

# ECE 4750 Computer Architecture, Fall 2016

## T01 Fundamental Processor Concepts

School of Electrical and Computer Engineering  
Cornell University

revision: 2016-09-05-15-52

<b>1</b>	<b>Instruction Set Architecture</b>	<b>2</b>
1.1.	IBM 360 Instruction Set Architecture . . . . .	4
1.2.	MIPS32 Instruction Set Architecture . . . . .	6
1.3.	Tiny RISC-V Instruction Set Architecture . . . . .	10
<b>2</b>	<b>Processor Functional-Level Model</b>	<b>14</b>
2.1.	Transactions and Steps . . . . .	14
2.2.	TinyRV1 Simple Assembly Example . . . . .	15
2.3.	TinyRV1 VVAdd Asm and C Program . . . . .	16
2.4.	TinyRV1 Mystery Asm and C Program . . . . .	17
<b>3</b>	<b>Processor/Laundry Analogy</b>	<b>18</b>
3.1.	Arch vs. $\mu$ Arch vs. VLSI Impl . . . . .	18
3.2.	Processor Microarchitectural Design Patterns . . . . .	19
3.3.	Transaction Diagrams . . . . .	20
<b>4</b>	<b>Analyzing Processor Performance</b>	<b>21</b>

## 1. Instruction Set Architecture

- By early 1960's, IBM had several incompatible lines of computers!
  - Defense : 701
  - Scientific : 704, 709, 7090, 7094
  - Business : 702, 705, 7080
  - Mid-Sized Business : 1400
  - Decimal Architectures : 7070, 7072, 7074
- Each system had its own:
  - Implementation and potentially even technology
  - Instruction set
  - I/O system and secondary storage (tapes, drums, disks)
  - Assemblers, compilers, libraries, etc
  - Application niche

- IBM 360 was the first line of machines to separate ISA from microarchitecture
  - Enabled same software to run on different current and future microarchitectures
  - Reduced impact of modifying the microarchitecture enabling rapid innovation in hardware

Application
Algorithm
Programming Language
Operating System
Instruction Set Architecture
Microarchitecture
Register-Transfer Level
Gate Level
Circuits
Devices
Physics

... the structure of a computer that a **machine language programmer** must understand to write a correct (**timing independent**) program for that machine.

