# Master Informatics Eng.

2017/18

*A.J.Proença*

### Data Parallelism 1 (*vector & SIMD extensions*)
#### (*most slides are borrowed*)

# Instruction and Data Streams

§7.6 SISD, MIMD, SIMD, SPMD, and Vector

|                        |          | Data Streams                   |                                   |
| ---------------------- | -------- | ------------------------------ | --------------------------------- |
|                        |          | Single                         | Multiple                          |
| Instruction Streams    | Single   | **SISD**: Intel Pentium 4      | **SIMD**: SSE instructions of x86 |
|                        | Multiple | **MISD**: No examples today    | **MIMD**: Intel Xeon e5345        |

- SPMD: Single Program Multiple Data
  - A parallel program on a MIMD computer
  - Conditional code for different processors

# Introduction

- SIMD architectures can exploit significant data-level parallelism for:
    - matrix-oriented <u>scientific computing</u>
    - media-oriented <u>image</u> and <u>sound</u> processing

- SIMD is more energy efficient than MIMD
    - only needs to fetch one instruction per data operation
    - makes SIMD attractive for personal mobile devices

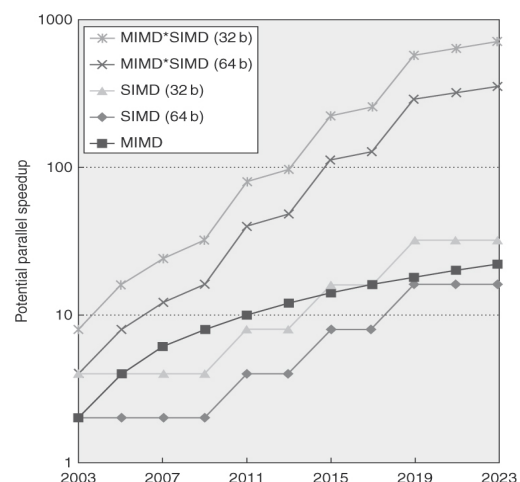- SIMD allows programmers to continue to think sequentially

**ΜΚ**®

---

# SIMD Parallelism

- Vector architectures *(slides 5 to 19)*
- SIMD & extensions *(slides 20 to 29 and next set)*
- Graphics Processor Units (GPUs) *(next set)*

- For x86 processors:
    - Expected grow:
      2 more cores/chip/year
    - SIMD width:
      2x every 4 years
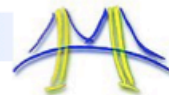    - Potential speedup:
      SIMD 2x that from MIMD!

# Vector Architectures

- ## Basic idea:
  - Read sets of data elements *(gather from memory)* into "vector registers"
  - Operate on those registers
  - Store/scatter the results back into memory

- ## Registers are controlled by the compiler
  - Used to hide memory latency
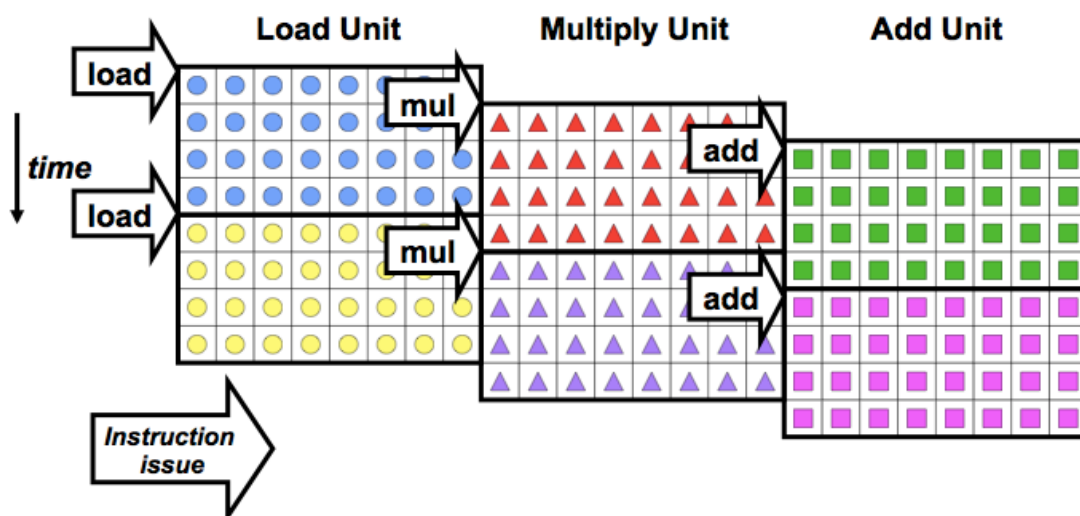  - Leverage memory bandwidth

