

درس معماری کامپیوتر پیشرفته

دکتر غلامی

موضوع: معماری کامپیوتر پتerson

۶ نمره ارائه

۱۴ نمره امتحان پایانی

۲ نمره ارائه روش یا ایده جدید

eh.gholami@gmail.com

کارایی: عبارتست از مقدار پردازش‌هایی که در واحد زمان انجام می‌شود.

$$\text{performance} = \frac{1}{\text{exe-time}}$$

کارایی و سرعت دو پارامتر جدا از هم هستند که در انتخاب سیستم‌هایی به عنوان اینکه کدام یک بهتر است به کار می‌روند.

$$P_{m1} = \frac{1}{x} \quad , \quad P_{m2} = \frac{1}{y} \quad \rightarrow \quad \frac{P_{m1}}{P_{m2}} = \frac{\frac{1}{x}}{\frac{1}{y}} = \frac{y}{x}$$

کارایی ماشین ۱

در بررسی کارایی بین دو ماشین مشاهده می‌گردد که نسبت کارایی آنها به یکدیگر با زمان نسبت عکس دارد.

زمان پاسخ: مدت زمان شروع تا خاتمه اجرای دستور از زمان پاسخ گویند.

واکثر - وابسته به حافظه
- رفتاری - اجرا
- وابسته به پردازنده

زمان اجرا: فقط شامل مدت زمان اجرای دستور می‌باشد. به زمان اجرا: cpu-time نیز گفته می‌شود.

زمان اجرا به موارد زیادی بستگی دارد مانند نوع مداری که برای اجرای آن دستور پیش‌بینی شده است.

زمان کاربری **user cpu-time** (معمولاً به صورت نخت افزاری است)
زمان سیستم **system cpu-time** (به صورت نخت افزاری یا نرم افزاری می‌باشد)
cpu-time (در معماری‌های پترن)

* زمان سیستم برای مشخص کردن اینکه کدام یک از cpuها وظیفه اجرای این دستور را بر عهده بگیرد در نظر گرفته می‌شود.

کاهش هر یک از دوزها فوق (زمان کاربری و زمان سیستم) می‌تواند باعث افزایش کارایی سیستم گردد.

معمولاً جهت بررسی کارایی یک سیستم، یک برنامه خاص را روی آن انجام می‌دهند. لذا برای بررسی کارایی

یک سیستم مدت زمان اجرای آن برنامه خاص را در نظر می‌گیرند.

$$\text{CPU time for a program} = \left(\text{Number of clock for a program} \right) \times \text{Time of each clock}$$

با توجه به آنکه زمان هر کلاک به فرکانس بستن دارد و عبارتست از $\text{clock time} = \frac{1}{f}$ لذا رابطه فوق را می توان به صورت زیر نیز بیان نمود:

$$\text{CPU time for a program} = \left(\text{Number of clock for a program} \right) \times \left(\frac{1}{f} \right)$$

clock rate

یا

clock cycle

مدت زمان یک کلاک

۹۳, ۱۳, ۱۴

CMSC 611: Advanced Computer Architecture

Performance

Some material adapted from Mohamed Younis, UMBC CMSC 611 Spr 2003 course slides.
Some material adapted from Hennessy & Patterson / © 2003 Elsevier Science

Integrated Circuits: Fueling Innovation

2

- Technology innovations over time

Year	Technology used in computers	Relative performance/unit cost
1951	Vacuum tube	1
1965	Transistor	35
1975	Integrated circuits	900
1995	Very large-scale integrated circuit	2,400,000

Advances of the IC technology affect H/W and S/W design philosophy

۳۱

کار آیی به معنای اجزای تراکنش های موفق است. یک تراکنش ممکن است بیس از یک دستور باشد. نکته ای که باید توجه شود آن است که ممکن است مقدار ۱۰۰ دستور مثلاً اجزای سرور که مربوط به تراکنش های متفاوت می باشند. در سنجش کار آیی هم آن است که چه مقدار از تراکنش ها به صورت صحیح انجام شده است.

Moore's Law

- Transistors double every year
 - Revised to every two years
 - Sometimes revised to performance instead of transistors

Year	Processor	Transistors
1971	4004	2,300
1974	8080	4,500
1978	8086	29,000
1989	486	1,200,000
1993	Pentium	3,100,000
2000	Pentium 4	42,000,000
2004	Itanium 2	592,000,000

(Source: Intel)

(Wikipedia - Wikipedia)

میزان توان مصرفی در بعضی جاها مانند یک web server هم نیست و به صورت نامحدود همواره در اختیار است در صورتی که مثلاً در یک سیستم هکتر محدودیت توان مصرفی وجود دارد.

Defining Performance

- Performance means different things to different people, therefore its assessment is subtle

Analogy from the airlines industry:

- How to measure performance for an airplane?
 - Cruising speed (How fast it gets to the destination)
 - Flight range (How far it can reach)
 - Passenger capacity (How many passengers it can carry)

Airplane	Passenger capacity	Cruising range (miles)	Cruising speed (m.p.h)	Passenger throughput (Passenger × m.p.h)
Boeing 777	375	4550	610	228,750
Boeing 747	470	4150	610	286,700
BAC/Jacq Concorde	132	4000	1350	178,200
Douglas DC-6-50	140	8720	544	79,424

Criteria of performance evaluation differs among users and designers

برای تعریف بازدهی یا throughput باید مشخص نماییم ماشین ما قادر است چه کاری را انجام دهد. مثلاً برای یک web server، انجام دستورات یا تراکنش ها به عنوان بازدهی آن تعریف نمی شود بلکه مقدار سرویس که به کاربران خود می دهد، بازدهی سیستم را مشخص می نماید. یا مثلاً سرور مربوط به آموزش یک دانشگاه براساس مقدار دانشجویانی که به صورت online می توانند سرویس دریافت نمایند به عنوان بازدهی تعریف می شود. معنی کارهاست که بازدهی یا throughput را تعریف کنند.

performance یا کارایی را به عنوان مقدار دستوراتی یا کارهایی که در واحد زمان انجام می‌شود و یا به صورت
 ۱. تعریف می‌نمایند ولی در این رابطه قید شده است که زمان اجرای یک دستور، یا یک برنامه یا
 یک دستور العمل در نظر گرفته می‌شود. زمانی که می‌خواهیم دو کامپیوتر را با یکدیگر مقایسه نماییم باید یک کار
 یکسان را برای مقایسه آنها در نظر بگیریم.

Performance Metrics

- Response (execution) time:
 - The time between the start and the completion of a task
 - Measures user perception of the system speed
 - Common in reactive and time critical systems
 - Single-user computer
- Throughput:
 - The total number of tasks done in a given time
 - Relevant to batch processing (billing, credit card processing)
 - Also many-user services (web servers)
- Power:
 - Power consumed or battery life
 - Especially relevant for mobile

در کامپیوتری که دستورات متفاوت انجام می‌دهد برای محاسبه کارایی کدام زمان را در نظر بگیریم. مشخص است
 که کمترین زمان، بهترین کارایی را می‌دهد و دستورات طولانی که دارای زمان اجرای طولانی‌تری هستند
 بهترین کارایی را مشخص می‌نمایند. در صورتی که می‌توانیم آنها را در نظر بگیریم باز هم مقدار دقیق کارایی مشخص
 نخواهد شد.

Response-time Metric

- Maximizing performance means minimizing response (execution) time

$$\text{Performance} = \frac{1}{\text{Execution time}}$$

و پارامترهایی که در سنجش کارایی یک سیستم مهم هستند:

- ۱- زمان پاسخ (Response Time): مدت زمان بین شروع و اتمام یک پردازش است.
- ۲- بازدهی (Throughput): تعداد کارها یا پردازش‌هایی که در یک مدت معین انجام می‌شود.
- ۳- توان مصرفی (power): مصرف توان یا طول عمر باتری

Response-time Metric

$$\text{Performance} = \frac{1}{\text{Execution time}}$$

- Performance of Processor P1 is better than P2 if
 - For a given work load L
 - P1 takes less time to execute L than P2

Performance (P_1) > Performance (P_2) w.r.t L

\Rightarrow Execution time (P_1, L) < Execution time (P_2, L)

Response-time Metric

$$\text{Performance} = \frac{1}{\text{Execution time}}$$

- Relative performance captures the performance ratio
 - For the same work load

$$\text{Speedup} = \frac{\text{CPU Performance } (P_2)}{\text{CPU Performance } (P_1)} = \frac{\text{Total execution time } (P_1)}{\text{Total execution time } (P_2)}$$

برای یک کار مشخص مقدار speedup را می توان از رابطه زیر بدست آورد:

$$\text{speedup} = \frac{\text{کارایی کامپیوتر } P_2}{\text{کارایی کامپیوتر } P_1} = \frac{\text{زمان اجرای کامل } P_1}{\text{زمان اجرای کامل } P_2}$$

در سمت افزار معمولاً به جای زمان اجرا از فرکانس کار استفاده می شود.

Designer's Performance Metrics

- Users and designers measure performance using different metrics
 - Users: quotable metrics (GHz)
 - Designers: program execution

CPU execution time for a program = CPU clock cycles for a program \times Clock cycle time

$$\frac{\text{CPU clock cycles for a program}}{\text{Clock rate}}$$

- Designer focuses on reducing the clock cycle time and the number of cycles per program
- Many techniques to decrease the number of clock cycles also increase the clock cycle time or the average number of cycles per instruction (CPI)

Example

A program runs in 10 seconds on a computer "A" with a 400 MHz clock. We desire a faster computer "B" that could run the program in 6 seconds. The designer has determined that a substantial increase in the clock speed is possible, however it would cause computer "B" to require 1.2 times as many clock cycles as computer "A". What should be the clock rate of computer "B"?

- Why would this happen?
 - Maybe one operation takes one full clock cycle at 400 MHz ($1/400 \text{ MHz} = 2.5 \text{ ns}$)
 - All the rest take less than one cycle
 - Split this operation across multiple cycles
 - Can now increase clock rate, but also increase total cycles

کارایی یا performance را از دوریدگاه می توان مورد بررسی قرار داد :

۱- از دید کاربرهایی که از این سیستم استفاده می کنند .

۲- از دید طراحان سیستم

* از دید کاربر سیستمی که دارای مقدار کم کار با لاتری است دارای کارایی بیشتری است در صورتی که از دید طراح ، زمان اجرای یک برنامه مهم است که سریعتر انجام شود .

مدت زمان هر کلاک \times تعداد کلاک های CPU برای هر برنامه = زمان اجرای CPU برای یک برنامه مشخص
 طراح می تواند با کم کردن زمان هر کلاک و نیز تعداد کلاک ها برای اجرای یک برنامه باعث افزایش کارایی سیستم شود.

11

Example

A program runs in 10 seconds on a computer "A" with a 400 MHz clock.
 We desire a faster computer "B" that could run the program in 6 seconds.
 The designer has determined that a substantial increase in the clock speed is possible, however it would cause computer "B" to require 1.2 times as many clock cycles as computer "A". What should be the clock rate of computer "B"?

CPU time (A) = $\frac{\text{CPU clock cycles}}{\text{Clock rate (A)}}$ 10 seconds = $\frac{\text{CPU clock cycles of program}}{400 \times 10^6 \text{ cycles/second}}$

CPU clock cycles of program = 10 seconds \times 400 \times 10⁶ cycles/second
 = 4000 \times 10⁶ cycles

To get the clock rate of the faster computer, we use the same formula

6 seconds = $\frac{1.2 \times \text{CPU clock cycles of program}}{\text{clock rate (B)}}$ = $\frac{1.2 \times 4000 \times 10^6 \text{ cycles}}{\text{clock rate (B)}}$

clock rate (B) = $\frac{1.2 \times 4000 \times 10^6 \text{ cycles}}{6 \text{ second}}$ = 800 \times 10⁶ cycles/second

مطابق برنامه از مجموعه ای از دستورات تشکیل شده است لذا ما می توانیم از روابط زیر بدست آوریم:

12

Calculation of CPU Time

CPU time = $\frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock cycle}}$

CPU time = Instruction count \times CPI \times Clock cycle time

Or

CPU time = $\frac{\text{تعداد دستورات} \times \text{تعداد کلاک برای هر دستور}}{\text{Clock rate}}$

Component of performance	Units of measure
CPU execution time for a program	Seconds for the program
Instruction count	Instructions executed for the program
Clock cycles per instructions (CPI)	Average number of clock cycles/instruction
Clock cycle time	Seconds per clock cycle

مثال: یک برنامه در مدت 10 ثانیه بر روی ماشین A که دارای کلاک 400 مگاهرتز است اجرا می شود. ما تصمیم داریم با یک کامپیوتر B که سریعتر است این برنامه را در مدت 6 ثانیه اجرا کنیم. در این صورت کامپیوتر B باید دارای چه فرکانس کلاکی باشد که بتواند این کار را انجام دهد؟ اگر بخواهیم تعداد کلاک های ماشین B را به اندازه 2 برابر بیشتر باشد در این صورت تعداد کلاک های ماشین B را بدست آوریم.

A: $10 \text{ s} = \frac{x}{400 \text{ MHz}} \rightarrow x = 4000 \times 10^6$ تعداد کلاک B = $1.2 \times x$

B: $6 \text{ s} = \frac{1.2 \times 4000 \times 10^6}{f_B} \rightarrow f_B = \frac{4000 \times 10^6}{6} \text{ Hz} = \frac{4000}{6} \text{ MHz} \approx 666.67 \text{ MHz}$

با توجه اینکه دید طراحی به زمان اجرای برنامه است لذا برای بهبود سیستم طراح می تواند تعداد کلاک ها را برای اجرای برنامه کمتر نماید. البته زمان کلاک نیز کمتر شود زمان اجرای برنامه کمتر می شود ولی نکته این که وجود دار تکنولوژی ساخت تراشه ها و گیت های محدودیت را برای کلاک ایجاد می کند

13

CPU Time (Cont.)

- CPU execution time can be measured by running the program
- The clock cycle is usually published by the manufacture
- Measuring the CPI and instruction count is not trivial
 - Instruction counts can be measured by: software profiling, using an architecture simulator, using hardware counters on some architecture
 - The CPI depends on many factors including: processor structure, memory system, the mix of instruction types and the implementation of these instructions

در مثال صفحه قبل مشاهده می شود که فرکانس کار ماشین B دو برابر ماشین A می باشد در صورتیکه زمان اجرای آن نصف می باشد یعنی نسبت $\frac{6}{10}$ را دارد.

14

CPU Time (Cont.)

- Designers sometimes uses the following formula:

$$\text{CPU clock cycles} = \sum_{i=1}^n \text{CPI}_i \times C_i$$

Where: C_i is the count of number of instructions of class i executed
 CPI_i is the average number of cycles per instruction for that instruction class
 n is the number of different instruction classes

تغییر تعداد کلاک وابسته به سازندگان است ولی برای اجرای یک دستور می توان با استفاده از تکنیک های تعداد کلاک را کاهش داد مثلاً از پردازش های موازی استفاده نمود. همچنین مقدار حافظه نیز در تعداد کلاک های مورد نیاز برای اجرای یک دستور تأثیر دارد.

طراحان سیستم گاهی اوقات برای بدست آوردن زمان CPU از رابطه زیر استفاده می نمایند:

$$\text{CPU clock cycles} = \sum_{i=1}^n \text{CPI}_i \times C_i$$

9

سؤال: فرض کنید برای دو مجموعه دستورالعمل دارای یک نوبت‌بازاری متفاوت داشته باشیم. بر روی ماشین A مدت زمان هر کلاک برابر 1 نانوثانیه و CPI برابر 2 است. در ماشین B مدت زمان هر کلاک 2 نانوثانیه و CPI برابر 1.2 است. کدامیک از این ماشین‌ها برنامه را سریعتر اجرا می‌کند و چه مقدار؟ هر دو ماشین دستورات یکسانی را برای برنامه اجرا کنند. در نظر بگیرید که تعداد دستورات برابر I باشد.

15

Example

Suppose we have two implementation of the same instruction set architecture. Machine "A" has a clock cycle time of 1 ns and a CPI of 2.0 for some program, and machine "B" has a clock cycle time of 2 ns and a CPI of 1.2 for the same program. Which machine is faster for this program and by how much?

Both machines execute the same instructions for the program. Assume the number of instructions is "I",

CPU clock cycles (A) = ~~I~~ × 2.0 CPU clock cycles (B) = ~~I~~ × 1.2

The CPU time required for each machine is as follows:

CPU time (A) = CPU clock cycles (A) × Clock cycle time (A)
 = ~~I~~ × 2.0 × 1 ns = 2 × I ns

CPU time (B) = CPU clock cycles (B) × Clock cycle time (B)
 = ~~I~~ × 1.2 × 2 ns = 2.4 × I ns

Therefore machine A will be faster by the following ratio:

$$\frac{\text{CPU Performance (A)}}{\text{CPU Performance (B)}} = \frac{\text{CPU time (B)}}{\text{CPU time (A)}} = \frac{2.4 \times I \text{ ns}}{2 \times I \text{ ns}} = 1.2$$

بنابراین ماشین A به مقدار ۱٫۲ برابر سریعتر از ماشین B خواهد بود.

16

Comparing Code Segments

A compiler designer is trying to decide between two code sequences for a particular machine. The hardware designers have supplied the following facts:

Instruction class	CPI for this instruction class
A	1
B	2
C	3

For a particular high-level language statement, the compiler writer is considering two code sequences that require the following instruction counts:

Code sequence	Instruction count for instruction class		
	A	B	C
1	2	1	3
2	4	1	1

Which code sequence executes the most instructions? Which will be faster? What is the CPI for each sequence?

