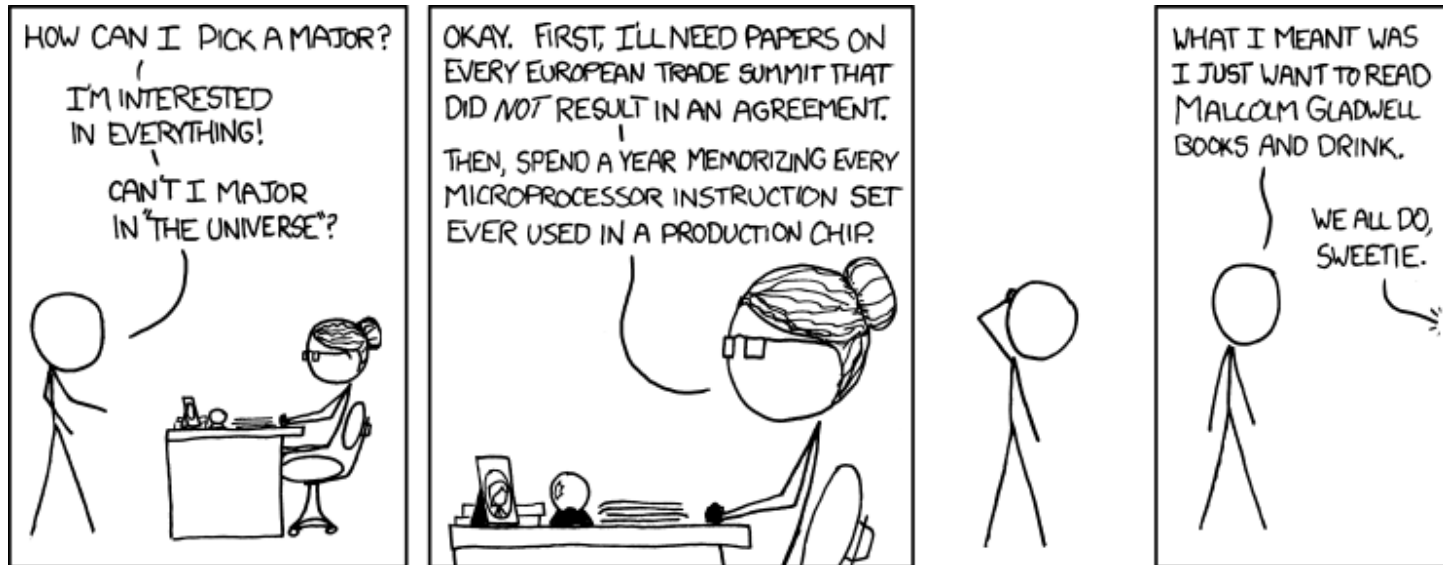


CS/ECE 752: Advanced Computer Architecture I



Prof. Matthew D. Sinclair

Source: XKCD

Hardware/Software Interfaces

Slide History/Attribution Diagram:

UW Madison
Hill, Sohi,
Smith, Wood

UPenn
Amir Roth,
Milo Martin

Various Universities
Asanovic, Falsafi, Hoe, Lipasti,
Shen, Smith, Vijaykumar

UW Madison
Hill, Sohi, Wood,
Sankaralingam, Sinclair

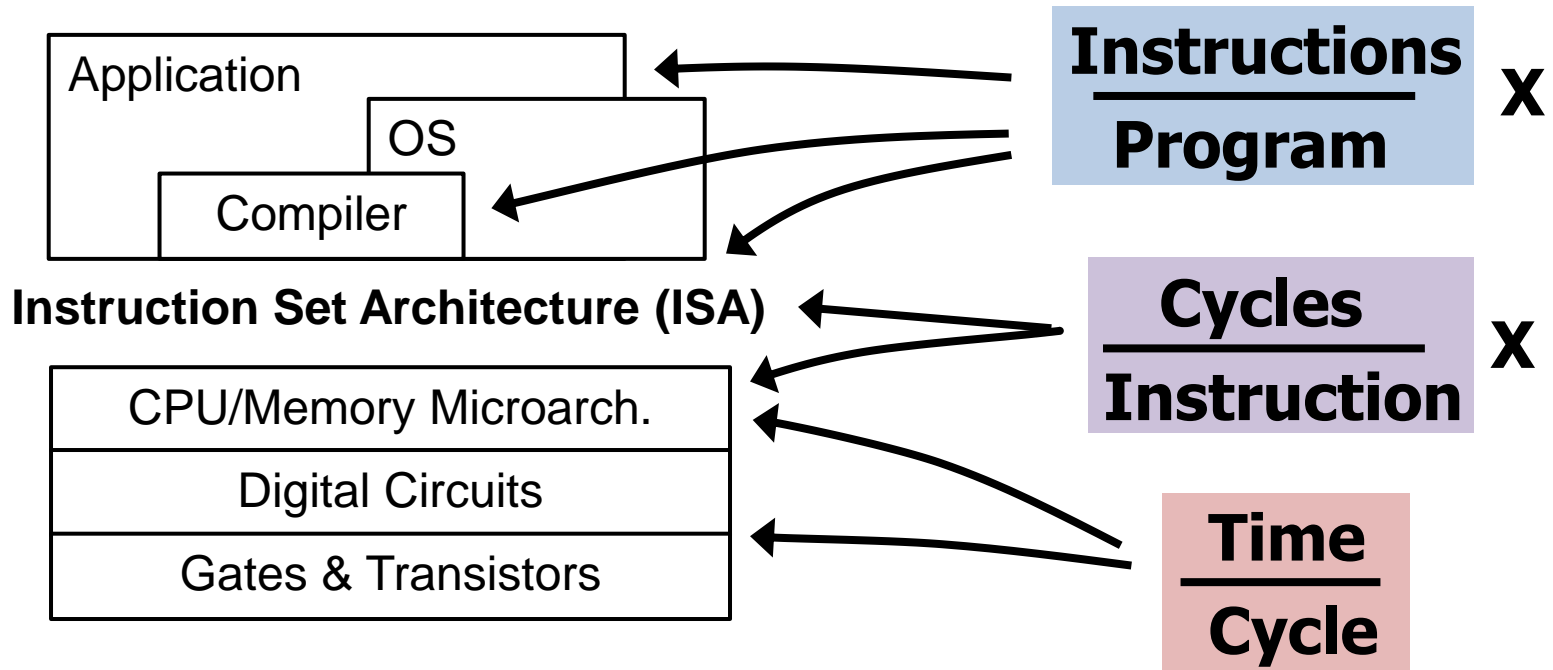
UCLA
Nowatzki

Announcements 9/15/24

- HW2 Due Saturday (one small update – see Piazza)
 - Posting FAQs on Piazza – let me know if useful
- HW1 grading hopefully done tonight
- Review 1 grading done
- Review 2 due Wednesday
- I'll work on being more concise
- Going to try and get through this unit in 1 lecture (lots of 252/552 review) to keep on schedule
 - Skip some slides, posted for reference
- Advanced Topic Winner: ML
 - Posted materials – see Canvas/class schedule
 - Advanced Branch Prediction will partially be folded into Dyn Sched
 - Other topics: feel free to email me

Computer System Layers

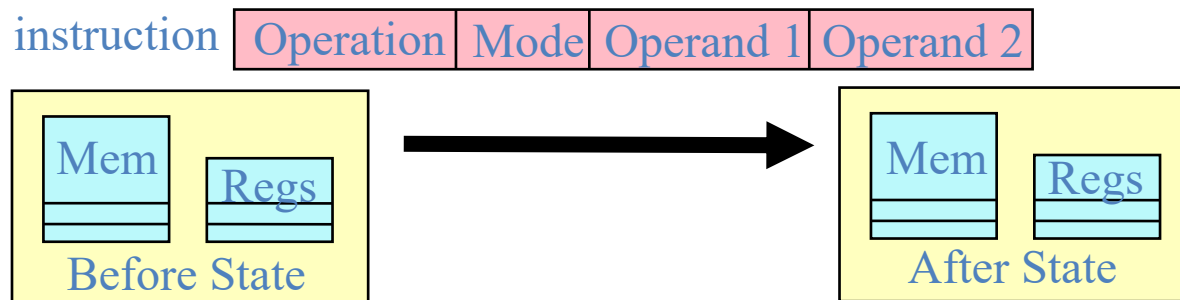
**Main
Focus
for
Today**



Basic Concepts / Review

What Is An ISA?

- The “**contract**” between software and hardware
 - **Functional definition** of operations, modes, and storage locations supported by hardware
 - **Precise description** of how to invoke, and access them
- Instruction Concept



- Execution Model: When to perform each instruction
- ISAs can hide some details about hardware choices:
 - How operations are implemented
 - Which operations are fast and which are slow and when
 - Which operations take more power and which take less

A Language Analogy for ISAs

- An ISA is analogous to a human language
 - **Allows communication**
 - Language: person to person
 - ISA: hardware to software
 - **Many common aspects**
 - Part of speech: verbs, nouns, adjectives, adverbs, etc.
 - Common operations: calculation, control/branch, memory
 - **Many different languages/ISAs, mostly they are similar...**
 - Sometimes fundamental differences
 - **Both evolve over time**
- Key differences: ISAs must be **unambiguous**
 - ISAs are **explicitly** engineered and extended

Program Compilation

```
int array[100], sum;
void array_sum() {
    for (int i=0; i<100;i++) {
        sum += array[i];
    }
}
```

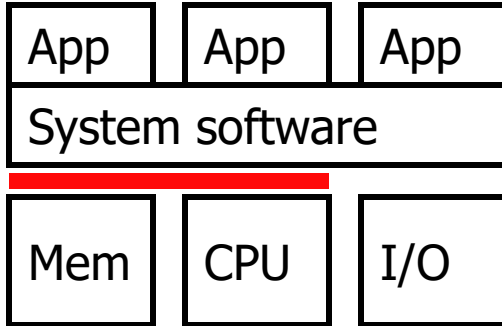
- **Program** written in a “high-level” programming language
 - C, C++, Java, C#
 - Hierarchical, structured control: loops, functions, conditionals
 - Hierarchical, structured data: scalars, arrays, pointers, structures
- **Compiler**: translates program to **assembly**
 - Parsing and straight-forward translation
 - Compiler also optimizes
 - Compiler itself another application ... who compiled compiler?

