# EECS 322 Computer Architecture

## Superpipline and the Cache



*Instructor: Francis G. Wolff*
*wolff@eecs.cwru.edu*
*Case Western Reserve University*
*This presentation uses powerpoint animation: please viewshow*
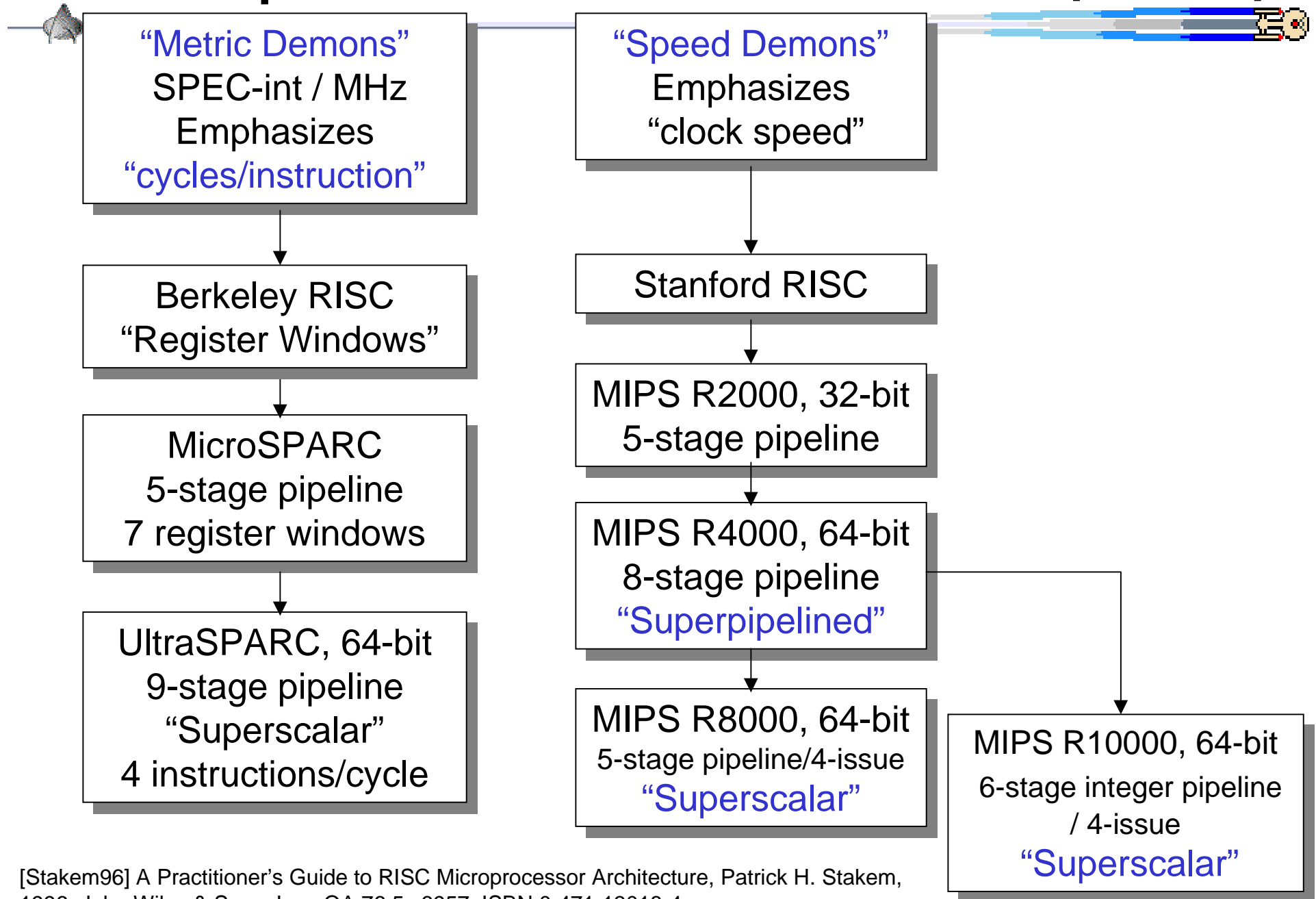
# Summary: Instruction Hazards

| Instruction Class | Integer Application | Floating-Point (FP) Application |
|---|---|---|
| Integer Arithmetic | 50% | 25% |
| FP-Arithmetic | 0% | 30% |
| Loads | 17% | 25% |
| Stores | 8% | 15% |
| Branches | 25% | 5% |

| | No-Forwarding | Forwarding | Hazard |
|---|---|---|---|
| R-Format | 1-3 | 1 | Data |
| Load | 1-3 | 1-2 | Data, Structural |
| Store | 1 | 1-2 | Structural |

| | No Delay Slot | Delay Slot | Hazard |
|---|---|---|---|
| Branch (decision is made in the ID stage) | 2 | 1 | Control |
| Branch (decision is made in the EX stage) | 3 | 1 | Control |
| Jump | 2 | 1 | |

**Structural Hazard:** Instruction & Data memory combined.

# RISC camps

**"Metric Demons"**
SPEC-int / MHz
Emphasizes
**"cycles/instruction"**

**"Speed Demons"**
Emphasizes
"clock speed"

Berkeley RISC
"Register Windows"

Stanford RISC

MicroSPARC
5-stage pipeline
7 register windows

MIPS R2000, 32-bit
5-stage pipeline

UltraSPARC, 64-bit
9-stage pipeline
"Superscalar"
4 instructions/cycle

MIPS R4000, 64-bit
8-stage pipeline
**"Superpipelined"**

MIPS R8000, 64-bit
5-stage pipeline/4-issue
**"Superscalar"**

MIPS R10000, 64-bit
6-stage integer pipeline
/ 4-issue
**"Superscalar"**

# Instruction Level Parallelism (ILP)

## Superpipelined scheme

• uses a longer pipeline with more stages to reduce cycle time

• simple dependencies: structural, data, control pipeline hazards.

• requires higher clock speeds

• require little additional logic to baseline processor

• Branches cause a latency of 3 internal clocks and loads a 2-cycle latency.

• However, superpipelining increases performance, because each stage can run at twice the system clock.
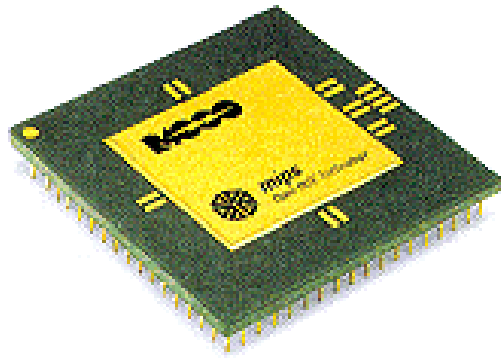
# Instruction Level Parallelism (ILP)

- Superscalar scheme

- multiple execution units by duplicating the functional units (ALUs)

- combinatorial dependance problem

  ✓ instructions can only be issued only if they are independent

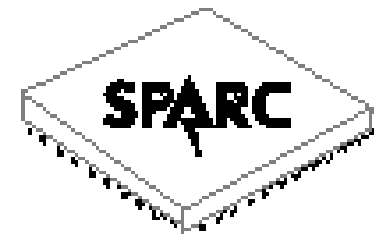- require sophisticated & complex logic (i.e. instruction Scheduler)

# R4400 processor

Both R4000 and R4400 use:

- True 64bit architecture
- On-chip TLB for virtual-to-physical address translation
- 8-stage superpipelined architecture
- large integrated caches

100MHz offers
85 MIPS, 16 Mflops,
70 SPECmarks

150MHz R4400 offers
136 MIPS, 24 Mflops,
100 SPECmarks.