Answer:

- Sequence 1: executes 2 + 1 + 2 = 5 instructions
- Sequence 2: executes 4 + 1 + 1 = 6 instructions

مثال فوق برای برنامه یکسانی بود که مقدار CPI به صورت میانگین داده شده بود. حال اگر بر روی ماشین های A و B دستوراتی اجرا شود که برای اجرای هر دستور مقدار کلاک متفاوتی لازم باشد در آن صورت از رابطه زیر استفاده می شود:

$$\text{cpu clock cycles} = \sum_{i=1}^n \text{CPI}_i \times C_i$$

17

Comparing Code Segments

Using the formula: $\text{CPU clock cycles} = \sum_{i=1}^n \text{CPI}_i \times C_i$

Sequence 1: $\text{CPU clock cycles} = (2 \times 1) + (1 \times 2) + (2 \times 3) = 10 \text{ cycles}$
 Sequence 2: $\text{CPU clock cycles} = (4 \times 1) + (1 \times 2) + (1 \times 3) = 9 \text{ cycles}$

☐ Therefore Sequence 2 is faster although it executes more instructions

Using the formula: $\text{CPI} = \frac{\text{CPU clock cycles}}{\text{Instruction count}}$

Sequence 1: $\text{CPI} = 10/5 = 2$
 Sequence 2: $\text{CPI} = 9/6 = 1.5$

☐ Since Sequence 2 takes fewer overall clock cycles but has more instructions it must have a lower CPI

پس بنابراین کارایی یک ماشین را به روش های مختلفی می توان بررسی نمود.

18

The Role of Performance

- Hardware performance is a key to the effectiveness of the entire system
- Performance has to be measured and compared to evaluate designs
- To optimize the performance, major affecting factors have to be known
- For different types of applications
 - different performance metrics may be appropriate
 - different aspects of a computer system may be most significant
- Instructions use and implementation, memory hierarchy and I/O handling are among the factors that affect the performance

اولین روش به صورت کارایی سخت افزاری است. Clock Rate به صورت سخت افزاری است. CPI هم سخت افزاری و هم نرم افزاری است زیرا CPI نسبت مقدار کلاک های CPU به مقدار دستورات است که مقدار کلاک ها سخت افزاری و مقدار دستورات نرم افزاری است و می توان آنها را با استفاده از روش های خاص تغییر داد. سخت افزار در کارایی بسیار مؤثر است و برای بهینه کردن کارایی می توان تغییراتی را در سخت افزار ایجاد کرد.

Calculation of CPU Time

$$\text{CPU time} = \frac{\text{Instruction count} \times \text{CPI}}{\text{Clock rate}}$$

	Instr. Count	CPI	Clock Rate
Program	X		
Compiler	X	X	
Instruction Set	X	X	
Organization		X	X
Technology			X

$$\text{CPU clock cycles} = \sum_{i=1}^n \text{CPI}_i \times C_i$$

Where: C_i is the count of number of instructions of class i executed
 CPI_i is the average number of cycles per instruction for that instruction class
 n is the number of different instruction classes

Important Equations (so far)

$$\text{Performance} = \frac{1}{\text{Execution time}}$$

$$\text{Speedup} = \frac{\text{Performance (B)}}{\text{Performance (A)}} = \frac{\text{Time (A)}}{\text{Time (B)}}$$

$$\text{CPU time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Cycle}}$$

$$\text{CPU clock cycles} = \sum_{i=1}^n \text{CPI}_i \times \text{Instructions}_i$$

Amdahl's Law

The performance enhancement possible with a given improvement is limited by the amount that the improved feature is used

Execution time after improvement =

$$\frac{\text{Execution time affected by the improvement}}{\text{Amount of improvement}}$$

+ Execution time unaffected

- A common theme in Hardware design is to *make the common case fast*
- Increasing the clock rate would not affect memory access time
- Using a floating point processing unit does not speed integer ALU operations

Example: Floating point instructions improved to run 2X; but only 34% of actual instructions are floating point

$$\text{Exec-Time}_{\text{new}} = \text{Exec-Time}_{\text{old}} \times (0.66 + .34/2) = 0.83 \times \text{Exec-Time}_{\text{old}}$$

$$\text{Speedup}_{\text{overall}} = \text{Exec-Time}_{\text{old}} / \text{Exec-Time}_{\text{new}} = 1/0.83 = 1.205$$

Ahmdal's Law Visually



Ahmdal's Law for Speedup

$$Time_{old} = Time_{old} * (Fraction_{unchanged} + Fraction_{enhanced})$$

$$Time_{new} = Time_{old} * \left(Fraction_{unchanged} + \frac{Fraction_{enhanced}}{Speedup_{enhanced}} \right)$$

$$\begin{aligned} Speedup_{overall} &= \frac{Time_{old}}{Time_{new}} = \frac{Time_{old}}{Time_{old} * \left(Fraction_{unchanged} + \frac{Fraction_{enhanced}}{Speedup_{enhanced}} \right)} \\ &= \frac{1}{Fraction_{unchanged} + \frac{Fraction_{enhanced}}{Speedup_{enhanced}}} \end{aligned}$$

$$Speedup_{overall} = \frac{1}{(1 - Fraction_{enhanced}) + \frac{Fraction_{enhanced}}{Speedup_{enhanced}}}$$

۹۳،۱۲،۲۱

pipeline حفظ داده تکلیف است که امکا اجرای سریعتر مجموعه دستورالعمل ها را فراهم کند. اجرای دستورات می تواند به صورت ترتیبی یا pipeline صورت گیرد و از این در حالت خارج نیست.

CMSC 611: Advanced Computer Architecture

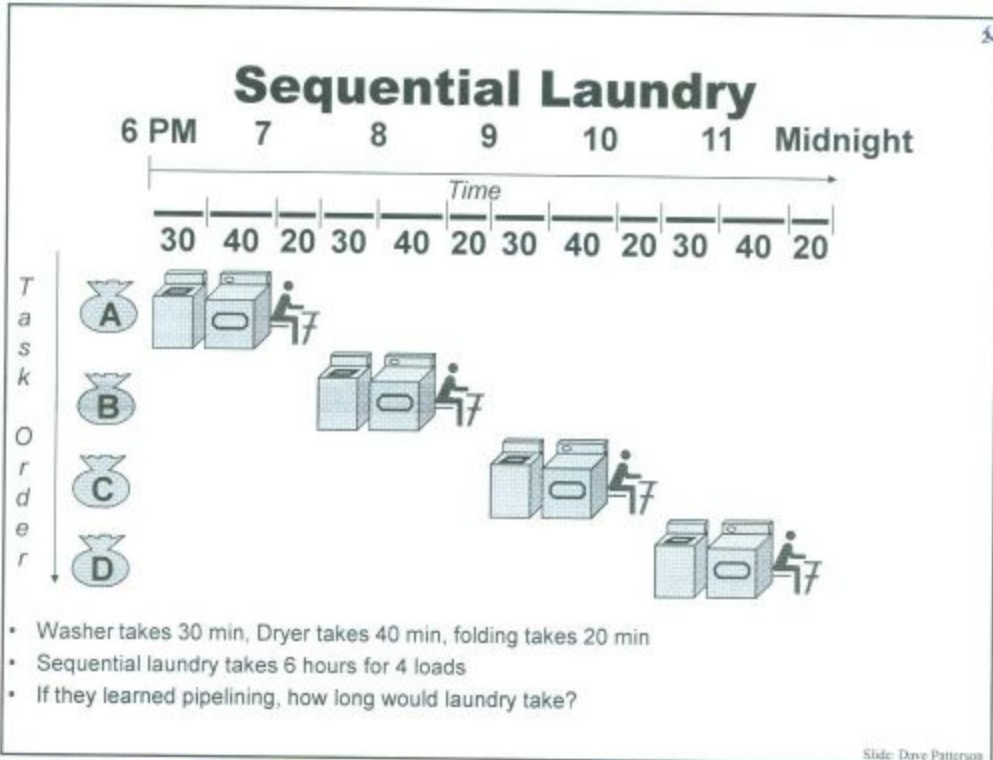
Pipelining

یکی از خصوصیات pipeline است که مدت اجرای یک دستور به تنهایی کاهش می یابد بلکه مجموع انجام کلیه کارها را کاهش می دهد یعنی سرعت اجرای مجموعه را کاهش می دهد. این روش برای سیستم های کامپیوتری نیز به همین صورت است.

Some material adapted from Mohamed Younis, UMBC CMSC 611 Sp2007 course slides
Some material adapted from Hennessy & Patterson / © 2003 Prentice Hall

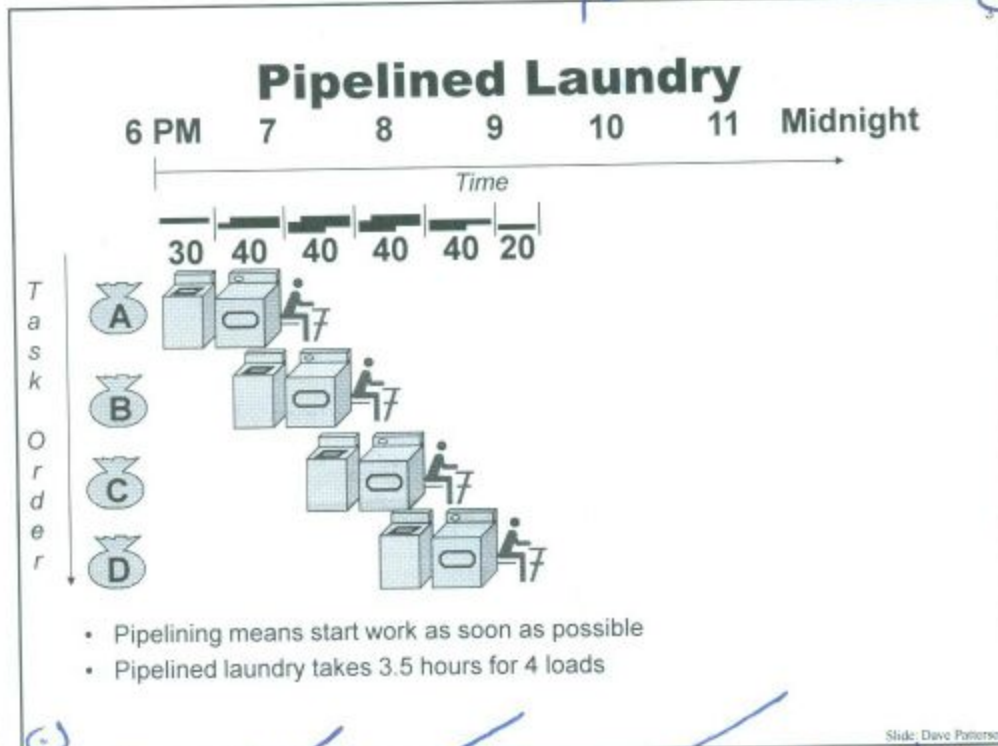
سوال زیر کار در این صورت مجموعه ای از ماشین، خشک کردن و اطو کردن ترتیب کرده است و می خواهد چهار بسته لباس را بسته، خشک کرده و اطو نماید. برای ماشین ۳۰ دقیقه، خشک کردن ۴۰ دقیقه و اطو کردن ۲۰ دقیقه

زمان نیاز است

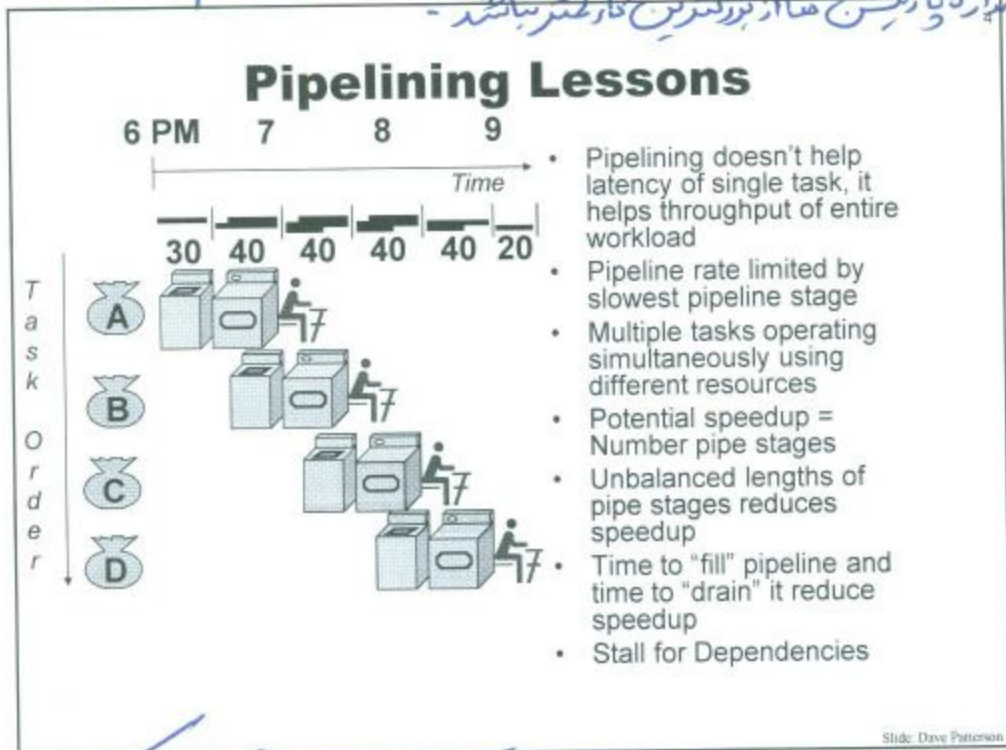


در سوال فوق اگر بخواهیم کارها را به صورت ترتیبی انجام دهیم محلاً برای هر بسته ۹۰ دقیقه زمان لازم است لذا برای چهار بسته به ۳۶۰ دقیقه زمان نیاز داریم در صورتیکه استفاده از pipeline این امکان را می دهد که همزمان از چند دستگاه استفاده نماییم که با استفاده از این روش مدت زمانی که طول می کشد تا یک بسته بسته شود، خشک شده و اطو شوند به سه ساعت و نیم یا ۲۱۰ دقیقه کاهش می یابد. البته باید توجه نمود برای آنکه بتوان از سیستم pipeline استفاده نمود باید امکان تغذیه منابع وجود داشته باشد. ۱۵

اولین قدم در pipeline کردن یک کار، تفکیک درست منابع آن کار است. pipeline به زمان انجام یک کار به تنهایی کمک نمی‌کند بلکه زمان اجرای کل کارها را کاهش می‌دهد و باعث ارتقای کارایی کل کار می‌شود. امکان آنکه بتوان در یک سیستم از pipeline استفاده نمود آن است که امکان استفاده از چند منبع بطور هم‌زمان را داشته باشیم.



در حالت استفاده از pipeline توصیه این نیست که منابع ممکن است دارای زمان اجرای یک نباشند لذا ما باید کمترین زمان را برای اجرای هر یک در نظر بگیریم مانند زمانی که می‌خواهیم کار در پارکینگ های مختلف قرار دهیم مجبوریم که این پارکینگ ها را بزرگترین کار کمتر نباشد.



هدف از به کار بردن pipeline امکان افزایش کارایی در سیستم می‌باشد. زمانی که pipeline استفاده می‌کنیم حداکثر میزان کارایی به میزان تعداد مراحل pipeline نسبت به حالت ترتیبی افزایش می‌یابد یعنی مثلاً اگر تعداد مراحل pipeline سه مرحله باشد در آن صورت میزان کارایی حداکثر در بهترین حالت سه برابر حالت ترتیبی خواهد شد. یا به طور کلی اگر تعداد مراحل n مرحله باشد، کارایی حداکثر n برابر حالت ترتیبی می‌شود.

یکی از عواملی که باعث می شود میزان سریع یا speed up به n نزده ، بالانس نبودن زمان اجرای مراحل pipeline خواهد بود لذا در pipeline کردن کارها ، مراحل را به گونه ای تعریف می کنیم که علاوه بر شرایط بالا از نظر مدت زمانی هم تقریباً یکسان باشند و تفاوتی با هم نداشته باشند.