Assembly & Machine Language

- **Assembly language**
 - Human-readable representation
- **Machine language**
 - Machine-readable representation
 - 1s and 0s (often displayed in “hex”)
- **Assembler**
 - Translates assembly to machine

<u>Machine code</u>	<u>Assembly code</u>
x9A00	CONST R5, #0
x9200	CONST R1, array
xD320	HICONST R1, array
x9464	CONST R2, sum
xD520	HICONST R2, sum
x6640	LDR R3, R1, #0
x6880	LDR R4, R2, #0
x18C4	ADD R4, R3, R4
x7880	STR R4, R2, #0
x1261	ADD R1, R1, #1
x1BA1	ADD R5, R5, #1
x2B64	CMPI R5, #100
x03F8	BRn array_sum_loop

Example Assembly Language & ISA



- **MIPS**: example of real ISA
 - 32/64-bit operations
 - 32-bit insns
 - 64 registers (32 integer, 32 FP)
 - ~100 different insns
 - Full OS support

Example code is MIPS, but
all ISAs are similar at some level

```
.data
array: .space 100
sum:   .word 0

.text
array_sum:
    li $5, 0
    la $1, array
    la $2, sum
array_sum_loop:
    lw $3, 0($1)
    lw $4, 0($2)
    add $4, $3, $4
    sw $4, 0($2)
    addi $1, $1, 1
    addi $5, $5, 1
    li $6, 100
    blt $5, $6, array_sum_loop
```

Question:

1. ISA Objectives?

2. ISA Considerations/Aspects?

What Makes a Good ISA?

- **Programmability**
 - Easy to express programs efficiently? (human or compiler)
- **Implementability**
 - Easy to design high-performance implementations?
 - More recently
 - Easy to design low-power implementations?
 - Easy to design high-reliability implementations?
 - Easy to design low-cost implementations?
- **Compatibility**
 - Easy to maintain programmability (implementability) as languages and programs evolve?
 - x86 (IA32) generations: 8086, 286, 386, 486, Pentium, PentiumII, PentiumIII, Pentium4,...

Programmability


- How easy is it to express programs efficiently?
 - For whom?
- Early: **human** *(Today quantum is here)*
 - Compilers were terrible, most code was hand-assembled
 - Want high-level coarse-grain instructions
 - As similar to high-level language as possible
- Last several decades: **compiler** *tension*
 - Optimizing compilers usually generate better code than you or I
 - Want low-level, simple, fine-grain instructions
 - Hard for the compiler to target higher-level more complex idioms

Human Programmability

- What makes an ISA easy for a human to program in?
 - Proximity to a high-level language (HLL)
 - Closing the “**semantic gap**”
 - Semantically heavy (CISC-like) insns that capture complete idioms
 - “Access array element”, “loop”, “procedure call”
 - Example: SPARC **save/restore**
 - Bad example: x86 **rep movsb** (copy string)
 - Ridiculous example: VAX **insque** (insert-into-queue)
 - “**Semantic clash**”: what if you have many high-level languages?

Compilers 101

- Compiler goals:
 - all correct programs execute correctly
 - most compiled programs execute fast
 - compile fast
 - provide support for debugging
- Use multiple phases to manage complexity
 - Lexical analysis (e.g., "+" means "add", "foobar" is an identifier)
 - Parsing (e.g., "x = a + b" means assign sum of variables a and b to x)
 - Generates intermediate representation
 - Optimization & code generation (transforms intermediate representation)
 - Procedure In-lining, Loop optimizations, Common sub-expression elimination, Jump optimization, Constant propagation, Register allocation, Strength reduction, Pipeline scheduling, Interprocedural analysis
 - Generation of assembly code



Which comes first?
Phase ordering
problem.

Compiler Programmability

- What makes an ISA easy for a compiler to program in?
 - Low level primitives from which solutions can be synthesized
 - Compilers good at breaking complex structures to simple ones
 - Requires decomposition
 - Not so good at combining simple structures into complex ones
 - Requires search, pattern matching (why AI is hard)
 - Easier to synthesize complex insns than to compare them
 - Rules of thumb
 - Regularity: “**principle of least astonishment**”
 - Orthogonality & composability

Implementability

- Every ISA can be implemented
 - Not every ISA can be implemented efficiently (at least easily)
- Classic high-performance implementation techniques
 - Pipelining, parallel execution, out-of-order execution (more later)
- Certain ISA features make the above difficult
 - Variable instruction lengths/formats: complicate decoding (x86)
 - Implicit state: complicates dynamic scheduling
 - Variable latencies: complicates scheduling
 - Difficult to interrupt instructions: complicate many things
 - A solution: High-performance x86 machines dynamically translate CISC instructions into internal micro-ops (e.g., RISC-ops)

Compatibility

- No-one buys new hardware... if it requires new software (traditionally)
 - IBM established ISA for mainframes; Intel for PCs
 - ISA must remain compatible, no matter what
 - x86 arguably one of the worst ISAs EVER, but survives
 - As does IBM's 360/370/390 (the *first* "ISA family")
 - This was traditionally true, but has been changing in recent years
- **Backward compatibility**
 - New processors must support old programs
 - Can't drop features, but can deprecate and emulate
 - Traditionally very important (except if you're ARM)
- **Forward (upward) compatibility** (usually lower perf.)
 - Old processors must support new programs (with software help)
 - New processors redefine only previously-illegal opcodes
 - Allow software to detect support for specific new instructions (CUDA)
 - Old processors emulate new instructions in low-level software

The Compatibility Trap

- Easy compatibility requires forethought
 - Temptation: use some ISA extension for 5% performance gain
 - Frequent outcome: gain diminishes, disappears, or turns to loss
 - Must continue to support gadget for eternity
- Example: register windows (SPARC)
 - Idea: Reduce register spills and fills for function calls by adding:
Implicit stack of registers, output registers -> input registers
 - Adds cost and complexity to out-of-order implementations of SPARC
- Example: branch delay slot (most RISCs)
 - Idea: Eliminates branch hazard in simple 5-stage pipeline by:
always executing the instruction after a branch
 - Complicates multi-instruction issue (superscalar)

The Compatibility *Trap Door*

- Compatibility's friends
 - **Trap**: instruction makes low-level "function call" to OS handler
 - **Nop**: "no operation" - instructions with no functional semantics
- Backward compatibility
 - Handle rarely used but hard to implement "legacy" opcodes
 - Define to trap in new implementation and emulate in software
 - Rid yourself of some ISA mistakes of the past
 - Problem: performance suffers for legacy codes
- Forward compatibility
 - Reserve sets of trap & nop opcodes (don't define uses)
 - Add ISA functionality by overloading traps
 - Release firmware patch to "add" to old implementation
 - Add ISA hints by overloading nops

Blocking the Compatibility Trap Door

- Temptation:
 - Define “unused” instruction fields as “don’t cares”
 - E.g., MIPS “shift length” field in an “add” instruction
 - Simplifies hardware logic needed to decode instructions
- Problem:
 - Can’t use “unused” values for new instructions
 - Same problem for special registers (e.g., Interrupt status register)
- Solution:
 - Define all bits (usually to be zero).
 - Undefined instructions are trapped

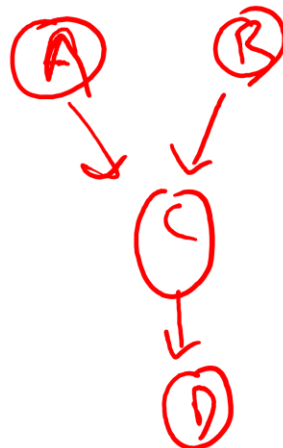
Aspects of ISAs

ISA Aspects

- Execution Model (Implicit vs Explicit Dependences)
 - Von Neumann vs Dataflow
- Implicit vs Explicit Parallelism Specification
 - Scalar vs Vector vs VLIW
- Format (Variable vs Fixed-width)
- Operand Model (where is temporary data stored)
 - Register, Stack, Accumulator
- Addressing Modes (how do you specify addresses)
 - Register-Indirect, Displacement, Index-base / *Base + Offset*
- Memory Abstractions (physical/virtual addr. space)
- Control (condition, branch/predication)
- Details: Register file size, Datatypes (float vs int, bitwidth), Instruction Complexity -- loops

Execution Model

- How to decide when to execute an instruction?
- Option 1: Execute them in order
 - Most high level imperative languages are anyways sequential...
 - Use original program order
- Option 2: Let dependences decide
 - Instructions trigger the execution of dependent instructions
 - Parallel execution "automatic" as instructions execute



] A & B can be parallel

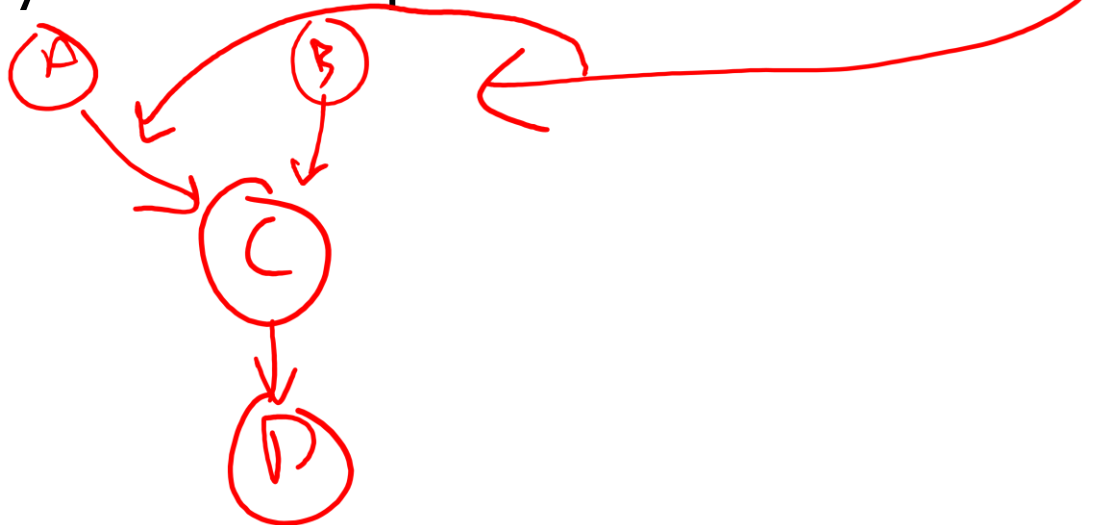
Von Neumann (sequential) Model

- Implicit model of modern commercial ISAs (ARM, x86, etc.)
 - Called von Neumann, but in ENIAC design before*
- Basic feature: the **program counter (PC)**
 - Defines **total order** on dynamic instruction
 - Next PC is PC++ unless insn says otherwise
 - Branch instructions change PC
 - Order and **named storage** define computation
 - Value flows from insn X to Y via storage A iff...
 - X names A as output, Y names A as input...
 - And Y after X in total order
 - (and no other instructions between X&Y outputs to A)

* Probably should be called Eckert/Mauchly Model ...