- ✦ UltraSPARC is a high-volume processor — the R4400 is a low volume processor
- ✦ UltraSPARC enjoys the installed base of the SPARC product line — the R4400 has a small installed base
- ✦ UltraSPARC is binary compatible with other SPARC processors — the R4400 has narrow application support
- ✦ UltraSPARC-I is network and multimedia ready — the R4400 does not provide multimedia support
- ✦ UltraSPARC provides high performance and high bandwidth — the R4400 provides comparatively low integer and floating point performance

Ref: http://sunsite.ics.forth.gr/sunsite/Sun/sun_microelectronics/UltraSparc/ultra_arch_versus.html

# μP Comparisons

| | MIPS R4400 | UltraSparc I |
|---|---|---|
| **Clock** | **250 MHz** | **167 MHz** |
| Bus speed | 50/66/75 | 83 Mhz |
| Pipeline | 8-stage | 9-stage |
| Superscalar | 1-issue | 4-issue |
| **Branch prediction** | **no** | **yes** |
| TLB | 48 even/odd | 64-Inst/64-Data |
| L1 I/D-cache | 16k/16k | 16k/16k |
| Associativity | 1-way (direct) | 2-way |
| L2 cache | 1 Mb | 0.5Mb |
| CMOS technology | 0.35μ | 0.5μ 4 layers |
| Fabrication Vendor | NEC, IDT, Toshiba | Fujitsu |
| **Year** | **1993** | **1995** |
| Voltage | 3.3 volts | 3.3 volts |
| **Transistors** | **2.2 million** | **3.8-5.2 million** |
| SpecInt92/fp92 | 175/178 | 250/350 |
| **SpecInt95/fp95** | **5.07/?** | **6.26/9.06** |
| Cost: | $1250 | $1395 |

Ref: http://sunsite.ics.forth.gr/sunsite/Sun/sun_microelectronics/UltraSparc/ultra_arch_versus.html
http://www.mips.com/Documentation/R4400_Overview.pdf
http://www.spec.org/osg/cpu95/results/res96q3/cpu95-960624-00962.html **and** http://www.eecs.umich.edu/~tnm/theses/mikeu.pdf

# R4000: no dynamic branch prediction

## R4000 Performance

- **Not ideal CPI of 1:**
  - **Load stalls** (1 or 2 clock cycles)
  - Branch stalls (2 cycles + unfilled slots)
  - **FP result stalls**: RAW data hazard (latency)
  - **FP structural stalls**: Not enough FP hardware (parallelism)



Legend: ■ Base ■ Load stalls ■ Branch stalls ■ FP result stalls ■ FP structural stalls

(Benchmarks: eqntott, espresso, gcc, li, doduc, nasa7, ora, spice2g6, su2cor, tomcatv)
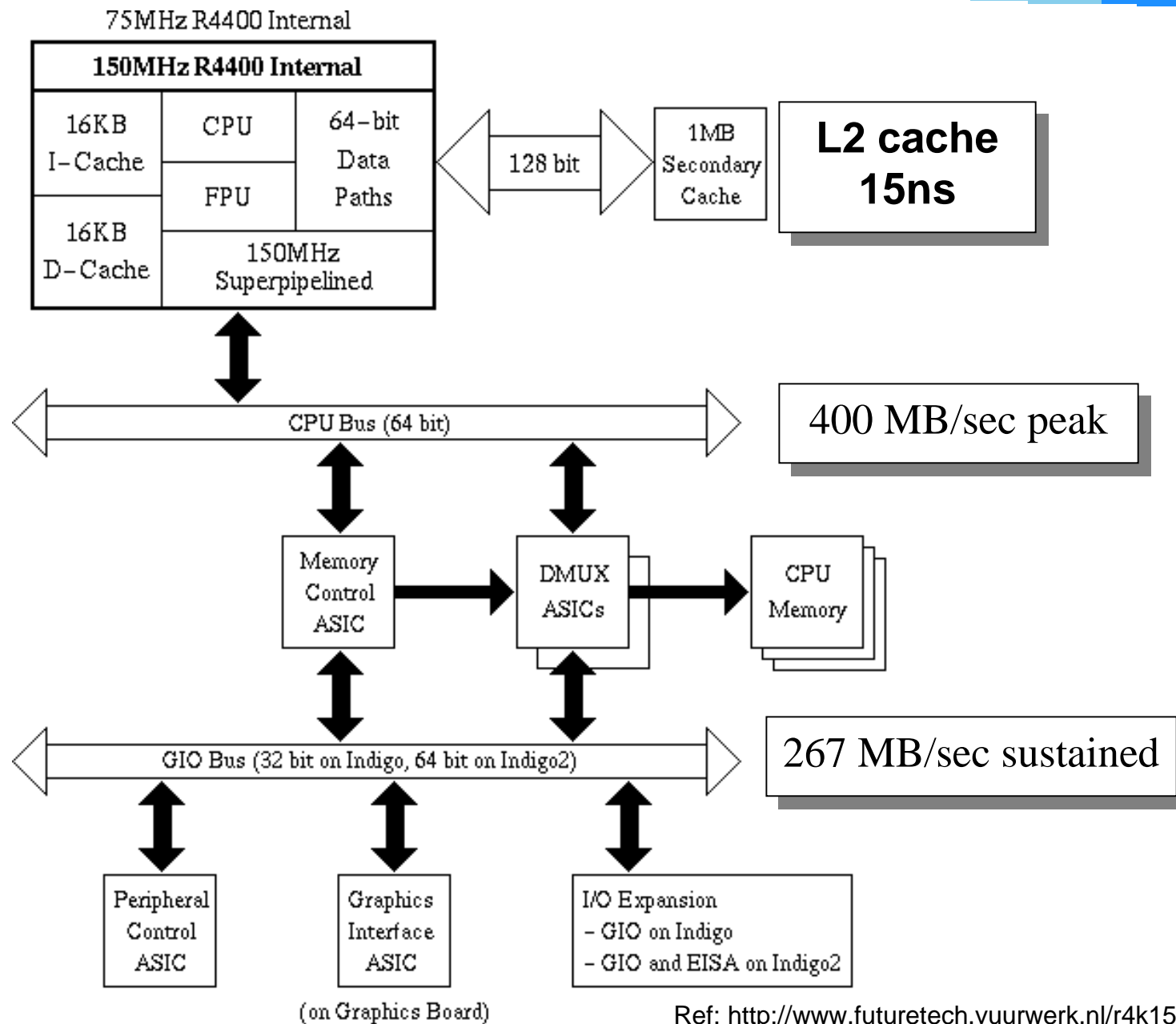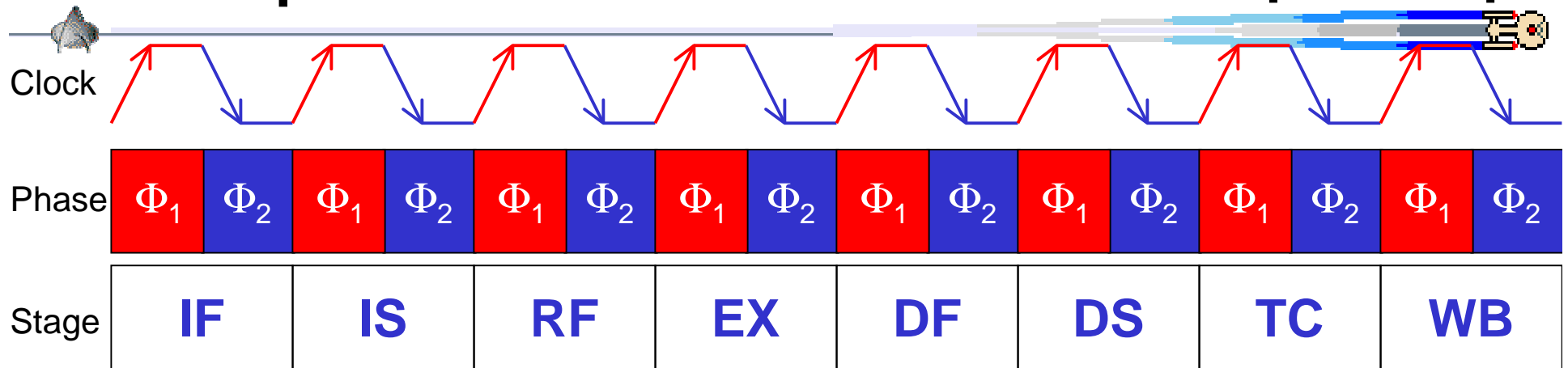
# Differences Between the MIPS R4400 and UltraSPARC-I