5

## Vector Instruction Parallelism

Can overlap execution of multiple vector instructions
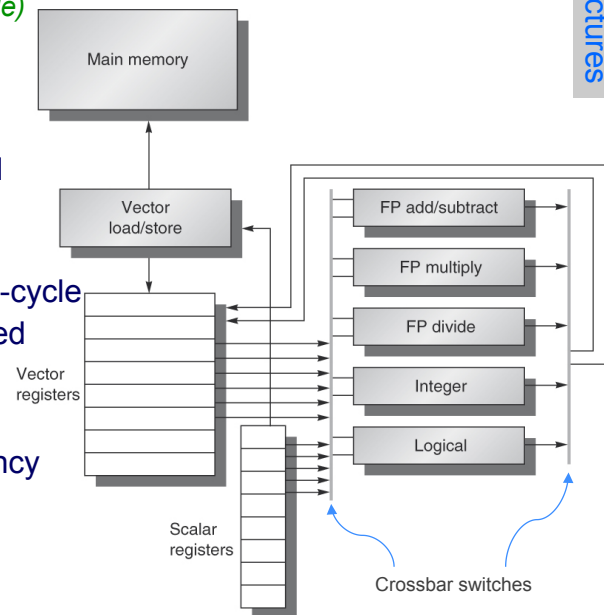- Consider machine with 32 elements per vector register and 8 lanes:



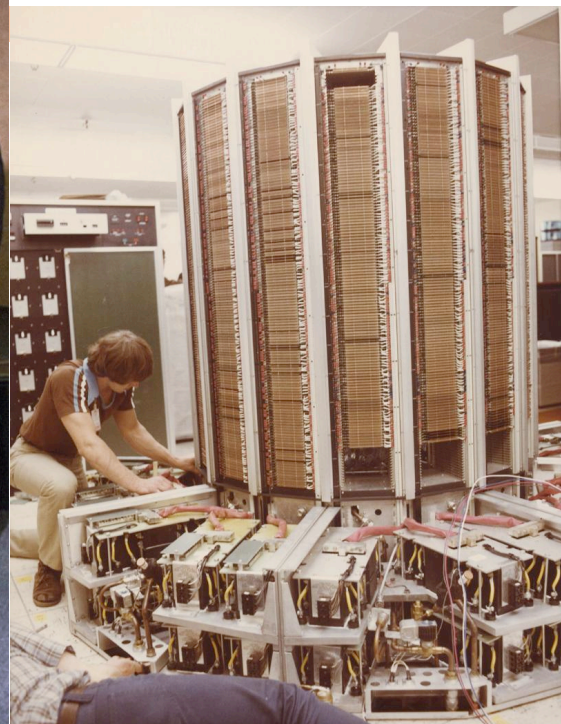Complete 24 operations/cycle while issuing 1 short instruction/cycle

8/19/2009          John Kubiatowicz          Parallel Architecture: 35

# VMIPS

- Example architecture:  VMIPS
  - Loosely based on Cray-1 *(next slide)*
  - Vector registers
    - Each register holds a 64-element, 64 bits/element vector
    - Register file has 16 read ports and 8 write ports
  - Vector functional units
    - Fully pipelined, new op each clock-cycle
    - Data & control hazards are detected
  - Vector load-store unit
    - Fully pipelined
    - 1 word/clock-cycle after initial latency
  - Scalar registers
    - 32 general-purpose registers
    - 32 floating-point registers

7



## Cray-1 Supercomputer
### (1976)

# VMIPS Instructions

- ADDVV.D: add two vectors
- ADDVS.D: add vector to a scalar
- LV/SV: vector load and vector store from address