Dataflow Model

- Dependences encoded explicitly in the ISA
 - E.g., $X \rightarrow Y$ is explicit (as opposed to implicit communication through general purpose registers/memory)
- Dataflow firing rule:
 - Instructions execute when all operands are ready
 - Upon completion, data values are forwarded to instructions which consume them.
- Lacks PC: Unnecessary because dependencies sufficient



Von Neumann/Dataflow Tradeoffs

- Von Neumann
 - Same sequential model as imperative programs
 - easy to reason about, easy to translate (not parallel)
 - Very easy to “pause” execution at any point in the program
 - Total ordering of instructions makes speculation recovery simpler
 - Total ordering of memory eases dependence/alias detection
- Dataflow
 - Multiple instructions may be ready simultaneously -> parallelism
 - Much harder to pause
 - Many instructions active – more state to save on context switch
 - Hard to reason about (debugging)
 - Related problem: unbounded state...
- One Solution: Hybrid!
 - Von Neumann Outside / Dataflow Inside
 - Bounded state, can still recover, still get parallelism
 - E.g., Out-of-order cores; dataflow processors: TRIPS/Wavescalar

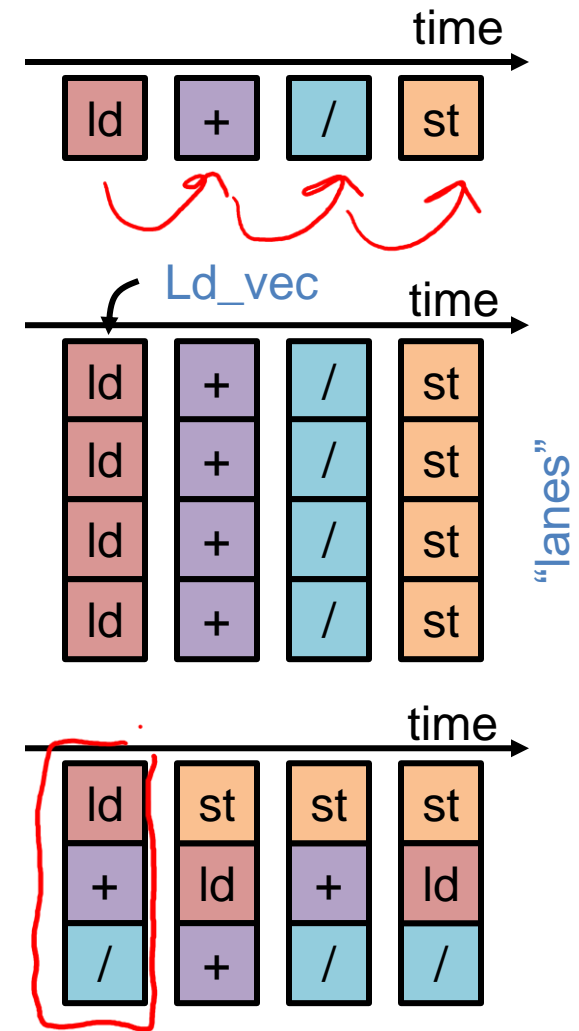
- Expressing *dependences* in the ISA (e.g. via dataflow) is one way to expose parallelism to microarchitecture
- Alternate way to expose parallelism:
 - **Express independence ?**

Implicit vs Explicit Parallelism

Implicit vs Explicit Parallelism

```
for i = 0 -> N  
  b[i] = (a[i]+1)/2
```

- Possibility 1: Scalar Instructions
 - No parallelism explicitly
 - Register Datatype: Scalar Types
- Possibility 2: Vector Instructions (GPU)
 - Registers: Wide Vectors (e.g. 2-128 elem)
 - Single instruction over multiple data (SIMD*)
 - Advantages:
 - Less instructions
 - Parallelism explicit (h/w simpler)
 - Issue: Need data-level parallelism...
- Possibility 3: VLIW: Very long inst. word
 - Instruction packets specify parallel insns
 - Registers: Scalar (more-or-less)
 - Advantage: Parallelism explicit (h/w simpler)
 - Issue: Need ILP and regular memory :)
- Worst part of 2&3: Many Compiler Challenges!



for ($i=0; i < N; i+=4$) {

$v_0 = \text{simd_ld}(a[i]);$

$v_1 = \text{simd_add}(v_0, 1);$

$v_2 = \text{simd_div}(v_1, 2);$

$\text{simd_st}(b[i], v_2);$

}

State-of-the-Art in Implicit vs Explicit

- Possible to combine above techniques in a single ISA, or only use them separately.
- Today there are basically two main approaches for CPUs:
 - Scalar ISA + vector extensions
 - Extensions useful for when lots of DLP *+ SVE*
 - E.g., X86 SSE/MMX/AVX, ARM Neon, IBM Altivec
 - “Pure” VLIW Processors
 - Useful for specialized tasks like signal processing
 - E.g., TI DSPs, also Intel IA64 (*Itanium*)
- Fancier SIMD: SIMT – Single Instruction Multiple Thread
 - Basic Idea: broadcast one instruction at a time to multiple threads, but each thread has its own logical context (so threads can diverge)
 - This is used by GPUs

ISA Aspects

- Execution Model (Implicit vs Explicit Dependences)
 - Von Neumann vs Dataflow
- Implicit vs Explicit Parallelism Specification
 - Scalar vs Vector vs VLIW
- Format (Variable vs Fixed-width)
- Operand Model (where is temporary data stored)
 - Register, Stack, Accumulator
- Addressing Modes (how do you specify addresses)
 - Register-Indirect, Displacement, Index-base
- Memory Abstractions (physical/virtual addr. space)
- Control (condition, branch/predication)
- Details: Register file size, Datatypes (float vs int, bitwidth), Instruction Complexity -- loops

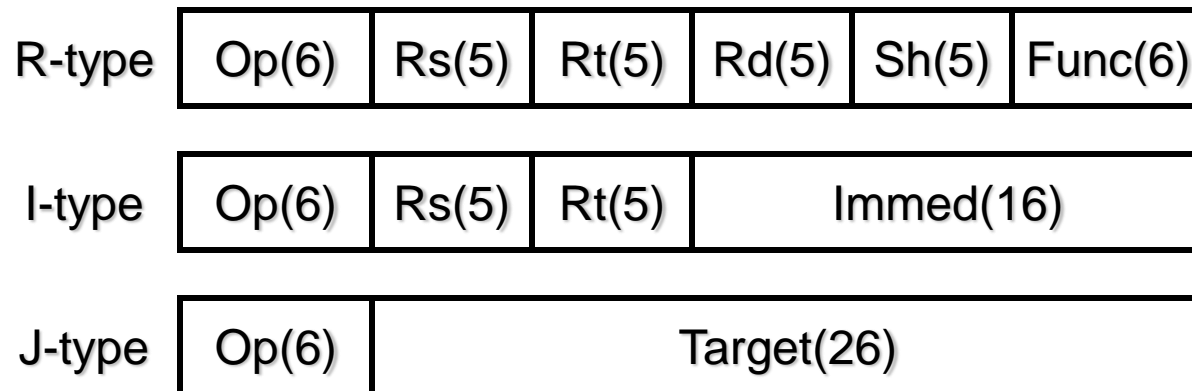
✓
✓
skip

Format

- **Length**
 - Fixed length
 - Most common is 32 bits
 - + Simple implementation: compute next PC using only PC
 - Code density: 32 bits to increment a register by 1?
 - x86 can do this in one 8-bit instruction
 - Variable length
 - Complex implementation
 - + Code density
 - Compromise: two lengths
 - MIPS16 or ARM's Thumb
- **Encoding**
 - A few simple encodings simplify decoder implementation

Simple Example: MIPS Format

- Length
 - 32-bits
- Encoding
 - 3 formats, simple encoding
 - Q: how many instructions can be encoded? A: 64? 127? 4096?



Not Simple Example: X86 Format

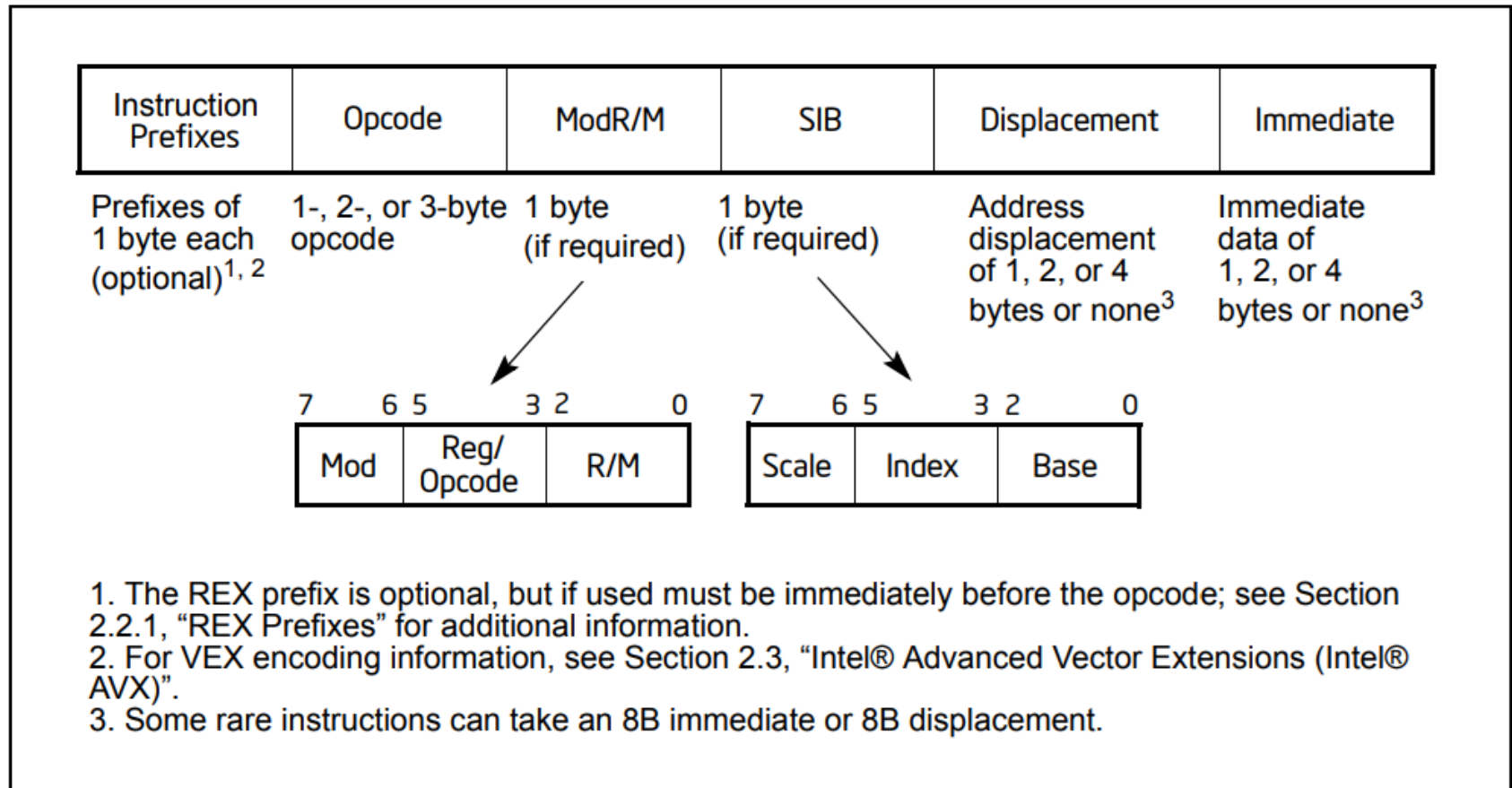


Figure 2-1. Intel 64 and IA-32 Architectures Instruction Format

Operand Model:

What to do with temporaries?