• The MIPS R4400 uses an 8-stage pipeline architecture, and is less efficient than the superscalar, pipelined UltraSPARC-I.

• Although it is an integrated processor, the R4400 requires several other modules in order to incorporate it into a system.

• External secondary caches (L2) must be designed around the processor, and multiprocessor and graphics support are not provided.

• The highly integrated UltraSPARC-I, utilizing on-chip caches, an advanced processor design and UPA architecture, requires little to complete its chip set, significantly easing its integration into systems.
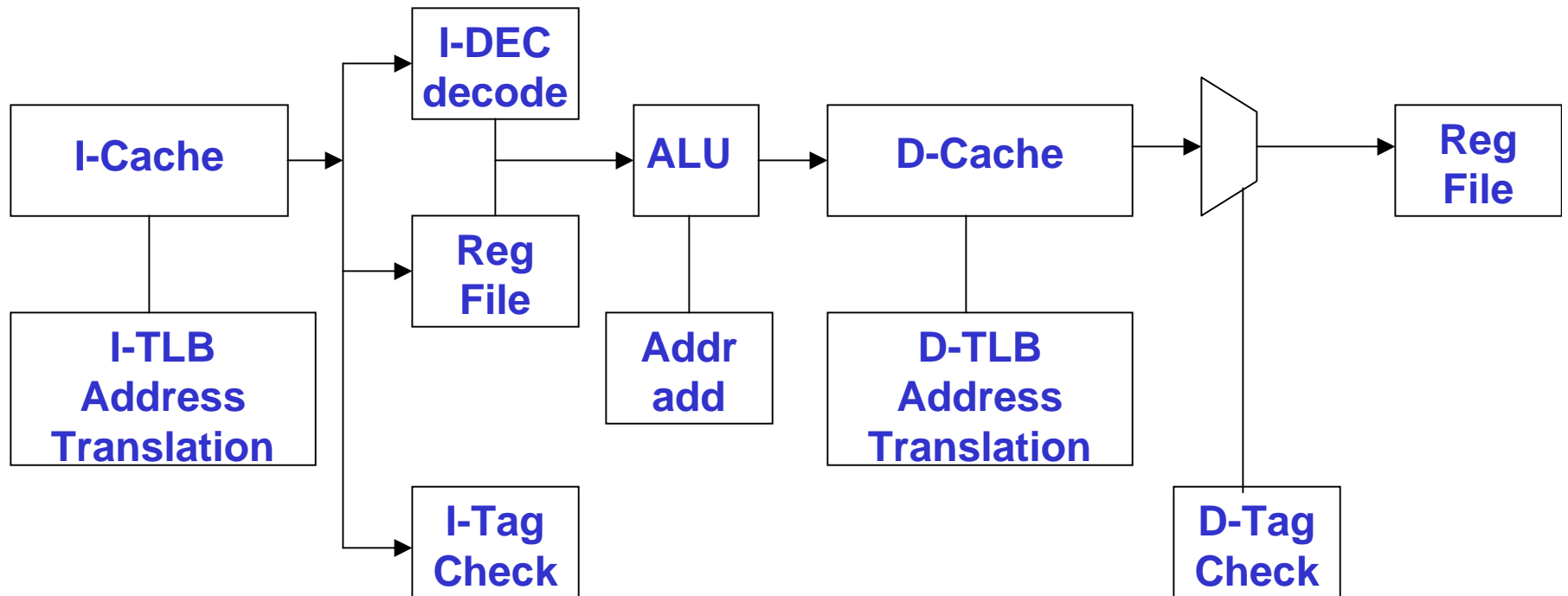
# R4400 Bus



75MHz R4400 Internal

| 150MHz R4400 Internal |
|---|

| 16KB I-Cache | CPU | 64-bit Data Paths |
| | FPU | |
| 16KB D-Cache | 150MHz Superpipelined | |

128 bit

1MB Secondary Cache

**L2 cache 15ns**

CPU Bus (64 bit)

400 MB/sec peak

Memory Control ASIC → DMUX ASICs → CPU Memory

GIO Bus (32 bit on Indigo, 64 bit on Indigo2)

267 MB/sec sustained

Peripheral Control ASIC

Graphics Interface ASIC

I/O Expansion
– GIO on Indigo
– GIO and EISA on Indigo2

(on Graphics Board)