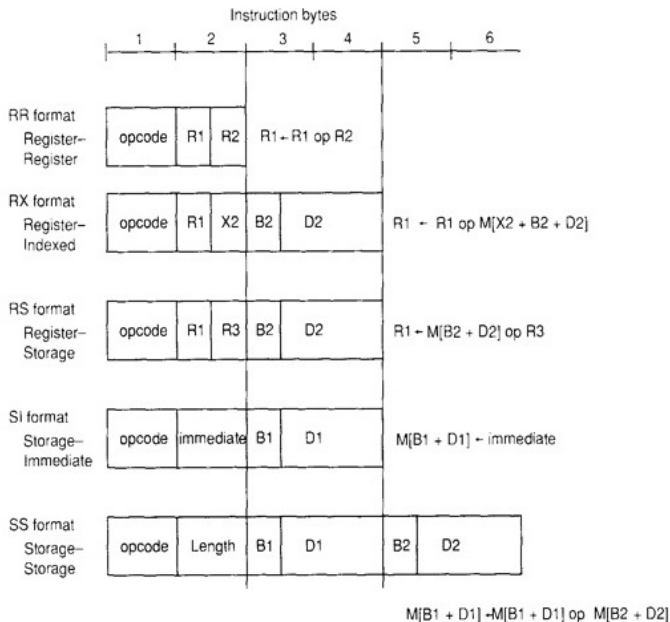
— Amdahl, Blaauw, Brooks, 1964

## ISA is the contract between software and hardware

- 1. \_\_\_\_\_
  - Representations for characters, integers, floating-point
  - Integer formats can be signed or unsigned
  - Floating-point formats can be single- or double-precision
  - Byte addresses can ordered within a word as either little- or big-endian
- 2. \_\_\_\_\_
  - Registers: general-purpose, floating-point, control
  - Memory: different addresses spaces for heap, stack, I/O
- 3. \_\_\_\_\_
  - Register: operand stored in registers
  - Immediate: operand is an immediate in the instruction
  - Direct: address of operand in memory is stored in instruction
  - Register Indirect: address of operand in memory is stored in register
  - Displacement: register indirect, addr is added to immediate
  - Autoincrement/decrement: register indirect, addr is automatically adj
  - PC-Relative: displacement is added to the program counter
- 4. \_\_\_\_\_
  - Integer and floating-point arithmetic instructions
  - Register and memory data movement instructions
  - Control transfer instructions
  - System control instructions
- 5. \_\_\_\_\_
  - Opcode, addresses of operands and destination, next instruction
  - Variable length vs. fixed length

## 1.1. IBM 360 Instruction Set Architecture

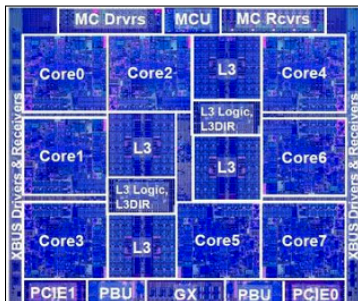
- How is data represented?
  - 8-bit bytes, 16-bit half-words, 32-bit words, 64-bit double-words
  - IBM 360 is why bytes are 8-bits long today!
- Where can data be stored?
  - $2^{24}$  32-bit memory locations
  - 16 general-purpose 32-bit registers and 4 floating-point 64-bit registers
  - Condition codes, control flags, program counter
- What operations can be done on data?
  - Large number of arithmetic, data movement, and control instructions



	Model 30	Model 70
Storage	8–64 KB	256–512 KB
Datapath	8-bit	64-bit
Circuit Delay	30 ns/level	5 ns/level
Local Store	Main store	Transistor registers
Control Store	Read only 1 $\mu$ s	Conventional circuits

- IBM 360 instruction set architecture completely hid the underlying technological differences between various models
- **Significant Milestone: The first true ISA designed as a portable hardware-software interface**
- IBM 360: 50 years later ...  
The zSeries z13 Microprocessor

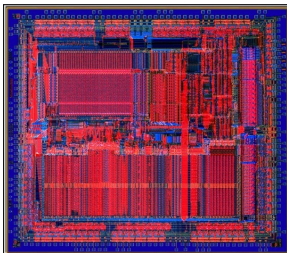
- 5 GHz in IBM 22 nm SOI
- 4B transistors in 678 mm<sup>2</sup>
- 17 metal layers
- $\approx$ 20K pads
- Eight cores per chip
- Aggressive out-of-order execution
- Four-level cache hierarchy
- On-chip 64MB eDRAM L3 cache
- Off-chip 480MB eDRAM L4 cache
- Can still run IBM 360 code!



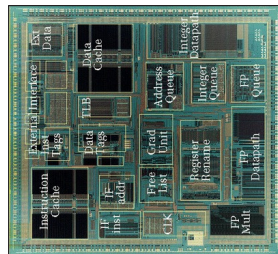
J. Warnock, et al., "22nm Next-Generation IBM System-Z Microprocessor,"  
Int'l Solid-State Circuits Conference, Feb. 2015.

## 1.2. MIPS32 Instruction Set Architecture

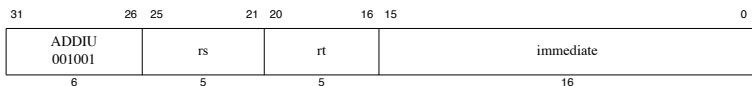
- How is data represented?
  - 8-bit bytes, 16-bit half-words, 32-bit words
  - 32-bit single-precision, 64-bit double-precision floating point
- Where can data be stored?
  - $2^{32}$  32-bit memory locations
  - 32 general-purpose 32-bit registers, 32 SP (16 DP) floating-point registers
  - FP status register, Program counter
- How can data be accessed?
  - Register, immediate, displacement
- What operations can be done on data?
  - Large number of arithmetic, data movement, and control instructions
- How are instructions encoded?
  - Fixed-length 32-bit instructions