Operand Model: Memory Only

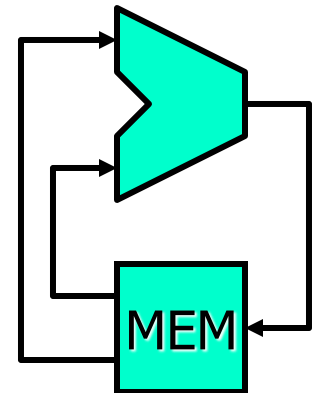
- Where (other than memory) can operands come from?
 - And how are they specified?
 - Example: $A = B + C$
 - Several options

- **Memory only**

`add B,C,A`

$\text{mem}[A] = \text{mem}[B] + \text{mem}[C]$

- Not practical



Operand Model: Accumulator

- **Accumulator**: implicit single element storage

load B

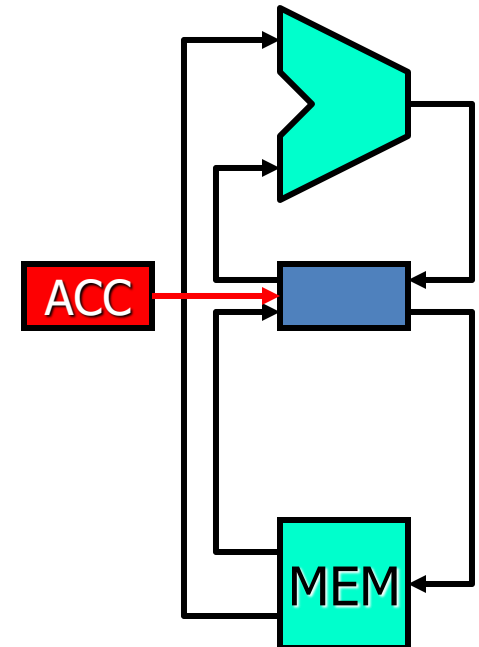
$ACC = mem[B]$

add C

$ACC = ACC + mem[C]$

store A

$mem[A] = ACC$



Operand Model: Stack

- **Stack:** TOS implicit in instructions

push B

$\text{stk}[\text{TOS}++] = \text{mem}[\text{B}]$

push C

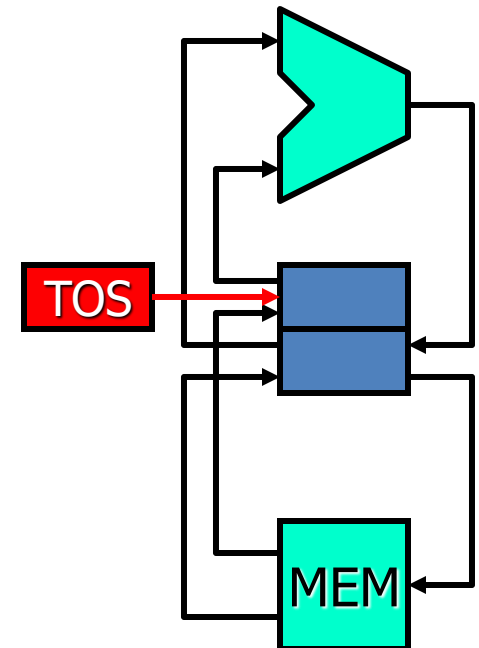
$\text{stk}[\text{TOS}++] = \text{mem}[\text{C}]$

add

$\text{stk}[\text{TOS}++] = \text{stk}[--\text{TOS}] + \text{stk}[--\text{TOS}]$

pop A

$\text{mem}[\text{A}] = \text{stk}[--\text{TOS}]$



Operand Model: Registers

- **General-purpose register:** multiple explicit accumulator

load B,R1

$R1 = \text{mem}[B]$

add C,R1

$R1 = R1 + \text{mem}[C]$

store R1,A

$\text{mem}[A] = R1$

- **Load-store:** GPR and only loads/stores access memory

load B,R1

$R1 = \text{mem}[B]$

load C,R2

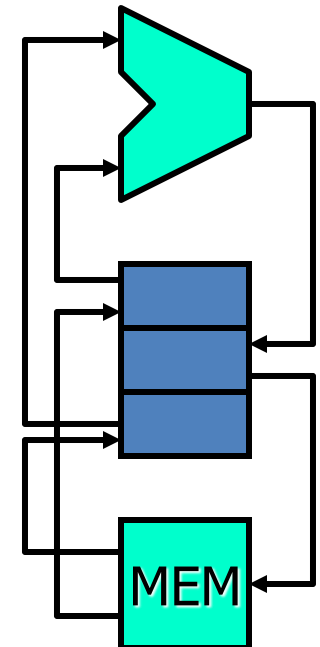
$R2 = \text{mem}[C]$

add R1,R2,R1

$R1 = R1 + R2$

store R1,A

$\text{mem}[A] = R1$



Operand Model Pros and Cons

- Metric I: **static code size**
 - Number of instructions needed to represent program, size of each
 - Want many implicit operands, high level instructions
 - Good → bad: accumulator, stack, GP-register, load-store
- Metric II: **data memory traffic**
 - Number of bytes move to and from memory (including cache)
 - Want as many long-lived operands in on-chip storage
 - Good → bad: load-store / GP-register, stack, accumulator
- Metric III: **cycles per instruction**
 - Want short (1 cycle?), little variability, few nearby dependences
 - Good → bad: load-store, GP-register, stack, accumulator
- Upshot: most new ISAs are load-store (or GP-register)

How Many Registers?

- Registers faster than memory, have as many as possible?
 - **No**
 - One reason registers are faster is that there are **fewer of them**
 - Smaller is faster in h/w (shorter wires, simpler decoder, etc.)
 - Another is that they are **directly addressed** (no address calc)
 - More of them, means larger specifiers
 - Fewer registers per instruction or indirect addressing
 - **Not everything can be put in registers**
 - Structures, arrays, anything pointed-to
 - Although compilers are getting better at putting more things in
 - More registers means **more saving/restoring**
- Upshot: trend to more registers: 8 (x86)→32 (MIPS) →128 (IA64)
 - 64-bit x86 has 16 64-bit integer and 16 128-bit FP registers

Virtual Address Size

- What is an n-bit processor?
 - Support memory size of 2^n
 - Alternative (wrong) definition: size of operators
- Virtual address size
 - Determines maximum size of addressable (usable) memory
 - Most high-perf devices are 64 bit (server/cell/laptop/desktop)
 - Most implementations limited to 40-50 bits
 - X86-64 uses 48-bit virtual address space
 - A pain to overcome too-small virtual address space
 - x86 evolution:
 - 12-bit (4004), 14-bit (8008), 16-bit (8086), 24-bit (80286),
 - 32-bit + protected memory (80386)
 - 64-bit (AMD's Opteron & Intel's EM64T Pentium4)

Size of Addressability

- What is the smallest unit of data an ISA can address
- Most architectures use 8-bits (byte-level)

Memory Addressing

- **Addressing mode:** way of specifying address
 - Used in memory-memory or load/store instructions in register ISA
- Examples
 - **Register-Indirect:** $R1 = \text{mem}[R2]$
 - **Displacement:** $R1 = \text{mem}[R2 + \text{immed}]$
 - **Index-base:** $R1 = \text{mem}[R2 + R3]$
 - **Memory-indirect:** $R1 = \text{mem}[\text{mem}[R2]]$
 - **Auto-increment:** $R1 = \text{mem}[R2], R2 = R2 + 1$
 - **Auto-indexing:** $R1 = \text{mem}[R2 + \text{immed}], R2 = R2 + \text{immed}$
 - **Scaled:** $R1 = \text{mem}[R2 + R3 * \text{immed1} + \text{immed2}]$
 - **PC-relative:** $R1 = \text{mem}[\text{PC} + \text{imm}]$
- What high-level program idioms are these used for?
- What implementation impact? What impact on insn count?

ISA Aspects

- Execution Model (Implicit vs Explicit Dependences)
 - Von Neumann vs Dataflow
- Implicit vs Explicit Parallelism Specification
 - Scalar vs Vector vs VLIW
- Format (Variable vs Fixed-width)
- Operand Model (where is temporary data stored)
 - Register, Stack, Accumulator
- Addressing Modes (how do you specify addresses)
 - Register-Indirect, Displacement, Index-base
- Memory Abstractions (physical/virtual addr. space)
- Control (condition, branch/predication)
- Details: Register file size, Datatypes (float vs int, bitwidth), Instruction Complexity -- loops

Control Instructions

- One issue: **testing for conditions**
 - Option I: **compare and branch insns**
`branch-less-than R1,10,target`
+ Simple, – two ALUs: one for condition, one for target address
if (R1 < 10) branch to target
 - Option II: **implicit condition codes**
 - `subtract R2,R1,10 // sets "negative" CC`
`branch-neg target`
+ Condition codes set "for free", – implicit dependence is tricky
 - Option III: **condition registers, separate branch insns**
`set-less-than R2,R1,10`
`branch-not-equal-zero R2,target`
– Additional instructions, + one ALU per, + explicit dependence

Example: MIPS Conditional Branches

- MIPS uses combination of options II/III
 - Compare 2 registers and branch: **beq, bne**
 - Equality and inequality only
 - + Don't need an adder for comparison
 - Compare 1 register to zero and branch: **bgtz, bgez, bltz, blez**
 - Greater/less than comparisons
 - + Don't need adder for comparison
 - Set explicit condition registers: **slt, sltu, slti, sltiu**, etc.
- Reasoning?
 - More than 80% of branches are (in)equalities or comparisons to 0
 - OK to take two insns to do remaining branches

Control Instructions II

- Another issue: **computing targets**
 - Option I: **PC-relative**
 - Position-independent within procedure
 - Used for branches and jumps within a procedure
 - Option II: **Absolute**
 - Position independent outside procedure
 - Used for procedure calls
 - Option III: **Indirect** (target found in register/memory)
 - Needed for jumping to dynamic targets
 - Used for returns, dynamic procedure calls, switches
- How far do you need to jump?
 - Typically not so far within a procedure (they don't get that big)
 - Further from one procedure to another

MIPS Control Instructions

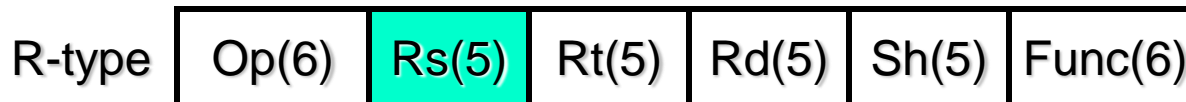
- MIPS uses all three
 - PC-relative conditional branches: **bne**, **beq**, **blez**, etc.
 - 16-bit relative offset, <0.1% branches need more



