# CS250P: Computer Systems Architecture
## Explicit Parallelism
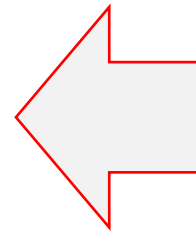
Sang-Woo Jun

Fall 2022

UCI

# Modern Processor Topics - Performance

❑ Transparent Performance Improvements
- o Pipelining, Caches
- o Superscalar, Out-of-Order, Branch Prediction, Speculation, …
- o Covered in CS250A and others

❑ Explicit Performance Improvements
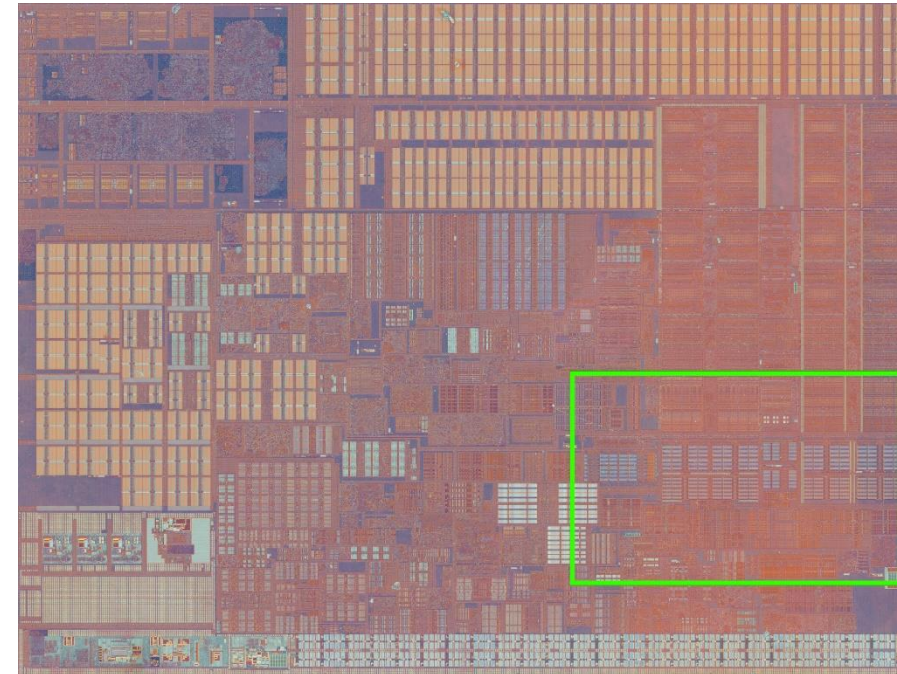- o SIMD extensions, AES extensions, …
- o …

# SIMD operations

❑ Single ISA instruction performs same computation on multiple data

❑ Typically implemented with special, wider registers

❑ Example operation:
  o Load 32 bytes from memory to special register X
  o Load 32 bytes from memory to special register Y
  o Perform addition between each 4-byte value in X and each 4 byte value in Y
  o Store the four results in special register Z     For i in (0 to 7): Z[i] = X[i] + Y[i];
  o Store Z to memory

❑ RISC-V SIMD extensions (P) is still being worked on (as of 2021)

# Example: Intel SIMD Extensions

❏ More transistors (Moore's law) but no faster clock, no more ILP…
  o More capabilities per processor has to be explicit!

❏ New instructions, new registers
  o Must be used explicitly by programmer or compiler!

❏ Introduced in phases/groups of functionality
  o SSE – SSE4 (1999 –2006)
    • 128 bit width operations
  o AVX, FMA, AVX2, AVX-512 (2008 – 2019)
    • 256 – 512 bit width operations
  o F16C, and more to come?



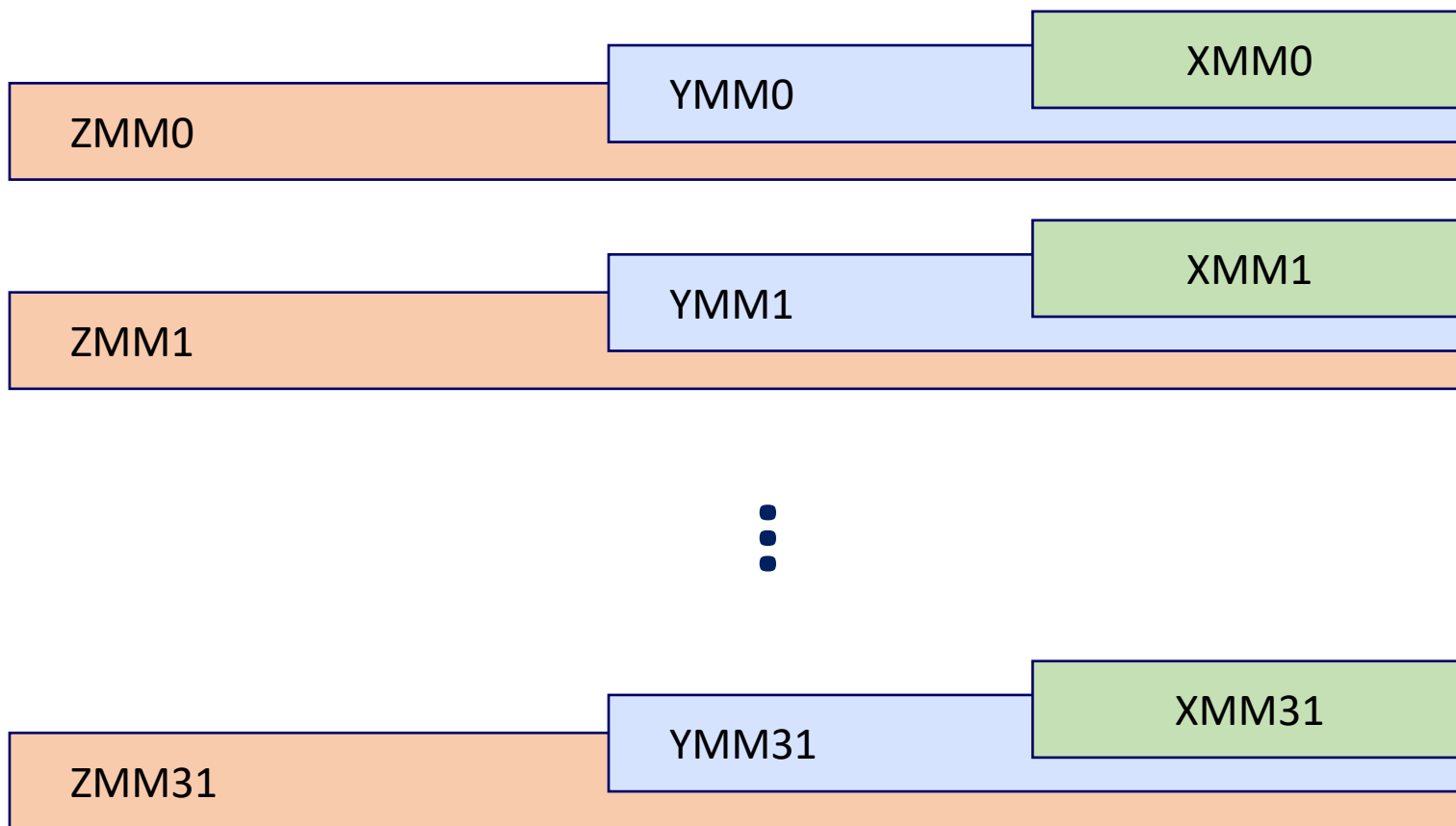Performance matters, "AVX-512 Mask Registers, Again." 2020

# Aside: Do I Have SIMD Capabilities?

❑ less /proc/cpuinfo

```
flags           : fpu vme de pse tsc msr pae mce cx8 apic sep mtrr pge mca cmov pat p
se36 clflush dts acpi mmx fxsr sse sse2 ss ht tm pbe syscall nx pdpe1gb rdtscp lm con
stant_tsc art arch_perfmon pebs bts rep_good nopl xtopology nonstop_tsc cpuid aperfmp
erf tsc_known_freq pni pclmulqdq dtes64 monitor ds_cpl vmx est tm2 ssse3 sdbg fma cx1
6 xtpr pdcm pcid sse4_1 sse4_2 x2apic movbe popcnt tsc_deadline_timer aes xsave avx f
16c rdrand lahf_lm abm 3dnowprefetch cpuid_fault epb invpcid_single pti ssbd ibrs ibp
b stibp tpr_shadow vnmi flexpriority ept vpid fsgsbase tsc_adjust bmi1 avx2 smep bmi2
 erms invpcid mpx rdseed adx smap clflushopt intel_pt xsaveopt xsavec xgetbv1 xsaves
dtherm ida arat pln pts hwp hwp_notify hwp_act_window hwp_epp flush_l1d
```