**MIPS R2K:** 1986, single-issue, in-order, off-chip caches, 2  $\mu\text{m}$ , 8–15 MHz, 110K transistors, 80  $\text{mm}^2$



**MIPS R10K:** 1996, quad-issue, out-of-order, on-chip caches, 0.35  $\mu\text{m}$ , 200 MHz, 6.8M transistors, 300  $\text{mm}^2$



**Format:** ADDIU *rt*, *rs*, *immediate*

MIPS32

**Purpose:** Add Immediate Unsigned Word

To add a constant to a 32-bit integer

**Description:**  $GPR[rt] \leftarrow GPR[rs] + immediate$

The 16-bit signed *immediate* is added to the 32-bit value in GPR *rs* and the 32-bit arithmetic result is placed into GPR *rt*.

No Integer Overflow exception occurs under any circumstances.

**Restrictions:**

None

**Operation:**

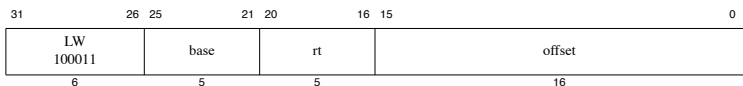
```
temp ← GPR[rs] + sign_extend(immediate)
GPR[rt] ← temp
```

**Exceptions:**

None

**Programming Notes:**

The term “unsigned” in the instruction name is a misnomer; this operation is 32-bit modulo arithmetic that does not trap on overflow. This instruction is appropriate for unsigned arithmetic, such as address arithmetic, or integer arithmetic environments that ignore overflow, such as C language arithmetic.



**Format:**  $LW\ rt,\ offset(base)$

MIPS32

**Purpose:** Load Word

To load a word from memory as a signed value

**Description:**  $GPR[rt] \leftarrow memory[GPR[base] + offset]$

The contents of the 32-bit word at the memory location specified by the aligned effective address are fetched, sign-extended to the GPR register length if necessary, and placed in GPR *rt*. The 16-bit signed *offset* is added to the contents of GPR *base* to form the effective address.

**Restrictions:**

The effective address must be naturally-aligned. If either of the 2 least-significant bits of the address is non-zero, an Address Error exception occurs.

**Operation:**

```

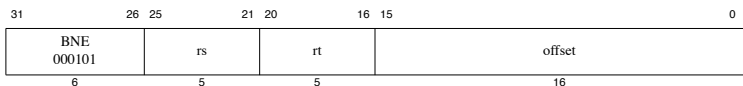
vAddr ← sign_extend(offset) + GPR[base]
if vAddr1..0 ≠ 02 then
    SignalException(AddressError)
endif
(pAddr, CCA) ← AddressTranslation (vAddr, DATA, LOAD)
memword ← LoadMemory (CCA, WORD, pAddr, vAddr, DATA)
GPR[rt] ← memword

```

**Exceptions:**

TLB Refill, TLB Invalid, Bus Error, Address Error, Watch





**Format:** BNE *rs*, *rt*, *offset*

MIPS32

**Purpose:** Branch on Not Equal

To compare GPRs then do a PC-relative conditional branch

**Description:** if  $GPR[rs] \neq GPR[rt]$  then branch

An 18-bit signed offset (the 16-bit *offset* field shifted left 2 bits) is added to the address of the instruction following the branch (not the branch itself), in the branch delay slot, to form a PC-relative effective target address.

If the contents of GPR *rs* and GPR *rt* are not equal, branch to the effective target address after the instruction in the delay slot is executed.

**Restrictions:**

Processor operation is **UNPREDICTABLE** if a branch, jump, ERET, DERET, or WAIT instruction is placed in the delay slot of a branch or jump.

**Operation:**

```

I:   target_offset ← sign_extend(offset || 02)
       condition ← (GPR[rs] ≠ GPR[rt])
I+1: if condition then
       PC ← PC + target_offset
       endif

```

**Exceptions:**

None

**Programming Notes:**

With the 18-bit signed instruction offset, the conditional branch range is  $\pm 128$  KBytes. Use jump (J) or jump register (JR) instructions to branch to addresses outside this range.