- Example: DAXPY *(Double-precision A x X Plus Y)*

```
L.D      F0,a          ; load scalar a
LV       V1,Rx         ; load vector X
MULVS.D  V2,V1,F0      ; vector-scalar multiply
LV       V3,Ry         ; load vector Y
ADDVV    V4,V2,V3      ; add
SV       Ry,V4         ; store the result
```

- Requires the execution of 6 instructions *versus* almost 600 for MIPS *(assuming DAXPY is operating on a vector with 64 elements)*
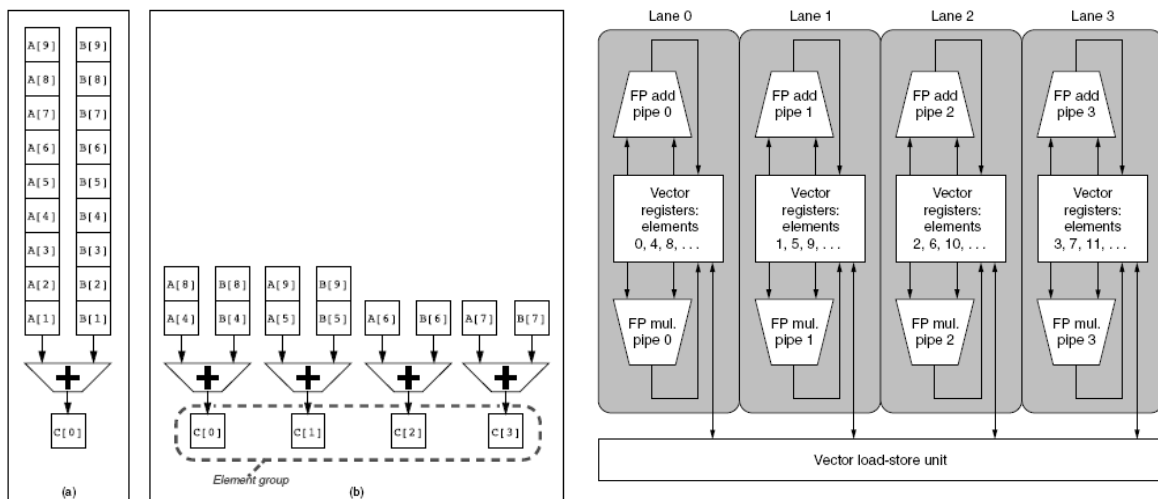
---

# Vector Execution Time

- Execution time depends on three factors:
    - Length of operand vectors
    - Structural hazards
    - Data dependencies

- VMIPS functional units consume one element per clock cycle
    - Execution time is approximately the vector length

- *Convoy*
    - Set of vector instructions that could potentially execute together in one unit of time, *chime*

# Challenges

- ## Start up time
    - Latency of vector functional unit
    - Assume the same as Cray-1
        - Floating-point add => 6 clock cycles
        - Floating-point multiply => 7 clock cycles
        - Floating-point divide => 20 clock cycles
        - Vector load => 12 clock cycles

- ## Improvements:
    - \> 1 element per clock cycle *(1)*
    - Non-64 wide vectors *(2)*
    - IF statements in vector code *(3)*
    - Memory system optimizations to support vector processors *(4)*
    - Multiple dimensional matrices *(5)*
    - Sparse matrices *(6)*
    - Programming a vector computer *(7)*