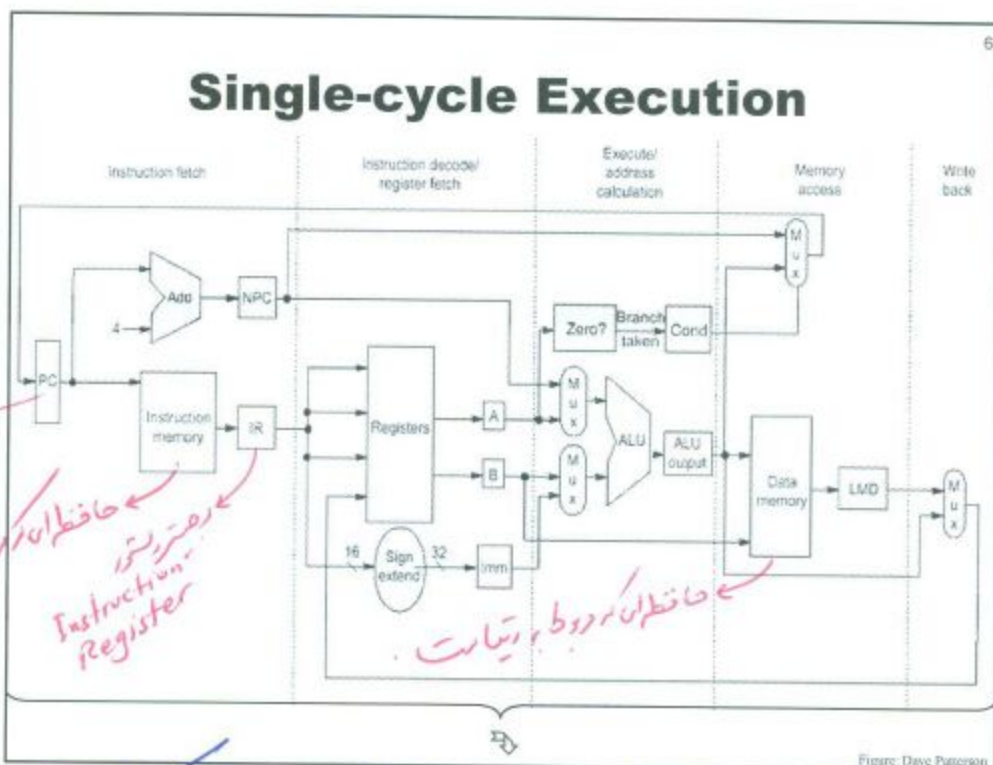
MIPS Instruction Set

- RISC characterized by the following features that simplify implementation:
 - All ALU operations apply only on registers
 - Memory is affected only by load and store
 - Instructions follow very few formats and typically are of the same size

	31	26	21	16	11	6	0
	op		rs	rt	rd	shamt	funct
	6 bits		5 bits	5 bits	5 bits	5 bits	6 bits

	31	26	21	16	0
	op		rs	rt	immediate
	6 bits		5 bits	5 bits	16 bits

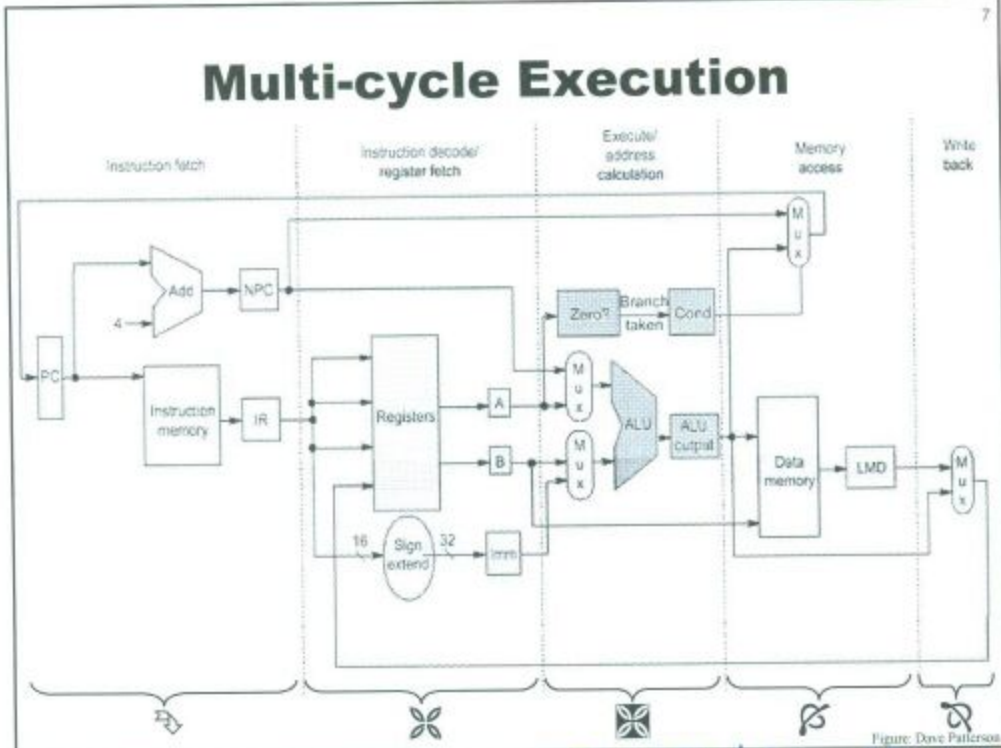
	31	26	0
	op		target address
	6 bits		26 bits



یکی دیگر از عوامل ، شروع کار کردن با pipeline و پایان آن است . برای آنکه کار شروع به انجام به صورت pipeline کند مدت زمانی طول می کشد که به آن زمان پر شدن یا fill گفته می شود . همچنین زمانی طول می کشد تا کل کارها انجام شود که اصطلاحاً به آن زمان خالی شدن pipeline یا drain می گویند . هر چه تعداد کارهایی که وارد pipeline می شود بیشتر باشد بهره برداری از pipeline بیشتر خواهد بود . در واقع در هر تعداد کارها بیشتر زمان پر شدن و خالی شدن بین کارهای بیشتری سرشکن می شود پس بنابراین زمان پر شدن و خالی شدن می تواند مقدار speed up را تغییر دهد .

۱۷

گاهی اوقات لازم است بعضی از کارها در بین کارهای دیگر به صورت تعلیق درآیند هر چه تعداد این کارها
 منطبق بیشتر شود باعث کاهش تاخیر یا speedup میگردد.
 توانایی speedup یا تسریع سیستم در زمانی که از pipeline استفاده میکنیم حداکثر به میزان مراحل انجام
 pipeline است.



در کامپیوترها RISC تعداد دستور و قالب دستور محدود هستند.

Multi-Cycle Implementation of MIPS

چون در دستورهای MIPS 4 بایت است پس دستور بعدی در 4 بایت بعد شروع می شود

رابطه با فصل که مقدار را در آن قرار می دهیم

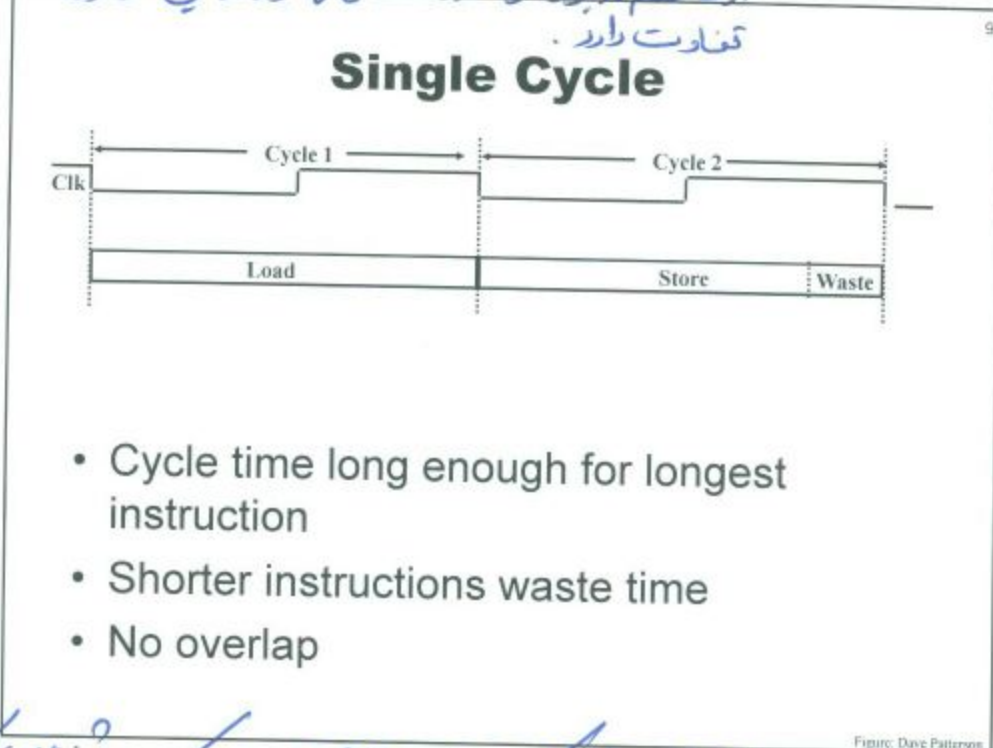
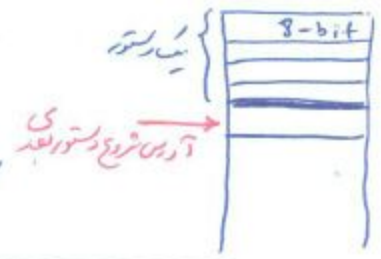
- Instruction fetch cycle (IF)**
 $IR \leftarrow Mem[PC]; \quad NPC \leftarrow PC + 4$
- Instruction decode/register fetch cycle (ID)**
 $A \leftarrow Regs[IR_{5..10}]; \quad B \leftarrow Regs[IR_{11..15}]; \quad Imm \leftarrow ((IR_{16})^{16} \# \# IR_{16..31})$
- Execution/effective address cycle (EX)**
 Memory ref: $ALUOutput \leftarrow A + Imm;$
 Reg-Reg ALU: $ALUOutput \leftarrow A \text{ func } B;$
 Reg-Imm ALU: $ALUOutput \leftarrow A \text{ op } Imm;$
 Branch: $ALUOutput \leftarrow NPC + Imm; \quad Cond \leftarrow (A \text{ op } 0)$
- Memory access/branch completion cycle (MEM)**
 Memory ref: $LMD \leftarrow Mem[ALUOutput] \text{ or } Mem(ALUOutput) \leftarrow B;$
 Branch: $\text{if (cond) } PC \leftarrow ALUOutput;$
- Write-back cycle (WB)**
 Reg-Reg ALU: $Regs[IR_{16..20}] \leftarrow ALUOutput;$
 Reg-Imm ALU: $Regs[IR_{11..15}] \leftarrow ALUOutput;$
 Load: $Regs[IR_{11..15}] \leftarrow LMD;$

فقط برای ریزترها استفاده می شوند

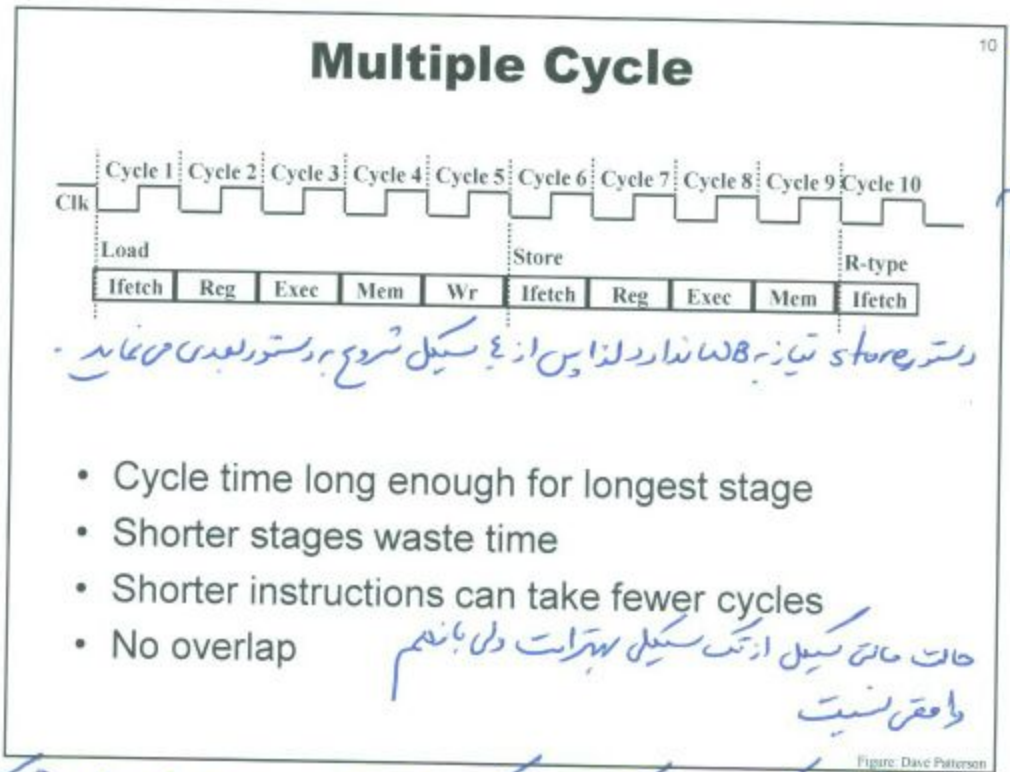
کامپیوترهای تریول
 Fetch
 واکش
 اجزای
 دیت اورد
 دستور
 آدرس حافظه
 حافظه
 * در کامپیوترهای تریول
 مراحل اجرای تک دستور
 چهار مرحله است.

در کتاب مانو مراحل اجرای تک دستور سه مرحله بود
 در کامپیوترهای تریول دو حافظه مجزا وجود دارد که یکی جهت دستورات و دیگری جهت داده به کار می رود.
 یکی از خصوصیات کامپیوترهای تریول آن است که امکان ایجاد pipeline را به ما می دهد.
 آدرس دستور همیشه در بایت PC قرار دارد.

حون هر دستور به ۴ بایت حافظه در کامپیوتر پیام پترون نیاز دارد پس دستور بعدی در آدرس چهارخانه بعد شروع خواهد شد.
 در کامپیوتر پیام مانده دستور یک بطن ۳۲ بیتی دار حافظه اشغال کند در صورتیکه در کامپیوتر پیام پترون هر دستور در چهارخانه یک بایتی در زیر هم قرار میگیرد لذا فضای اضافه حافظه که برای هر دستور اشغال می شود در این دو مدل انتظار نیست افزای با هم تفاوت دارد.

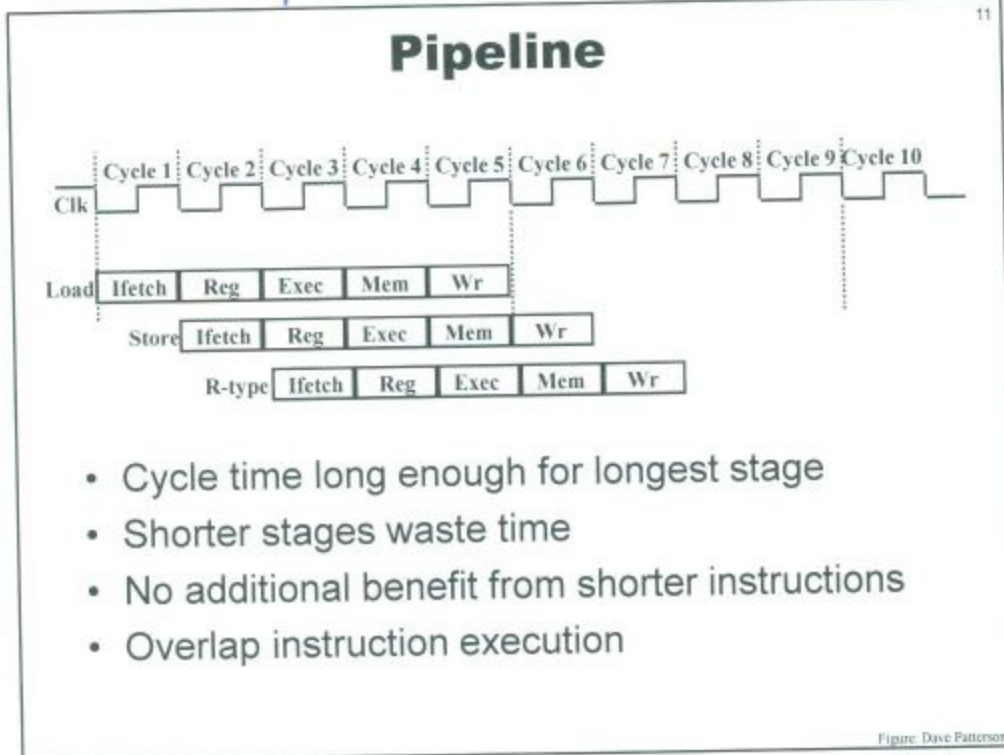


کارها یک بخت افتاده می تواند به چند بخش تقسیم شود که برای اجرای هر بخش مجبوریم یک سری سینکلیشن کنیم تا ایجا رعایت واکش -
 دی نگردن -
 اجرا -
 دسترسی حافظه -
 write back



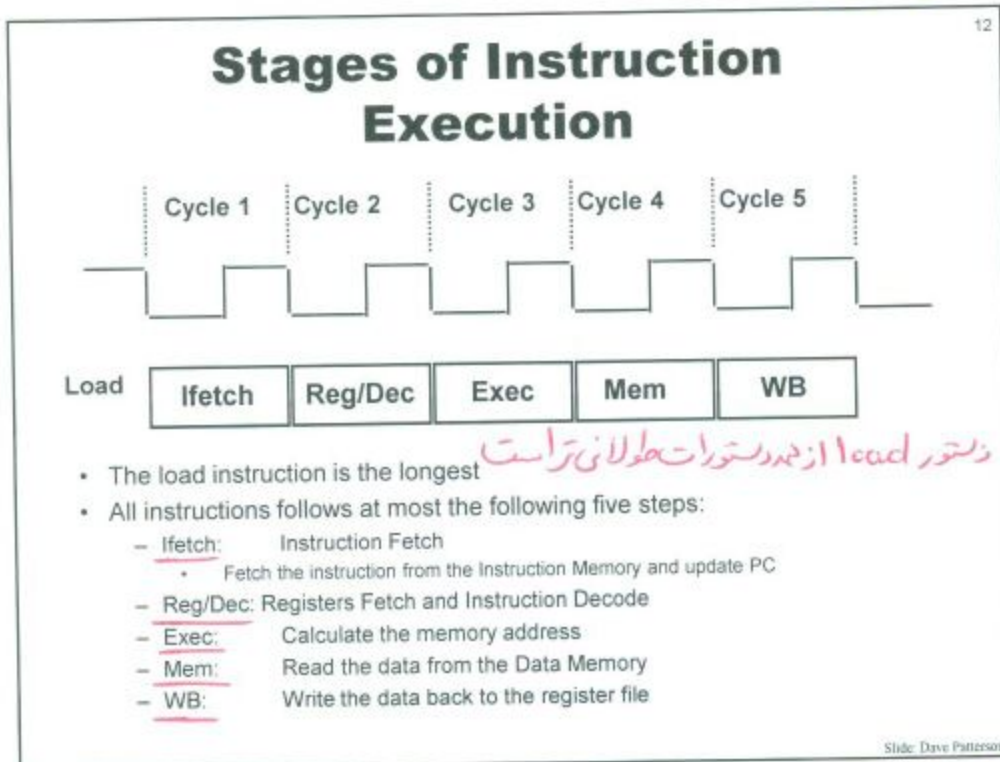
برض از کامپیوترها به صورت تک سیکلی عمل می کنند بطوریکه برای اجرای هر دستور مجبور هستیم تک سیکل را فعال کنیم که این سیکل باید به اندازه مدت زمان لازم برای اجرای آن دستور فعال باشد. در صورتیکه مدت زمان این سیکل به پایان رسید، کامپیوتر ما در صورتی موفق است که بتواند در مدت زمان آن سیکل دستورات را به پایان برساند.
 در کامپیوترها دستورات متفاوت با مدت زمان اجرای متفاوت وجود دارد لذا در کامپیوترهای تک سیکلی مدت زمان کلاک باید طوری باشد که کمترین دستورات بتوانند در یک سیکل اجرا شوند لذا برای بعضی از دستورات که سریعتر هستند مقدار آنلاف زمان خواهیم داشت. بنا بر این پیاده سازی تک سیکلی باعث هدر رفتن زمان برای دستوراتی که به زمان کمتری نیاز دارند می شود.