### 1.3. Tiny RISC-V Instruction Set Architecture

- RISC-V instruction set architecture
  - Brand new free, open instruction set architecture
  - Significant excitement around RISC-V hardware/software ecosystem
  - Helping to energize “open-source hardware”
  - Specifically designed to encourage subsetting and extension
  - Link to official ISA manual on course webpage
- Tiny RISC-V instruction set architecture
  - Subset we use in this course
  - Small enough for teaching, powerful enough for running real C programs
  - How is data represented? \_\_\_\_\_
  - Where can data be stored? \_\_\_\_\_
  - How can data be accessed? \_\_\_\_\_
  - What ops can be done on data? \_\_\_\_\_
  - How are inst encoded? \_\_\_\_\_
  - <http://www.csl.cornell.edu/courses/ece4750/handouts>
- **TinyRV1:** Small subset suitable for lecture, homeworks, exams
  - \_\_\_\_\_
  - \_\_\_\_\_
  - \_\_\_\_\_
- **TinyRV2:** Subset suitable for lab assignments and capable of executing simple C programs without an operating system
  - add, addi, sub, mul, and, andi, or, ori, xor, xori
  - slt, slti, sltu, sltiu
  - sra, srai, srl, srli, sll, slli
  - lui, auipc, lw, sw
  - jal, jalr, beq, bne, blt, bge, bltu, bgeu
  - csrr, csrwr

## TinyRV1 instruction assembly, semantics, and encoding

### ADD

add rd, rs1, rs2

$R[rd] \leftarrow R[rs1] + R[rs2]$

$PC \leftarrow PC + 4$

31	25 24	20 19	15 14 12 11	7 6	0
0000000	rs2	rs1	000	rd	0110011

### ADDI

addi rd, rs1, imm

$R[rd] \leftarrow R[rs1] + \text{sext}(imm)$

$PC \leftarrow PC + 4$

31	20 19	15 14 12 11	7 6	0	
imm		rs1	000	rd	0010011

### MUL

mul rd, rs1, rs2

$R[rd] \leftarrow R[rs1] \times R[rs2]$

$PC \leftarrow PC + 4$

31	25 24	20 19	15 14 12 11	7 6	0
0000001	rs2	rs1	000	rd	0110011

### LW

lw rd, imm(rs1)

$R[rd] \leftarrow M[R[rs1] + \text{sext}(imm)]$

$PC \leftarrow PC + 4$

31	20 19	15 14 12 11	7 6	0	
imm		rs1	010	rd	0000011

### SW

sw rs2, imm(rs1)