# Intel SIMD Registers (AVX-512)

| | | |
|---|---|---|
| | | XMM0 |
| | YMM0 | |
| ZMM0 | | |

| | | |
|---|---|---|
| | | XMM1 |
| | YMM1 | |
| ZMM1 | | |

⋮

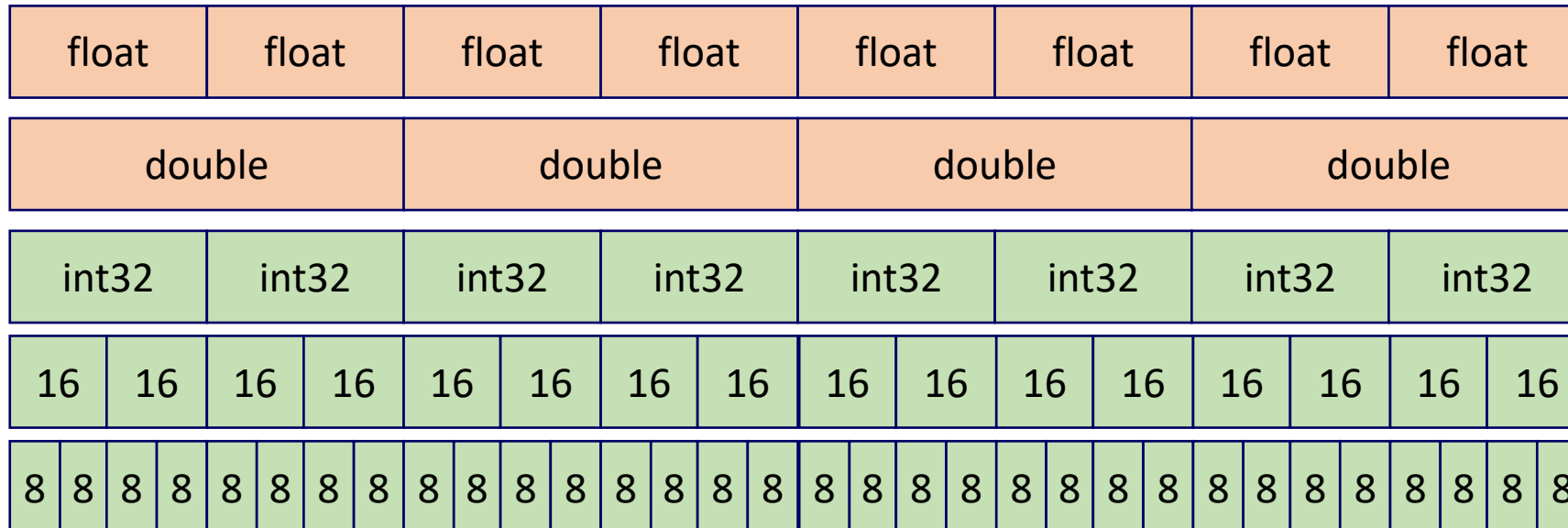| | | |
|---|---|---|
| | | XMM31 |
| | YMM31 | |
| ZMM31 | | |

❑ XMM0 – XMM15
  o 128-bit registers
  o SSE

❑ YMM0 – YMM15
  o 256-bit registers
  o AVX, AVX2

❑ ZMM0 – ZMM31
  o 512-bit registers
  o AVX-512

# SSE/AVX Data Types

| 255 | | | | | | | 0 |
|-----|---|---|---|---|---|---|---|
| | | | YMM0 | | | | |

| float | float | float | float | float | float | float | float |
|-------|-------|-------|-------|-------|-------|-------|-------|

| double | double | double | double |
|--------|--------|--------|--------|

| int32 | int32 | int32 | int32 | int32 | int32 | int32 | int32 |
|-------|-------|-------|-------|-------|-------|-------|-------|

| 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|

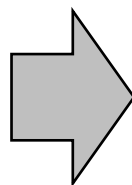| 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Operation on
32 8-bit values
in one instruction!

# Compiler Automatic Vectorization

❑ In gcc, flags "-O3 -mavx -mavx2" attempts automatic vectorization

❑ Works pretty well for simple loops

```
int a[256], b[256], c[256];
void foo () {
    for (int i=0; i<256; i++) a[i] = b[i] * c[i];
}
```

```
.L2:
        vmovdqa xmm1, XMMWORD PTR b[rax]
        add     rax, 16
        vpmulld xmm0, xmm1, XMMWORD PTR c[rax-16]
        vmovaps XMMWORD PTR a[rax-16], xmm0
        cmp     rax, 1024
        jne     .L2
```

Generated using GCC explorer: https://gcc.godbolt.org/

❑ But not for anything complex

   o E.g., naïve bubblesort code not parallelized at all

# Intel SIMD Intrinsics

❑ Use C functions instead of inline assembly to call AVX instructions

❑ Compiler manages registers, etc

❑ Intel Intrinsics Guide

    o https://software.intel.com/sites/landingpage/IntrinsicsGuide

    o One of my most-visited pages…

```
e.g.,
__m256 a, b, c;
__m256 d = _mm256_fmadd_ps(a, b, c); // d[i] = a[i]*b[i]+c[i] for i = 0 …7
```

# Intrinsic Naming Convention

❑ _mm**<width>**_**[function]**_**[type]**

  o E.g., _mm256_fmadd_ps :
    perform fmadd (floating point multiply-add) on
    256 bits of
    packed single-precision floating point values (8 of them)

| Width | Prefix |
|---|---|
| 128 | _mm_ |
| 256 | _mm256_ |
| 512 | _mm512_ |

| Type | Postfix |
|---|---|
| Single precision | _ps |
| Double precision | _pd |
| Packed signed integer | _epiNNN (e.g., epi256) |
| Packed unsigned integer | _epuNNN (e.g., epu256) |
| Scalar integer | _siNNN (e.g., si256) |

