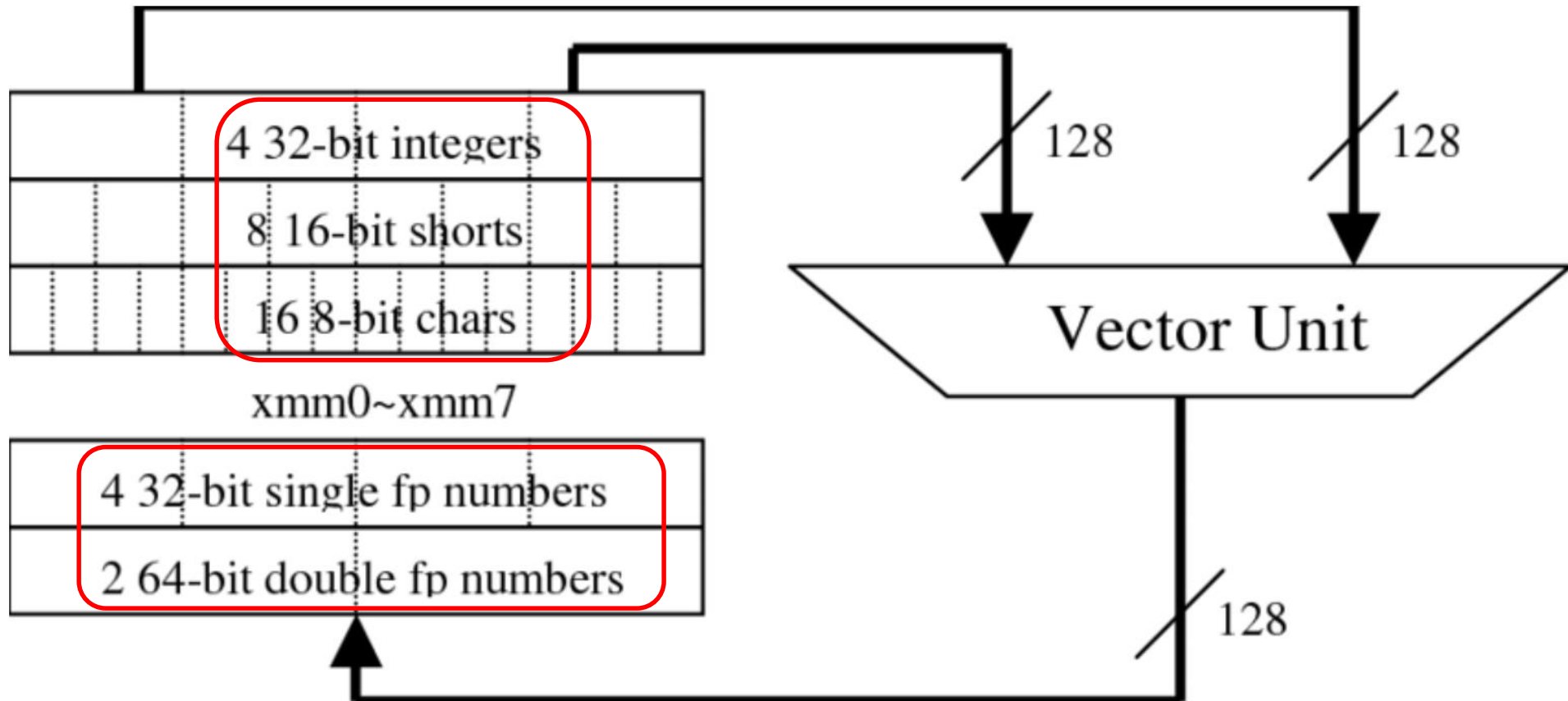


4.3 SIMD Instruction Set Extensions for Multimedia

SIMD Extensions

- Many media applications operate on **narrower data types** than the 32-bit word size.
 - Graphics: **8-bit color**
 - Audio samples: **8-16 bits**
- By partitioning carry chains within a **256-bit adder**, a processor could perform **simultaneous operations on short vectors**
 - 32 * **8-bit** operands
 - 16 * **16-bit** operands
 - 8 * **32-bit** operands

Example: 128-bit adder



SIMD Extensions vs. Vector Arch.

- **Limitations, compared to vector instructions:**
 - **Fix # of data operands** encoded into op code
 - led to the addition of hundreds of instructions
 - **No sophisticated addressing modes**
 - Not support strided, scatter-gather
 - **No mask registers.**
 - Not support conditional branching

Typical SIMD multimedia support

- Summary of typical SIMD multimedia support for 256-bit-wide operations

| Instruction category | Operands | | | | |
|-----------------------|---|--------|---------|---------|------------------------------|
| Unsigned add/subtract | Thirty-two | 8-bit, | sixteen | 16-bit, | eight 32-bit, or four 64-bit |
| Maximum/minimum | Thirty-two | 8-bit, | sixteen | 16-bit, | eight 32-bit, or four 64-bit |
| Average | Thirty-two | 8-bit, | sixteen | 16-bit, | eight 32-bit, or four 64-bit |
| Shift right/left | Thirty-two | 8-bit, | sixteen | 16-bit, | eight 32-bit, or four 64-bit |
| Floating point | Sixteen 16-bit, eight 32-bit, four 64-bit, or two 128-bit | | | | |

SIMD Implementations

- **Intel 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)**
 - Four 64-bit integer/fp ops
- **AVX-512 (2017)**
 - Eight 64-bit integer/fp ops