$M[R[rs1] + \text{sext}(imm)] \leftarrow R[rs2]$

$PC \leftarrow PC + 4$

31	25 24	20 19	15 14 12 11	7 6	0
imm	rs2	rs1	010	imm	0100011

$imm = \{ inst[31:25], inst[11:7] \}$

### JAL

jal rd, imm

$R[rd] \leftarrow PC + 4$

$PC \leftarrow PC + \text{sext}(imm)$

31	12 11	7 6	0
imm		rd	1101111

$imm = \{ inst[31], inst[19:12], inst[20], inst[30:21], 0 \}$

### JR

jr rs1

$PC \leftarrow R[rs1]$

31	20 19	15 14 12 11	7 6	0	
000000000000		rs1	000	00000	1100111

### BNE

bne rs1, rs2, imm

if ( $R[rs1] \neq R[rs2]$ )  $PC \leftarrow PC + \text{sext}(imm)$

else  $PC \leftarrow PC + 4$

31	25 24	20 19	15 14 12 11	7 6	0
imm	rs2	rs1	001	imm	1100011

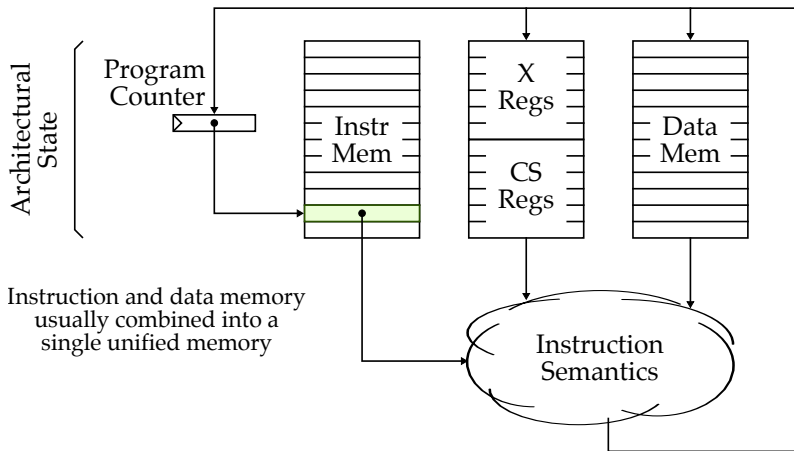
$imm = \{ inst[31], inst[7], inst[30:25], inst[11:8], 0 \}$