Ref: http://www.futuretech.vuurwerk.nl/r4k150upgrade.html

# R4000 Pipeline

**Clock**

**Phase** | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$

**Stage** | IF | IS | RF | EX | DF | DS | TC | WB
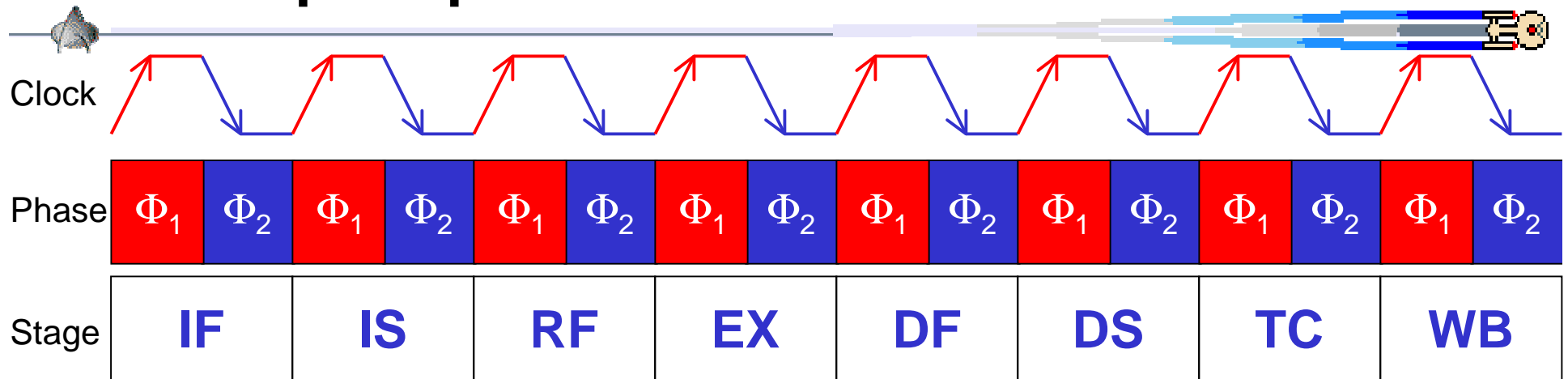
- IF - Instruction Fetch, first half
- IS - Instruction fetch, Second half
- RF - Register Fetch
- EX - Execution (branch compare)
- DF - Data Fetch, first half
- DS - Data fetch, Second half
- TC - Tag Check
- WB - Write Back

# R4400 SuperPipeline

| Stage | | | | | | | | |
|---|---|---|---|---|---|---|---|---|

**Clock**

**Phase**

| $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ | $\Phi_1$ | $\Phi_2$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**Stage**

| IF | IS | RF | EX | DF | DS | TC | WB |
|---|---|---|---|---|---|---|---|

I-Cache

I-DEC decode

Reg File

I-Tag Check

I-TLB Address Translation

ALU

Addr add

D-Cache

D-TLB Address Translation

D-Tag Check

Reg File

# R4000 Pipeline stages: IF & IS

IF -  instruction Fetch, First half

• PC: Branch logic selects an instruction address and

• instruction catch fetch begins

• I-TLB: instruction translation lookaside buffer begins the virtual-to-physical address translation

IS -  instruction fetch, second half

• Complete instruction catch fetch and

• the virtual-to-physical address translation

# R4000 Pipeline stages: RF

## RF - register fetch

- I-DEC: instruction decoder decodes the instruction and checks for interlock conditions

- instruction cache tag is checked against the page frame number (PFN) obtained from the ITLB.

- Any required operands are fetched from the register file

# R4000 Pipeline stages: EX

EX - execution

• Register-to-register instructions: The ALU performs arithmetic or logical operation

• Load & Store instructions: the ALU calculates the data virtual address (i.e. offset + base register).

• Branch instructions: The ALU determines whether the branch condition is true & calculates the virtual target address.

# R4000 Pipeline stages: DF & DS

## DF - data fetch, first half

• **Register-to-Register:** No operations are performed during DF, DS, and TC stages

• **Load & Store instructions:** The data cache fetch and the data virtual-to-physical translation begins

• **Branch instructions:** address translation and TLB update begins

## DS - data fetch, second half

• **Load & Store:** completion of data cache fetch & data virtual-to-physical translation. The shifter aligns data to its word or doubleword boundary

• **branch:** completion of instruction address translation and TLB update

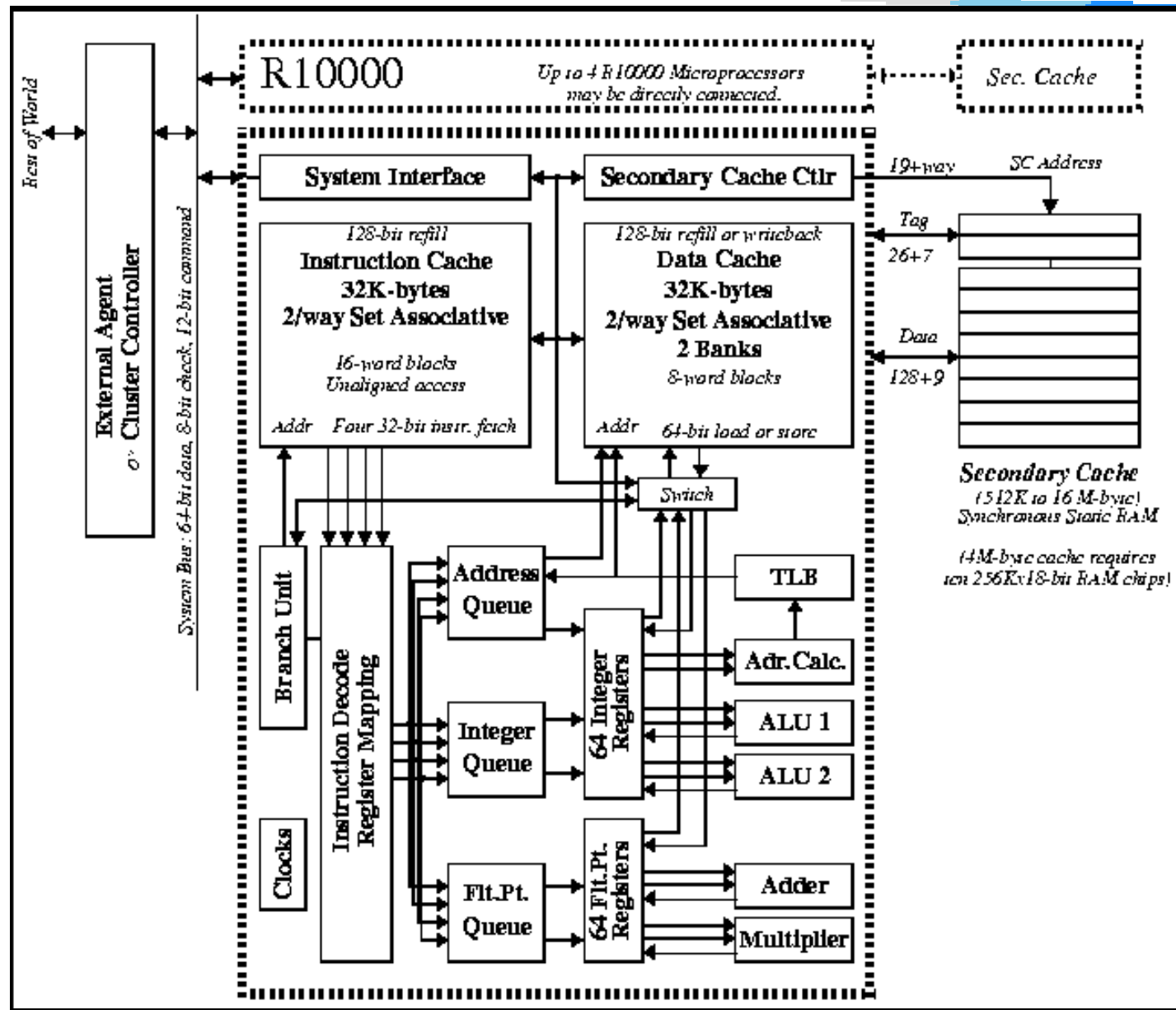# R4000 Pipeline stages: TC & WB

## TC - Tag check

• **Load & Store instructions:** the cache performs the tag check.