---
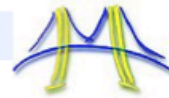
# Multiple Lanes *(1)*

- Element *n* of vector register *A* is "hardwired" to element *n* of vector register *B*
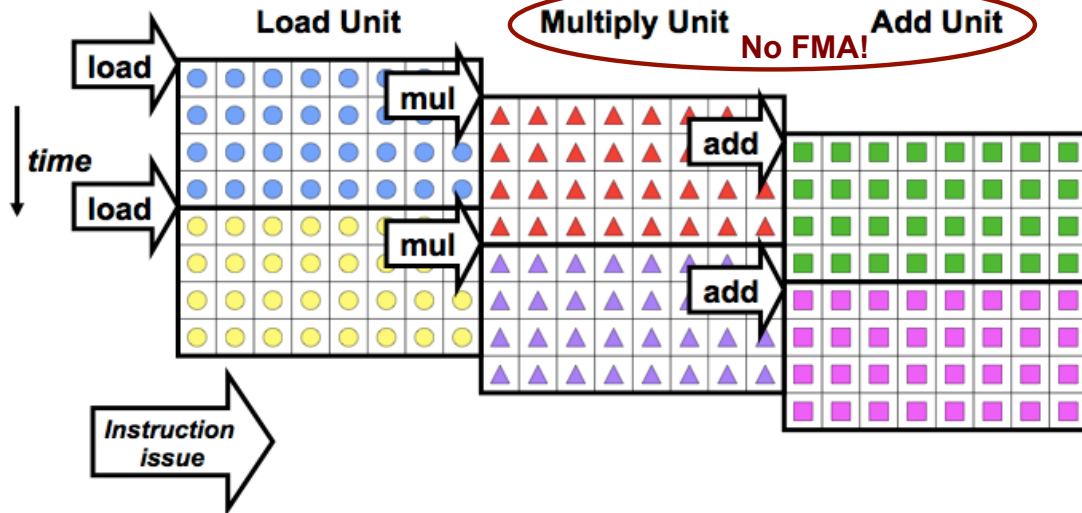    - Allows for multiple hardware lanes

## Vector Instruction Parallelism

Can overlap execution of multiple vector instructions
- Consider machine with 32 elements per vector register and 8 lanes

Complete 24 operations/cycle while issuing 1 short instruction/cycle

# Vector Length Register *(2)*

- Handling vector length not known at compile time
- Use Vector Length Register (VLR)
- Use strip mining for vectors over the maximum length:

```
low = 0;
VL = (n % MVL); /*find odd-size piece using modulo op % */
for (j = 0; j <= (n/MVL); j=j+1) { /*outer loop*/
   for (i = low; i < (low+VL); i=i+1) /*runs for length VL*/
    Y[i] = a * X[i] + Y[i] ; /*main operation*/
   low = low + VL; /*start of next vector*/
   VL = MVL; /*reset the length to maximum vector length*/
}
```

# Vector Mask Registers *(3)*

- Handling IF statements in Vector Loops:

```
for (i = 0; i < 64; i=i+1)
    if (X[i] != 0)
        X[i] = X[i] – Y[i];
```

- Use vector mask register to "disable" elements:

```
LV        V1,Rx        ;load vector X into V1
LV        V2,Ry        ;load vector Y
L.D       F0,#0        ;load FP zero into F0
SNEVS.D   V1,F0        ;sets VM(i) to 1 if V1(i)!=F0
SUBVV.D   V1,V1,V2     ;subtract under vector mask
SV        Rx,V1        ;store the result in X
```

- GFLOPS rate decreases!

# Memory Banks *(4)*

- Memory system must be designed to support high bandwidth for vector loads and stores

- Spread accesses across multiple banks
  - Control bank addresses independently
  - Load or store non sequential words
  - Support multiple vector processors sharing the same memory

- Example (Cray T932, 1996; Ford acquired 1 out of 13, $39M):
  - 32 processors, each generating 4 loads and 2 stores per cycle
  - Processor cycle time is 2.167 ns, SRAM cycle time is 15 ns
  - How many memory banks needed?

# Stride *(5)*

- Handling **multidimensional arrays** in Vector Architectures:

```
for (i = 0; i < 100; i=i+1) {
    for (j = 0; j < 100; j=j+1) {
    A[i][j] = 0.0;
    for (k = 0; k < 100; k=k+1)
            A[i][j] = A[i][j] + B[i][k] * D[k][j];
    }
}
```

- Must vectorize multiplication of rows of B with columns of D
- Use *non-unit stride* (in VMIPS: load/store vector with stride)
- Bank conflict (stall) occurs when the same bank is hit faster than bank busy time:
  - #banks / Least_Common_Multiple (stride, #banks) < bank busy time