اگر کامپیوتری به صورت یک سیکلی باشد نمی توان آن را به صورت pipeline استفاده نمود بلکه باید ابتدا آنرا مرحله‌ای نموده و سپس به صورت pipeline استفاده نماییم. مدت زمان سیکل‌ها را باید به گونه‌ای انتخاب نماییم که کمترین دستورات قابل اجرا باشند لذا در اجرای دستورات به صورت pipeline، حداکثر دستور را مثلاً load که نیاز به پنج سیکل دارد را انتخاب می‌نماییم.



- Cycle time long enough for longest stage
- Shorter stages waste time
- No additional benefit from shorter instructions
- Overlap instruction execution

مدت زمان اجرا }
مدت زمان بین اجرا }



مدت زمان اجرا :

مدت زمان شروع اولین دستور تا شروع آخرین دستور را گوئیم که هر چه این مقدار کمتر باشد عملاً مثل این است که زمان پاسخ کوتاه‌تر است زیرا ما داریم زمان پاسخ برابر مجموع زمان انتظار و زمان اجرا می‌باشد. در pipeline زمان اجرا تغییر نمی‌کند در صورتی که زمان انتظار کاهش می‌یابد.

کارایی یک pipeline به چند مورد بستگی دارد :
 - هر چه تعداد دستورهای بیشتر شود کارایی بیشتر می شود ولی زمان اجرای یک دستور تغییر نمی کند .
 - بران بر فرض از دستورات مستقل ممکن است زمان اجرا نسبت به حالت ترتیبی بیشتر شود .

Instruction Pipelining

13

- Start handling next instruction while the current instruction is in progress
- Feasible when different devices at different stages

Time between instructions_{pipelined} = $\frac{\text{Time between instructions}_{\text{nonpipelined}}}{\text{Number of pipe stages}}$

Pipelining improves performance by increasing instruction throughput

Example of Instruction Pipelining

14

Time between first & fourth instructions is $3 \times 8 = 24 \text{ ns}$

Time between first & fourth instructions is $3 \times 2 = 6 \text{ ns}$

Ideal and upper bound for speedup is number of stages in the pipeline

در حالت استفاده از pipeline مشاهده می شود که زمان بین شروع اولین دستور و شروع دستور چهارم فقط ۶ نانوثانیه است .

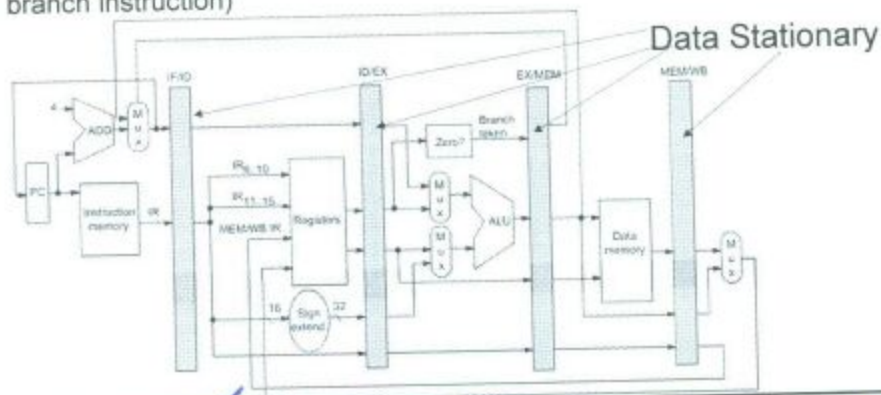
Pipeline Performance

- Pipeline increases the instruction throughput
 - not execution time of an individual instruction
- An individual instruction can be **slower**:
 - Additional pipeline control
 - Imbalance among pipeline stages
- Suppose we execute 100 instructions:
 - Single Cycle Machine
 - $45 \text{ ns/cycle} \times 1 \text{ CPI} \times 100 \text{ inst} = 4500 \text{ ns}$
 - Multi-cycle Machine
 - $10 \text{ ns/cycle} \times 4.2 \text{ CPI (due to inst mix)} \times 100 \text{ inst} = 4200 \text{ ns}$
 - Ideal 5 stages pipelined machine
 - $10 \text{ ns/cycle} \times (1 \text{ CPI} \times 100 \text{ inst} + 4 \text{ cycle drain}) = 1040 \text{ ns}$
- Lose performance due to fill and drain

در مثال فوق مقایسه ای بین زمان اجرا برای کامپیوترهای تک سیکلی، چند سیکلی و pipeline صورت گرفته است.

Pipeline Datapath

- Every stage must be completed in one clock cycle to avoid stalls
- Values must be latched to ensure correct execution of instructions
- The PC multiplexer has moved to the IF stage to prevent two instructions from updating the PC simultaneously (in case of branch instruction)



در پیاده سازی pipeline باید هر مرحله بدون داشتن وقفه در یک سیکل اجرا شود یعنی مراحل اجرای pipeline باید به صورت پشت سرهم انجام شود. همچنین دستورالعمل باید پشت سرهم وارد شوند.
 1 در اجرای pipeline باید از یک سری حافظه های موقت بنام latch استفاده کنیم تا از اجرای صحیح دستورالعمل مطمئن شویم. مقدار latch ها باید برابر n-1 باشد که n تعداد مراحل اجرای pipeline می باشد یعنی اگر 5 مرحله برای اجرای pipeline داریم به مقدار 4 عدد latch نیاز داریم.

علیرغم استفاده از حافظه‌های latch با زحم ممکن است مخاطره‌هایی وجود داشته باشد که نگذارند به کارایی لازم دست یابی داشته باشیم که باید شناسایی شده و برطرف گردند. این مخاطرات را می‌توان به صورت زیر دسته بندی نمود:

انواع Hazard: ساختاری - ساختاری
 Data hazard - دیتا
 Control hazard - کنترل

17

Pipeline Stage Interface

Stage	Any Instruction		
IF	$IF/ID\ IR \leftarrow M[EX/PC]$ $IF/ID\ NPC\ PC \leftarrow (if\ (EX/MEM\ opcode == branch) \&\ EX/MEM\ cond) \{ EX/MEM\ ALUOutput\} \text{ else } \{ PC + 4 \}$		
ID	$ID/EX\ A \leftarrow Regs[IF/ID\ IR_{2..4}]$ $ID/EX\ B \leftarrow Regs[IF/ID\ IR_{5..7}]$ $ID/EX\ NPC \leftarrow IF/ID\ NPC$ $ID/EX\ IR \leftarrow IF/ID\ IR$ $ID/EX\ Imm \leftarrow (IF/ID\ IR_{12}) \ll 20 \text{ or } IF/ID\ IR_{13..15}$		
	ALU	Load or Store	Branch
EX	$EX/MEM\ IR \leftarrow ID/EX\ IR$ $EX/MEM\ ALUOutput \leftarrow ID/EX\ A\ Func\ ID/EX\ B$ Or $EX/MEM\ ALUOutput \leftarrow ID/EX\ A\ op\ ID/EX\ Imm$ $EX/MEM\ cond \leftarrow 0$	$EX/MEM\ IR \leftarrow ID/EX\ IR$ $EX/MEM\ ALUOutput \leftarrow ID/EX\ A + ID/EX\ Imm$ $EX/MEM\ cond \leftarrow 0$ $EX/MEM\ S \leftarrow ID/EX\ B$	$EX/MEM\ ALUOutput \leftarrow ID/EX\ NPC + ID/EX\ Imm$ $EX/MEM\ cond \leftarrow (ID/EX\ A\ op\ 0)$
MEM	$MEM/WB\ IR \leftarrow EX/MEM\ IR$ $MEM/WB\ ALUOutput \leftarrow EX/MEM\ ALUOutput$	$MEM/WB\ IR \leftarrow EX/MEM\ IR$ $MEM/WB\ LMD \leftarrow Mem[EX/MEM\ ALUOutput]$ Or $Mem[EX/MEM\ ALUOutput] \leftarrow EX/MEM\ S$	
WB	$Regs[MEM/WB\ IR_{2..4}] \leftarrow EX/MEM\ ALUOutput$ Or $Regs[MEM/WB\ IR_{5..7}] \leftarrow MEM/WB\ ALUOutput$	For load only: $Regs[MEM/WB\ IR_{2..4}] \leftarrow MEM/WB\ LMD$	

Hazardها یا مخاطره‌ها همیشه یا منتظر ماندن قبل رفع شدن هستند ولی این کار باعث کاهش کارایی می‌گردد لذا برای افزایش سرعت باید با استفاده از روش‌های دیگری این مخاطره‌ها را برطرف نمائیم

18

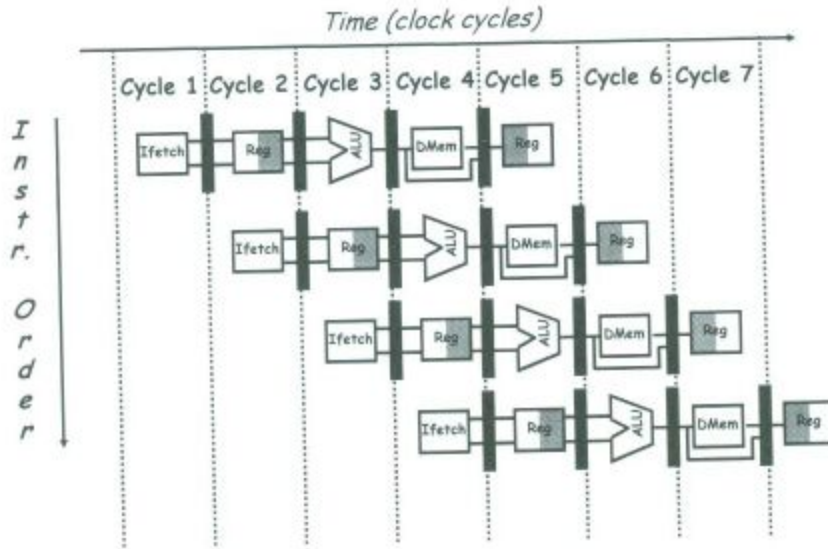
Pipeline Hazards

- Cases that affect instruction execution semantics and thus need to be detected and corrected
- Hazards types
 - Structural hazard: attempt to use a resource two different ways at same time
 - Single memory for instruction and data
 - Data hazard: attempt to use item before it is ready
 - Instruction depends on result of prior instruction still in the pipeline
 - Control hazard: attempt to make a decision before condition is evaluated
 - branch instructions
- Hazards can always be resolved by waiting

همانقدر گفته شد pipeline اجرا برداشتن دسته‌ای از دستورات را با هم سریعتر می‌کند. اولین کاره که بتوانیم برناسه راه صورت pipeline اجرا نمود آن است که بتوان آن را مرحله بندی کرد. Hazardهای ساختاری مربوط به سخت افزار می‌باشد. توسط سخت افزار باید بتوان دسته‌های مختلف دستورات را مرحله بندی نمود یعنی در واقع سخت افزار باید امکان تفکیک مراحل را فراهم نماید. زمانی که دستوراتی را به صورت pipeline اجرا می‌نمائیم در خیلی از مواقع ممکن است دستورات به هم وابسته باشند و نتوان آنها را تفکیک نمود.

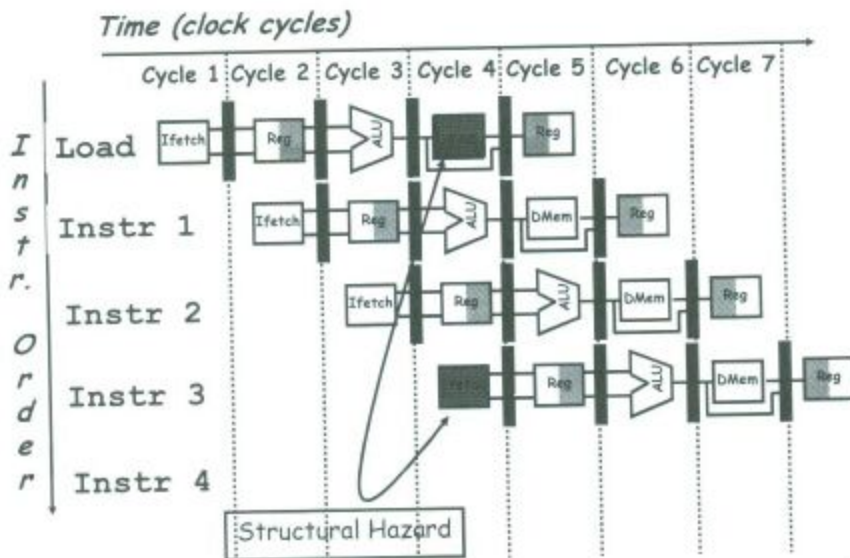
۲۳

Visualizing Pipelining



Slide: David Culler

Example: One Memory Port/Structural Hazard



Slide: David Culler

در مثال فوق هم می خواهیم در حافظه بنویسیم و هم در دستورواکش می خواهیم از حافظه بخوانیم که این امر باعث ایجاد خطا می شود. به این نوع فرط بهره استفاده یا سمت افزاری گفته می شود. برای رفع این اشکال یا باید با ایجاد یک صیاب منتظر بمانیم یعنی کیفیت در اجرای دستورات بعدی ایجاد نمی کنیم یا از سمت افزارهای رید استفاده کنیم.

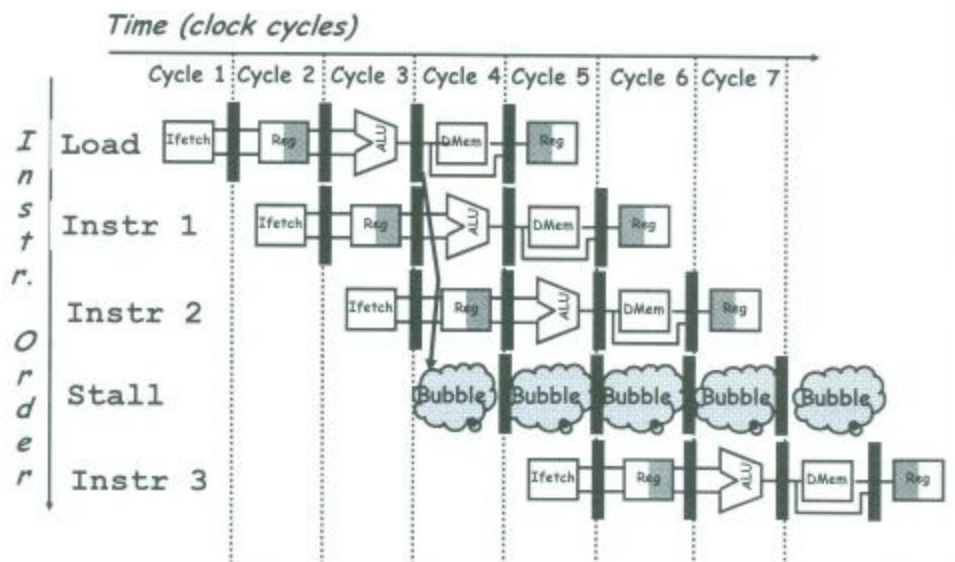
if $A < B$
 $r_1 = r_2 + r_3$
 else
 $r_1 = r_4 + r_5$

مخاطره‌های ساختاری را فقط به صورت سخت‌افزاری می‌توان برطرف نمود.
 دو راهکار مختلف برای مقابله با مخاطره‌های ساختاری وجود دارد:
 ۱- منتظر بمانیم یعنی عملاً آن مرحله را نمی‌توانیم تکمیل کرد یعنی سخت‌افزار را تا تکمیل مراحل دانش‌نامه باسدین بناچار باید تاخیر را تحمل نماییم.
 ۲- سخت‌افزار را تغییر دهیم.

Resolving Structural Hazards

1. Wait
 - Must detect the hazard
 - Easier with uniform ISA
 - Must have mechanism to stall
 - Easier with uniform pipeline organization
2. Throw more hardware at the problem
 - Use instruction & data cache rather than direct access to memory

Detecting and Resolving Structural Hazard



$$r_1 = r_2 + r_3$$

$$r_2 = r_1 + r_4$$

وابستگی داده



در این مثال عمل می‌کند کردن دستور دوم و اجرای دستور اول فترک
 من خواهد انجام شود در صورتیکه زمانی می‌توان دستور بعدی را اجرا
 نمود که دستور قبلی کاملاً اجرا شده باشد. این نوع مخاطره را
 مخاطره دیتا می‌نامند.