Not all permutations exist! Check guide

# Example: Vertical Vector Instructions

❑ Add/Subtract/Multiply
  ○ _mm256_add/sub/mul/div_ps/pd/epi
    • Mul only supported for epi32/epu32/ps/pd
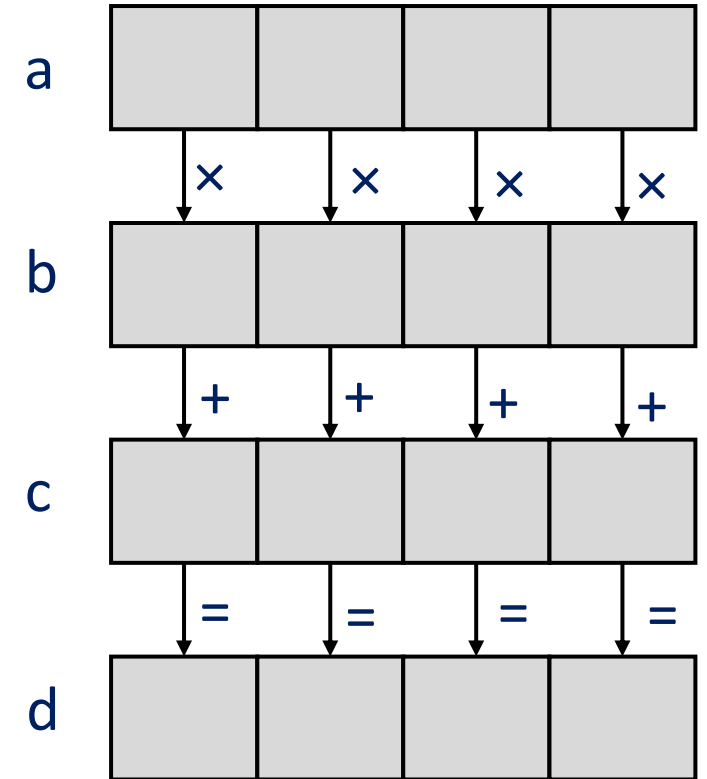    • Div only supported for ps/pd
    • Consult the guide!
❑ Max/Min/GreaterThan/Equals
❑ Sqrt, Reciprocal, Shift, etc...
❑ FMA (Fused Multiply-Add)
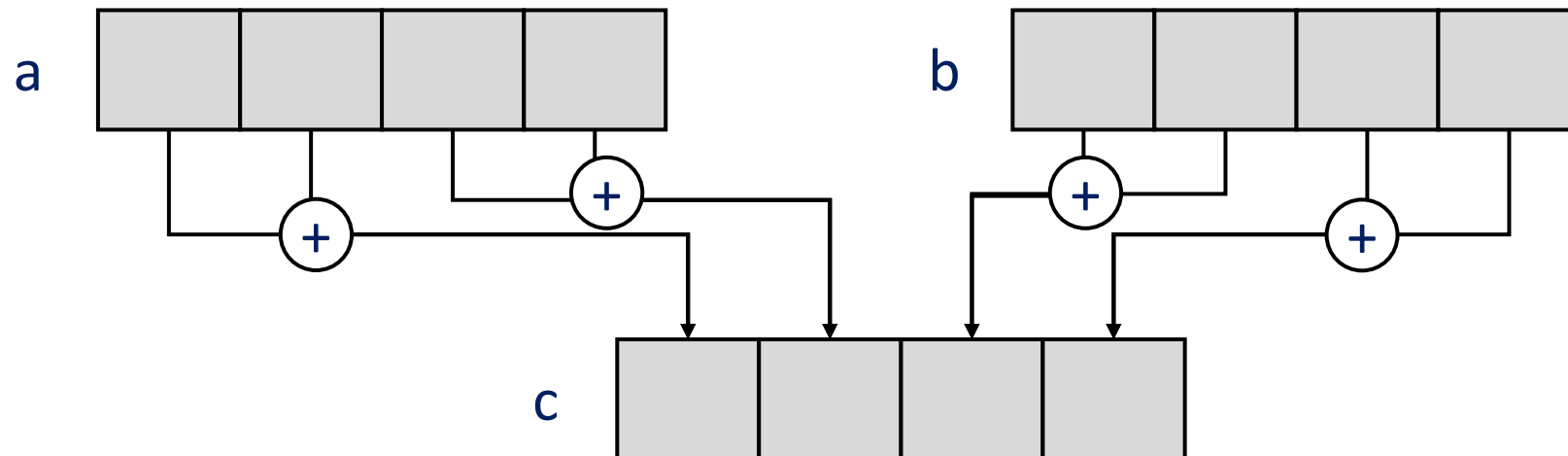  ○ (a*b)+c, -(a*b)-c, -(a*b)+c, and other permutations!
  ○ Consult the guide!
❑ ...

a

× × × ×

b

+ + + +

c

= = = =

d

__m256 a, b, c;
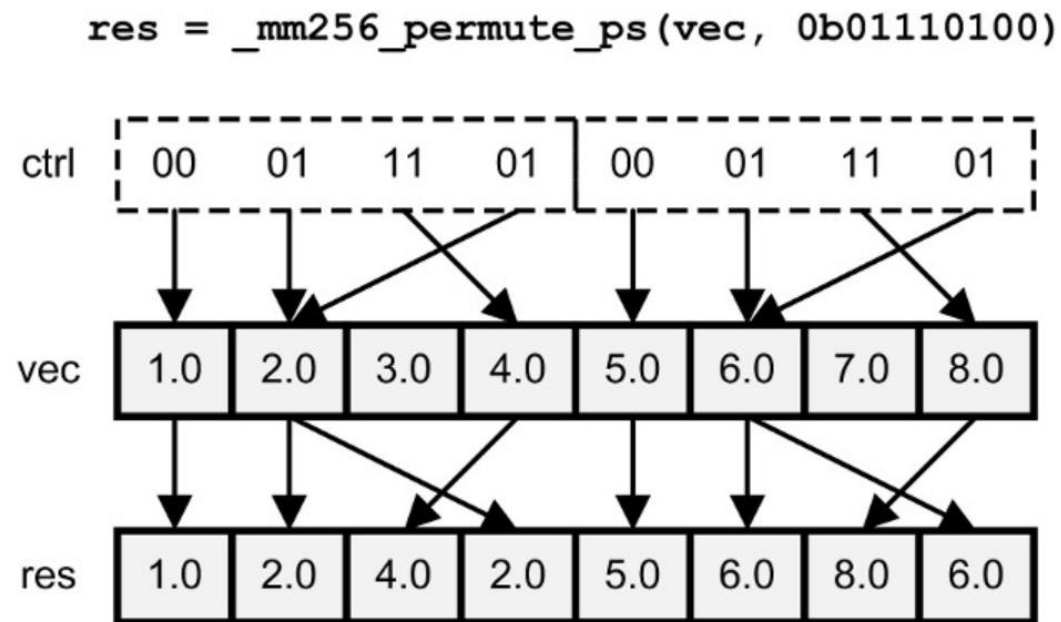__m256 d = _mm256_fmadd_pd(a, b, c);

# Horizontal Vector Instructions

❑ Horizontal add/subtraction
  o Adds adjacent pairs of values
  o E.g., __m256d _mm256_hadd_pd (__m256d a, __m256d b)

# Shuffling/Permutation

❑ Within 128-bit lanes
  o _mm256_shuffle_ps/pd/… (a,b, imm8)
  o _mm256_permute_ps/pd
  o _mm256_permutevar_ps/…

❑ Across 128-bit lanes
  o _mm256_permute2x128/4x64 : Uses 8 bit control
  o _mm256_permutevar8x32/… : Uses 256 bit control

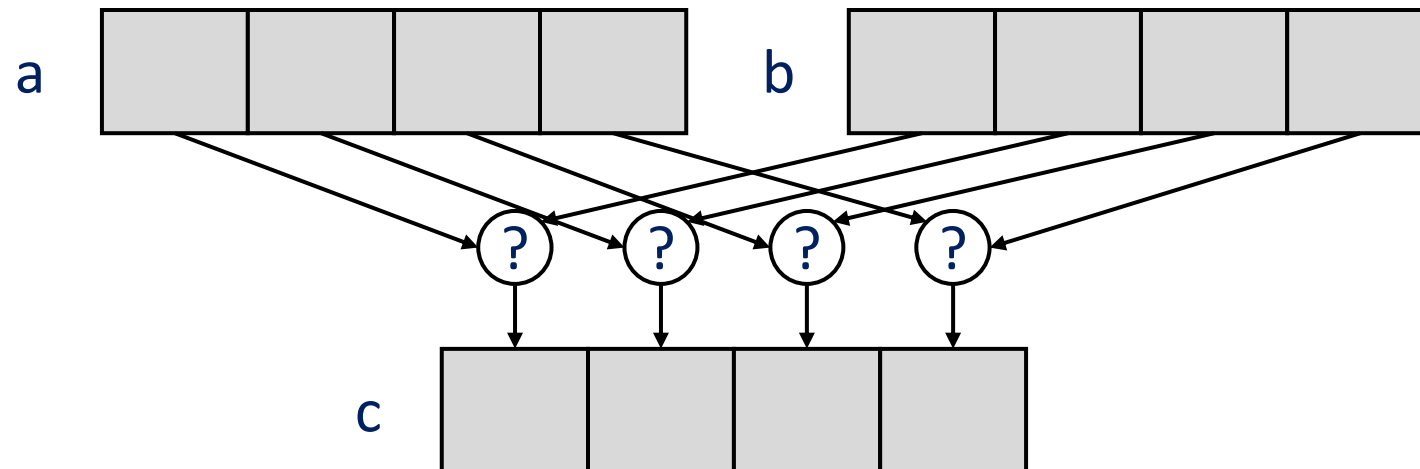❑ Not all type permutations exist for each type, but variables can be cast back and forth between types



```
res = _mm256_permute_ps(vec, 0b01110100)
```