# Scatter-Gather *(6)*

- Handling **sparse matrices** in Vector Architectures:

```
for (i = 0; i < n; i=i+1)
    A[K[i]] = A[K[i]] + C[M[i]];
```

- Use index vector:

```
LV      Vk, Rk          ;load K
LVI     Va, (Ra+Vk)     ;load A[K[]]
LV      Vm, Rm          ;load M
LVI     Vc, (Rc+Vm)     ;load C[M[]]
ADDVV.D Va, Va, Vc      ;add them
SVI     (Ra+Vk), Va     ;store A[K[]]
```

# Vector Programming *(7)*

- Compilers are a key element to give hints on whether a code section will vectorize or not

- Check if loop iterations have data dependencies, otherwise vectorization is compromised

- Vector Architectures have a too high cost, but simpler variants are currently available on off-the-shelf devices; however:
  - most do not support non-unit stride => care must be taken in the design of data structures
  - same applies for gather-scatter...

19

---

# SIMD Extensions

- Media applications operate on data types narrower than the native word size
  - Intel SIMD Ext started with 64-bit wide vectors and grew to wider vectors and more capabilities
  - Current AVX generation is 512-bit wide



```
double *x, *y, *z;
for (i=0; i<n; i++)    z[i] = x[i] + y[i];
```
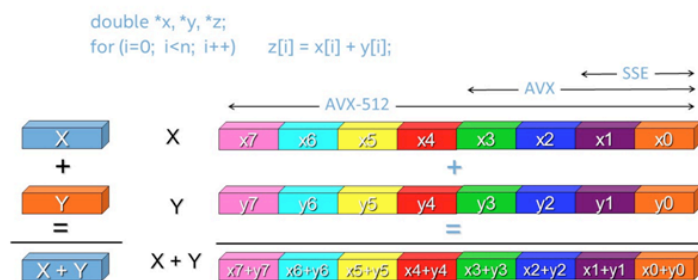
Figure 1   Scalar and vectorized loop versions with Intel® SSE, AVX and AVX-512.

- Limitations, compared to vector architectures **(before AVX...)**:
  - Number of data operands encoded into op code
  - No sophisticated addressing modes (strided, scatter-gather)
  - No mask registers

20

# Example SIMD Code

- Example DAXPY *(in MIPS SIMD)*:

```
L.D     F0,a       ;load scalar a
MOV     F1, F0     ;copy a into F1 for SIMD MUL
MOV     F2, F0     ;copy a into F2 for SIMD MUL
MOV     F3, F0     ;copy a into F3 for SIMD MUL
DADDIU  R4,Rx,#512 ;last address to load
Loop:
L.4D    F4,0[Rx]   ;load X[i], X[i+1], X[i+2], X[i+3]
MUL.4D  F4,F4,F0   ;a×X[i],a×X[i+1],a×X[i+2],a×X[i+3]
L.4D    F8,0[Ry]   ;load Y[i], Y[i+1], Y[i+2], Y[i+3]
ADD.4D  F8,F8,F4   ;a×X[i]+Y[i], ..., a×X[i+3]+Y[i+3]
S.4D    0[Ry],F8   ;store into Y[i],Y[i+1],Y[i+2],Y[i+3]
DADDIU  Rx,Rx,#32  ;increment index to X
DADDIU  Ry,Ry,#32  ;increment index to Y
DSUBU   R20,R4,Rx  ;compute bound
BNEZ    R20,Loop   ;check if done
```

---

# SIMD Implementations

- Intel implementations:
  - MMX (1996)
    - Eight 8-bit integer ops or four 16-bit integer ops
  - Streaming SIMD Extensions (SSE) (1999)
    - Eight 16-bit integer ops
    - Four 32-bit integer/fp ops or two 64-bit integer/fp ops
  - Advanced Vector eXtensions (AVX) (2010...)
    - Eight 32-bit fp ops or Four 64-bit fp ops (integers in AVX-2)
    - 512-bits wide in AVX-512 (and also in Larrabee & Phi-KNC)
  - Operands **must / should be in consecutive and aligned** memory locations

- AMD Zen/Epyc (Opteron follow-up): with AVX-2

- ARM v8 (64-bit) implementations (next...)