23

Stalls & Pipeline Performance

✓ Pipelining Speedup = $\frac{\text{Average instruction time unpipelined}}{\text{Average instruction time pipelined}}$

$$= \frac{\text{CPI unpipelined} \times \text{Clock cycle unpipelined}}{\text{CPI pipelined} \times \text{Clock cycle pipelined}}$$

Ideal CPI pipelined = 1

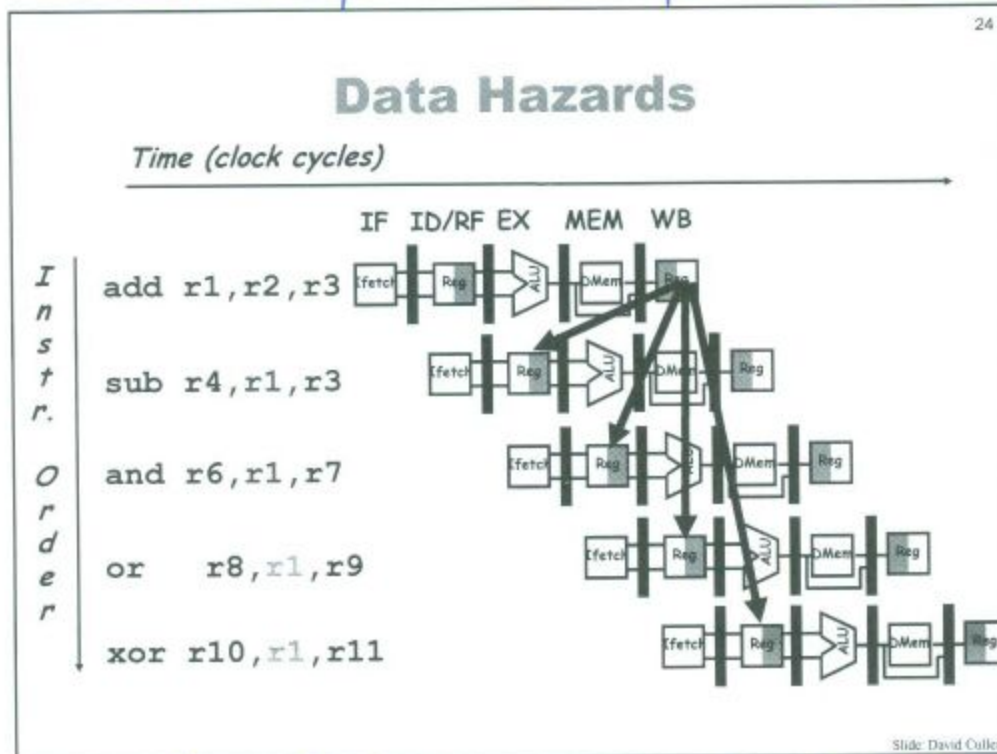
CPI pipelined = Ideal CPI + Pipeline stall cycles per instruction
 = 1 + Pipeline stall cycles per instruction

$$\text{Speedup} = \frac{\text{CPI unpipelined}}{1 + \text{Pipeline stall cycles per instruction}} \times \frac{\text{Clock cycle unpipelined}}{\text{Clock cycle pipelined}}$$

Assuming all pipeline stages are balanced

$$\text{Speedup} = \frac{\text{Pipeline depth}}{1 + \text{Pipeline stall cycles per instruction}}$$

مخاطره‌های دیتا را هم می‌توان به صورت سخت‌افزاری و هم به صورت نرم‌افزاری برطرف نمود.
 مخاطرات دیتا زمانی بوجود می‌آید که بخواهیم از داده‌ای استفاده کنیم قبل از آنکه آن داده بدست آید.



در شکل فوق مشاهده می‌رود ما می‌خواهیم از داده‌ای استفاده کنیم که دو سیکل بعد در ثبت نوشته می‌شود
 پس عملاً مخاطره دیتا به وجود می‌آید. در دستور سوم هم می‌خواهیم از داده‌ای استفاده کنیم که در سیکل
 بعدی تهیه می‌شود. اما مشاهده می‌شود که از دستور چهارم به بعد مشکل بوجود نمی‌آید. همانگونه که مشاهده می‌شود در
 دستور اول در نیمه اول WB محل ثبت داده انجام می‌شود لذا در دستور چهارم در نیمه دوم می‌توان از ثبت برای داده استفاده
 نمود و در این حالت اشکالی بوجود نمی‌آید. قوی‌ترین صورت آید WB کردن و دی‌کد کردن نیز در یک سیکل انجام می‌شود و آن
 صورت در اجرای دستور چهارم هم مشاهده می‌توانیم می‌آید.

۲۴

- سه نوع مخاطره رتیباً وجود دارد:
- ۱- خواندن بعد از نوشتن
 - ۲- نوشتن بعد از خواندن
 - ۳- نوشتن بعد از نوشتن

25

Three Generic Data Hazards

- Read After Write (RAW)
Instr_j tries to read operand before Instr_i writes it

```

I: add r1, r2, r3
  ↘
J: sub r4, r1, r3

```

- Caused by a "Data Dependence" (in compiler nomenclature). This hazard results from an actual need for communication.

Slide: David Culler

با استفاده از روش های نرم افزاری و سخت افزاری یا با دیدن وقوع مخاطرات رتیباً جلوگیری می‌کنیم یا آنرا به وقوع پیوسته آنهارا تسخیر می‌کنیم و به طرفت می‌کنیم.

26

Three Generic Data Hazards

- Write After Read (WAR)
Instr_j writes operand before Instr_i reads it

```

I: sub r4, r1, r3
  ↘
J: add r1, r2, r3
  ↘
K: mul r6, r1, r7

```


- Called an "anti-dependence" in compilers.
 - This results from reuse of the name "r1".
- Can't happen in MIPS 5 stage pipeline because:
 - All instructions take 5 stages, and
 - Reads are always in stage 2, and
 - Writes are always in stage 5

Slide: David Culler

Three Generic Data Hazards

- Write After Write (WAW)
Instr_j writes operand before Instr_i writes it.

```
    I: mul r1, r4, r3  
    J: add r1, r2, r3  
    K: sub r6, r1, r7
```



- Called an "output dependence" in compilers
 - This also results from the reuse of name "r1".
- Can't happen in MIPS 5 stage pipeline:
 - All instructions take 5 stages, and
 - Writes are always in stage 5
- Do see WAR and WAW in more complicated pipes

۹۴, ۱, ۲۰

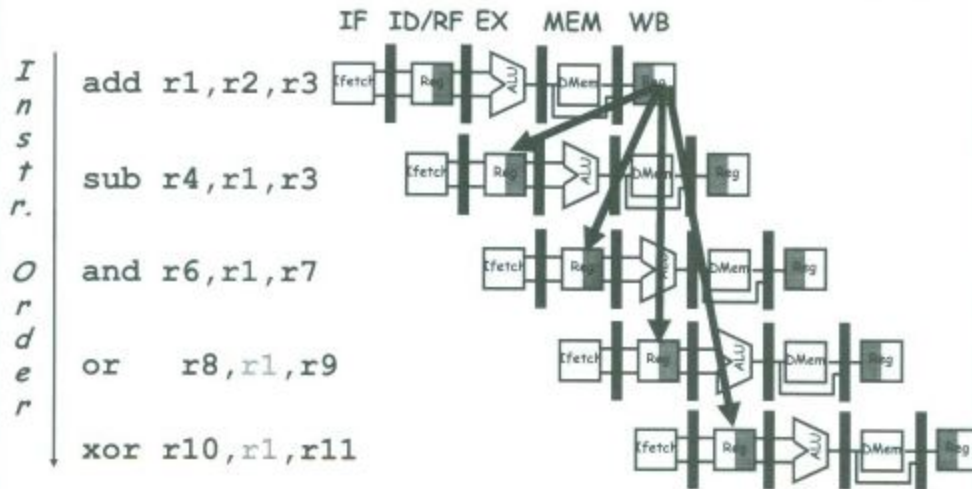
CMSC 611: Advanced Computer Architecture

Pipelining

Some material adapted from Mohamed Younis, UMBC CMSC 611 Spr 2003 course slides
Some material adapted from Hennessy & Patterson / © 2003 Elsevier Science

Data Hazards

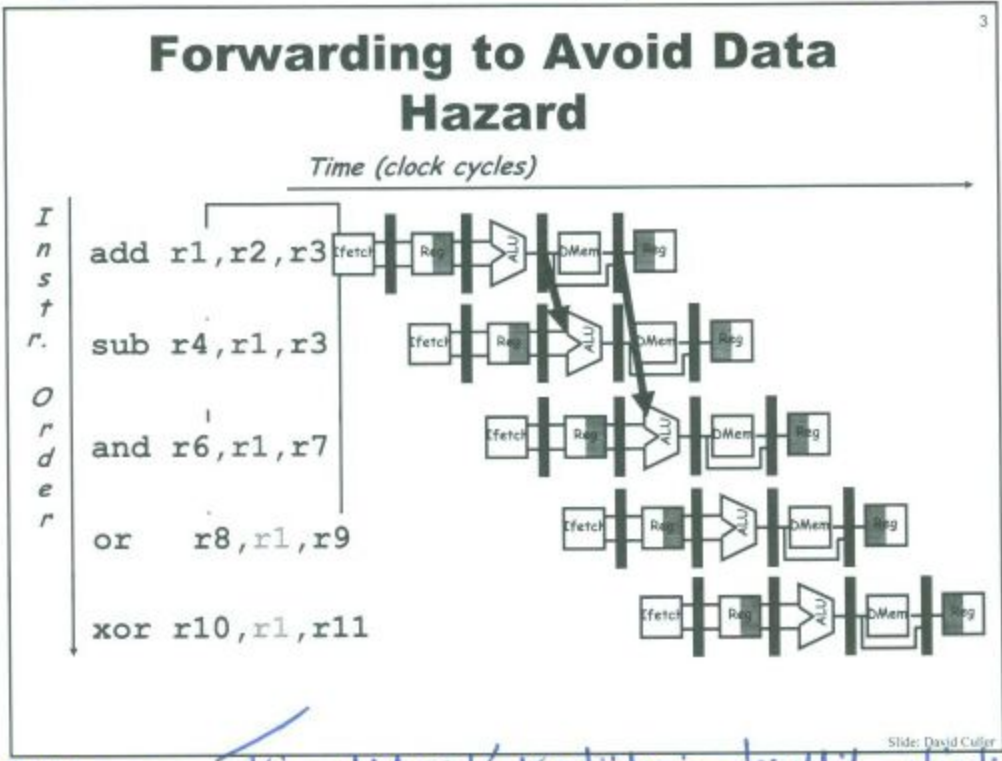
Time (clock cycles)



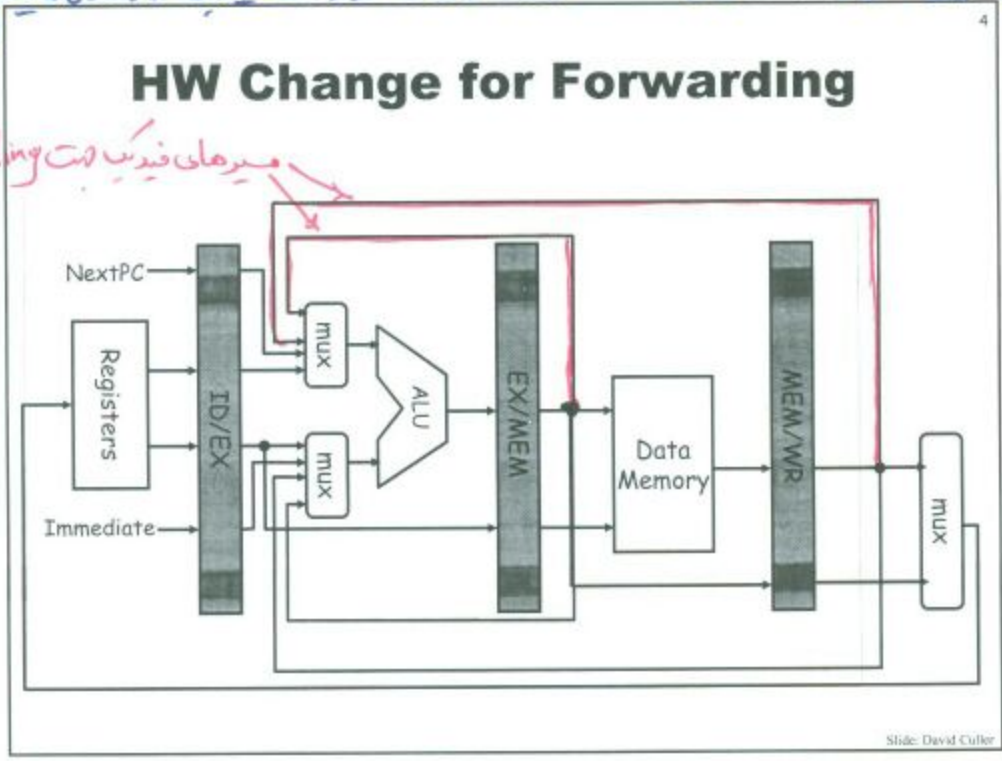
Slide: David Culler

همانگونه گفته شد برای رفع مخاطره دیتا هیزی از روش‌ها انتقاد است که این کار باعث ایجاد تأخیر می‌گردد لذا باید سعی شود این مشکل بدون ایجاد تأخیر برطرف شود. برای رفع اشکال مربوط به مخاطره دیتا هیزی از راه‌ها استفاده از forwarding است.

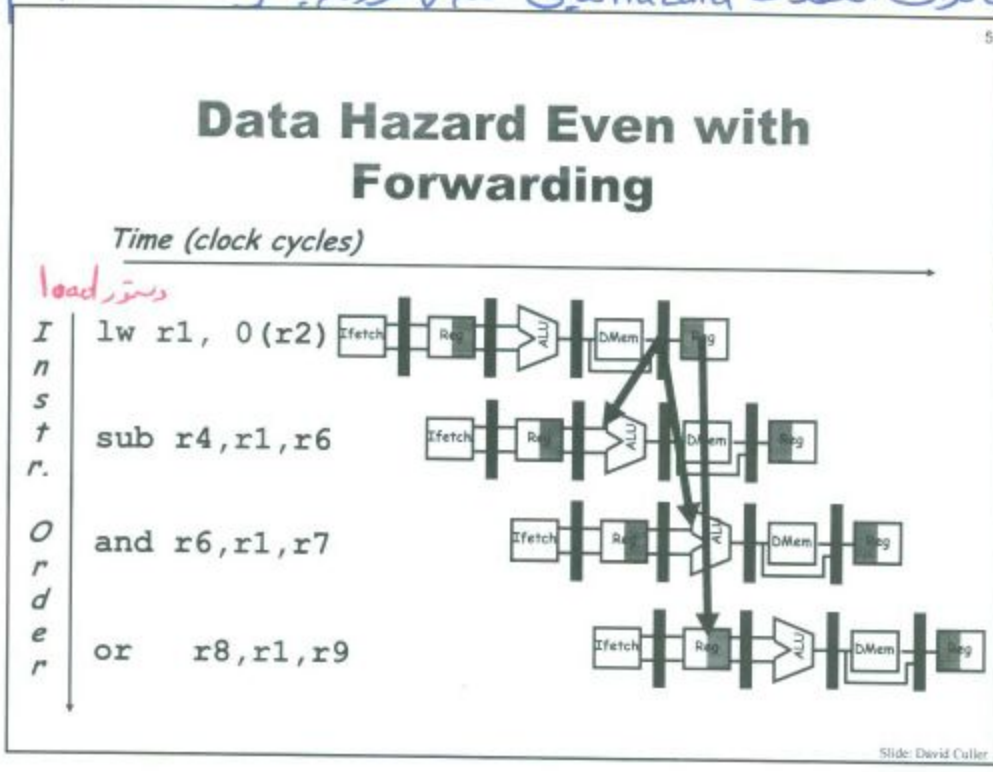
latch مانند یک نویسن است که به طور موقت اطلاعات هر مرحله را ذخیره می کند. latch این امکان را می دهد که اطلاعات دستور اول را ذخیره نموده و اطلاعات مربوط به دستور دوم را در ثبات های خودش بنویسد. در این حالت ALU در دستور دوم قرار است از اطلاعات قبلی استفاده نماید که هنوز اطلاعات در آن ثبت نگردیده است ولی به صورت یک نویسن در latch مربوط به ALU قرار دارد و ALU می تواند از این اطلاعات مربوط



به latch خودش استفاده نماید. لذا از نظر سخت افزاری امکانی را فراهم می نمایم که دوری ALU به خروجی latch خودش منتقل گردد که به این عمل Forwarding می گویند. در دستور سوم نیز می توان از Forward مربوط به حافظه استفاده نمود. این مسیرها را از طریق سخت افزاری و توسط مسیرهای فیدبک ایجاد می نمایم.



گاهی اوقات حتی با استفاده از Forwarding نیز Hazard متبایخ خواهد داشت. آن در دستور load و Store وجود داشته باشد علی‌رغم استفاده از Forwarding باز هم مضطره دارد به وجود برآید. چون در این حالت یک مرحله عقب تر هستیم پس در دستور دوم نمی‌توان از خروجی latch مربوط به حافظه استفاده نمود. به این مضطرات اصطلاحاً Hazardهایی گفته می‌شود که با Forwarding هم نمی‌توان آنها را از بین



برود. برای رفع این Hazardها حتماً باید یک بازه زمانی منتظر بمانیم

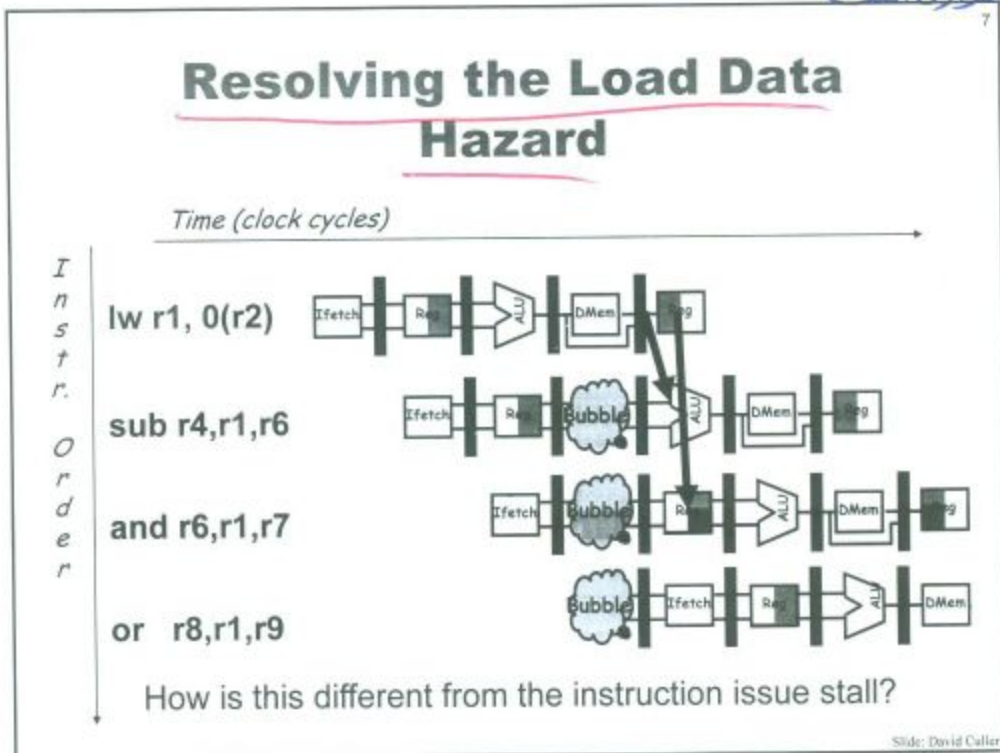
Resolving Load Hazards

- Adding hardware? How? Where?
- Detection?
- Compilation techniques?
- What is the cost of load delays?