| ctrl | 00 | 01 | 11 | 01 | 00 | 01 | 11 | 01 |

| vec | 1.0 | 2.0 | 3.0 | 4.0 | 5.0 | 6.0 | 7.0 | 8.0 |

| res | 1.0 | 2.0 | 4.0 | 2.0 | 5.0 | 6.0 | 8.0 | 6.0 |

Matt Scarpino, "Crunching Numbers with AVX and AVX2," 2016

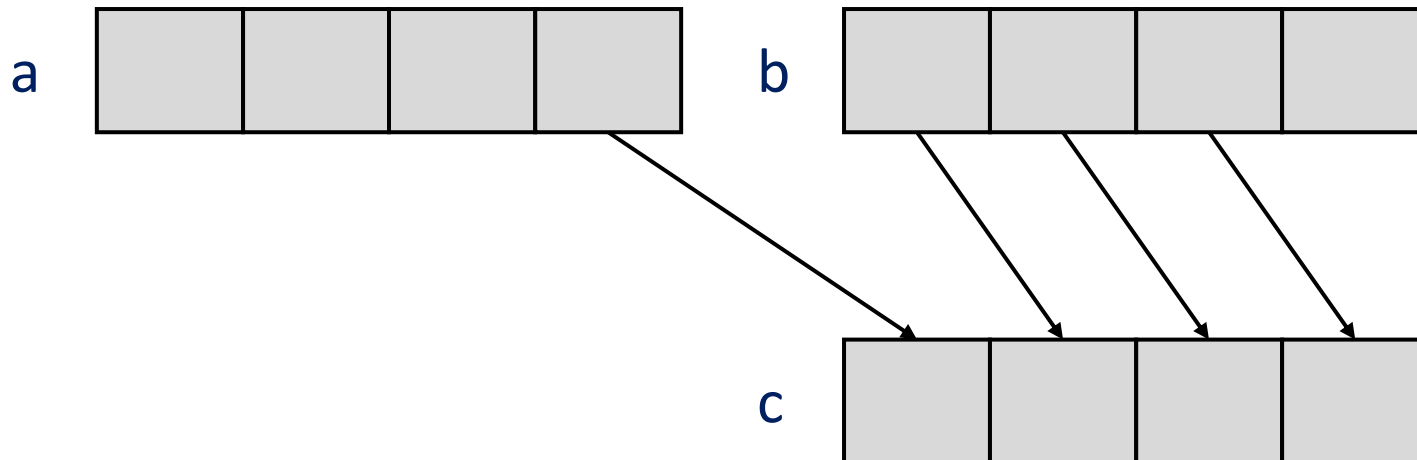# Blend

❑ Merges two vectors using a control
  - ○ _mm256_blend_... : Uses 8 bit control
    - • e.g., _mm256_blend_epi32
  - ○ _mm256_blendv_... : Uses 256 bit control
    - • e.g., _mm256_blendv_epi8

# Alignr

❑ Right-shifts concatenated value of two registers, by byte
- o Often used to implement circular shift by using two same register inputs
- o _mm256_alignr_epi8 (a, b, count)

Example of 64-bit values being shifted by 8

# Helper Instructions

- ❏ Cast
  - o __mm256i <-> __mm256, etc…
  - o Syntactic sugar -- does not spend cycles
- ❏ Convert
  - o 4 floats <-> 4 doubles, etc…
- ❏ Movemask
  - o __mm256 mask to -> int imm8
- ❏ And many more…

# Case Study: Matrix Multiplication

❑ Remember simply transposing matrix B brought 6x performance
   o At that point, we are bottlenecked by single-thread processing performance
   o Adding SIMD gets us more!
   o After this we are again bottlenecked by memory, but that is for another time

A $\qquad$ B $\qquad$ VS $\qquad$ A $\qquad$ B$^T$

63.19 seconds

10.89 seconds
2.20 seconds
(6x performance)
(29x performance!)

# Case Study: Sorting

❑ Important, fundamental application!

❑ Can be parallelized via divide-and-conquer

❑ How can SIMD help?

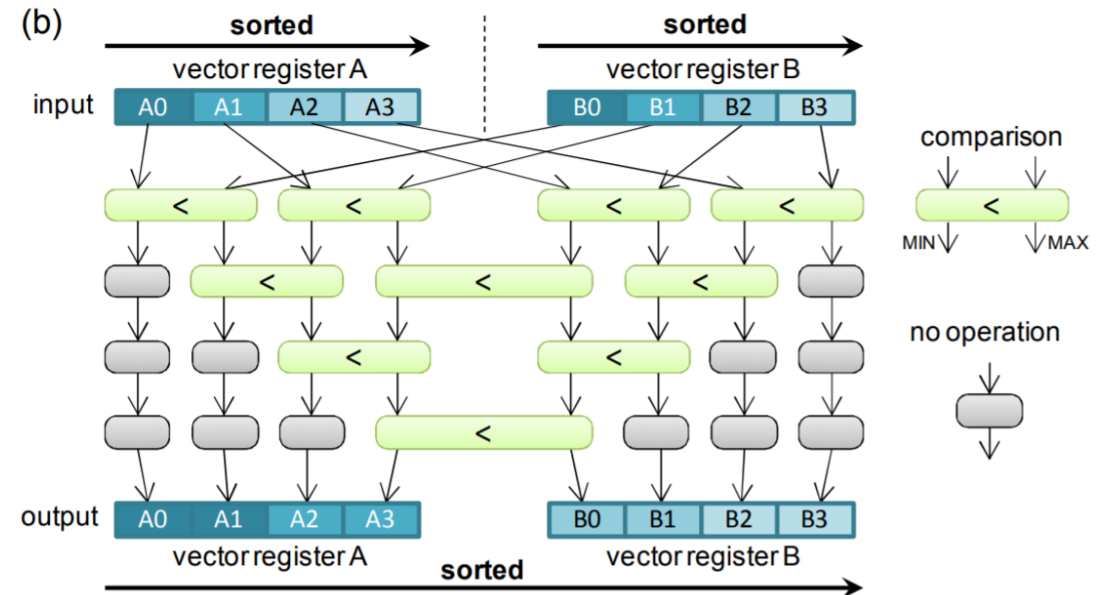# The Two Register Merge

❑ **Sort units of two pre-sorted registers, K elements**
- o minv = A, maxv = B

- o // Repeat K times
  - minv = min(minv,maxv)
  - maxv = max(minv,maxv)
  - // circular shift one value down
  - minv = alignr(minv, minv, sizeof(int))



Inoue et.al., "SIMD- and Cache-Friendly Algorithm for Sorting an Array of Structures," VLDB 2015

# SIMD And Merge Sort

☐ Hierarchically merged sorted subsections

☐ Using the SIMD merger for sorting

    o vector_merge is the two-register sorter from before

```
aPos = bPos = outPos = 0;
vMin = va[aPos++];
vMax = vb[bPos++];
while (aPos < aEnd && bPos < bEnd) {
  /*  merge vMin and vMax */
  vector_merge(vMin, vMax);

  /*  store the smaller vector as output*/
  vMergedArray[outPos++] = vMin;

  /*  load next vector and advance pointer    */
  /*  a[aPos*4] is first element of va[aPos] */
  /*  and b[bPos*4] is that of vb[bPos]       */
  if (a[aPos*4] < b[bPos*4])
      vMin = va[aPos++];
  else
      vMin = vb[bPos++];
}
```
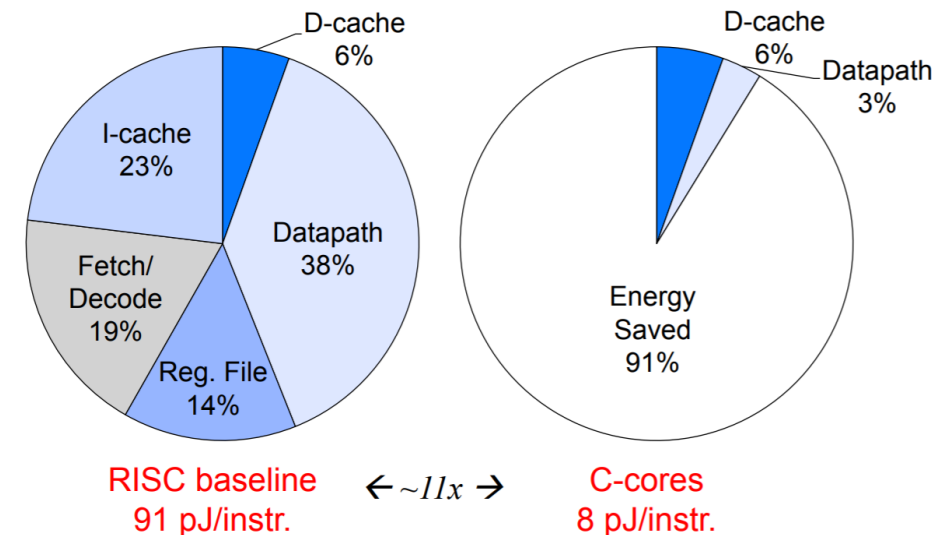
Inoue et.al., "SIMD- and Cache-Friendly Algorithm for Sorting an Array of Structures," VLDB 2015

# Topic Under Active Research!

❑ Papers being written about…

    o Architecture-optimized matrix transposition

    o Register-level sorting algorithm

    o Merge-sort

    o … and more!