• **Hit or Miss:** physical address from TLB is checked against the tag check to determine if there is a hit or miss.
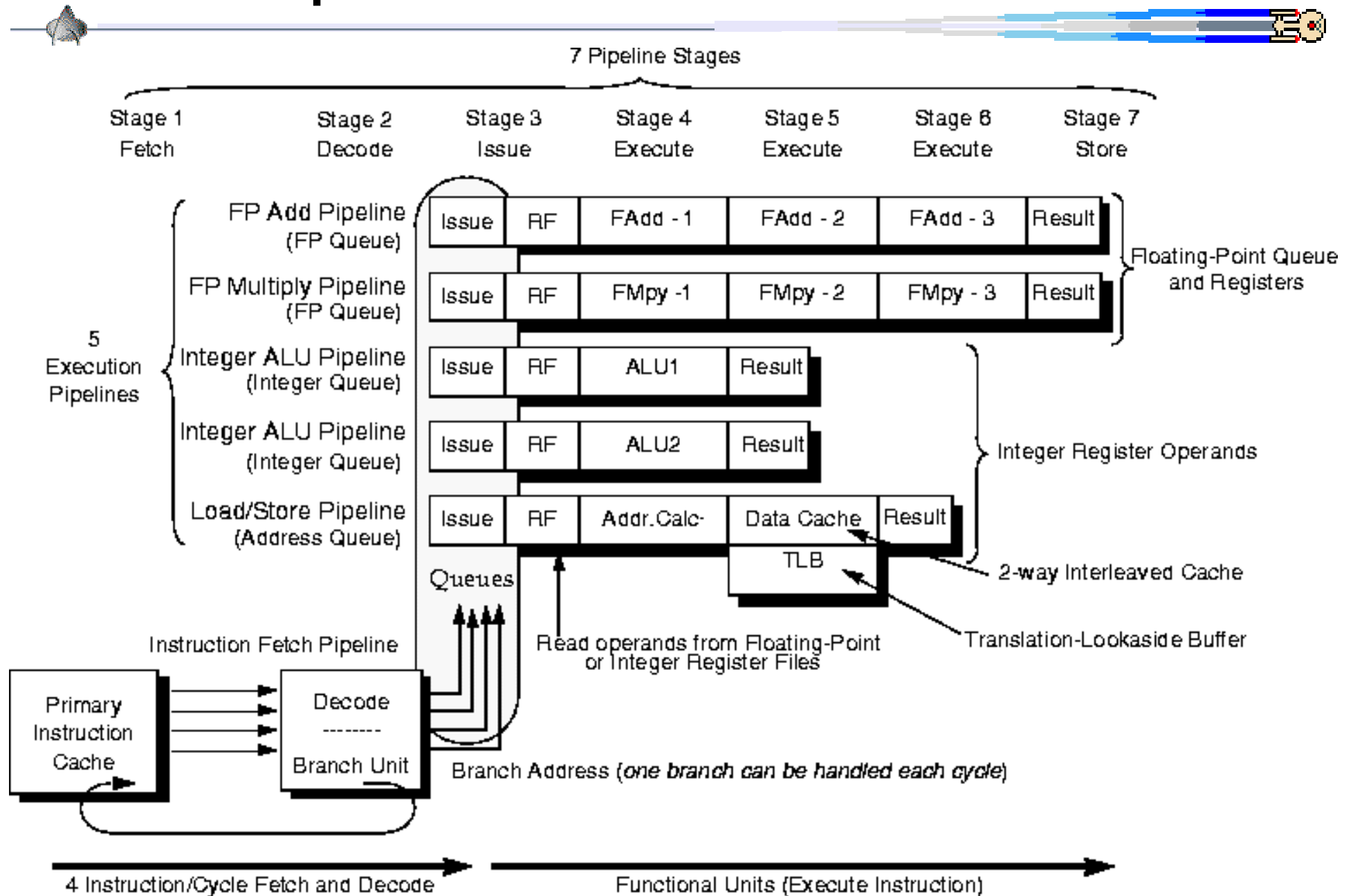
## WB - write back

• **Register-to-register & load:** the instruction result is written back to the register file

• **Branch:** no operation
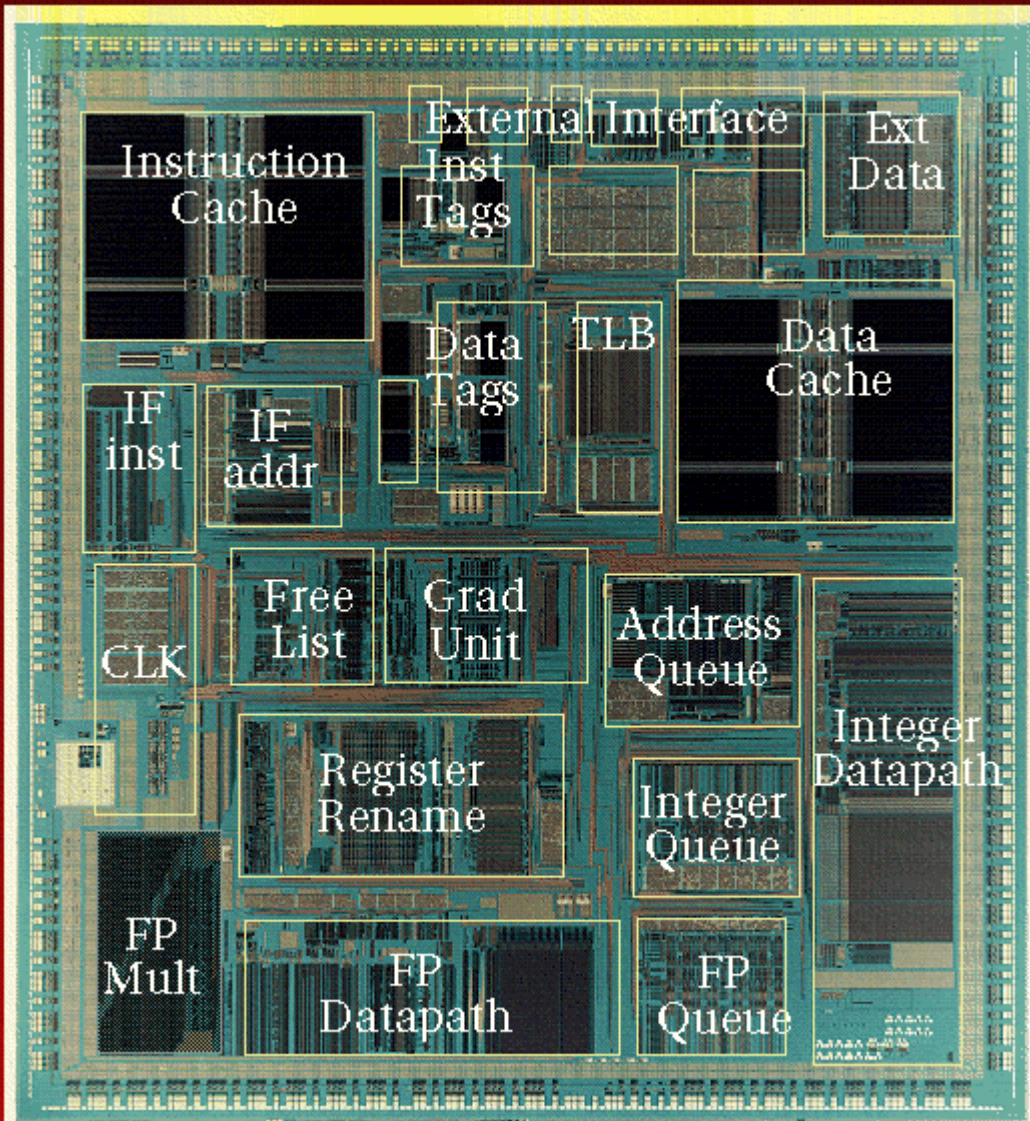
# R10000 superscalar architecture