Slide: David Culler

برای رفع Hazardهای مربوط به دستور load که با Forwarding هم نمی‌توان آنها را برطرف نمود یا باید سخت‌افزارها را دگرگونی اضافه کنیم یا ترتیب اجرای دستورات را تغییر دهیم. در هر صورت باید توجه نماییم که در اجرای هر یک از این روش‌ها باید بررسی نماییم آیا روش به کار برده شده از نظر هزینه اجرای آن به صرفه است یا بهتر است که منتظر بمانیم. در هر صورت یک کامپایلر یا مدیری باید وجود داشته باشد که ابتدا تشخیص دهد Hazard وجود دارد و سپس تشخیص دهد از کدام روش برای برطرف نمودن آن استفاده نماید.

تاکنون ما گفتیم برای رفع Hazard می‌توانیم Wait ایجاد نموده و اجازه ورود کل مراحل دستور را به Pipeline ندهیم ولی برای دستوری مانند load یا store در صورتی که مضامی این رخ دهد می‌توانیم جای آنکه کل مراحل دستورات را با ایجاد حباب در حالت انتظار قرار دهیم فقط همان مرحله‌ای را که ایجاد مضامی می‌نماید را منتظر بمانیم. حسن این روش آن است که در ورود دستورات وقفه‌ای ایجاد نمی‌شود. این روش یک روش سخت افزاری است.



Software Scheduling to Avoid Load Hazards

Try producing fast code for

$$a = b + c;$$

$$d = e - f;$$

assuming a, b, c, d, e, and f in memory.

Slow code:	Fast code:
LW Rb,b	LW Rb,b
LW Rc,c	LW Rc,c
ADD Ra,Rb,Rc	ADD Ra,Rb,Rc
SW a,Ra	LW Rf,f
LW Re,e	SW a,Ra
LW Rf,f	SUB Rd,Re,Rf
SUB Rd,Re,Rf	SW d,Rd
SW d,Rd	

Slide: David Callier

این دستورات از روش‌های برطرف نمودن Hazardهای مربوط به دیتا، استفاده از روش‌های جایابی دستورات است. در مثال فوق با توجه به آنکه دستور LW Re,e هیچ وابستگی به Rc ندارد لذا با جایابی این دو دستور می‌توان Hazard مربوط به دیتا را برطرف نمود. این روش در جاهایی استفاده می‌شود که از پردازنده‌های برای اجرای pipeline دستورات مشخص استفاده می‌شود که در آن موقع به صرفه است دستورات را به گونه‌ای جایابی نمود که کمترین Hazard مربوط به دیتا وجود داشته باشد. برطرف نمودن Hazardهای دیتا به این روش اصلاً الگوریتم ثابتی ندارد و در خیلی از جاها امکان جایابی دستورات وجود ندارد. ۲۱

Instruction Set Connection

- What is exposed about this organizational hazard in the instruction set?
- k cycle delay?
 - bad, CPI is not part of ISA
- k instruction slot delay
 - load should not be followed by use of the value in the next k instructions
- Nothing, but code can reduce run-time delays
- MIPS did the transformation in the assembler

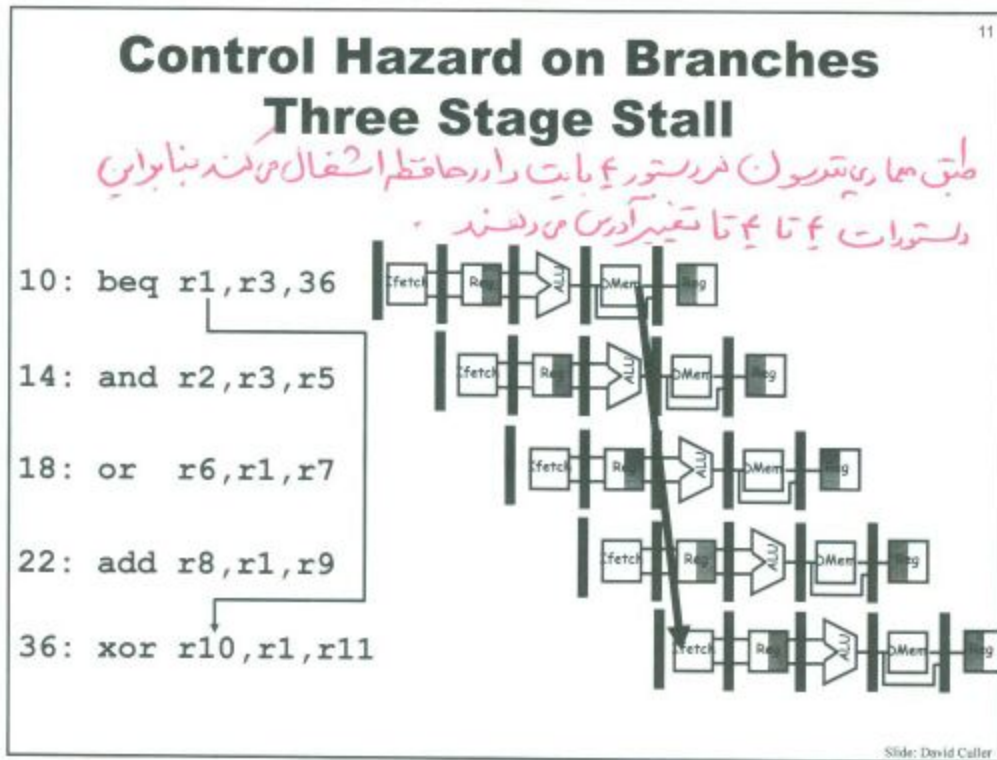
Slide: David Culler

Pipeline Hazards

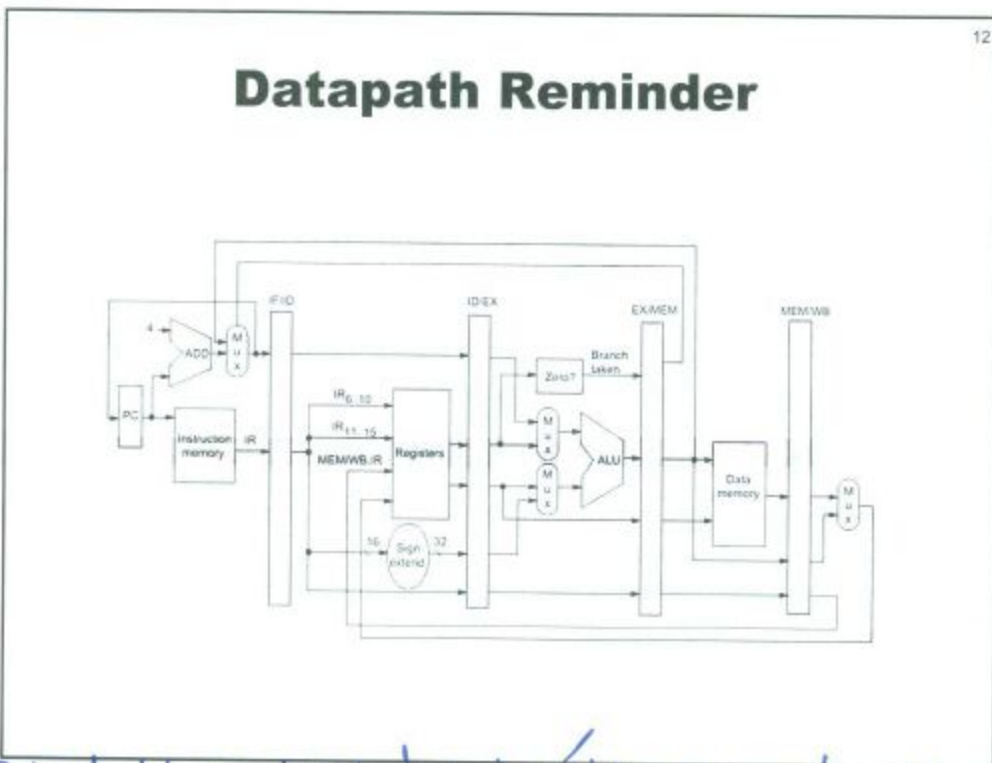
- Cases that affect instruction execution semantics and thus need to be detected and corrected
- Hazards types
 - Structural hazard: attempt to use a resource two different ways at same time
 - Single memory for instruction and data
 - Data hazard: attempt to use item before it is ready
 - Instruction depends on result of prior instruction still in the pipeline
 - Control hazard: attempt to make a decision before condition is evaluated
 - branch instructions
- Hazards can always be resolved by waiting

Hazard های کنترلی مخاطراتی هستند در دستورات شرطی یا اشعاب (پوش) با آنها برخورد می کنیم. این نوع مخاطرات را می توان فقط با استفاده از روش های سخت افزاری و نرم افزاری بیهوش نمود و نمی توان آنها را از بین برد.

در مثال زیر بیان شده است اگر ۲۱ برابر ۱۳ باشد به خانه ۳۶ حافظه پرش نماید در غیر این صورت دستور بعدی اجرا شود. این مظاهره مربوط به Hazard کنترل خواهد بود. اگر شرط را اجرا نکنیم و دستورات را به ترتیب اجرا نمائیم در آن صورت اگر شرط برقرار بوده باشد سه دستور ۱۴، ۱۸، ۲۲ نباید اجرا می شده اند و اگر تغییری ایجاد شده باشد باید دوباره آنرا به حالت اول برگردانیم و اگر پرش نمائیم و شرط برقرار نباشد در



۲ صورت اجرای سه دستور ۱۴، ۱۸، ۲۲ را حذف کرده ایم.



ساده ترین راه حل برای برطرف نمودن مظاهره کنترلی منتظر ماندن است یا اجرای دستور شرطی که این روش بدترین روش است. روش دوم تخمین است یعنی پیش بینی کنیم که شرط درست است یا نادرست است. روش دیگر تخمین آماره یا پیش بینی آماره است یعنی برای آن نمونه برداری های انجام شده تخمین بزنیم یعنی به صورت هوشمندانه و براس تجربیه و کارهای قبلی که انجام شده، پیش بینی را انجام دهیم. روش دیگر بستن به نوع دستور دارد مثلاً برای دستور for بهترین است که پرش را انجام ندهیم در صورتیکه برای دستور if، else بهترین است که پرش نمائیم.

۳۴

اگر دستورات شرطی بیش از ۳۰٪ نیاز دارند. اگر ۳۰٪ دستورات به صورت شرطی باشند و هر دستوری که سیکل نیاز داشته باشد، برای رفع مخاطره کنترلی دو کار باید انجام داد:

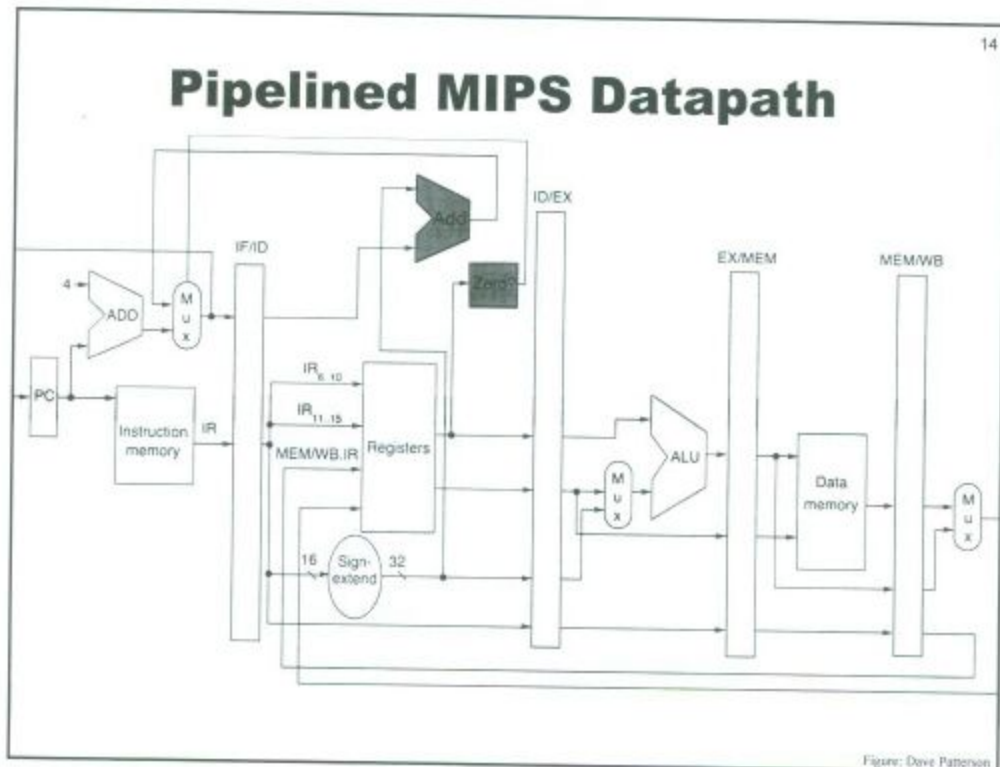
- ۱- تشخیص دهیم که شرط برقرار است یا نه که این کار با حوک کردن پریم ثبات انجام می شود.
- ۲- آدرس مورد جایی که قرار است پرش به آنجا انجام شود را بدست آوریم.

13

Example: Branch Stall Impact

- If 30% branch, 3-cycle stall significant!
- Two part solution:
 - Determine branch taken or not sooner, AND
 - Compute taken branch address earlier
- MIPS branch tests if register = 0 or $\neq 0$
- MIPS Solution:
 - Move Zero test to ID/RF stage
 - Adder to calculate new PC in ID/RF stage
 - 1 clock cycle penalty for branch versus 3

Slide: David Callar



چون هر دستوری که سیکل انجام می شود ندارد صورت اشتباه تشخیص دادن شرط بصورت دو سیکل به عقب برگردیم که به این کار جریمه شدن یا پینالتی گفته می شود. هر چه تعداد سیکل های مربوط به اجرای یک دستور بیشتر شود در صورت جریمه شدن تعداد سیکل بیشتری باید به عقب برگشت تا تمام

برای برخورد با Hazard های کنترلی چهار روش به صورت زیر ارائه گردیده است :

- ۱- منتظر بمانیم تا ببینیم شرط درست است یا نه
- ۲- با فرض برقرار نبودن شرط جلو رویم (معمولاً ۴۷٪ دستوراً بر اساس معماری سیستم پتربون شرطها را برقرار است)
- ۳- با فرض برقرار بودن شرط پیش رویم (معمولاً طبق آمار در ۵۳٪ دستوراً، شرطها برقرار هستند)

15

Four Branch Hazard Alternatives

1. Stall until branch direction is clear
2. Predict Branch Not Taken
 - Execute successor instructions in sequence
 - "Squash" instructions in pipeline if branch taken
 - Advantage of late pipeline state update
 - 47% MIPS branches not taken on average
 - PC+4 already calculated, so use it to get next instruction
3. Predict Branch Taken
 - 53% MIPS branches taken on average
 - But haven't calculated branch target address in MIPS
 - MIPS still incurs 1 cycle branch penalty
 - Other machines: branch target known before outcome

Slide: David Callier

۴- بر اساس حالت های قبلی در دستورات قبلی از این تخمین استفاده می کنند، یعنی مثلاً اگر PC بود از فرض برقرار نبودن شرط پیش می رود.

16

Four Branch Hazard Alternatives

4. Delayed Branch
 - Define branch to take place AFTER a following instruction
 - branch instruction
 - sequential successor₁
 - sequential successor₂
 -
 - sequential successor_n
 - } Branch delay of length n
 -
 - branch target if taken
 - 1 slot delay allows proper decision and branch target address in 5 stage pipeline
 - MIPS uses this

Slide: David Callier

زاینده ای که بر این بدست آوردن speedup مربوط به pipeline استفاده شود به صورت زیر است:

$$\text{pipeline speedup} = \frac{\text{Pipeline depth}}{1 + \text{pipeline stall CPI}}$$

طول پایپ لاین (چند مرحله است)
 مقدار مازانی که باید منتظر بمانیم تا دستورات pipeline کامل شود.

Branch-Delay Scheduling Requirements

Scheduling Strategy	Requirements	Improves performance when?
(a) From before	Branch must not depend on the rescheduled instructions	Always
(b) From target	Must be OK to execute rescheduled instructions if branch is not taken. May need to duplicate instructions.	When branch is taken. May enlarge programs if instructions are duplicated.
(c) From fall through	Must be okay to execute instructions if branch is taken.	When branch is not taken.

- Limitation on delayed-branch scheduling arise from:
 - Restrictions on instructions scheduled into the delay slots
 - Ability to predict at compile-time whether a branch is likely to be taken
- May have to fill with a no-op instruction
 - Average 30% wasted
- Additional PC is needed to allow safe operation in case of interrupts (more on this later)

Example: Evaluating Branch Alternatives

$$\text{Pipeline speedup} = \frac{\text{Pipeline depth}}{1 + \text{Pipeline stall CPI}}$$

$$= \frac{\text{Pipeline depth}}{1 + \text{Branch frequency} \times \text{Branch penalty}}$$

Assume:
 14% Conditional & Unconditional
 65% Taken; 52% Delay slots not usefully filled

Scheduling Scheme	Branch Penalty	CPI	Pipeline Speedup	Speedup vs stall
Stall pipeline	3.00	1.42	3.52	1.00
Predict taken	1.00	1.14	4.39	1.25
Predict not taken	1.00	1.09	4.58	1.30
Delayed branch	0.52	1.07	4.66	1.32

Slide: David Culler

اگر بخواهیم بر اساس تخمین جلو رویم خواهیم داشت:

$$\text{pipeline speedup} = \frac{\text{Pipeline depth}}{1 + \text{Branch frequency} \times \text{Branch penalty}}$$

چند بار مجبور هستیم به دلیل این تخمین نادرست زود آید به عقب برگردیم (تقداری که جای است که اگر دستورات pipeline را به اشتباه جلو رفتیم باید برگردیم)

چند درصد دستورات به صورت اشتباه هستند

Multiprocessors

سیستم های چند پردازنده ای

اولین ایده برای استفاده از سیستم های چند پردازنده ای این بود که به جای استفاده از یک پردازنده قوی با حجم پردازش بالا، از چند پردازنده کوچکتر با در کنار هم قرار دادن آنها از توان پردازش مجموع آنها برای رسیدن به هدف استفاده کنیم.