❑ Good find can accelerate your application kernel Nx

# Processor Microarchitectural Effects on Power Efficiency

❑ The majority of power consumption of a CPU is not from the ALU
  o Cache management, data movement, decoding, and other infrastructure
  o Adding a few more ALUs should not impact power consumption

❑ Indeed, 4X performance via AVX does not add 4X power consumption
  o From i7 4770K measurements:
  o Idle: 40 W
  o Under load : 117 W
  o Under AVX load : 128 W



D-cache 6%
I-cache 23%
Datapath 38%
Fetch/ Decode 19%
Reg. File 14%

D-cache 6%
Datapath 3%
Energy Saved 91%

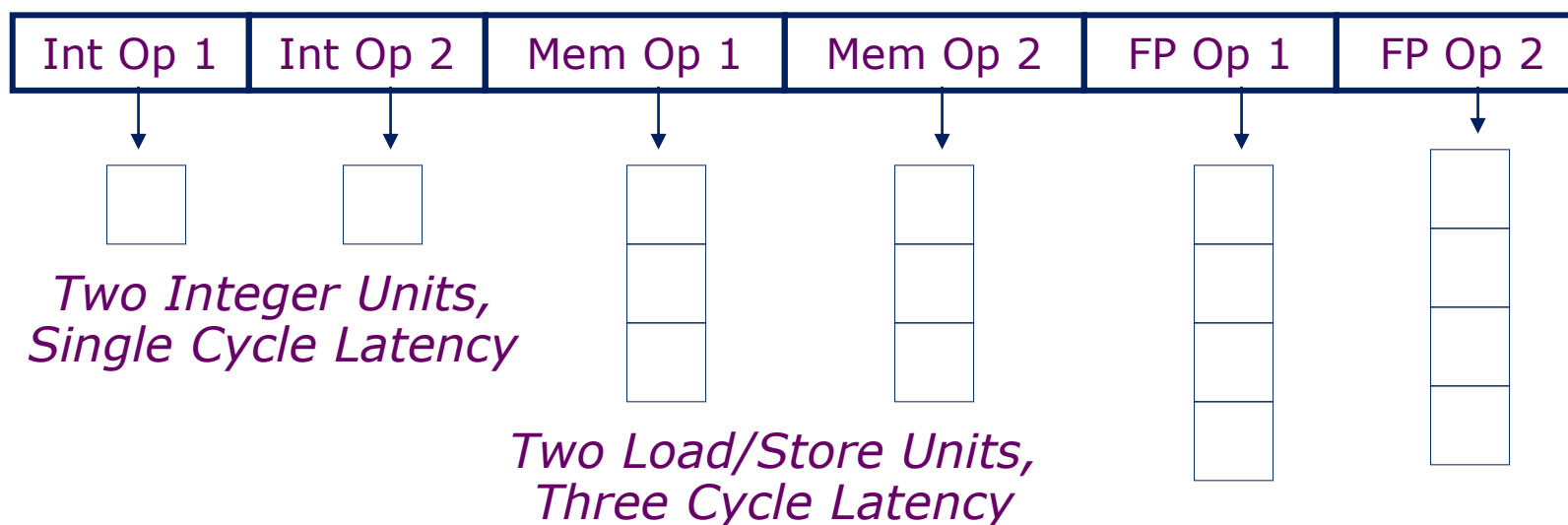RISC baseline 91 pJ/instr.  ← ~11x →  C-cores 8 pJ/instr.

# Very Long Instruction Word (VLIW)

❑ Superscalar does not change the ISA

   o Complicates hardware in charge of detecting dependencies!

❑ What if we changed the ISA, and made the compiler manage ILP?
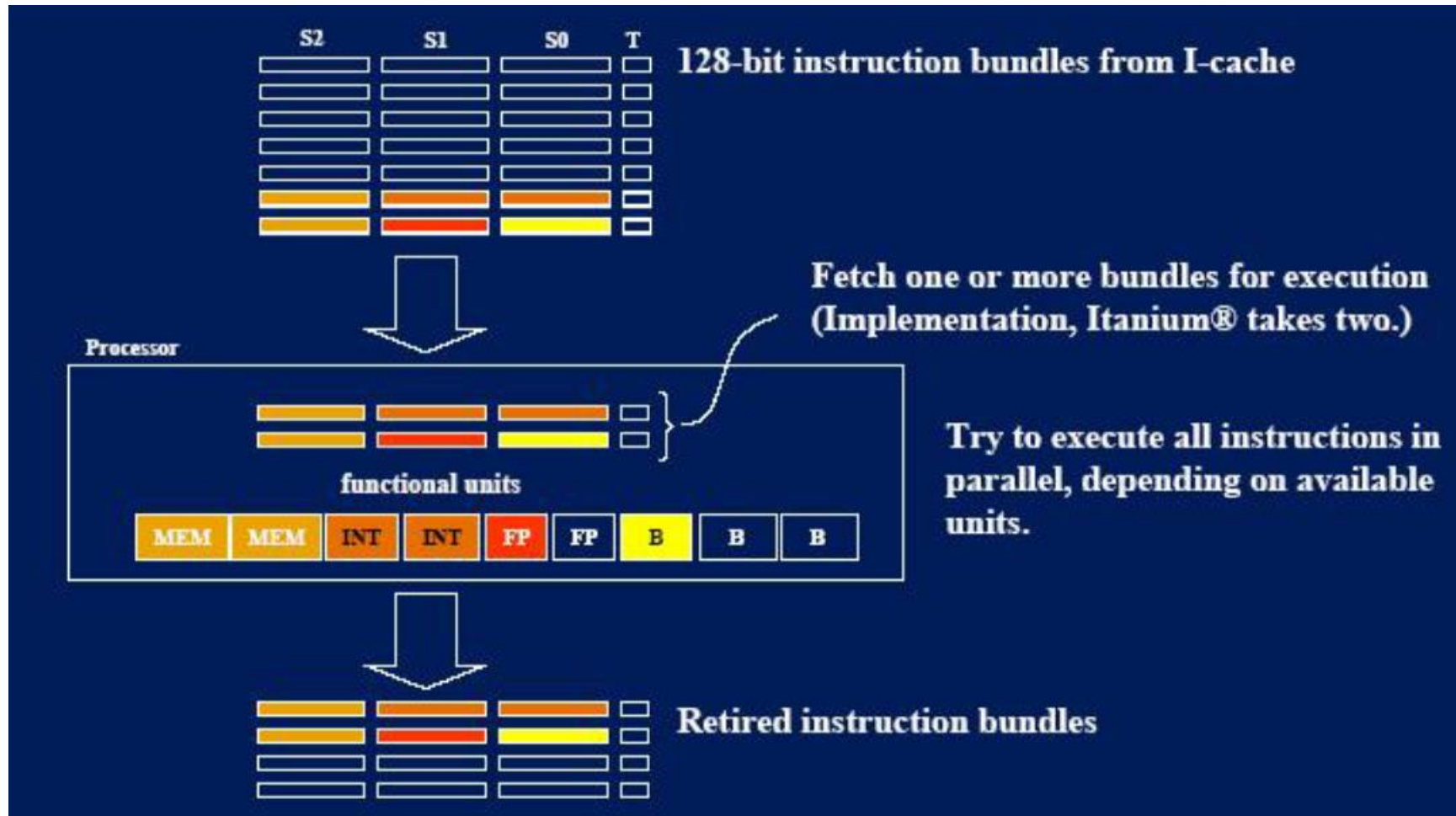

❑ Not in x86/RISC-V/ARM/…

   o Sometimes as accelerator extensions!

   o (RISC-V "V" extension)

# Very Long Instruction Word (VLIW)

❑ Multiple instructions packaged into a Very Long Instruction
  o Sometimes "bundle"