# R10000 - superscalar



7 Pipeline Stages

| Stage 1 Fetch | Stage 2 Decode | Stage 3 Issue | Stage 4 Execute | Stage 5 Execute | Stage 6 Execute | Stage 7 Store |

FP Add Pipeline (FP Queue) — Issue | RF | FAdd - 1 | FAdd - 2 | FAdd - 3 | Result

FP Multiply Pipeline (FP Queue) — Issue | RF | FMpy -1 | FMpy - 2 | FMpy - 3 | Result

Floating-Point Queue and Registers

5 Execution Pipelines

Integer ALU Pipeline (Integer Queue) — Issue | RF | ALU1 | Result

Integer ALU Pipeline (Integer Queue) — Issue | RF | ALU2 | Result

Load/Store Pipeline (Address Queue) — Issue | RF | Addr.Calc- | Data Cache | Result

TLB

Integer Register Operands

2-way Interleaved Cache

Translation-Lookaside Buffer

Queues

Read operands from Floating-Point or Integer Register Files

Instruction Fetch Pipeline

Primary Instruction Cache

Decode
--------
Branch Unit

Branch Address (*one branch can be handled each cycle*)

4 Instruction/Cycle Fetch and Decode

Functional Units (Execute Instruction)

Ref: http://www.sgi.com/processors/r10k/manual/T5.HW.Ch01.intro_AFrame_16.gif

# R10000 die



R10K die size 16.6mm x 17.9mm

Instruction Cache
External Interface
Inst Tags
Ext Data
IF inst
IF addr
Data Tags
TLB
Data Cache
CLK
Free List
Grad Unit
Address Queue
Register Rename
Integer Queue
Integer Datapath
FP Mult
FP Datapath
FP Queue

## R10000

SPECint95 base 14.1
SPECint95 peak 14.7 SPECfp95 base 22.6
SPECfp95 peak 24.5

200 MHz Clock

I/D-cache: 32k/32k

TLB: 64 entries
Virtual Page Sizes: 16k-16M

0.35µ 4-layer CMOS technology

17 mm x18 mm chip

contains about **6.7** million transistors

including about 4.4 million transistors in its primary caches.

# Principle of Locality

- **Principle of Locality**
  states that programs access a relatively small portion of their address space at <u>any instance of time</u>

- <u>**Two types of locality**</u>

  - <u>**Temporal locality**</u> **(locality in time)**
    If an item is referenced, then
    <u>the same</u> item will tend to be referenced soon
    "the tendency to reuse recently accessed data items"

  - <u>**Spatial locality**</u> **(locality in space)**
    If an item is referenced, then
    <u>nearby</u> items will be referenced soon
    "the tendency to reference nearby data items"

**Figure 7.2**

# Cache Example

**Time 1: Hit: in cache**

Processor

**Time 3: deliver to CPU**

**Time 1: Miss**

Data are transferred

**Time 2: fetch from lower level into cache**

**Hit time** = Time 1          **Miss penalty** = Time 2 + Time 3

# Modern Systems: Pentium Pro and PowerPC



| Characteristic | Intel Pentium Pro | PowerPC 604 |
|---|---|---|
| Cache organization | Split instruction and data caches | Split intruction and data caches |
| Cache size | 8 KB each for instructions/data | 16 KB each for instructions/data |
| Cache associativity | Four-way set associative | Four-way set associative |
| Replacement | Approximated LRU replacement | LRU replacement |
| Block size | 32 bytes | 32 bytes |
| Write policy | Write-back | Write-back or write-through |

# Cache Terminology

**A hit** if the data requested by the CPU is in the upper level

**Hit rate** or **Hit ratio**
        is the fraction of accesses found in the upper level

**Hit time**
        is the time required to access data in the upper level
        = <detection time for hit or miss> + <hit access time>

**A miss** if the data is not found in the upper level

**Miss rate** or **(1 − hit rate)**
        is the fraction of accesses <u>not</u> found in the upper level

**Miss penalty**
        is the time required to access data in the lower level
        = <lower access time>+<reload processor time>

# Direct Mapped Cache

- **Direct Mapped:** assign the cache location based on the address of the word in memory

- **cache_address = memory_address** modulo **cache_size;**

Cache

| 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |

Memory

| 00001 | 00101 | 01001 | 01101 | 10001 | 10101 | 11001 | 11101 |

**Observe there is a Many-to-1 memory to cache relationship**

# Direct Mapped Cache: Data Structure
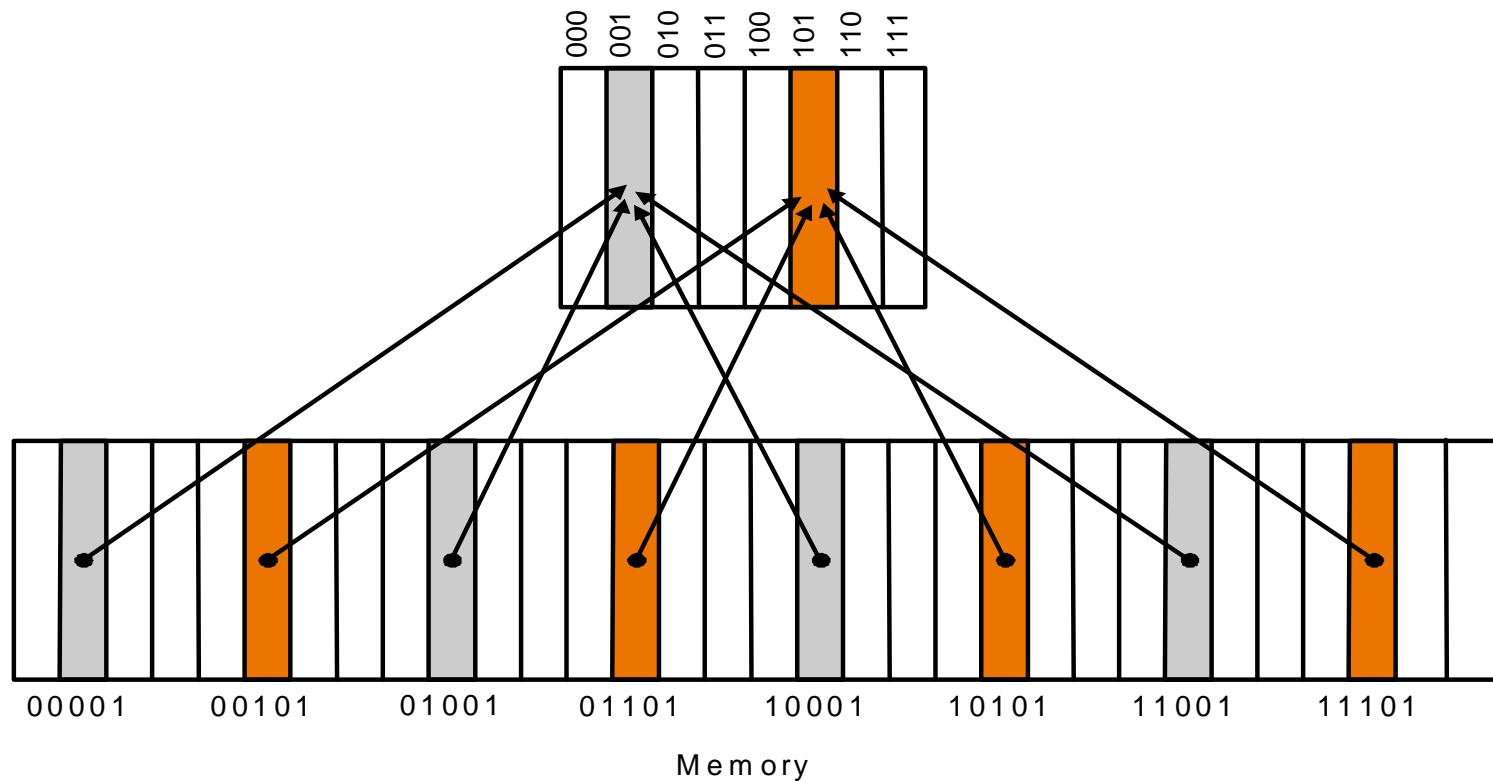
There is a **Many-to-1 relationship** between memory and cache