پردازنده قوی و بزرگ \equiv استفاده از چندین پردازنده کوچک در جهت رسیدن به توان پردازش بالا

این پردازنده ها باید بتوانند دستورات را به صورت موازی انجام دهند. در pipeline نیز دستورات همزمان اجرا نمی شوند بلکه روند اجرای دستورات سریعتر بود.

یکی از ویژگی های که سیستم ها را به صورت چند پردازنده ای استفاده می نمایند، قابلیت دسترسی پذیری یا Availability آنهاست. در زمانیکه از چند پردازنده استفاده می کنیم در اثر افزایش کار افتادن یکی از آنها سیستم هر چند بتواند کمتر و ولی به کار خود ادامه خواهد داد در صورتیکه اگر از یک پردازنده قوی استفاده می کنیم، در صورت بروز اشکال در آن، سیستم Down یا متوقف خواهد شد. پس بهترین خصوصیت استفاده از سیستم های چند پردازنده ای، خصوصیتی بنام قابلیت دسترسی پذیری است. لازم به ذکر است مدیریت سیستم های چند پردازنده به نسبت تک پردازنده مشکل تر است.

لازمه استفاده از حداکثر کارایی یک سیستم چند پردازنده ای، موازی سازی در سطح برنامه های می باشد که می باید توسط آن اجرا شوند. یعنی در واقع زمانی که می خواهیم از سیستم چند پردازنده ای استفاده می کنیم باید بتوانیم دستوراتی را که قرار است توسط این پردازنده اجرا شوند به صورت موازی اجرا می کنیم و گرنه نمی توانیم از حداکثر توان پردازنده ها استفاده کنیم. اگر بتوانیم دستورات خود را ۱۰۰٪ موازی سازی می کنیم آنگاه می توانیم از ۱۰۰٪ توان سیستم چند پردازنده ای خود استفاده می کنیم.

در سطح سخت افزار ما با برنامه سروکار نداریم بلکه با دستورات سروکار داریم. در یک لول بالاتر در سطح سیستم عامل با دو قسمت سروکار داریم پس موازی سازی می تواند در سطح انجام شود.

job (کاره که به سیستم عامل می دهیم تا سیستم عامل آن را مدیریت نموده و اجرا نماید)

process (هر job از تعدادی پروسس یا پردازش تشکیل شده است)

پس بنا بر این هر چه بتوانیم طرح‌های بیشتری را به صورت موازی اجرا نماییم، اجرای ما سریعتر خواهد بود.

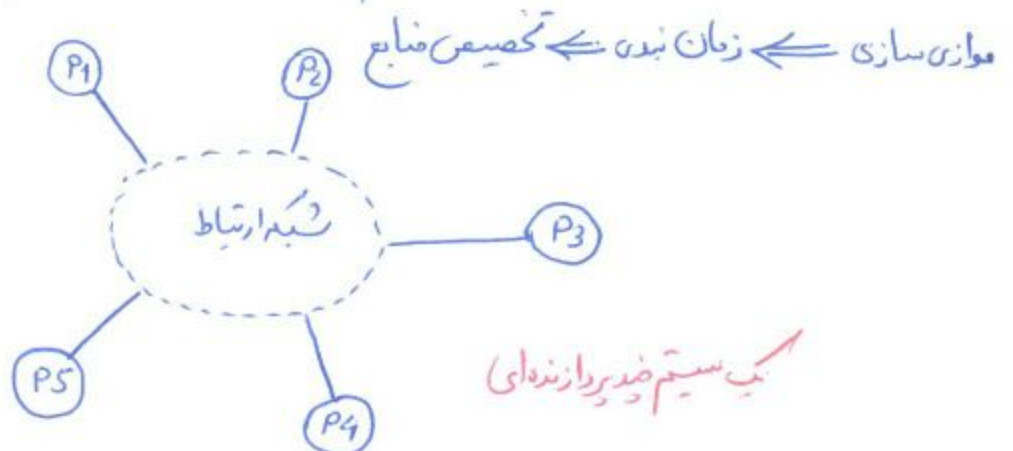
مقیاس پذیری:

در صورتیکه روش را که برای اجرای یک کار ارائه می‌کنیم با بزرگتر شدن کار باز هم قابل اجرا باشد گفته می‌شود که این روش مقیاس پذیر است. یکی از ویژگی‌های سیستم‌های ضد پردازنده‌ای مقیاس پذیری است اما هدف مقیاس پذیری انجام نمی‌شود.

* مقیاس پذیری در صورت مقیاس پذیری موازی سازی نرم افزارها و برنامه‌های قابل اجرا به روش یک سیستم ضد پردازنده، برای یک سیستم ضد پردازنده امکان پذیر است. به عبارتی دیگر اگر برنامه بتواند خود را با هر تعداد از پردازنده‌ها تطبیق دهد، سیستم ضد پردازنده‌ای دارای ویژگی مقیاس پذیری است. اکثر سیستم‌های ضد پردازنده‌ای موفق سیستم‌های دستگیر به منظورهای خاصی طراحی شده‌اند زیرا سیستم‌های ضد پردازنده‌ای کاملاً وابسته به نرم افزارهای هستند که بر روی آنها اجرا می‌شوند یعنی کارایی یک سیستم ضد پردازنده کاملاً وابسته به موازی سازی برنامه‌هایی است که روی آن اجرا می‌شوند.

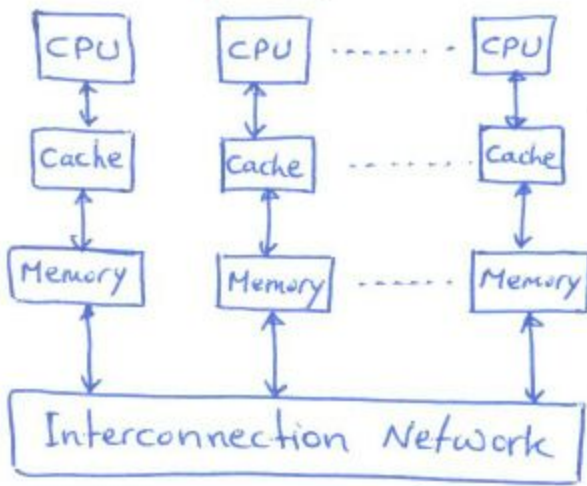
پس از مطالعه موازی سازی بهترین مدل‌ها که در استفاده از سیستم‌های ضد پردازنده‌ای مطرح است مطالعه زمان بندی است یعنی اینکه اگر قرار است مثلاً دستورالعمل توسط ۱۰۰ پردازنده به صورت موازی انجام شوند، کدام یک از دستورات جهت اجرا بر روی پردازنده‌ها قرار داده شوند یا زمانیکه از پردازنده‌هایی با توان‌های اجرای متفاوت استفاده می‌کنیم، کدام قسمت از برنامه را بر روی کدام یک از پردازنده‌ها اجرا نماییم تا کارایی سیستم بالاتر باشد.

پس از زمان بندی باید تخصیص منابع انجام دهیم یعنی ممکن است منابع متفاوتی وجود داشته باشد، حال اینکه کدام منبع را در چه زمانی در اختیار کدام پردازنده قرار دهیم که می‌تواند در کارایی سیستم مؤثر باشد.



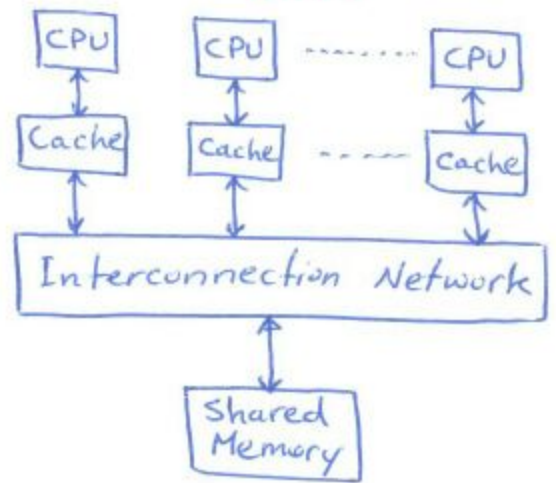
در سیستم‌های چند پردازنده‌ای، استفاده از حافظه می‌تواند به دو صورت زیر باشد. توجه شود که در هر دو روش، همه پردازنده‌ها دارای حافظه اختصاصی مرتبط خود بنام Cache می‌باشند.

Clusters Memory Multiprocessor (CMP)



سنتز کلاسیک (متبنی بر تبادل پیام)

Shared Memory Multiprocessor (SMP)

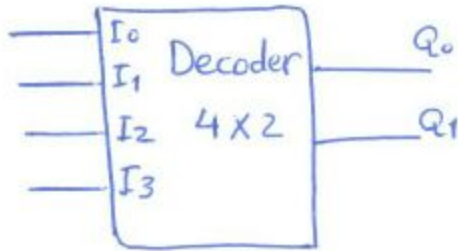


در حالت استفاده از حافظه مشترک، فضای حافظه به طور منطقی بین تمام پردازنده‌ها تقسیم می‌شود. با توجه به آنکه حافظه فقط دارای یک آدرس دهی می‌باشد لذا در یک زمان فقط یکی از پردازنده‌ها می‌تواند به آن دسترسی داشته باشد یعنی بر روی آدرس فقط می‌توان آدرس یک خانه را قرار داد که از آن خوانده یا نوشته شود. پس در هر زمان فقط یکی از پردازنده‌ها به حافظه دسترسی دارد.

در روش cluster که بر اساس سیستم متبنی بر تبادل پیام می‌باشد، هر پردازنده علاوه بر حافظه Cache دارای یک حافظه اختصاصی بوده که دستورات را از آنجا می‌خواند و پس از اجرای دستورات، اطلاعات را مجدداً در آنجا ذخیره می‌کند. از طریق یک شبکه اتصال داخلی، پردازنده‌ها به حافظه‌های گنبدی دسترسی داشته و می‌توانند از اطلاعات ذخیره شده در آنجا استفاده نمایند. این روش هم‌استفاده از چند کامپیوتر مستقل برای رسیدن به یک توان پردازشی بالاتر می‌باشد که از طریق یک شبکه ارتباطی به یکدیگر متصل شده‌اند.

یکی از روش‌های دسترسی به حافظه مشترک یا منابع به صورت Random می‌باشد. در این حالت اولویت بندی جهت دسترسی لحاظ نمی‌گردد. روش دیگر به صورت ترتیبی است که در این حالت از یک دیکتور معمولی می‌توان استفاده نمود. روش دیگر استفاده از Encoder یا اولویت می‌باشد. مشخص است که روش اولویت بندی گام آبی سیستم را بالاتر برد ولی دسترسی به این روش به پیچیده‌تر است.

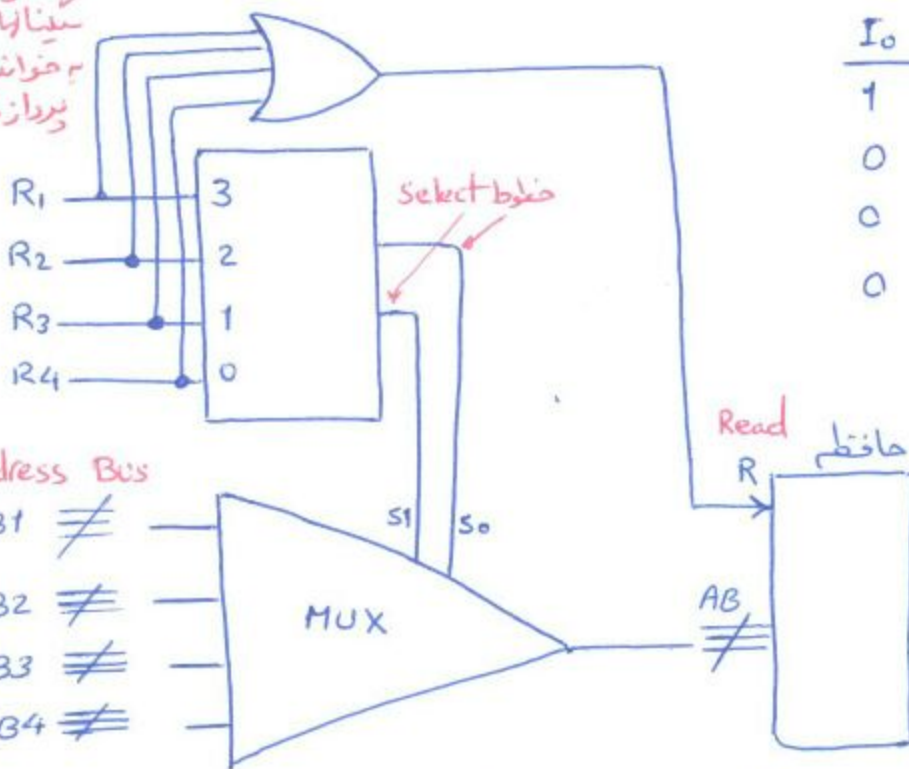
روش دیگر دسترسی به حافظه به این صورت است که اگر مثلاً تمام پردازنده‌ها همزمان نیاز به دسترسی به حافظه داشته باشند، اولویت یا پردازنده شماره پائین‌تر باشد.
یکی از روش‌های سخت افزاری به صورت مدار زیر می‌باشد. در این روش می‌توان حافظه را به عنوان یک منبع ایجاد وقفه نیز در نظر گرفت.



I_0	I_1	I_2	I_3	Q_0	Q_1
1	0	0	0	0	0
0	1	0	0	0	1
0	0	1	0	1	0
0	0	0	1	1	1

اشکالی که Decoder معمولی دارد آن است که اگر همزمان دو ورودی فعال وجود داشته باشد جواب غیر دلخواه می‌دهد.

کمیترها مربوط به خواندن هر پردازنده



I_0	I_1	I_2	I_3	Q_0	Q_1
1	X	X	X	1	1
0	1	X	X	1	0
0	0	1	X	0	1
0	0	0	1	0	0

Encoder با اولویت

با استفاده از این روش سخت افزاری، می‌توان دسترسی به حافظه با اولویت ایجاد نمود.

Uniform Memory Access : UMA

زمان دسترسی به حافظه در تمام پردازنده‌ها یکسان و یکنواخت است. برنامه نویسی آن به لحاظ زمان نبدی ساده است. دسترسی به حافظه به صورت سخت افزاری نمی‌باشد بلکه به صورت Random است.

دسترسی به حافظه در SMP

Non-Uniform Memory Access : NUMA

پردازنده‌ها در دسترسی به حافظه نسبت به یکدیگر دارای اولویت می‌باشند که زمان دسترسی به حافظه را برای آنها متفاوت می‌کند. برنامه نویسی آن سخت‌تر است.

در سطح دسترسی به حافظه به صورت یکسان و یکنواخت باشد امکان برخورد و خطا بیشتر است لذا امکان گسترش پذیری در NUMA بیشتر است.

این نوع ساختارها را ساختارهای مبتنی بر پیام می‌گویند زیرا پردازنده مستقیماً به حافظه دسترسی ندارد بلکه هر اطلاعاتی را که نیاز داشته باشد با ارسال یک پیام، حافظه یک کپی از آن را برای پردازنده ارسال می‌کند. نکته‌ای که باید توجه نمود آن است که مدیریت سیستم‌هایی که از حافظه‌های مجزا استفاده می‌نمایند دشوارتر از حافظه‌های مشترک است.

در حالتی که درها (Multi cores) چندین پردازنده بر روی یک تراشه قرار گرفته‌اند و تکنولوژی آن از حالتی پروسورهای جدیدتر است ولی منطق همان منطق است.

NOC = Network on chip

کپی‌های از پردازنده‌ها که روی یک chip قرار گرفته‌اند.

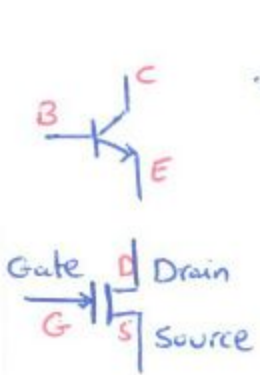
تفاوتی که بین NOC ها با حالتی که درها وجود دارند آن است که در NOC ها از چندین پردازنده مختلف که جهت استفاده کارهای مختلف می‌باشند بر روی یک chip استفاده می‌گردد. NOC ها برای کاربردهای خاص استفاده می‌شوند مثلاً برای کاربردهای مخابراتی یا تصویری. در NOC موازی سازی را می‌توان ساده‌تر

انجام داد.

* حالتی که درها از یک حافظه مشترک استفاده می‌کنند.

کارت گرافیک نیز به صورت یک پردازنده است که جهت پردازش تصویر به صورت یک چند پردازنده عمل می‌کنند. لذا در کامپیوترها امکان ترکیب حالتی که درها با حالتی پروسورها وجود دارد که در این حالت نیز امکان استفاده از حافظه‌های مجزا یا مشترک وجود دارد.

حافظه ها



بین ترانزیستورهای که در مدل اول مورد استفاده قرار گرفتند از ترانزیستورهای BJT بودند. این ترانزیستورها دارای سرعت بالایی بودند ولی مصرف توان بسیار زیادی نیز داشتند که به همین علت باعث شد تا از نوع دیگری از ترانزیستورها به نام Mos fet استفاده کرد. ترانزیستورهای Mos fet در حالت استاتیف توان مصرفی صفر دارد و تنها توان مصرفی آن در زمان تغییر حالت قطع و وصل و برعکس می باشد.