❑ Each execution operation slot has a fixed function (ALU, Mem, FP, etc)

❑ Compiler's responsibility to create efficient instructions
  o Inter-slot dependency is not checked by hardware!

| Int Op 1 | Int Op 2 | Mem Op 1 | Mem Op 2 | FP Op 1 | FP Op 2 |
|----------|----------|----------|----------|---------|---------|

*Two Integer Units,*
*Single Cycle Latency*

*Two Load/Store Units,*
*Three Cycle Latency*

# Intel Explicitly Parallel Instruction Computing (EPIC, Itanium)

# VLIW Characteristics

❑ Very good performance for computation-intensive code

❑ Very bad performance for code with many dependencies/hazards!

    o  Much more sensitive to hazards than single-issue pipelines

    o  Example: short loops

```
for (i=0; i<N; i++)
    B[i] = A[i] + C;
```

*Compile*

```
loop:  ld f1, 0(r1)
       add r1, 8
       fadd f2, f0, f1
       sd f2, 0(r2)
       add r2, 8
       bne r1, r3, loop
```

*Schedule*

loop:

| Int1 | Int 2 | M1 | M2 | FP+ | FPx |
|------|-------|----|----|-----|-----|
| add r1 |     | ld |    |     |     |
|      |       |    |    |     |     |
|      |       |    |    |     |     |
|      |       |    |    | fadd |    |
|      |       |    |    |     |     |
|      |       |    |    |     |     |
| add r2 | bne | sd |    |     |     |
|      |       |    |    |     |     |

How many FP ops/cycle?
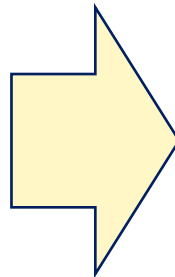
1 fadd / 8 cycles = 0.125

# Compiler's job is important!

❑ e.g., Loop unrolling to keep execution units busy

```
for (i=0; i<N; i++)
    B[i] = A[i] + C;
```

Unroll inner loop to perform 4 iterations at once

```
for (i=0; i<N; i+=4)
{
    B[i]   = A[i] + C;
    B[i+1] = A[i+1] + C;
    B[i+2] = A[i+2] + C;
    B[i+3] = A[i+3] + C;
}
```

*Unroll 4 ways*

```
loop:  ld f1, 0(r1)
       ld f2, 8(r1)
       ld f3, 16(r1)
       ld f4, 24(r1)
       add r1, 32
       fadd f5, f0, f1
       fadd f6, f0, f2
       fadd f7, f0, f3
       fadd f8, f0, f4
       sd f5, 0(r2)
       sd f6, 8(r2)
       sd f7, 16(r2)
       sd f8, 24(r2)
       add r2, 32
       bne r1, r3, loop
```

*Schedule* →

| Int1 | Int 2 | M1 | M2 | FP+ | FPx |
|------|-------|------|------|--------|-----|
|      |       | ld f1 |    |        |     |
|      |       | ld f2 |    |        |     |
|      |       | ld f3 |    |        |     |
| add r1 |     | ld f4 |    | fadd f5 |     |
|      |       |      |    | fadd f6 |     |
|      |       |      |    | fadd f7 |     |
|      |       |      |    | fadd f8 |     |
|      |       | sd f5 |   |        |     |
|      |       | sd f6 |   |        |     |
|      |       | sd f7 |   |        |     |
| add r2 | bne | sd f8 |   |        |     |
|      |       |      |    |        |     |
|      |       |      |    |        |     |

loop:

How many FLOPS/cycle?

4 fadds / 11 cycles = 0.36

# Issues with VLIW

❑ Execution unit configurations change across models
  o How many Integer units, how many float units, neural units …?
  o Cannot be binary compatible across models!
    • Unless hardware provides an abstraction layer…?
    • But that would add scheduler overhead, undermining VLIW (Itanium tried a good balance)

❑ Dependency/hazards difficult for compiler to manage
  o Too many slots end up empty (low performance, large binary)


❑ But when it works well, it works remarkably well
  o e.g., Scientific computing
  o That's why it is often resurrected as potential solution (Itanium, ATi TeraScale, …)

# CS250P: Computer Systems Architecture
## Out of Order Processing

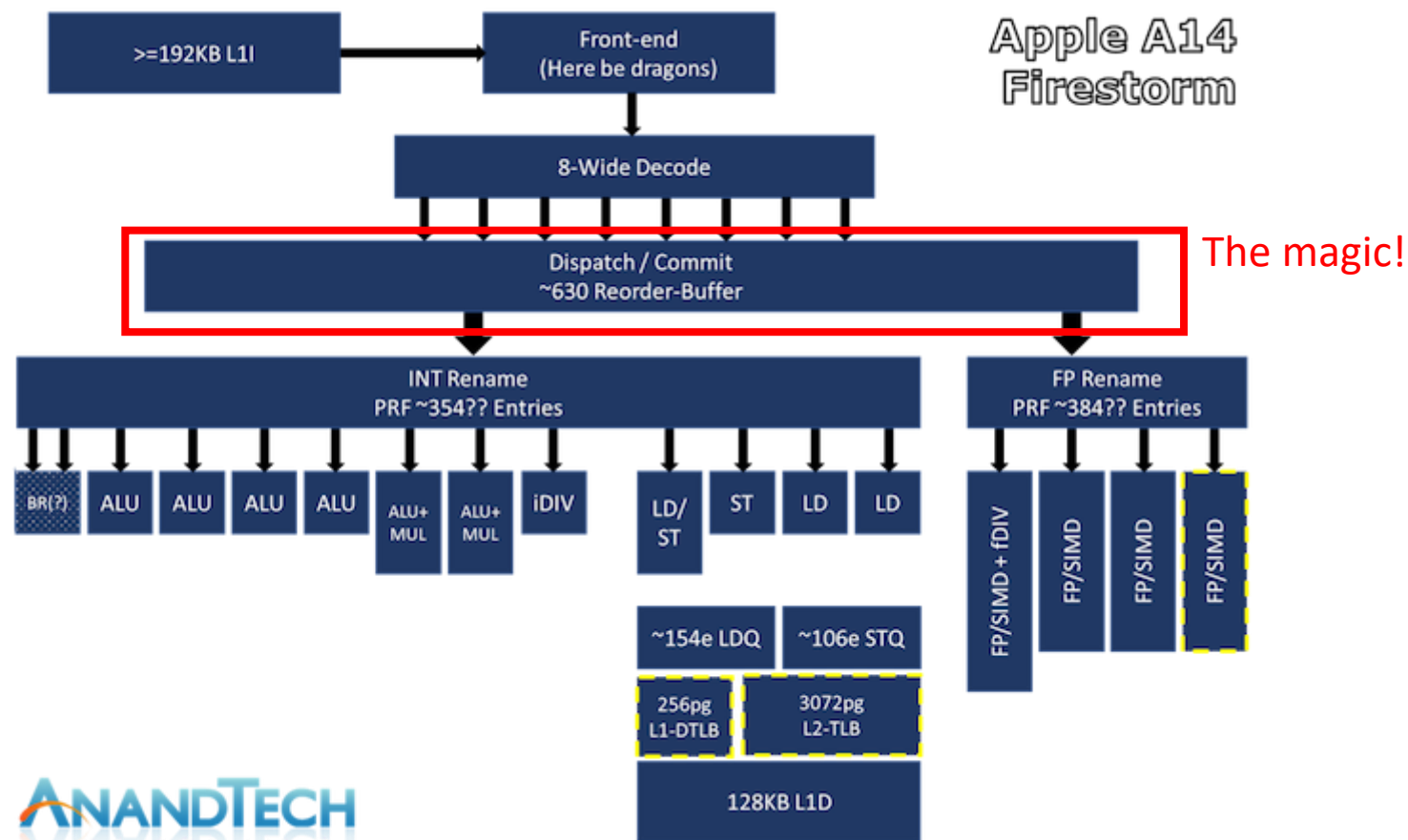Sang-Woo Jun

Fall 2022

# Back to Transparent Parallelism

❑ Explicit parallelism is not as popular as transparent
  o Everyone wants performance for free!