Intel SIMD ISA Evolution

SIMD extensions on top of x86/x87

64b
SIMD

128b
SIMD

256b
SIMD

512b
SIMD



RISC-V SIMD code for DAXPY

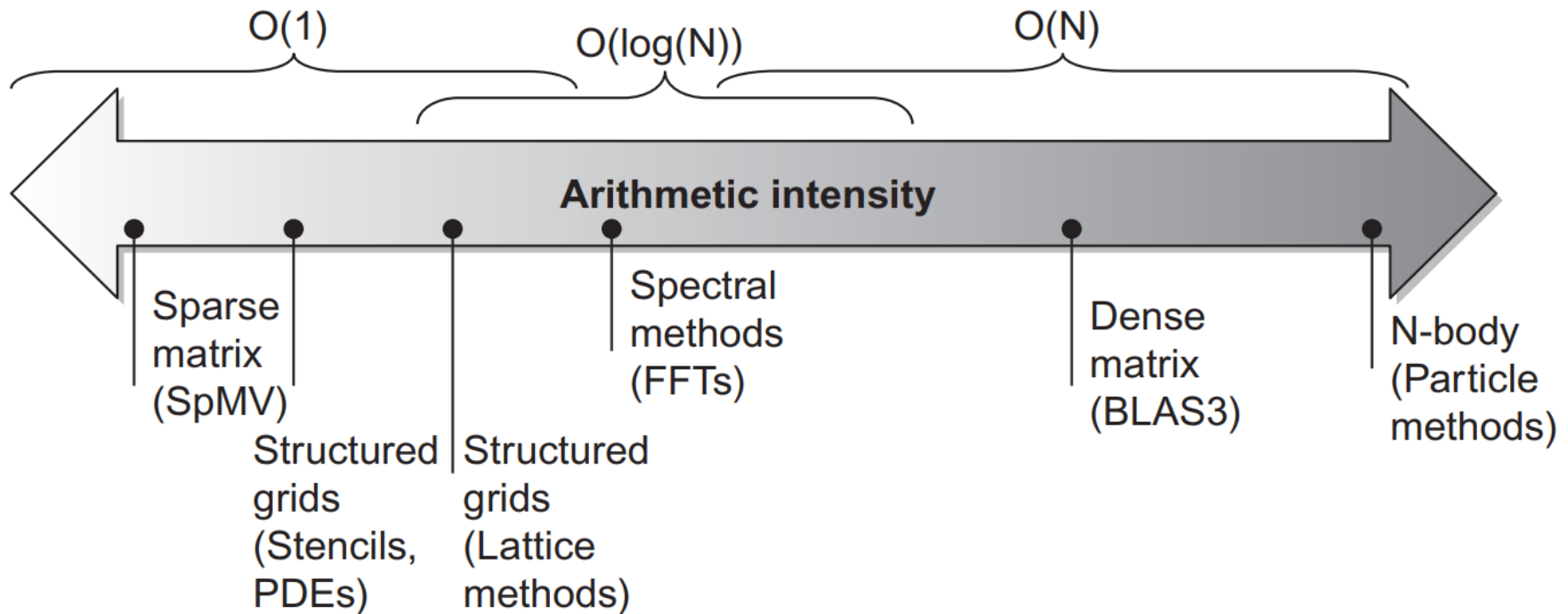
| | | | |
|-------|-----------------------|----------------------------|---------------------------------------|
| | <code>fld</code> | <code>f0, a</code> | <code>#Load scalar a</code> |
| | <code>splat.4D</code> | <code>f0, f0</code> | <code>#Make 4 copies of a</code> |
| | <code>addi</code> | <code>x28, x5, #256</code> | <code>#Last address to load</code> |
| Loop: | <code>fld.4D</code> | <code>f1, 0(x5)</code> | <code>#Load X[i] ... X[i+3]</code> |
| | <code>fmul.4D</code> | <code>f1, f1, f0</code> | <code>#a × X[i] ... a × X[i+3]</code> |
| | <code>fld.4D</code> | <code>f2, 0(x6)</code> | <code>#Load Y[i] ... Y[i+3]</code> |
| | <code>fadd.4D</code> | <code>f2, f2, f1</code> | <code># a × X[i] + Y[i] ...</code> |
| | | | <code># a × X[i+3] + Y[i+3]</code> |
| | <code>fsd.4D</code> | <code>f2, 0(x6)</code> | <code>#Store Y[i] ... Y[i+3]</code> |
| | <code>addi</code> | <code>x5, x5, #32</code> | <code>#Increment index to X</code> |
| | <code>addi</code> | <code>x6, x6, #32</code> | <code>#Increment index to Y</code> |
| | <code>bne</code> | <code>x28, x5, Loop</code> | <code>#Check if done</code> |

NOTE: We add the suffix “4D” on instructions that operate on four double-precision operands at once

Roofline Performance Model

Arithmetic intensity

- **Arithmetic intensity:** the ratio of floating-point operations per byte of memory accessed.



Roofline Performance Model

- One intuitive way to **compare variations of SIMD architectures** is the **Roofline model**
- **Basic idea:**
 - Plot **peak floating-point throughput** as a function of **arithmetic intensity**
 - Ties together **floating-point performance** and **memory performance** for a target machine

Examples

Attainable GFLOPs/s = Min(Peak Memory BW

× Arithmetic Intensity, Peak Floating – Point Perf.)