<b>Optional Multiply-Divide Instruction Extension: RVM</b>				
Category	Name	Fmt	RV32M (Multiply-Divide)	+RV(64,128)
<b>Multiply</b>	Multiply	R	MUL rd,rs1,rs2	MUL{W D} rd,rs1,rs2
	Multiply upper Half	R	MULH rd,rs1,rs2	
	Multiply Half Sign/Uns	R	MULHSU rd,rs1,rs2	
	Multiply upper Half Uns	R	MULHU rd,rs1,rs2	
<b>Divide</b>	DIVide	R	DIV rd,rs1,rs2	DIV{W D} rd,rs1,rs2
	DIVide Unsigned	R	DIVU rd,rs1,rs2	
<b>Remainder</b>	REMAinder	R	REM rd,rs1,rs2	REM{W D} rd,rs1,rs2
	REMAinder Unsigned	R	REMU rd,rs1,rs2	REMU{W D} rd,rs1,rs2
<b>Optional Atomic Instruction Extension: RVA</b>				
Category	Name	Fmt	RV32A (Atomic)	+RV(64,128)
<b>Load</b>	Load Reserved	R	LR.W rd,rs1	LR.{D Q} rd,rs1
<b>Store</b>	Store Conditional	R	SC.W rd,rs1,rs2	SC.{D Q} rd,rs1,rs2
<b>Swap</b>	SWAP	R	AMOSWAP.W rd,rs1,rs2	AMOSWAP.{D Q} rd,rs1,rs2
<b>Add</b>	ADD	R	AMOADD.W rd,rs1,rs2	AMOADD.{D Q} rd,rs1,rs2
<b>Logical</b>	XOR	R	AMOXOR.W rd,rs1,rs2	AMOXOR.{D Q} rd,rs1,rs2
	AND	R	AMOAND.W rd,rs1,rs2	AMOAND.{D Q} rd,rs1,rs2
	OR	R	AMOOR.W rd,rs1,rs2	AMOOR.{D Q} rd,rs1,rs2
<b>Min/Max</b>	MINimum	R	AMOMIN.W rd,rs1,rs2	AMOMIN.{D Q} rd,rs1,rs2
	MAXimum	R	AMOMAX.W rd,rs1,rs2	AMOMAX.{D Q} rd,rs1,rs2
	MINimum Unsigned	R	AMOMINU.W rd,rs1,rs2	AMOMINU.{D Q} rd,rs1,rs2
	MAXimum Unsigned	R	AMOMAXU.W rd,rs1,rs2	AMOMAXU.{D Q} rd,rs1,rs2
<b>Three Optional Floating-Point Instruction Extensions: RVF, RVD, &amp; RVQ</b>				
Category	Name	Fmt	RV32{F D Q} (HP/SP,DP,QP FP Pt)	+RV(64,128)
<b>Move</b>	Move from Integer	R	FMV.{H S}.X rd,rs1	FMV.{D Q}.X rd,rs1
	Move to Integer	R	FMV.X.{H S} rd,rs1	FMV.X.{D Q} rd,rs1
<b>Convert</b>	Convert from Int	R	FCVTE.{H S D Q}.W rd,rs1	FCVTE.{H S D Q}.L{T} rd,rs1
	Convert from Int Unsigned	R	FCVTE.{H S D Q}.WU rd,rs1	FCVTE.{H S D Q}.L{T}U rd,rs1
	Convert to Int	R	FCVTE.W.{H S D Q} rd,rs1	FCVTE.L{T}.{H S D Q} rd,rs1
	Convert to Int Unsigned	R	FCVTE.WU.{H S D Q} rd,rs1	FCVTE.L{T}U.{H S D Q} rd,rs1
<b>Load</b>	Load	I	FL{W,D,Q} rd,rs1,imm	
<b>Store</b>	Store	S	FS{W,D,Q} rs1,rs2,imm	
<b>Arithmetic</b>	ADD	R	FADD.{S D Q} rd,rs1,rs2	x0 zero --- Hard-wired zero
	SUBtract	R	FSUB.{S D Q} rd,rs1,rs2	x1 ra Caller Return address
	MULTiply	R	FMUL.{S D Q} rd,rs1,rs2	x2 sp Callee Stack pointer
	DIVide	R	FDIV.{S D Q} rd,rs1,rs2	x3 gp --- Global pointer
	SQuare RoOt	R	FSQRT.{S D Q} rd,rs1	x4 tp --- Thread pointer
<b>Mul-Add</b>	Multiply-ADD	R	FMAADD.{S D Q} rd,rs1,rs2,rs3	x5-7 t0-2 Caller Temporaries
	Multiply-SUBtract	R	FMSUB.{S D Q} rd,rs1,rs2,rs3	x8 s0/fp Callee Saved register/frame pointer
	Negative Multiply-SUBtract	R	FNMSUB.{S D Q} rd,rs1,rs2,rs3	x9 s1 Callee Saved register
	Negative Multiply-ADD	R	FNMADD.{S D Q} rd,rs1,rs2,rs3	x10-11 a0-1 Caller Function arguments/return values
<b>Sign Inject</b>	SIGN source	R	FSGNJ.{S D Q} rd,rs1,rs2	x12-17 a2-7 Caller Function arguments
	Negative SIGN source	R	FSGNJN.{S D Q} rd,rs1,rs2	x18-27 s2-11 Callee Saved registers
	Xor SIGN source	R	FSGNJX.{S D Q} rd,rs1,rs2	x28-31 t3-t6 Caller Temporaries
<b>Min/Max</b>	MINimum	R	FMIN.{S D Q} rd,rs1,rs2	f0-7 ft0-7 Caller FP temporaries
	MAXimum	R	FMAX.{S D Q} rd,rs1,rs2	f8-9 fa0-1 Callee FP saved registers
<b>Compare</b>	Compare Float =	R	FEQ.{S D Q} rd,rs1,rs2	f10-11 fa0-1 Caller FP arguments/return values
	Compare Float <	R	FLT.{S D Q} rd,rs1,rs2	f12-17 fa2-7 Caller FP arguments
	Compare Float ≤	R	FLE.{S D Q} rd,rs1,rs2	f18-27 fs2-11 Callee FP saved registers
<b>Categorization</b>	Classify Type	R	FCLASS.{S D Q} rd,rs1	f28-31 ft8-11 Caller FP temporaries
<b>Configuration</b>	Read Status	R	FRCSR rd	
	Read Rounding Mode	R	FRRM rd	
	Read Flags	R	FRFLAGS rd	
	Swap Status Reg	R	FSCSR rd,rs1	
	Swap Rounding Mode	R	FSRM rd,rs1	
	Swap Flags	R	FSFLAGS rd,rs1	
	Swap Rounding Mode Imm	I	FSRMI rd,imm	
Swap Flags Imm	I	FSFLAGSI rd,imm		

RISC-V calling convention and five optional extensions: 10 multiply-divide instructions (RV32M); 11 optional atomic instructions (RV32A); and 25 floating-point instructions each for single-, double-, and quadruple-precision (RV32F, RV32D, RV32Q). The latter add registers f0-f31, whose width matches the widest precision, and a floating-point control and status register fcsr. Each larger address adds some instructions: 4 for RVM, 11 for RVA, and 6 each for RVF/D/Q. Using regex notation, {} means set, so L{D|Q} is both LD and LQ. See riscv.org. (8/21/15 revision)