❑ Can we keep execution slots busy, using backwards-compatible single-thread instruction streams?

| Int Op 1 | Int Op 2 | Mem Op 1 | Mem Op 2 | FP Op 1 | FP Op 2 |
|----------|----------|----------|----------|---------|---------|

*Two Integer Units,*
*Single Cycle Latency*

*Two Load/Store Units,*
*Three Cycle Latency*

# Skylake-X Microarchitecture (2019)

# Apple M1 Microarchitecture (2020)

# OoO: Determining dependencies

```
(1)    r8 ← 20
(2)    r1 ← addr1
(3)    r2 ← addr2
(4)    r3 ← addr3

LOOP:
(5)    r4 ← MEM[r1]
(6)    r1 ← r1 + 4
(7)    r5 ← MEM[r2]
(8)    r2 ← r2 + 4
(9)    r6 ← r4 + r5
(10)   MEM[r3] ← r6
(11)   r3 ← r3 + 4
(12)   r8 ← r8 - 1
(13)   bnz r8, LOOP
```

Instruction dependence graph
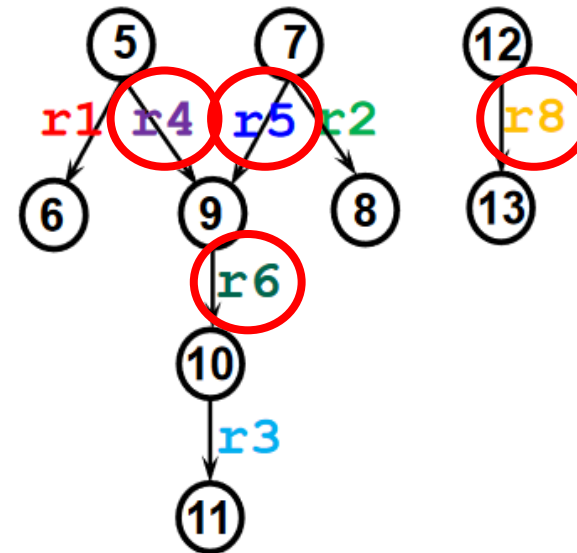
# Data dependency types: RaW

❑ Read-after-Write
   o A "true" dependency
   o We must wait until r5's value is materialized... No other choice

```
(7)    r5 ← MEM[r2]
(9)    r6 ← r4 + r5
```

```
LOOP:
    (5)    r4 ← MEM[r1]
    (6)    r1 ← r1 + 4
    (7)    r5 ← MEM[r2]
    (8)    r2 ← r2 + 4
    (9)    r6 ← r4 + r5
    (10) MEM[r3] ← r6
    (11) r3 ← r3 + 4
    (12) r8 ← r8 - 1
    (13) bnz r8, LOOP
```

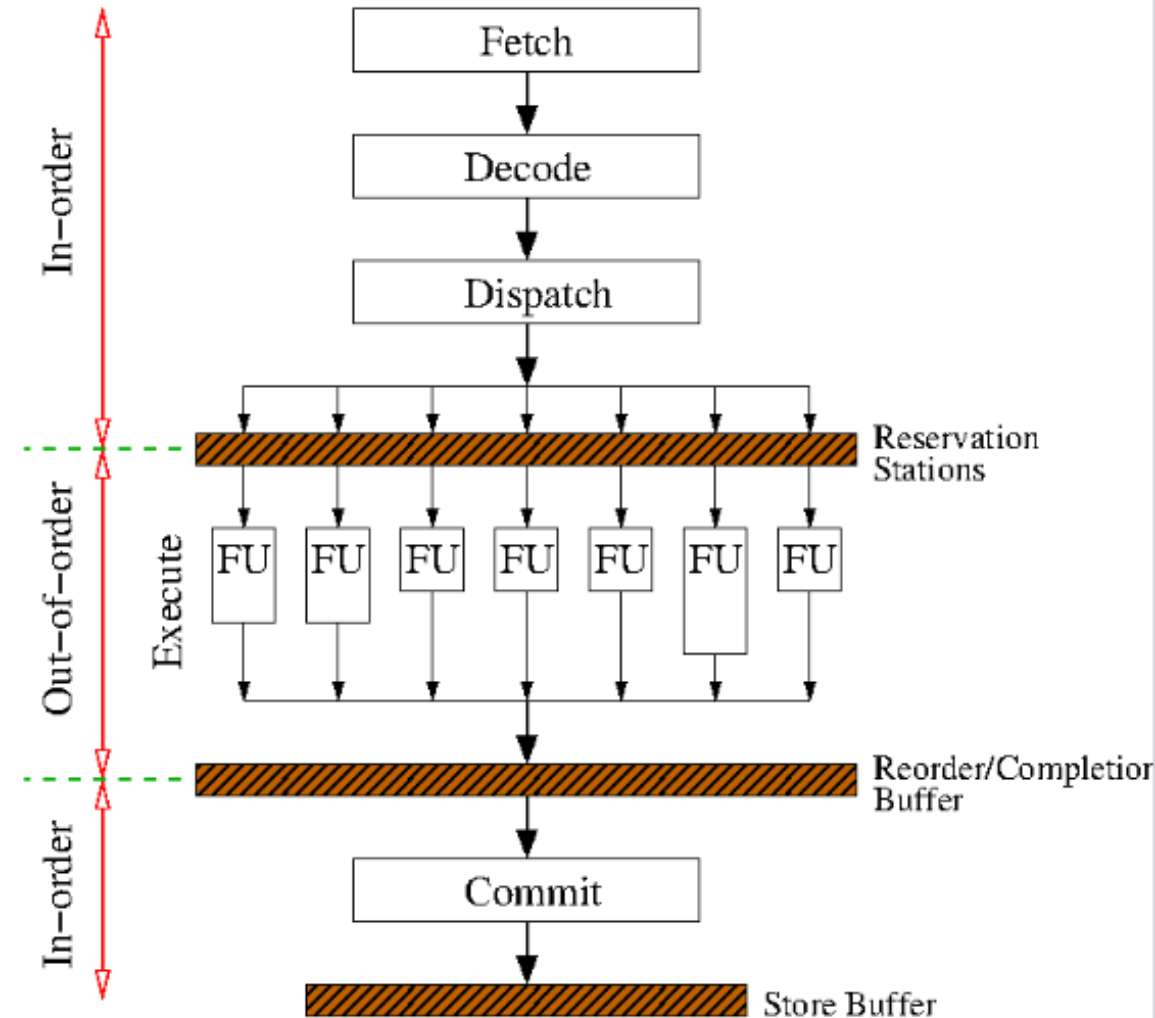# OoO managing dependencies

❑ Looks like dispatch+Commit stages added to VLIW

- o Instructions wait at "reservation stations"
- o Listens to forwarding paths
  - "Is my input operand being written to"
- o Forwarded to FU when ready
  - Out of order

# OoO managing dependencies

❑ Arithmetic can happen OoO, BUT Commits should happen In-Order!

  o Register writes, memory updates, etc

❑ Decoded instructions line up at Reorder Buffer(RoB)

  o Wait until execute results available

  o Wait until branch mispredict ruled out

  o Commits in order of insertion

# Many topics we won't go into today!

❑ Effectively matching available operands to waiting instructions
  o Looping over instructions is too slow
  o N-to-N broadcast is too expensive (slow clocks!)
  o Tomasulo's algorithm!

❑ Precise interrupts become complicated
  o Things are executing OoO, when a breakpoint happens, how do we line things back up for debugging?

# Just one more topic: Register renaming

❑ Not all dependencies are RaW. Some can be resolved!

    ○ Write-after-Read (WaR)

- e.g., 5->6
- "Anti-dependence": r1's value clobbered after (6)
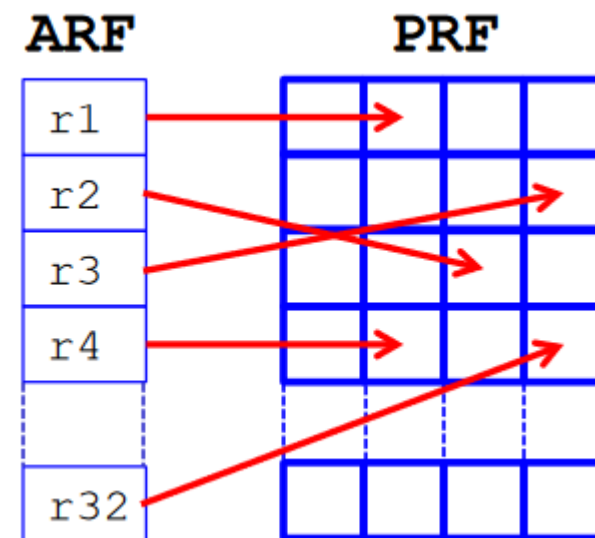- If we had used "r9" instead, no dependency!