- Absolute jumps unconditional jumps: **j**
 - 26-bit offset



- Indirect jumps: **jr**

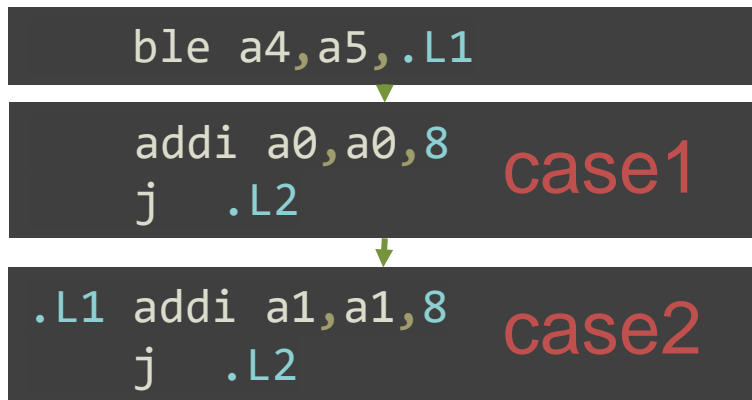


Control Instructions III

- Another issue: support for procedure calls?
 - Link (remember) address of calling insn + 4 so we can return to it
- MIPS
 - jump-and-link: `jalr Rs, Rd`
 - $Rd = PC + 8$
 - $PC = Rs$
- X86
 - Call
 - target can be displacement, absolute, absolute indirect...
 - Put return address on stack (memory pointed to by stack pointer)
 - Ret
 - Return to location on head of stack

Control Instructions IV

- Unpredictable branches are bad – can we eliminate?
- Sure: Predication (ARM, NVIDIA)
 - Predicate registers can be set by compare instructions
 - Instructions can be “predicated” by these registers
 - Turns control dependence into data dependence (similar to dataflow ISA, but on a small scale)
 - Example: `if(a4 < a5) {a0++;} else {a1++;}`



pred reg.

```
set_le p1, a4, a5
p1==false: addi a0,a0,8
p1==true: addi a1,a1,8
j .L2
```

- What's the tradeoff, when would you want to use this?
- Common feature (x86 has limited support, ARM general support)

ISA Aspects

- Execution Model (Implicit vs Explicit Dependences)
 - Von Neumann vs Dataflow
- Implicit vs Explicit Parallelism Specification
 - Scalar vs Vector vs VLIW
- Format (Variable vs Fixed-width)
- Operand Model (where is temporary data stored)
 - Register, Stack, Accumulator
- Addressing Modes (how do you specify addresses)
 - Register-Indirect, Displacement, Index-base
- Memory Abstractions (physical/virtual addr. space)
- Control (condition, branch/predication)
- Details: Register file size, Datatypes (float vs int, bitwidth), Instruction Complexity -- loops

Economic Concerns: License-ability?

- X86:
 - Very limited and expensive
 - Licensed to handful of companies: AMD
 - AMD X86-64 extension licensed to Intel :)

- ARM:
 - Easy to obtain license
 - License fee: \$1M - \$10M
 - Royalty: 0.5 - 2% of chip selling price

IP	Royalty
ARM7/9/11	1.0% - 1.5%
ARM Cortex A-series	1.5% - 2.0%
ARMv8 Based Cortex A-series	2.0% +
Mali GPU	0.75% - 1.25%
Physical IP Package (POP)	0.5%

Source: anandtech.com (2013)

- RISC-V:
 - Fully open source – no fees

The RISC vs. CISC Debate

RISC and CISC

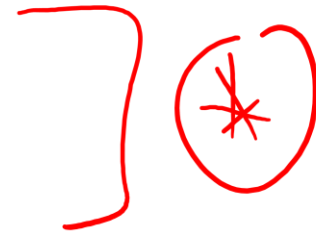
- **RISC**: reduced-instruction set computer
 - Coined by Patterson in early 80's
 - RISC-I (Patterson), MIPS (Hennessy), IBM 801 (Cocke)
 - Examples: PowerPC, ARM, SPARC, Alpha, PA-RISC
- **CISC**: complex-instruction set computer
 - Term didn't exist before "RISC"
 - Examples: x86, VAX, Motorola 68000, etc.
- Philosophical war (one of several) started in mid 1980's
 - RISC "won" the technology battles
 - CISC won the high-end commercial war (1990s to today)
 - Compatibility a stronger force than anyone (but Intel) thought
 - RISC won the embedded computing war

The Context

- Pre 1980
 - Bad compilers (so assembly written by hand)
 - Complex, high-level ISAs (easier to write assembly)
 - Slow multi-chip implementations
- Late 1970s/Early 1980s
 - Moore's Law makes single-chip microprocessor possible...
 - **...but only for small, simple ISAs**
 - Performance advantage of this "integration" was compelling
 - Compilers had to get involved in a big way
- **RISC manifesto**: create ISAs that...
 - **Simplify single-chip implementation**
 - **Facilitate optimizing compilation**

The RISC Design Tenets

- **Single-cycle execution** *(ignoring FP/vector)*
 - CISC: many multicycle operations
- **Hardwired control**
 - CISC: microcoded multi-cycle operations
- **Load/store architecture**
 - CISC: register-memory and memory-memory
- **Few memory addressing modes**
 - CISC: many modes
- **Fixed-length instruction format**
 - CISC: many formats and lengths
- **Reliance on compiler optimizations**
 - CISC: hand assemble to get good performance
- **Many registers** (compilers are better at using them)
 - CISC: few registers



CISCs and RISCs

- The CISCs: x86, VAX (**V**irtual **A**ddress e**X**tension to PDP-11)
 - Variable length instructions: 1-321 bytes!!!
 - 14 registers + PC + stack-pointer + condition codes
 - Data sizes: 8, 16, 32, 64, 128 bit, decimal, string
 - Memory-memory instructions for all data sizes
 - Special insns: `crc`, `insque`, `polyf`, and a cast of hundreds
 - x86: "Difficult to explain and impossible to love"
- The RISCs: MIPS, PA-RISC, SPARC, PowerPC, Alpha, ARM, RISC-V
 - 32-bit instructions
 - 32 integer registers, 32 floating point registers, load-store
 - 64-bit virtual address space
 - Few addressing modes
 - Why so many basically similar ISAs? Everyone wanted their own

The Debate

- RISC argument
 - CISC is fundamentally handicapped
 - For a given technology, RISC implementation will be better (faster)
 - Current technology enables single-chip RISC
 - When it enables single-chip CISC, RISC will be pipelined
 - When it enables pipelined CISC, RISC will have caches
 - When it enables CISC with caches, RISC will have next thing...
- CISC rebuttal
 - CISC flaws not fundamental, can be fixed with more transistors
 - Moore's Law will narrow the RISC/CISC gap (true)
 - Good pipeline: RISC = 100K transistors, CISC = 300K
 - By 1995: 2M+ transistors had evened playing field
 - Software costs dominate, **compatibility** is paramount

RISC vs CISC Performance Argument

$$\frac{\text{Time}}{\text{Program}} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Cycles}}{\text{Instruction}} \times \frac{\text{Time}}{\text{Cycle}}$$

- **CISC** (Complex Instruction Set Computing)
 - Reduce “instructions/program” with “complex” instructions
 - But tends to increase CPI or clock period
 - Easy for assembly-level programmers, good code density
- **RISC** (Reduced Instruction Set Computing)
 - Improve “cycles/instruction” with many single-cycle instructions
 - Increases “instruction/program”, but hopefully not as much
 - **Help from smart compiler**
 - Perhaps improve clock cycle time (seconds/cycle)
 - **via aggressive implementation allowed by simpler insn**

Current Winner (Server/Desktop): x86

- x86 was first 16-bit chip by ~2 years
 - IBM put it into its PCs because there was no competing choice
 - Rest is historical inertia and “financial feedback”
- Moore’s law has helped Intel in a big way
 - Many engineering problems can be solved with more transistors
- Complex architecture due to evolution over time
 - Typical of many older ISAs, e.g. IBM 360/370/390
 - Started as 16-bit microprocessor (later, 32-bits)
 - Upward compatible from 8080 (accumulator-based)

Current Winner (mobile/embed.): ARM

- ARM (Advanced RISC Machine)
 - First ARM chip in mid-1980s (from Acorn Computer Ltd).
 - Over 10 billion units sold (75% of 32/64-bit CPUs)
 - Low-power and embedded devices (phones, for example)
- 32-bit RISC ISA
 - 16 registers
 - Many addressing modes (for example, auto increment)
 - Condition codes, each instruction can be predicated
- Multiple compatible implementations
 - Intel's X-scale (was DEC's)
 - Others: Freescale (was Motorola), IBM, Texas Instruments, Nintendo, STMicroelectronics, Samsung, Sharp, Philips, etc.
- "Thumb" 16-bit wide instructions
 - Increase code density

Intel's Compatibility Trick: RISC Inside

- 1993: Intel wanted out-of-order execution in Pentium Pro
 - Hard to do with a coarse grain ISA like x86
- Solution? Translate x86 to RISC **μops** in hardware
 - `push $eax`
becomes (we think, uops are proprietary)
`store $eax, -4($esp)`
`addi $esp, $esp, -4`
 - + Processor maintains **x86 ISA externally for compatibility**
 - + But executes **RISC μISA internally for implementability**
 - Given translator, x86 almost as easy to implement as RISC
 - Intel implemented out-of-order before any RISC company
 - Also, OoO also benefits x86 more (because ISA limits compiler)
 - Idea co-opted by other x86 companies: AMD and Transmeta
 - Different **μops** for different designs
 - **Not part of the ISA specification**, not publicly disclosed

Potential Micro-op Scheme (1 of 2)

- Most instructions are a **single** micro-op
 - Add, xor, compare, branch, etc.
 - Loads example: `mov -4(%rax), %ebx`
 - Stores example: `mov %ebx, -4(%rax)`
- Each memory operation adds a micro-op
 - `"addl -4(%rax), %ebx"` is two micro-ops (load, add)
 - `"addl %ebx, -4(%rax)"` is three micro-ops (load, add, store)
- What about address generation?
 - Simple address generation is generally part of single micro-op
 - Sometime store addresses are calculated separately
 - More complicated (scaled addressing) might be separate micro-op

Potential Micro-op Scheme (2 of 2)

- Function call (CALL) – 4 uops
 - Get program counter, store program counter to stack, adjust stack pointer, unconditional jump to function start
- Return from function (RET) – 3 uops
 - Adjust stack pointer, load return address from stack, jump to return address
- Other operations
 - String manipulations instructions
 - For example STOS is around six micro-ops, etc.
- Again, this is just a basic idea, the exact micro-ops are specific to each chip

Power Struggles: Revisiting the RISC vs. CISC Debate on Contemporary ARM and x86 Architectures

- Goal & Position?
 - Determine does RISC vs CISC matter? (no, just microarchitecture)
- Why now, what are the challenges?
 - Factors independent of ISA and microarchitecture: technology, memories, OS, compiler backend, workloads.
- Methodology?
 - Native execution for performance/uarch events, Physically measure power (Wattsup), deal with technology scaling analytically, emulation/simulation for instruction mix
- Interesting Findings?