# Registers for vector processing in Intel 64



GPR — x87 (FP) — Program Counter — SSE/SSE2

63  32 31  0
RAX  EAX
RCX  ECX
RDX  EDX
RBX  EBX
RSP  ESP
RBP  EBP
RSI  ESI
RDI  EDI
R8
R9
R10
R11
R12
R13
R14
R15

79  0
x87 (FP)

63  32 31  0
EIP
Program Counter

127  0
XMM0
XMM7
XMM8
XMM15
SSE/SSE2

Original x86/IA32
Added by x86-64

Kmask k0..k7   zmm0..zmm31   ymm0..ymm15   xmm0..xmm7
AVX512   AVX   SSE   SSE – IA64

High amounts of compute need large amounts of state to compensate for memory BW
AVX512 has 8x state compared to SSE (commensurate with its 8x flops level)

*AJProença, Sistemas de Computação, UMinho, 2015/16*

**Advanced Vector eXtensions, AVX 1.0**

8x floats
4x doubles
16x bytes
8x 16-bit shorts
4x 32-bit integers
2x 64-bit integers
1x 128-bit integer

8x floats w/FMA3
4x doubles w/FMA3
32x bytes
16x 16-bit shorts
8x 32-bit integers
4x 64-bit integers
2x 128-bit integers

**AVX 2.0**

**Next: AVX-512 (Skylake...)**    23

---

# Vector & FP register sizes in ARMv8 (64-bit)



| Unused | D31 |
| Unused | S31 |
| Unused | H31 |
| Register V31 | |

127   64 63   32 31   16 15   0

...

DP register — Unused | D0
SP register — Unused | S0
HP register — Unused | H0
Vector register — Register V0

127   64 63   32 31   16 15   0
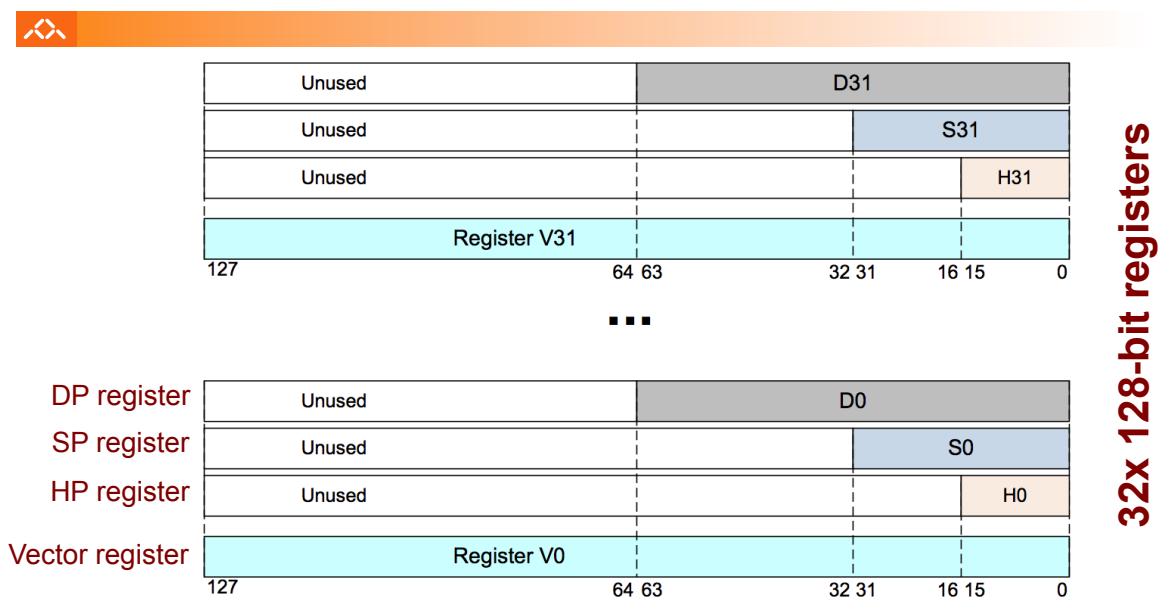
**32x 128-bit registers**

**Figure 4-10 Arrangement of floating-point values**

— Note —
16-bit floating-point is supported, but only as a format to be converted from or to. It is not supported for data processing operations.

