autovectorization	23	
A2.1.1 Compiler options and pragmas	23	
A2.1.2 Generating the Autovectorization diagnostics	23	
A2.2 Calling conventions	24	
Chapter A3 SVE Instruction Set		
A3.1 Overview	26	
A3.1.1 Constructive and destructive instructions	26	
A3.1.2 Predication of constructive and destructive instructions	26	
A3.1.3 Move prefix	27	
A3.2 Element size and type conversion instructions	28	
A3.2.1 Element size	28	
A3.2.2 Element size conversions	28	
A3.2.3 Type conversions	28	
A3.3 Loads and Stores	29	
A3.3.1 Contiguous loads and stores	29	
ARM-DAI-0548	Copyright © 2019-2020 Arm Limited or its affiliates. All rights reserved.	iv
REL-01	Non-confidential	

SVE Resources

<http://developer.arm.com/hpc>

- **Porting and Optimizing Guides**
 - For SVE: <https://developer.arm.com/docs/101726/0110>
 - For Arm in general: <https://developer.arm.com/docs/101725/0110>
- **The SVE Specification and Arm Instruction Reference**
 - [Arm Architecture Reference Manual Supplement, SVE for ARMv8-A](#)
 - <https://developer.arm.com/docs/ddi0596/i/a64-sve-instructions-alphabetic-order>
- **ACLE References and Examples**
 - ACLE for SVE: <https://developer.arm.com/docs/100987/latest>
 - Worked examples: [A Sneak Peek Into SVE and VLA Programming](#)
 - Optimized machine learning: [Arm SVE and Applications to Machine Learning](#)

arm

Thank You

Danke

Merci

谢谢

ありがとう

Gracias

Kiitos

감사합니다

धन्यवाद

شكرًا

תודה