Crux of argument

- Performance differs a lot....

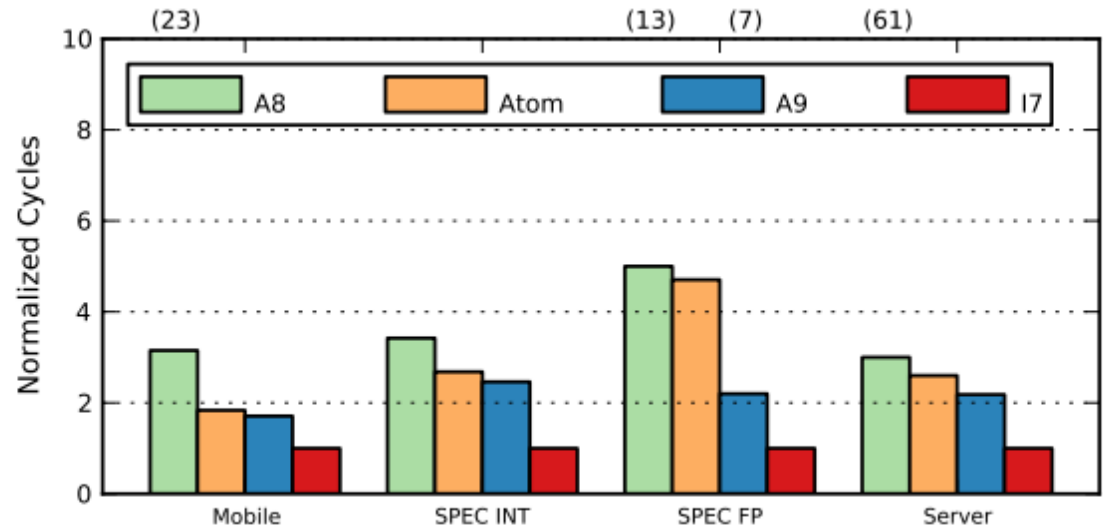


Figure 3. Cycle Count Normalized to i7.

- And so does the power....

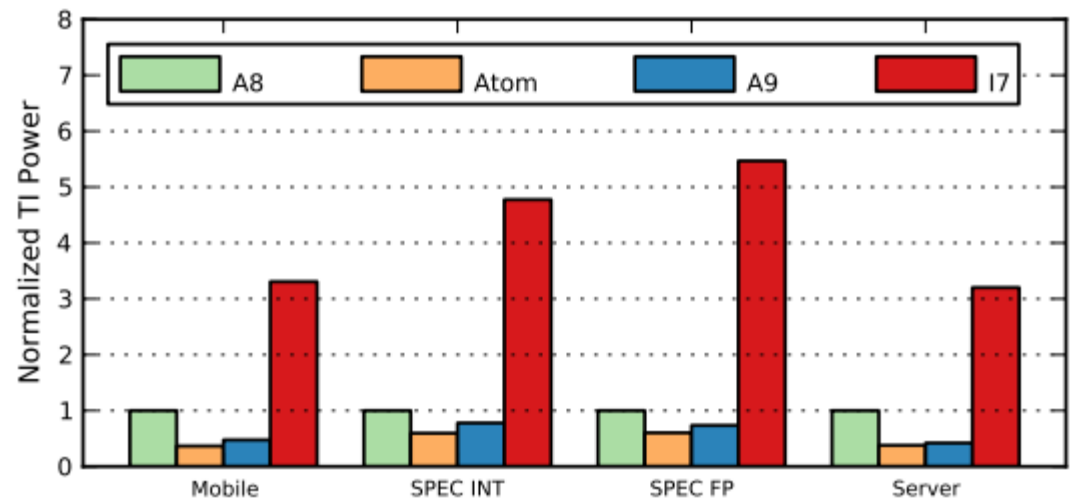


Figure 12. Tech. Independent Avg. Power Normalized to A8.

ISA doesn't seem to be playing a role...

- # Macro-ops and micro-ops are similar

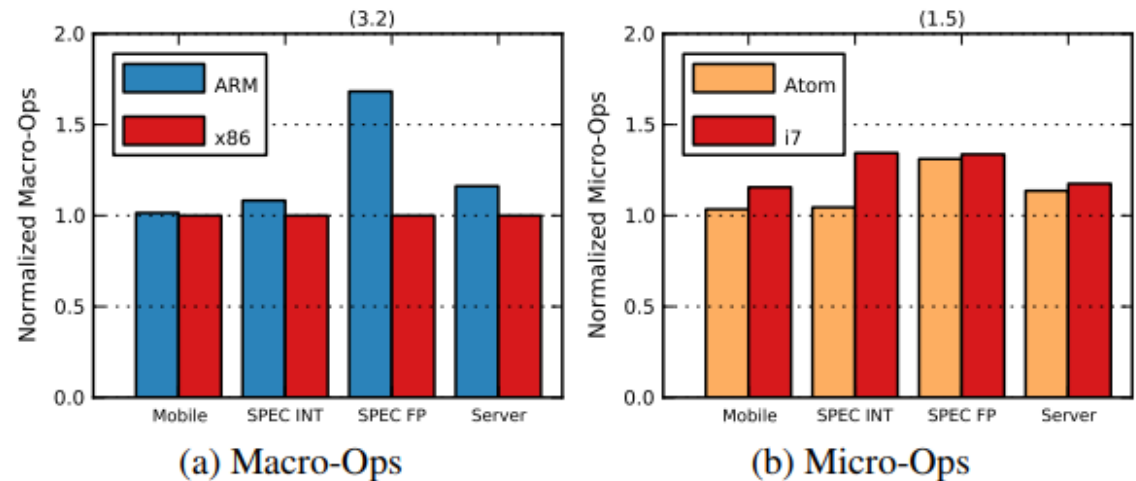


Figure 4. Instructions Normalized to i7 macro-ops.

- Instruction mix is pretty similar

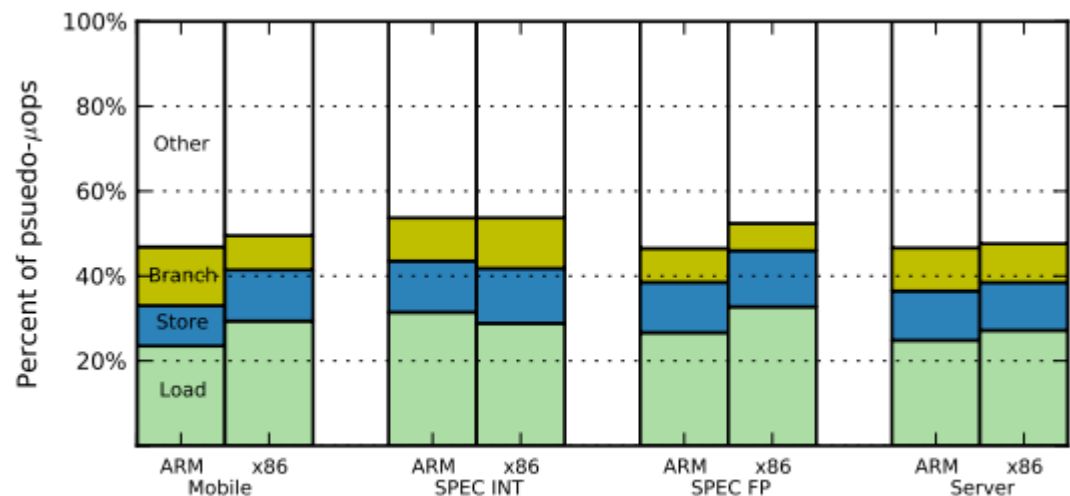


Figure 5. Instruction Mix (Performance Counters).

But microarchitecture plays a large role:

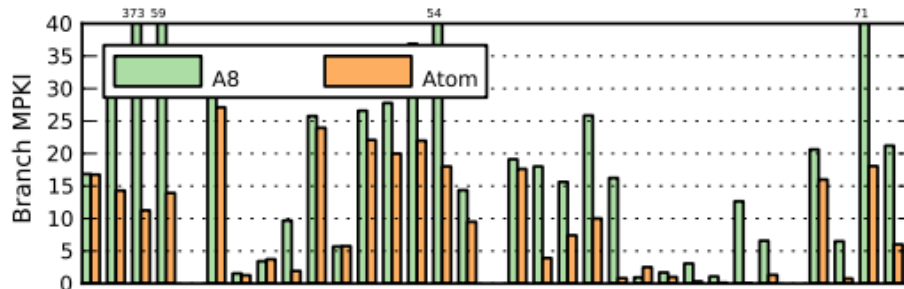


Figure 8. Branch Misses per 1000 ARM Instructions.