*AJProença, Adv*    24

# NEON FP registers in ARMv8 (64-bit)

| | 127 | 112 111 | 96 95 | 80 79 | 64 63 | 48 47 | 32 31 | 16 15 | 0 |
|---|---|---|---|---|---|---|---|---|---|
| **128-bit NEON register** | | | | | | | | | |

| | 127 | | | | 64 63 | | | | 0 |
|---|---|---|---|---|---|---|---|---|---|
| **2 x 64-bit lanes** | | | 1 | | | | 0 | | |

| | 127 | | 96 95 | | 64 63 | | 32 31 | | 0 |
|---|---|---|---|---|---|---|---|---|---|
| **4 x 32-bit lanes** | | 3 | | 2 | | 1 | | 0 | |

| | 127 | 112 111 | 96 95 | 80 79 | 64 63 | 48 47 | 32 31 | 16 15 | 0 |
|---|---|---|---|---|---|---|---|---|---|
| **8 x 16-bit lanes** | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 | |

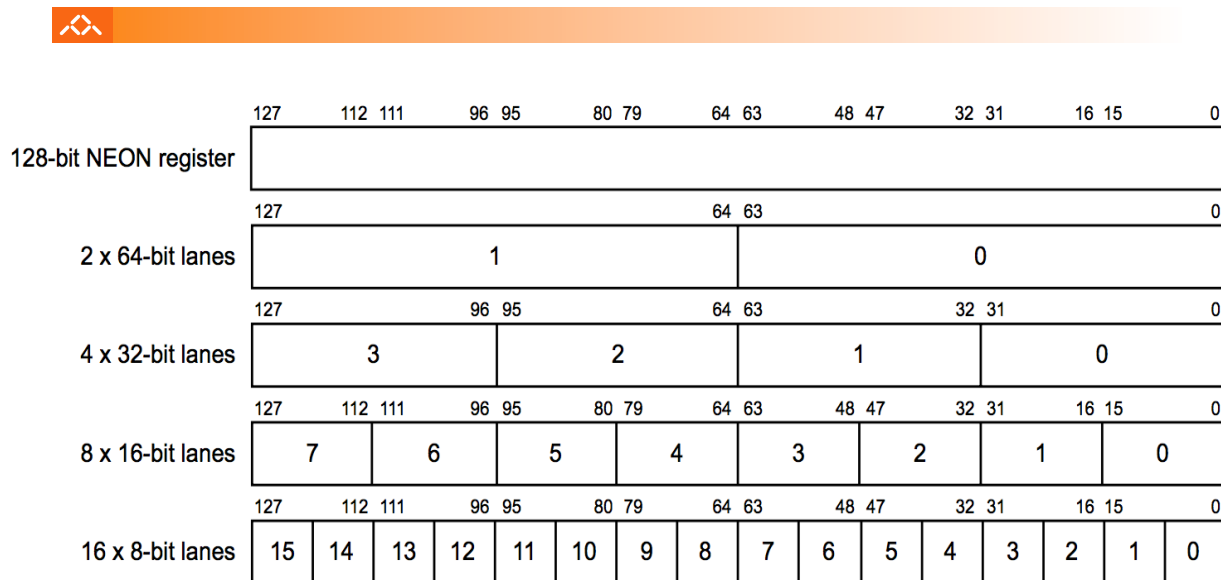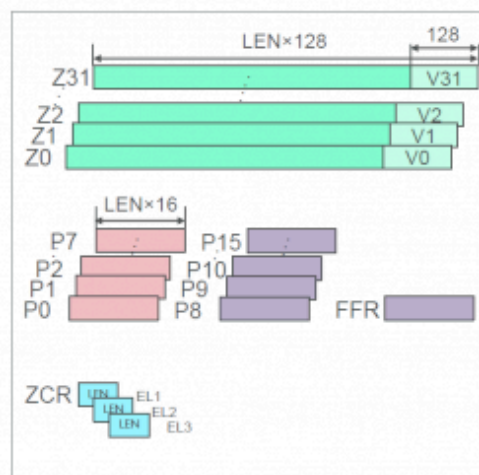| | 127 | 112 111 | 96 95 | 80 79 | 64 63 | 48 47 | 32 31 | 16 15 | 0 |
|---|---|---|---|---|---|---|---|---|---|
| **16 x 8-bit lanes** | 15 14 13 12 | 11 10 9 8 | 7 6 5 4 | 3 2 1 0 | | | | | |