How do we know whether the data in the cache corresponds to the requested word?

## tags
- contain the address information required to identify whether a word in the cache corresponds to the requested word.

- tags need only to contain the upper portion of the memory address (often referred to as a page address)

## valid bit
- indicates whether an entry contains a valid address

# Direct Mapped Cache: Temporal Example

Figure 7.6

lw     $1,**10 110** ($0)      **Miss: valid**      lw     $1,22($0)

lw     $2,**11 010** ($0)      **Miss: valid**      lw     $2,26($0)

lw     $3,**10 110** ($0)      **Hit!**      lw     $3,22($0)

| Index | Valid | Tag | Data |
|-------|-------|-----|------|
| 000 | N | | |
| 001 | N | | |
| 010 | Y | 11 | Memory[11010] |
| 011 | N | | |
| 100 | N | | |
| 101 | N | | |
| 110 | Y | 10 | Memory[10110] |
| 111 | N | | |

# Direct Mapped Cache: Worst case, always miss!

Figure 7.6

lw     $1,**10 110** ($0)

lw     $2,**11 110** ($0)

lw     $3,**00 110** ($0)

**Miss: valid**

**Miss: tag**

**Miss: tag**

lw     $1,22($0)

lw     $2,30($0)

lw     $3,6($0)

| Index | Valid | Tag | Data |
|-------|-------|-----|------|
| 000 | N | | |
| 001 | N | | |
| 010 | N | | |
| 011 | N | | |
| 100 | N | | |
| 101 | N | | |
| 110 | Y | 00 | Memory[00110] |
| 111 | N | | |

# Direct Mapped Cache: Mips Architecture

Figure 7.7

**Tag**

**Index**

31 30 · · · 13 12 11 · · · 2 1 0

| | |
|---|---|
| | Byte offset |

20

10

Tag

Index

**Hit**

**Data**

| Index | Valid | Tag | Data |
|---|---|---|---|
| 0 | | | |
| 1 | | | |
| 2 | | | |
| · · · | | | |
| · · · | | | |
| · · · | | | |
| 1021 | | | |
| 1022 | | | |
| 1023 | | | |

20

32

=

**Compare Tags**

# Direct Mapped Cache

- **Direct Mapped:** assign the cache location based on the address of the word in memory

- **cache_address = memory_address % cache_size;**



**Observe there is a Many-to-1 memory to cache relationship**

# Direct Mapped Cache: Mips Architecture

Figure 7.7

**Tag**

**Index**

31 30 · · · 13 12 11 · · · 2 1 0

| | | | Byte offset |
|---|---|---|---|

20

10

Tag

Index

**Hit**

**Data**

| Index | Valid | Tag | Data |
|---|---|---|---|
| 0 | | | |
| 1 | | | |
| 2 | | | |
| · · · | | | |
| · · · | | | |
| 1021 | | | |
| 1022 | | | |
| 1023 | | | |

20

32

=

**Compare Tags**

# Bits in a Direct Mapped Cache

How many total bits are required for a direct mapped cache
with 64KB (= $2^{16}$ KiloBytes) of data
and one word (=32 bit) blocks
assuming a 32 bit <u>byte</u> memory address?

Cache index width = $\log_2$ words
$$= \log_2 2^{16}/4 = \log_2 2^{14} \text{ words} = 14 \text{ bits}$$

Block address width = <byte address width> − $\log_2$ word
$$= 32 - 2 = 30 \text{ bits}$$

Tag size = <block address width> − <cache index width>
$$= 30 - 14 = 16 \text{ bits}$$

Cache block size = <valid size>+<tag size>+<block data size>
$$= 1 \text{ bit} + 16 \text{ bits} + 32 \text{ bits} = 49 \text{ bits}$$

Total size = <Cache word size> $\times$ <Cache block size>
$$= 2^{14} \text{ words} \times 49 \text{ bits} = 784 \times 2^{10} = 784 \text{ Kbits} = 98 \text{ KB}$$
$$= 98 \text{ KB}/64 \text{ KB} = 1.5 \text{ times overhead}$$

# The DECStation 3100 cache

## write-through cache

Always write the data into both the
cache and memory and then wait for memory.

## DECStation uses a write-through cache

- 128 KB total cache size (=32K words)
  - = 64 KB instruction cache (=16K words)
  - + 64 KB data cache (=16K words)

- 10 processor clock cycles to write to memory

In a gcc benchmark, 13% of the instructions are stores.

- Thus, CPI of 1.2 becomes 1.2+13%x10 = 2.5
- Reduces the performance by more than a factor of 2!

# Cache schemes

**write-through cache**
Always write the data into both the cache and memory and then wait for memory.

**write buffer**
write data into cache and write buffer.
If write buffer full processor must stall.

No amount of buffering can help
if writes are being generated faster
than the memory system can accept them.

**write-back cache**
Write data into the cache block and
only write to memory when block is modified
but complex to implement in hardware.

Chip Area    Speed

# Hits vs. Misses

- **Read hits**
  - ✓ this is what we want!

- **Read misses**
  - ✓ stall the CPU, fetch block from memory, deliver to cache, and restart.

- **Write hits**
  - ✓ write-through: can replace data in cache and memory.
  - ✓ write-buffer: write data into cache and buffer.
  - ✓ write-back: write the data only into the cache.

- **Write misses**
  - ✓ read the entire block into the cache, then write the word.