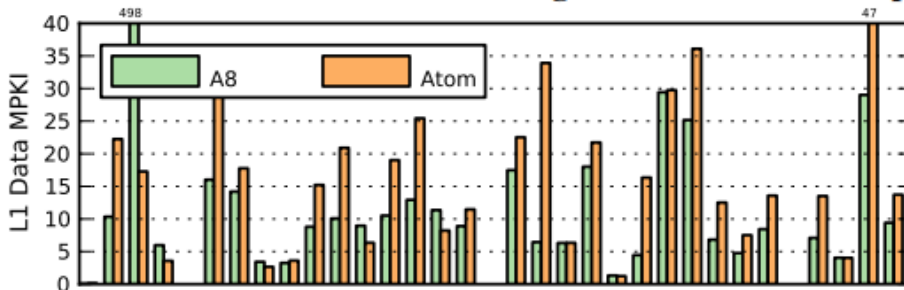
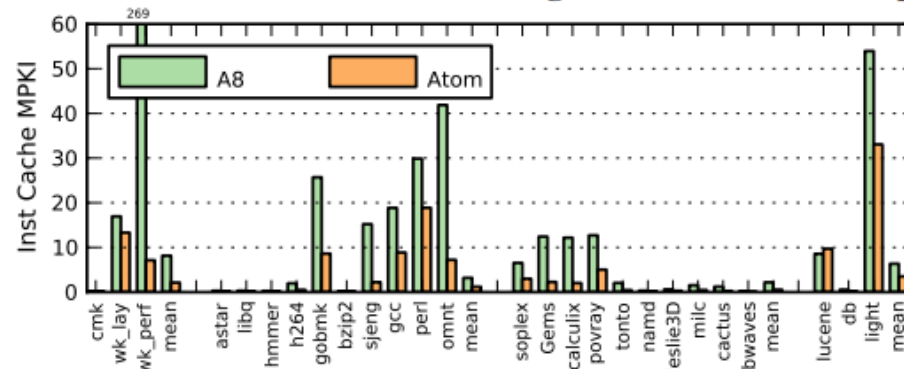
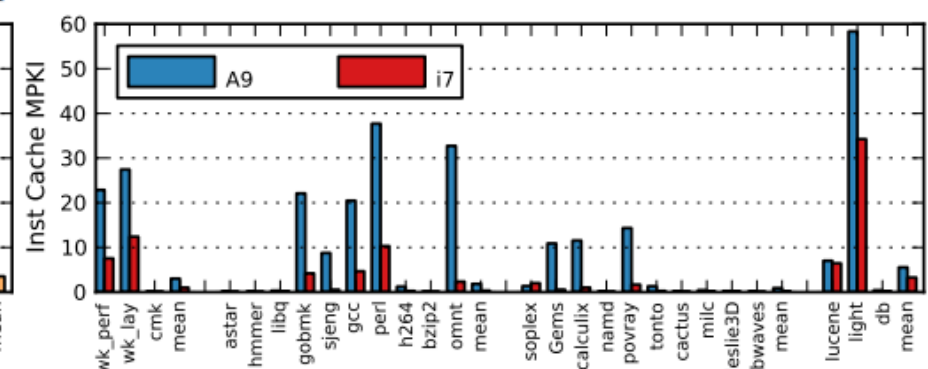
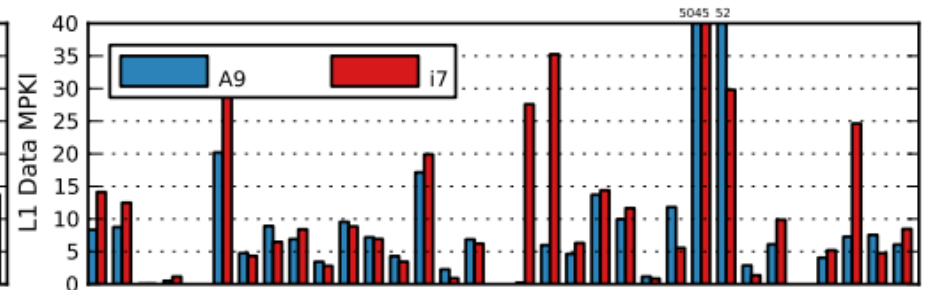
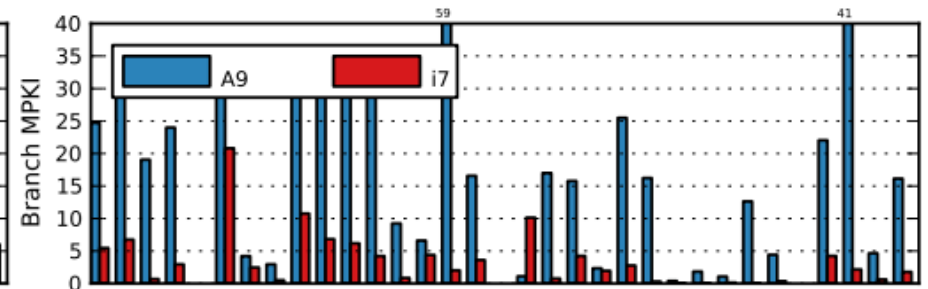


Figure 9. Data L1 Misses per 1000 ARM Instructions.



(a) In-Order



(b) Out-of-Order

At least one difference...

- Code density

Table 6. Instruction Size Summary.

		(a) Binary Size (MB)		(b) Instruction Length (B)	
		ARM	x86	ARM	x86
Mobile	Minimum	0.02	0.02	4.0	2.4
	Average	0.95	0.87	4.0	3.3
	Maximum	1.30	1.42	4.0	3.7
Desktop INT	Minimum	0.53	0.65	4.0	2.7
	Average	1.47	1.46	4.0	3.1
	Maximum	3.88	4.05	4.0	3.5
Desktop FP	Minimum	0.66	0.74	4.0	2.6
	Average	1.70	1.73	4.0	3.4
	Maximum	4.75	5.24	4.0	6.4
Server	Minimum	0.12	0.18	4.0	2.5
	Average	0.39	0.59	4.0	3.2
	Maximum	0.47	1.00	4.0	3.7

Power Struggles

- Any opinions on validity of results?
 - One opinion: Nobody was disputing in the first place that microarchitecture makes a big difference...
 - Microarchitectures are similar only in terms of high-level parameters...
- Could you do the same study on a simulator?
 - Advantage: Microarchitectures can be made to be identical (just change gem5 ISA between ARM & RISC)
 - Disadvantage: Simulator validity/bugs become a big problem. :)

Will ISAs (and ISA research) continue to be important?

Redux: Are ISAs Important?

- Does “quality” of ISA actually matter?
 - Not for performance (mostly)
 - Mostly comes as a design complexity issue
 - Insn/program: everything is compiled, compilers are good
 - Cycles/insn and seconds/cycle: μ ISA, many other tricks
 - What about power efficiency? *Maybe*
 - ARMs are most power efficient today...
 - ...but Intel is moving x86 that way (e.g., Intel’s Atom)
 - **Open question: can x86 be as power efficient as ARM?**
- Does “nastiness” of ISA matter?
 - Mostly no, only compiler writers and hardware designers see it
- Even compatibility is not what it used to be
 - Software emulation
 - **Open question: will “ARM compatibility” be the next x86?**

Summary

- What is an ISA?
 - A functional contract
- What makes a good ISA
 - {Programm|Implement|Compat}-ability
 - Enables high-performance
 - At least doesn't get in the way
- ISAs similar in many high-level ways
 - But many design choices in details
 - Two “philosophies”: CISC/RISC
 - Difference is blurring
 - Other distinctions more important, dataflow, vector, etc...
- Compatibility is a powerful force
 - Tricks: binary translation, μ ISAs mitigation

Hidden Bonus Slides

Ultimate Compatibility Trick

- Support old ISA by...
 - ...having a simple processor for that ISA somewhere in the system
 - How first Itanium supported x86 code
 - x86 processor (comparable to Pentium) on chip
 - How PlayStation2 supported PlayStation games
 - Used PlayStation processor for I/O chip **& emulation**

More X86 RISC Tricks

- x86 code is becoming more “RISC-like”
 - In 32-bit to 64-bit transition, x86 made two key changes:
 - Double number of registers, better function calling conventions
 - More registers (can pass parameters too), fewer **pushes/pops**
 - Result? Fewer complicated instructions
 - Moved from $\sim 1.6 \mu\text{ops} / \text{x86 insn}$ to $\sim 1.1 \mu\text{ops} / \text{x86 insn}$
- More recent: “macro-op fusion” and “micro-op fusion”
 - Intel’s recent processors fuse certain instruction pairs
 - Macro-op fusion: fuses “compare” and “branch” instructions
 - Micro-op fusion: fuses load/add pairs, fuses store “address” & “data”

Post-RISC: VLIW and EPIC

- ISAs explicitly targeted for multiple-issue (superscalar) cores
 - VLIW: Very Long Insn Word
 - Later rebranded as “EPIC”: Explicitly Parallel Insn Computing
- Intel/HP IA64 (Itanium): 2000
 - EPIC: 128-bit 3-operation bundles
 - 128 64-bit registers
 - + Some neat features: Full predication, explicit cache control
 - Predication: every instruction is conditional (to avoid branches)
 - But lots of difficult to use baggage as well: software speculation
 - Every new ISA feature suggested in last two decades
 - Relies on younger (less mature) compiler technology
 - Commercially dead – last Itanium in 2017

ISA Research

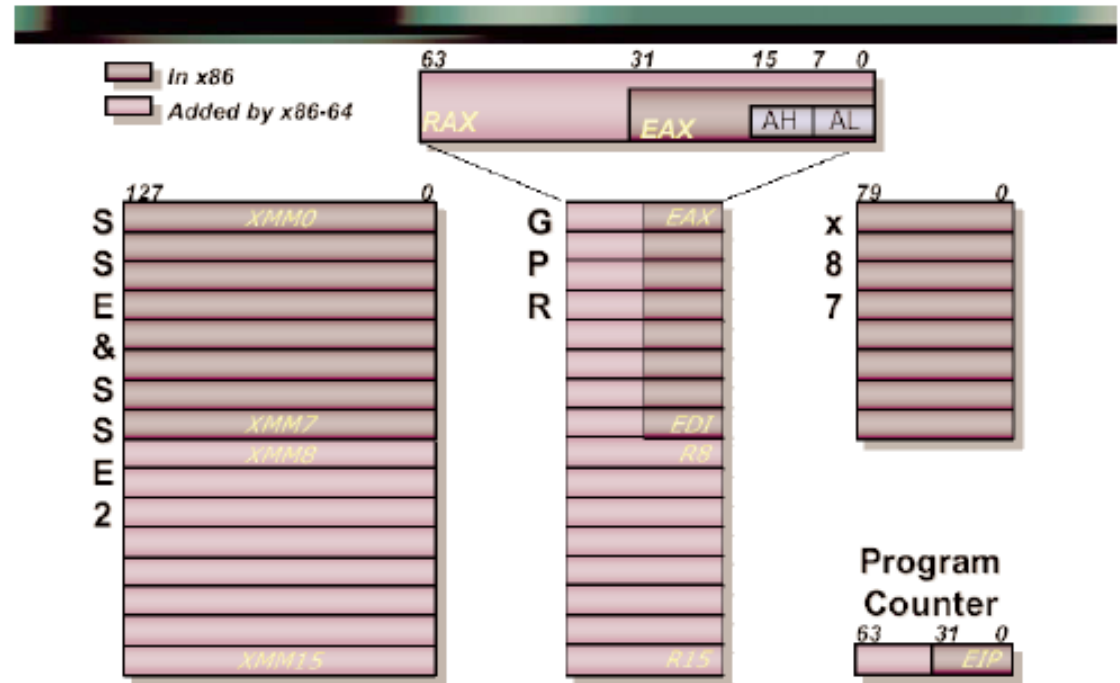
- Compatibility killed ISA research for a while
 - But binary translation/emulation has revived it
 - ... and so did the end of Moore's law
- Some ISA-related projects
 - General Purpose Dataflow
 - “WaveScalar” [Washington], “TRIPS EDGE” [Texas]
 - Explicit dataflow ISAs (von Neumann alternatives)
 - DySER [Wisconsin], CCA, BERET [Michigan]
 - (arguably) Dataflow extensions to von Neumann ISA
 - Specialized computers:
 - Stream-Dataflow [me]
 - Explicit memory movement + reconfigurable CGRA

x86: Registers

- 4 arithmetic,
 - 4 address,
 - 4 segment,
 - 2 control
-
- Accumulator
 - AH, AL (8 bits)
 - AX (16 bits)
 - EAX (32 bits)
 - RAX (64 bits)

blech!