**Figure 7-1 Divisions of the V register**

# ARMv8-A *Scalable Vector Extension (SVE)*
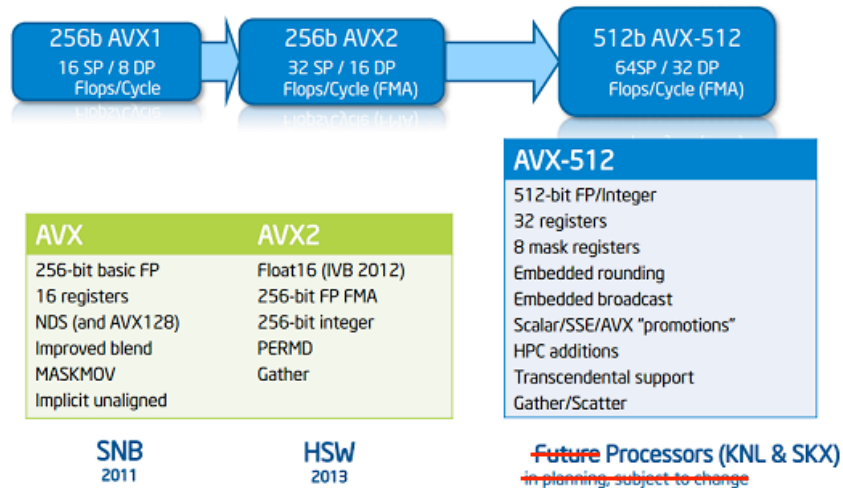
## SVE architectural state

- Scalable vector registers
  - Z0-Z31   extending NEON's V0-V31
    - DP & SP floating-point
    - 64, 32, 16 & 8-bit integer

- Scalable predicate registers
  - P0-P7    lane masks for ld/st/arith
  - P8-P15   for predicate manipulation
  - FFR      *first fault register*

- Scalable vector control registers
  - ZCR_ELx  vector length (LEN=1..16)
  - Exception / privilege level EL1 to EL3

LEN×128   128
Z31   V31
Z2   V2
Z1   V1
Z0   V0

LEN×16
P7   P15
P2   P10
P1   P9
P0   P8   FFR

ZCR   LEN EL1, LEN EL2, LEN EL3

5  © ARM 2016

**ARM**

# Intel evolution to the AVX-512

# Intel SIMD ISA evolution

# The AVX-512 across Intel devices



The 12 levels of AVX512 in Intel processors
2017-06-01
according to Intel SDE 8.40

CannonLake

SkyLake Xeon
Core-X

Knights Mill

Knights
Landing

AVX512VBMI
AVX512IFMA

AVX512BW
AVX512DQ
AVX512VL

AVX512F
AVX512CD

AVX512ER
AVX512PF

AvX512_4FMAPS
AVX512_4VNNIW
AVX512VPOPCNTDQ

AVX-512PF → **Prefetch**: Gather and Scatter Prefetch
AVX-512ER → **E**xponential and **R**eciprocal Instructions

| Intel AVX-512 Instruction Types | |
|---|---|
| AVX-512-F | AVX-512 Foundation Instructions |
| AVX-512-VL | Vector Length Orthogonality : ability to operate on sub-512 vector sizes |
| AVX-512-BW | 512-bit Byte/Word support |
| AVX-512-DQ | Additional D/Q/SP/DP instructions (converts, transcendental support, etc.) |
| AVX-512-CD | Conflict Detect : used in vectorizing loops with potential address conflicts |

# Additional features in Intel x86



from Marat Dukhan, 2014