    ○ Write-after-Write (WaW)

- e.g., 7->7 across loop iteration
- (7) does not read from "r5"…
- If each loop iteration used a different reg, (r9, r10, r11,…) no dependency!

```
(5)   r4 ← MEM[r1]
(6)   r1 ← r1 + 4
(7)   r5 ← MEM[r2]
(8)   r2 ← r2 + 4
(9)   r6 ← r4 + r5
(10)  MEM[r3] ← r6
(11)  r3 ← r3 + 4
(12)  r8 ← r8 - 1
(13)  bnz r8, LOOP
```

# OoO: Register renaming

❑ Two different concepts of registers
  o "Architectural Registers": Conceptually defined in ISA, software abstraction
    • "RISC-V has 32 registers in the register file"
  o "Physical Registers": Larger number of registers actually in silicon
    • Scheduler dynamically renames registers to an empty slot in the physical register file
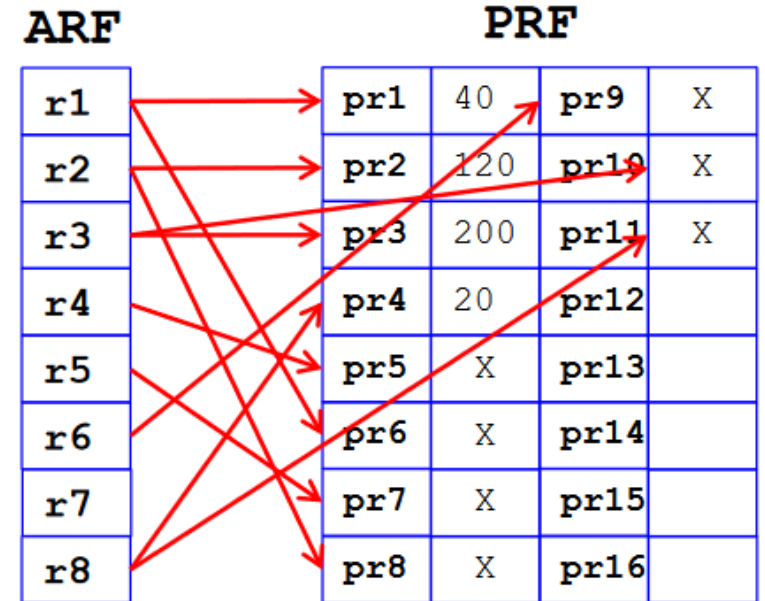    • When WaW or WaR dependencies are discovered

# Register renaming: Previous example...



Original code:
→ (5)   r4 ← MEM[r1]
→ (6)   r1 ← r1 + 4
→ (7)   r5 ← MEM[r2]
→ (8)   r2 ← r2 + 4
→ (9)   r6 ← r4 + r5
→ (10)  MEM[r3] ← r6
→ (11)  r3 ← r3 + 4
→ (12)  r8 ← r8 - 1
→ (13)  bnz r8, LOOP

Post-decode(as seen by RS):
(5)   pr5  ← MEM[40]
(6)   pr6  ← 40 + 4
(7)   pr7  ← MEM[120]
(8)   pr8  ← 120 + 4
(9)   pr9  ← pr5 + pr7
(10)  MEM[200] ← pr9
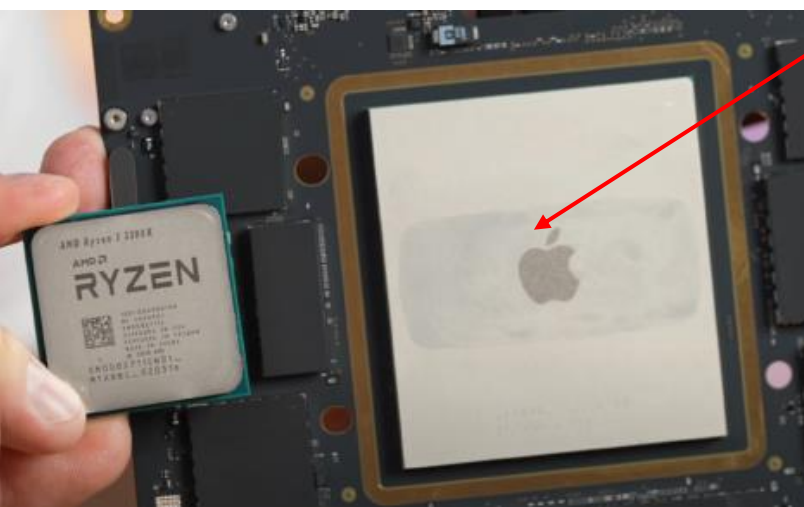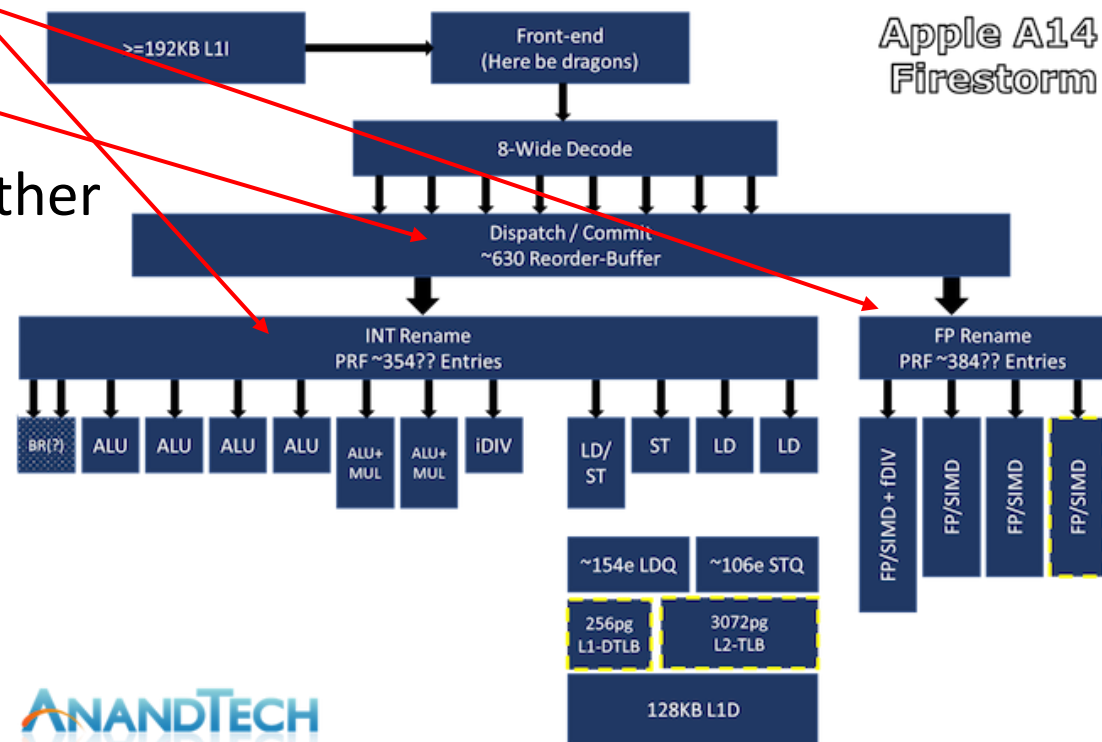(11)  pr10 ← 200 + 4
(12)  pr11 ← 20 - 1
(13)  bnz  pr11, LOOP

ARF
r1
r2
r3
r4
r5
r6
r7
r8

PRF

| pr1 | 40 | pr9 | X |
| pr2 | 120 | pr10 | X |
| pr3 | 200 | pr11 | X |
| pr4 | 20 | pr12 | |
| pr5 | X | pr13 | |
| pr6 | X | pr14 | |
| pr7 | X | pr15 | |
| pr8 | X | pr16 | |

# Back to timely example: Apple M1

❑ Really good single-thread performance!

❑ How?

RISC!

- "8-wide decoder" […] "16 execution units (per core)"
- "(Estimated) 630-deep out-of-order"
- "Unified memory architecture"
- Hardware/software optimized for each other

M1 Ultra
Image source: wccftech