x86-64 Programmer's Model



August 2002

x86-64 ISA

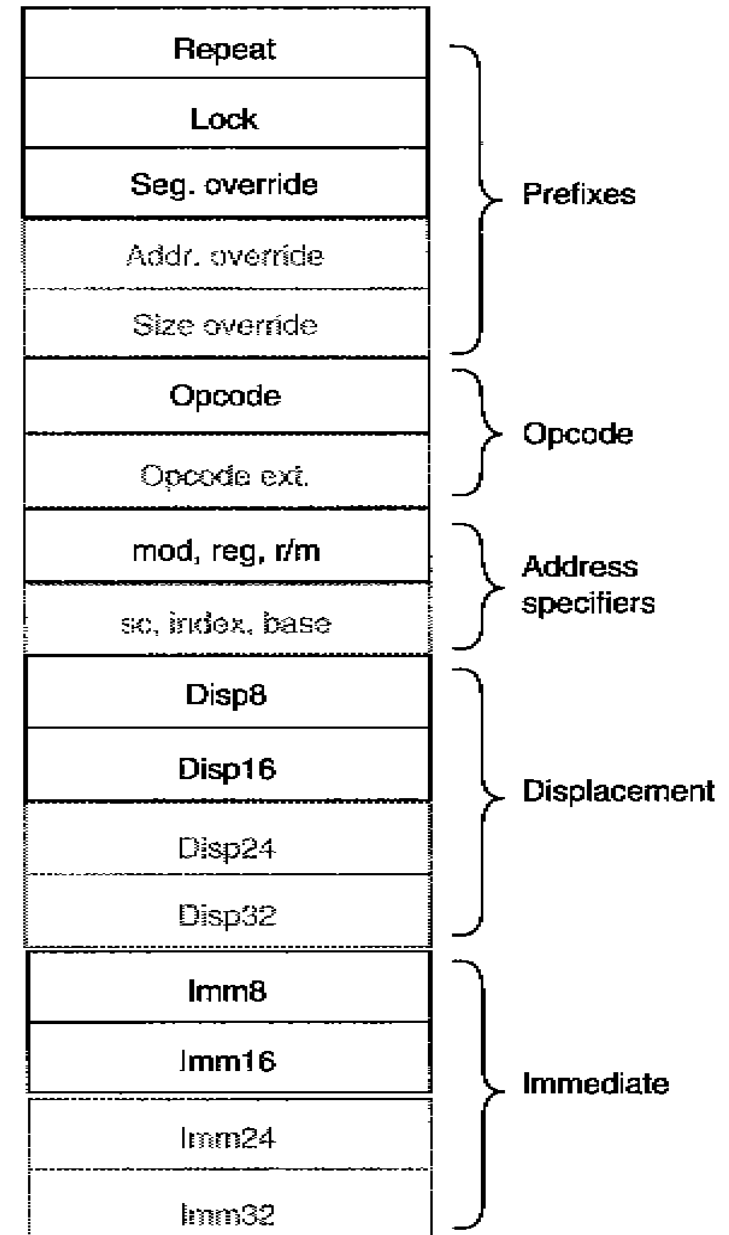
HotChips 14

9

x86 Details

- Seven address modes
 - Absolute
 - Register indirect
 - Based
 - Indexed
 - Based indexed with displacement
 - Based with scaled indexed
 - Based with scaled indexed and displacement

yak!

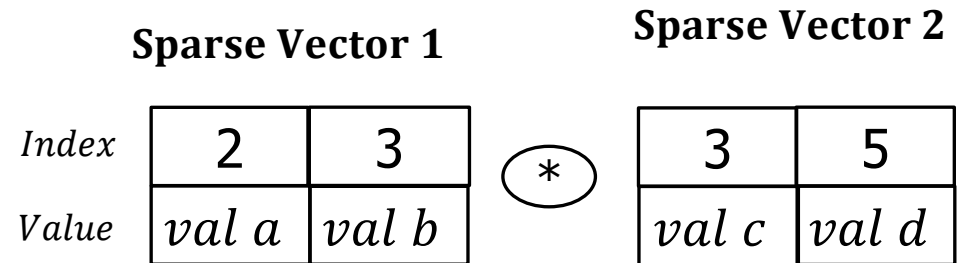


Many instruction formats

Example Program: Sparse Vector Multiply

```
typedef struct node {
    int index;
    float val;
} node;

float func(node* n1, node* n2)
{
    float total=0;
    while(n1 && n2) {
        if(n1->index==n2->index) {
            total+=n1->val*n2->val;
        } else if (n1->index >
                    n2->index) {
            n1++;
        } else {
            n2++;
        }
    }
}
```



Match with 3 -> Perform *val b * val c*

Von Neumann VS. Dataflow Program

Von Neumann

101b0: beqz a0,101c8 <func+0x2e>

101b2: beqz a1,101c8 <func+0x2e>

101b4: lw a4,0(a0)

101b6: lw a5,0(a1)

101b8: beq a4,a5,101a0 <func+0x6>

101a0: flw fa5,4(a0) **Ind1==ind2**

101a4: flw fa4,4(a1)

101a8: fmul.s fa5,fa5,fa4

101ac: fadd.s fa0,fa0,fa5

101bc: ble a4,a5,101c4 <func+0x2a>

101c0: addi a0,a0,8 **Ind1<ind2**

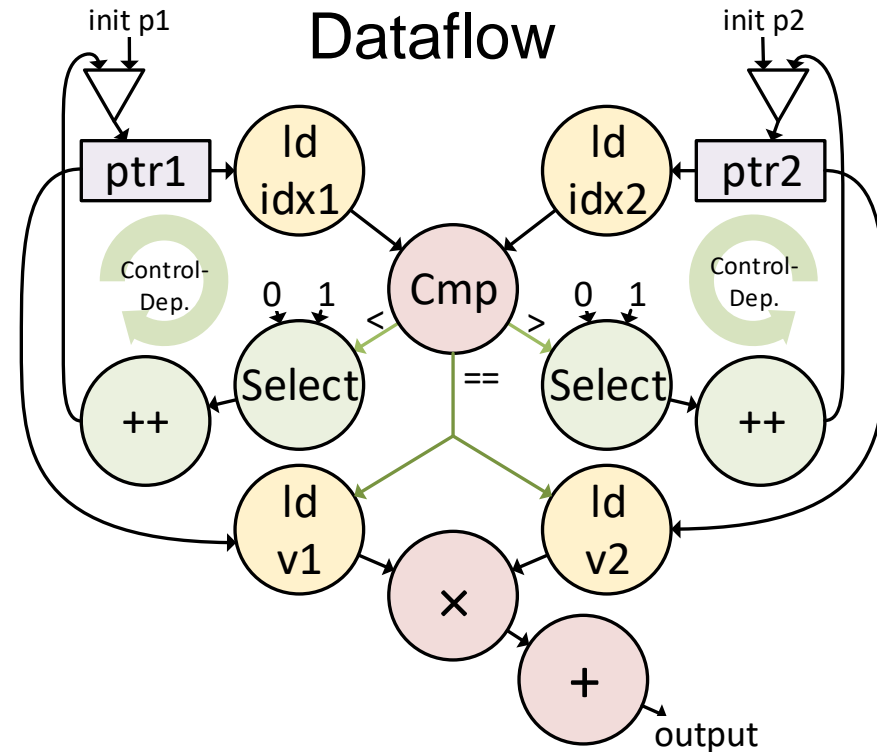
101c2: j 101b0 <func+0x16>

101c4: addi a1,a1,8 **Ind2<ind1**

101c6: j 101b0 <func+0x16>

Exit

Dataflow



Sparse Vector 1

Index

2

3

Value

val a

val b

Sparse Vector 2

3

5

val c

val d

*

Match index 3 -> Perform *val b* * *val c*

Virtual ISAs

- Move portability concerns to the software level!
- Java and C# use an ISA-like interface
 - JavaVM uses a stack-based bytecode
 - C# has the CLR (common language runtime)
 - NVIDIA's "PTX"
 - Higher-level than machine ISA
 - Design for translation (not direct execution)
- Goals:
 - Portability (abstract away the actual hardware)
 - Target for high-level compiler (one per language)
 - Source for low-level translator (one per ISA)
 - Flexibility over time
- May allow ISA research to overcome compatibility "gorilla"
 - But Intel wants x86 to be the winning "virtual ISA"

Transmeta's Take: Code Morphing

- **Code morphing**: x86 translation performed in software
 - Crusoe/Astro are x86 emulators, no actual x86 hardware anywhere
 - Only “code morphing” translation software written in native ISA
 - Native ISA is invisible to applications, even BIOS
 - Different Crusoe versions have (slightly) different ISAs: can't tell
- How was it done?
 - Code morphing software resides in boot ROM
 - On startup boot ROM hijacks 16MB of main memory
 - Translator loaded into 512KB, rest is **translation cache**
 - Software starts running in **interpreter** mode
 - Interpreter profiles to find “hot” regions: procedures, loops
 - Hot region compiled to native, optimized, cached
 - Gradually, more and more of application starts running native

Translation and Virtual ISAs

- New compatibility interface: ISA + translation software
 - Less necessary that processor ISA be compatible
 - As long as some combination of ISA + software translation layer is
- Kinds
 - **Binary-translation**: transform static image, run native
 - **Emulation**: unmodified image, interpret each dynamic insn
 - Typically optimized with just-in-time (JIT) compilation
- Examples: FX!32 (x86 on Alpha), Rosetta (PowerPC on x86)
- Downside: performance overheads
 - Performance overheads reasonable (many recent advances)

Adventures in ISA Missteps (funny as in *uh-oh*)

Motivation: Function calls incur many overhead instructions to save/restore registers.

Idea: Specialize communication between functions by exposing function->function dataflow to ISA?

Register Windows

- **Register windows**: hardware activation records
 - Sun SPARC (from the RISC I)
 - 32 integer registers divided into: 8 global, 8 local, 8 input, 8 output
 - Explicit **save/restore** instructions
 - Global registers fixed
 - **save**: inputs “pushed”, outputs → inputs, locals zeroed
 - **restore**: locals zeroed, inputs → outputs, inputs “popped”
 - Hardware stack provides few (8) on-chip register frames
 - Spilled-to/filled-from memory on over/under flow
- + Automatic parameter passing, caller-saved registers
- + No memory traffic on shallow (<8 deep) call graphs
- Hidden memory operations (some restores fast, others slow)
- A nightmare for register renaming (more later)

Motivation: Pipelined processor stalls after branch instructions because it takes time to resolve the branch.

Idea: Play with program semantics to eliminate?

Delay Slot

```
ble  a4,a5,101c4  
addi a0,a0,8  
j     101b0
```

This instruction
always
executes

- **Branch Delay Slot:**
 - Instruction after branch always executes
 - Possible to have multiple! (or have load delay slots)
 - Example ISAs: MIPS, PA-RISC, ETRAX CRIS, SuperH, and SPARC
- Good
 - Helps prevent pipeline bubbles/hazards (more on that later)
- Bad
 - Inelegant: Exposes the pipelined nature of the h/w to the ISA
 - Totally irrelevant for OOO execution – headache to deal with
 - Extra compiler/assembler work
- Later RISC ISAs do not include: PowerPC, ARM, Alpha, and RISC-V