## 2. Processor Functional-Level Model



### 2.1. Transactions and Steps

- We can think of each instruction as a **transaction**
- Executing a transaction involves a sequence of **steps**

---

add   addi   mul   lw   sw   jal   jr   bne
---

---

Fetch Instruction

---

Decode Instruction

---

Read Registers

---

Register Arithmetic

---

Read Memory

---

Write Memory

---

Write Registers

---

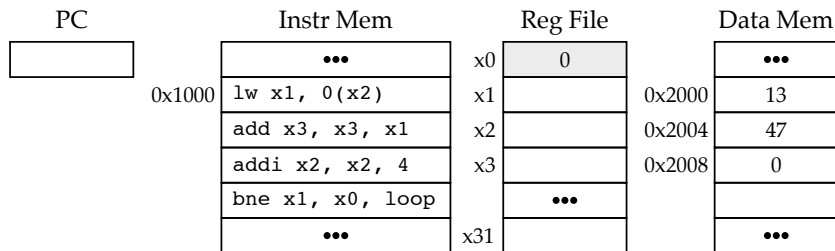
Update PC

---

## 2.2. TinyRV1 Simple Assembly Example

Static Asm Sequence	Instruction Semantics
loop: lw x1, 0(x2)	
add x3, x3, x1	
addi x2, x2, 4	
bne x1, x0, loop	

### Worksheet illustrating processor functional-level model



### Table illustrating processor functional-level model

PC	Dynamic Asm Sequence	x1	x2	x3
	lw x1, 0(x2)			
	add x3, x3, x1			
	addi x2, x2, 4			
	bne x1, x0, loop			
	lw x1, 0(x2)			
	add x3, x3, x1			





## 2.4. TinyRV1 Mystery Assembly and C Program

What is the C code corresponding to the TinyRV1 assembly shown below? Assume assembly implements a function.

---

---

---

---

---

---

---

---

---

---

```
addi x5, x0, 0

loop:
  lw   x6, 0(x12)
  bne  x6, x14, foo
  addi x10, x5, 0
  jr   x1

foo:
  addi x12, x12, 4
  addi x5, x5, 1
  bne  x5, x13, loop

addi x10, x0, -1
jr   x1
```

### 3. Processor/Laundry Analogy

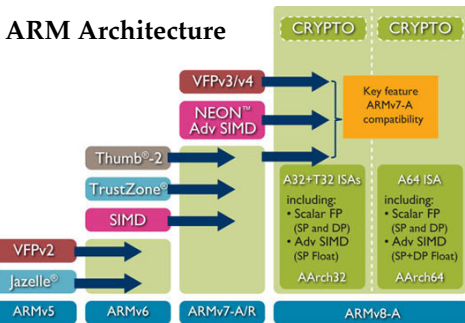
#### • Processor

- Instructions are “transactions” that execute on a processor
- Architecture: defines the hardware/software interface
- Microarchitecture: how hardware executes sequence of instructions

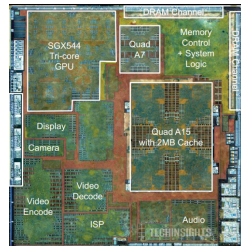
#### • Laundry

- Cleaning a load of laundry is a “transaction”
- Architecture: high-level specification, dirty clothes in, clean clothes out
- Microarchitecture: how laundry room actually processes multiple loads

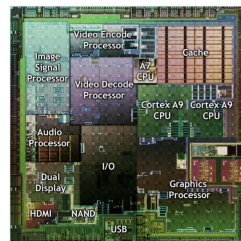
### 3.1. Arch vs. $\mu$ Arch vs. VLSI Impl