# The DECStation 3100 miss rates

**Figure 7.9**

- A split instruction and data cache increases the bandwidth

| Benchmark Program | gcc | spice |
|---|---|---|
| Instruction miss rate | 6.1% | 1.2% |
| Data miss rate | 2.1% | 1.3% |
| Effective split miss rate | 5.4% | 1.2% |
| Combined miss rate | 4.8% | |

**Why a lower miss rate?**

**Numerical programs tend to consist of a lot of small program loops**

**split cache has slightly worse miss rate**

# Spatial Locality

- **Temporal only cache**
  cache block contains only one word (No spatial locality).

- **Spatial locality**
  Cache block contains multiple words.

  - When a miss occurs, then fetch multiple words.

  - Advantage
    Hit ratio increases because there is a high probability that the adjacent words will be needed shortly.

  - Disadvantage
    Miss penalty increases with block size

Figure 7.10

# Spatial Locality: 64 KB cache, 4 words

- **64KB cache using four-word (16-byte word)**
- **16 bit tag, 12 bit index, 2 bit block offset, 2 bit byte offset.**

- **Use split caches because there is more spatial locality in code:**

| Program Block size | gcc =1 | gcc =4 | spice =1 | spice =4 |
|---|---|---|---|---|
| Instruction miss rate | 6.1% | 2.0% | 1.2% | 0.3% |
| Data miss rate | 2.1% | 1.7% | 1.3% | 0.6% |
| Effective split miss rate | 5.4% | 1.9% | 1.2% | 0.4% |
| Combined miss rate | 4.8% | 4.8% | | |

**Temporal only split cache: has slightly worse miss rate**

**Spatial split cache: has lower miss rate**

# Cache Block size Performance

Figure 7.12

- **Increasing the block size tends to decrease miss rate:**



Block size (bytes)

- 1 KB
- 8 KB
- 16 KB
- 64 KB
- 256 KB

# Designing the Memory System

Figure 7.13

- **Make reading multiple words easier by using banks of memory**



a. One-word-wide memory organization
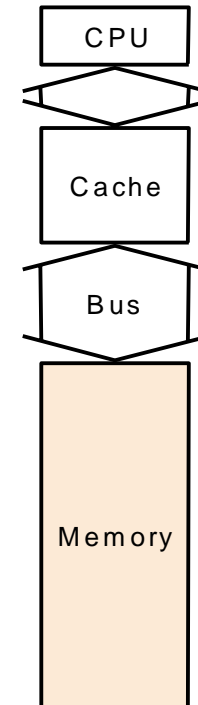
b. Wide memory organization

c. Interleaved memory organization

# 1-word-wide memory organization

Figure 7.13

**Suppose we have a system as follows**

- **1-word-wide memory organization**
- **1 cycle to send the address**
- **15 cycles to access DRAM**
- **1 cycle to send a word of data**

CPU

Cache

Bus

Memory

a. One-word-wide memory organization

**If we have a cache block of 4 words**

**Then the miss penalty is**
**=(1 address send) + 4×(15 DRAM reads)+4×(1 data send)**
**= 65 clocks per block read**

**Thus the number of bytes transferred per clock cycle**
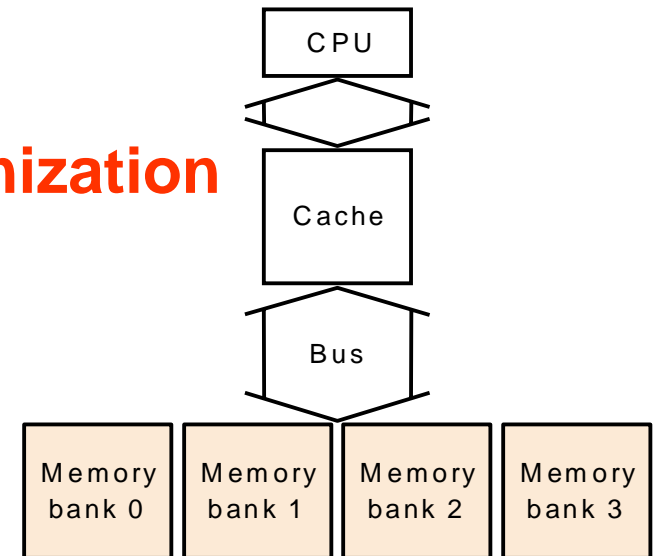**= 4 bytes/word x 4 words/65 clocks = 0.25 bytes/clock**

# Interleaved memory organization

Figure 7.13

**Suppose we have a system as follows**

- **4-bank memory <u>interleaving</u> organization**
- **1 cycle to send the address**
- **15 cycles to access DRAM**
- **1 cycle to send a word of data**

```
          CPU
           |
        <  >
           |
         Cache
           |
        <  >
           |
          Bus
           |
   ┌─────┬─────┬─────┬─────┐
 Memory  Memory  Memory  Memory
 bank 0  bank 1  bank 2  bank 3
```

c. Interleaved memory organization

**If we have a cache block of 4 words**

**Then the miss penalty is**
**= (1 address send) + 1×(15 DRAM reads)+ 4×(1 data send)**
**= 20 clocks per block read**

**Thus the number of bytes transferred per clock cycle**
**= 4 bytes/word x 4 words/17 clocks = 0.80 bytes/clock**
**we improved from 0.25 to 0.80 bytes/clock!**

# Wide bus: 4-word-wide memory organization Figure 7.13

**Suppose we have a system as follows**

- **4-word-wide memory organization**
- **1 cycle to send the address**
- **15 cycles to access DRAM**
- **1 cycle to send a word of data**



b. Wide memory organization

**If we have a cache block of 4 words**

**Then the miss penalty is**
**= (1 address send) + 1×(15 DRAM reads)+ 1×(1 data send)**
**= 17 clocks per block read**

**Thus the number of bytes transferred per clock cycle**
**= 4 bytes/word x 4 words/17 clocks = 0.94 bytes/clock**
**we improved from 0.25 to 0.80 to 0.94 bytes/clock!**

# Memory organizations

Figure 7.13

**One word wide memory organization**          Chip Area    Speed

    **Advantage**

        **Easy to implement, low hardware overhead**

    **Disadvantage**

        **Slow: 0.25 bytes/clock transfer rate**

**Interleave memory organization**

    **Advantage**

        **Better: 0.80 bytes/clock transfer rate**

        **Banks are valuable on writes: independently**

    **Disadvantage**

    **more complex bus hardware**

**Wide memory organization**

    **Advantage**

        **Fastest: 0.94 bytes/clock transfer rate**

    **Disadvantage**

        **Wider bus and increase in cache access time**

Suppose we have a processor with

| | |
|---|---|
| CPI | = 1.0 |
| Clock Rate | = 500 Mhz = 2 ns |
| L1 Cache Miss rate | = 5% |
| DRAM | = 200 ns |

**How mach faster will the machine will be if we add a**

| | |
|---|---|
| L2 Cache | = 20 ns (hit time = miss penalty) |
| L1 Cache Miss rate | = 2% |

$$\text{L to M Miss Penalty} = \frac{200\,\text{ns}}{2\,\text{ns per Clock Cycle}} = 100\,\text{Clock Cycles}$$

$$\text{L1 to L2 Miss Penalty} = \frac{20\,\text{ns}}{2\,\text{ns per Clock Cycle}} = 10\,\text{Clock Cycles}$$