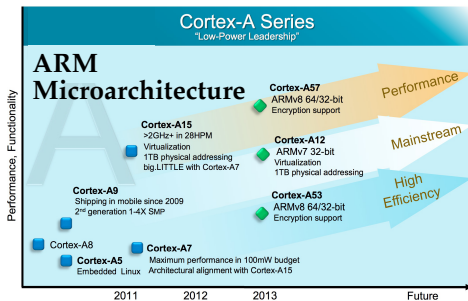
#### ARM VLSI Implementation



Samsung Exynos Octa

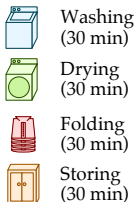


NVIDIA Tegra 2

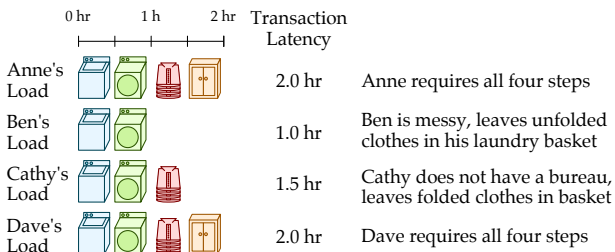


## 3.2. Processor Microarchitectural Design Patterns

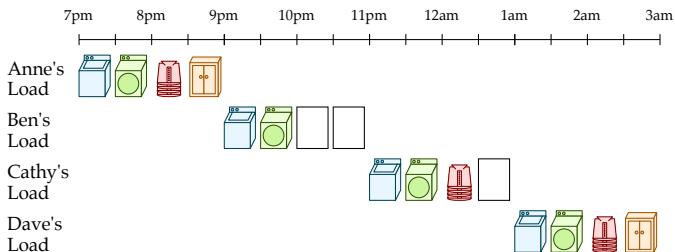
### Transaction Steps



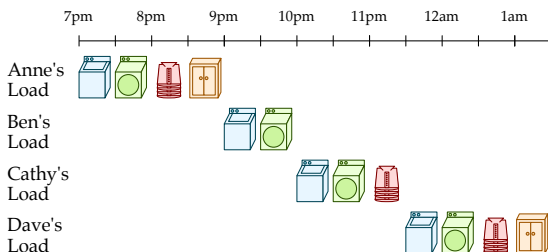
### Four Types of Transactions



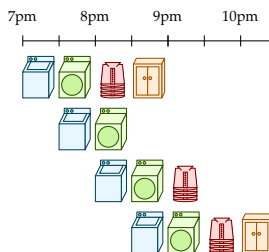
### Fixed Time Slot Laundry (Single-Cycle Processors)



### Variable Time Slot Laundry (FSM Processors)



### Pipelined Laundry





## 4. Analyzing Processor Performance

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

- Instructions / program depends on source code, compiler, ISA
- Avg cycles / instruction (CPI) depends on ISA, microarchitecture
- Time / cycle depends upon microarchitecture and implementation

Using our first-order equation for processor performance and a functional-level model, the execution time is just the number of dynamic instructions.

Microarchitecture	CPI	Cycle Time
Single-Cycle Processor	1	long
FSM Processor	>1	short
Pipelined Processor	≈1	short



Students often confuse “Cycle Time” with the execution time of a sequence of transactions measured in cycles. “Cycle Time” is the clock period or the inverse of the clock frequency.

## Estimating dynamic instruction count

Estimate the dynamic instruction count for the vector-vector add example assuming  $n$  is 64?

```
loop:
    lw    x5, 0(x13)
    lw    x6, 0(x14)
    add   x7, x5, x6
    sw    x7, 0(x12)
    addi  x13, x12, 4
    addi  x14, x14, 4
    addi  x12, x12, 4
    addi  x15, x15, -1
    bne   x15, x0, loop
    jr    x1
```

Estimate the dynamic instruction count for the mystery program assuming  $n$  is 64 and that we find a match on the final element.

```
    addi  x5, x0, 0
loop:
    lw    x6, 0(x12)
    bne   x6, x14, foo
    addi  x10, x5, 0
    jr    x1
foo:
    addi  x12, x12, 4
    addi  x5, x5, 1
    bne   x5, x13, loop
    addi  x10, x0, -1
    jr    x1
```