در ترانزیستورهای BJT معیار ولتاژ بود در حالی که در حالت قطع و وصل بودند در حالت وصل اینتریکتور وصل شده و اختلاف ولتاژ بین آنها صفر می شود در حالت قطع، اختلاف ولتاژ بین اینتریکتور برابر ولتاژ منبع تغذیه می باشد. در Mos fet ها معیار جریان می باشد بدین صورت که در حالت وصل، یک جریان بین درین و سورس ایجاد می شود که در حالت قطع مقدار این جریان تقریباً صفر است.

در Mos fet ها در حالت وصل دو حالت داریم: تغذیه / اشباع

تغذیه: در این حالت جریان عبوری بین درین و سورس بر اساس جریان گیت کنترل می شود و درجه جریان گیت بیشتر شود، جریان عبوری بین درین و سورس نیز بیشتر می شود.

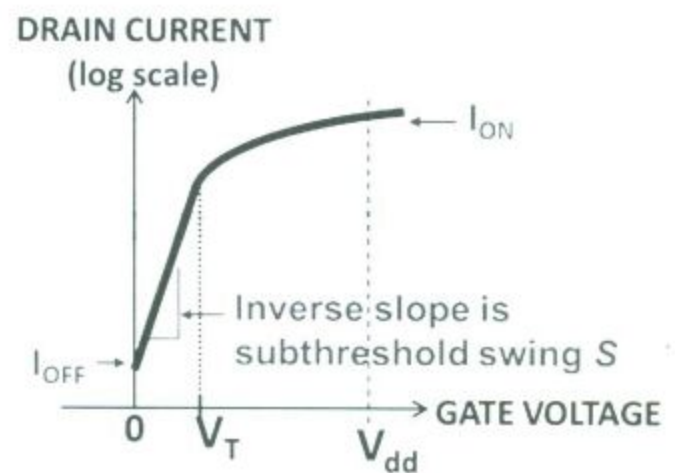
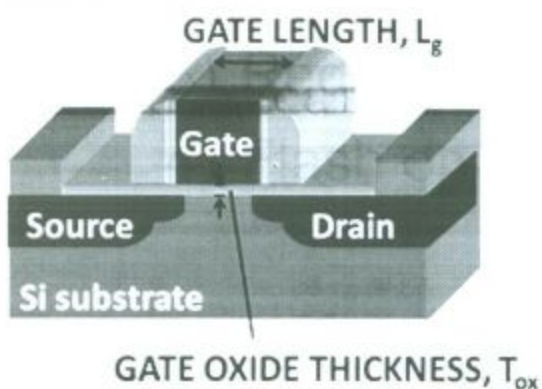
اشباع: در این حالت مقدار جریان عبوری بین درین و سورس مقدار ثابتی بوده و درجه جریان گیت را بیشتر نماییم، این جریان تغییری نمی کند و در واقع به حالت اشباع رسیده است.

درجه می توانیم در یک فضای مشخص تعداد ترانزیستور بیشتری را بر روی یک تراشه ایجاد نماییم در واقع مقدار گیت بیشتری خواهیم داشت و درجه تعداد گیت بیشتری داشته باشیم می توان محاسبات بیشتری را انجام دهیم و درجه در یک فضای مشخص بتوانیم محاسبات بیشتری را انجام دهیم، توان پردازش ما بیشتر خواهد بود. برای رسیدن به توان پردازش بالا رگین از کارها، استفاده از قطعات سیستم های پردازش به صورت موازی می باشد.

etimes

- 1) IBM Fabrication
- 2) Renesas: big/little
- 3) Reducing AMD power with gating
- 4) Flash on DDR4

Metal Oxide Semiconductor
Field-Effect Transistor:



DARPA UPSIDE

[August 17, 2012]

Digital Processors Limited by Power; What's The Upside?

(Targeted News Service Via Acquire Media NewsEdge) WASHINGTON, Aug. 14 -- The U.S. Department of Defense's Defense Advanced Research Projects Agency issued the following news release: Today's Defense missions rely on a massive amount of sensor data collected by intelligence, surveillance and reconnaissance (ISR) platforms. Not only has the volume of sensor data increased exponentially, there has also been a dramatic increase in the complexity of analysis required for applications such as target identification and tracking. The digital processors used for ISR data analysis are limited by power requirements, potentially limiting the speed and type of data analysis that can be done. A new, ultra-low power processing method may enable faster, mission critical analysis of ISR data.

Chapter Outline

- Introduction
- Parallel Computers
- Shared-Memory Programming
- Synchronization
- Cache Coherence

Introduction

- Growth in data-intensive applications.
 - Data bases, file servers, ...
- Growing interest in servers, server performance.
- Increasing desktop performance less important.
 - Outside of graphics
- Improved understanding in how to use multiprocessors effectively.
 - Especially servers where significant natural TLP
- Advantage of leveraging design investment by replication => CMPs or Multicores.
 - Rather than unique design

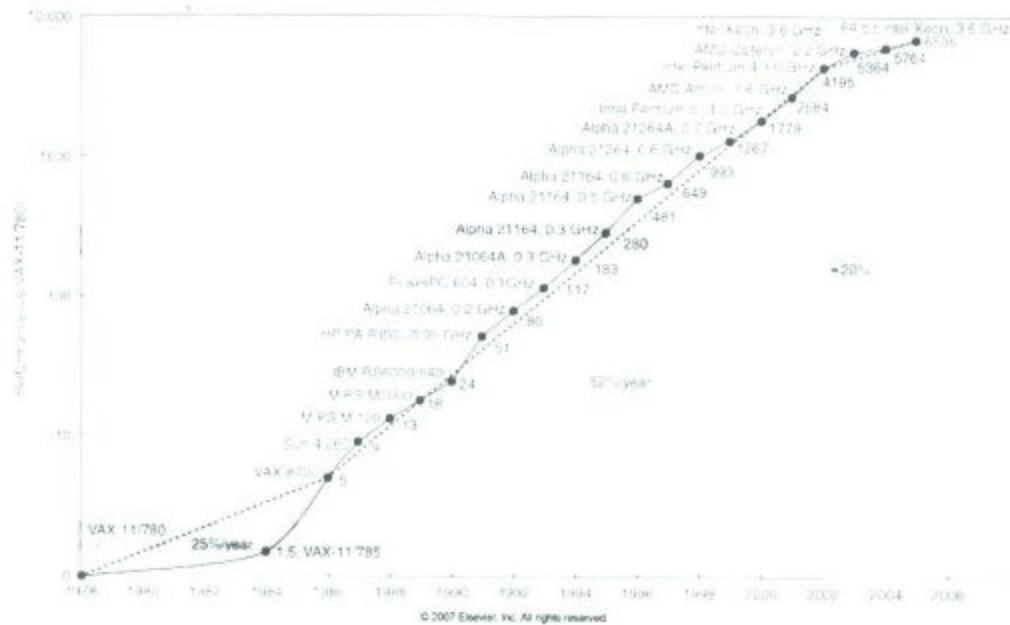
عواملی که باعث نیاز به پردازنده‌های سریع می‌باشند:

- ۱- اولین چیزی که باعث نیاز به پردازنده‌های سریع بود اینلیکین‌ها یا برنامه‌های مبتنی بر دیتا بود. مانند دیتا سنترها یا اینلیکین‌هایی که نیازمند پردازش و جستجوی سریع هستند.
 - ۲- علاقه بیشتر به افزایش بهره‌برداری از سرورهای سریع.
 - ۳- در کامپیوترهای شخصی در زمانی که از گرافیک و پردازش تصویر استفاده می‌کنیم نیاز به پردازنده‌های سریع‌تر می‌باشد.
- همه این عوامل باعث شد تا ما سیستمی را بنام سیستم‌های حالتی پرسر یا چند پردازنده معرفی کنیم که بهترین مزایای آن این است که الزامات ما را که همان انجام پردازش سریع می‌باشد را فراهم می‌نماید.

محدودیت‌هایی که برای کارایی پردازنده وجود دارد:

- ۱- بیچیدگی: هر قدر تراشه‌های بیشتری ایجاد کنیم به لحاظ کنترل و مدارات ALU پیچیده‌تر خواهند بود و به همین دلیل پشتیبانی کردن از چنین ساختاری دشوار است.
- ۲- تعدادسیم‌های ورودی در chip: تکنولوژی ساخت chip تا یک حدی به ما اجازه می‌دهد که تعدادسیم‌های ورودی در chip را افزایش دهیم.
- ۳- خنک کردن: بسته به هر پردازنده که می‌تواند مشکلاتی را در ایجاد تبادل حرارت ایجاد کند. هر چه تراشه‌های بیشتری در chip به کار رفته باشد، پردازش‌های بیشتری را انجام می‌دهد و به همین نسبت حرارت بیشتری نیز ایجاد می‌شود که خنک کردن آن به صورت یک معضل می‌باشد.
- ۴- کارایی پردازنده‌ها در استفاده از حافظه: زیرا حافظه‌ها کمترین عنصر در سیستم‌ها می‌باشند لذا حتی اگر پردازنده‌ای بسیار سریع ایجاد کنیم بدلیل استفاده از دیتاهایی که در حافظه‌های کند وجود دارند، سرعت پردازنده نمی‌تواند از یک حد مشخص بیشتر شود و همان‌طور که گفته شد، رسیدن سرعت حافظه‌ها چندین برابر کندتر از رسیدن سرعت پردازنده‌ها بوده است لذا سرعت حافظه باعث ایجاد محدودیت می‌گردد.

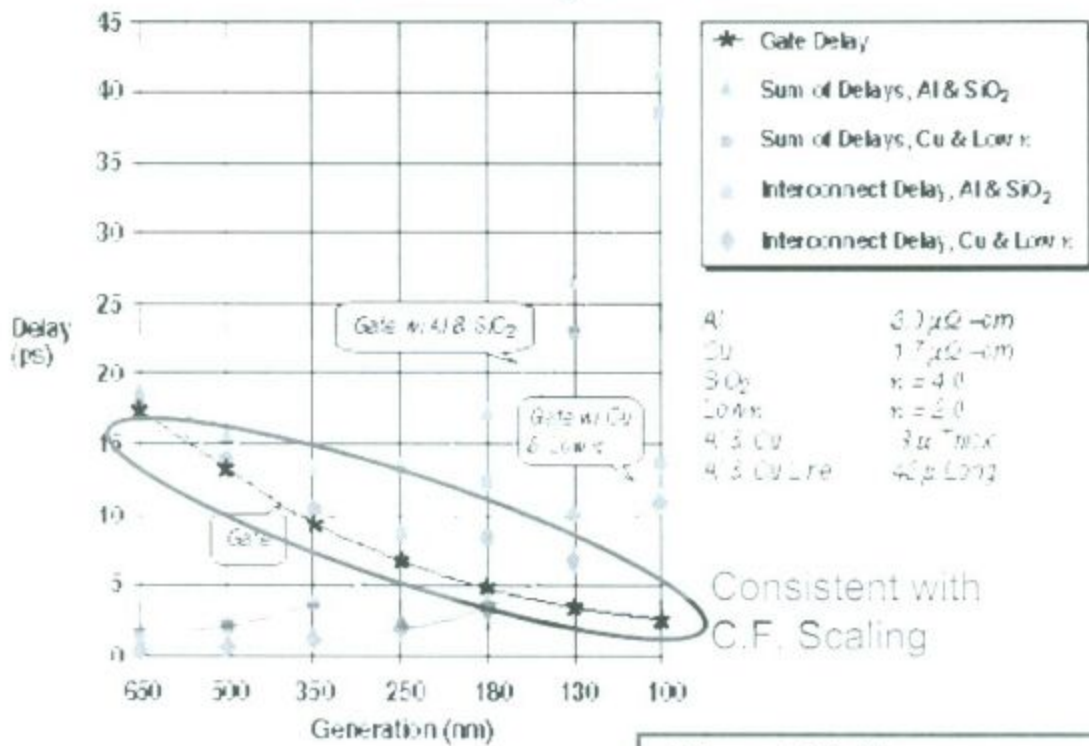
Another Reason



Limit to Processor Performance

- Complexity of exploiting ILP.
 - Difficult to support large instruction window and number of in-flight instructions.
- On-chip wires are becoming slower than logic gates.
 - Only a fraction of the die will be reachable within a single clock cycle.
- Cooling and packaging will be a real challenge due heat release.
- Memory and processor performance gap will continue to be a challenge.

Gate Delay Trends

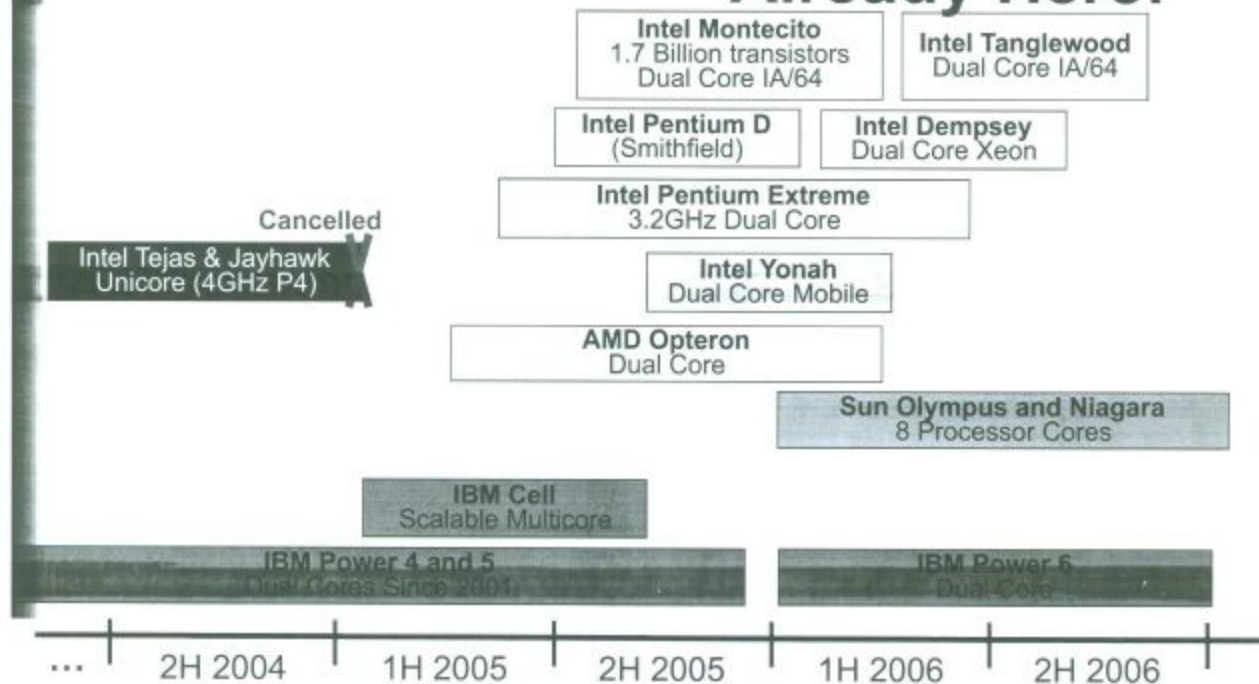


Each technology generation, gate delay reduced about 30% (src: ITRS '01)

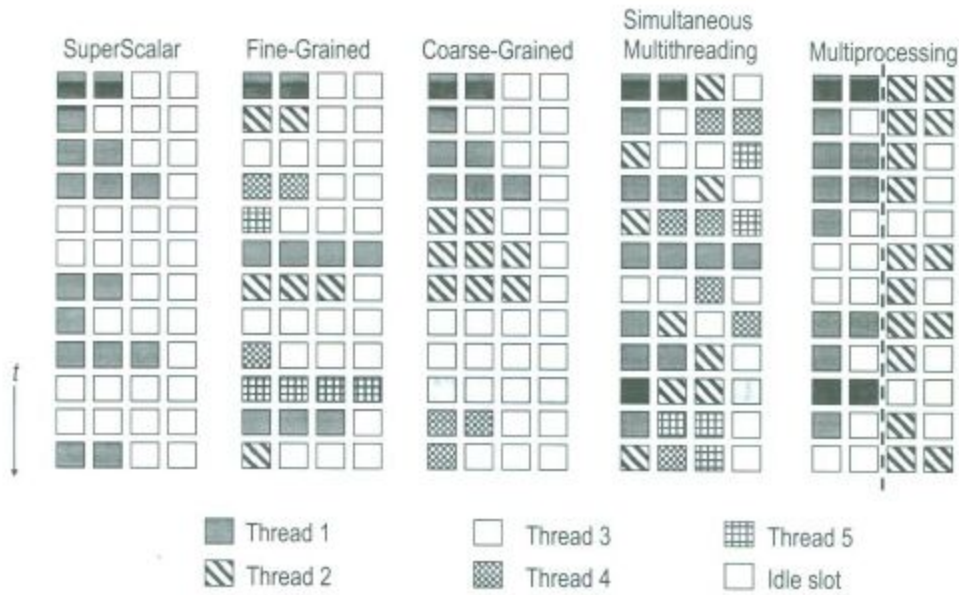
$$T_d = kCV/I$$

$$= kCV/(V_{dd} - V_t)^\alpha$$

Multicores are Coming Already Here!



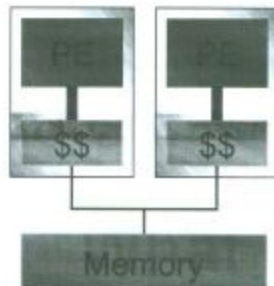
SS, MT, SMT, & MP



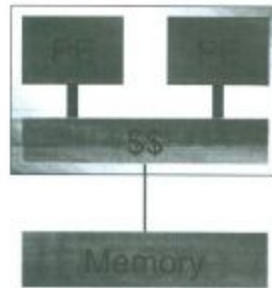
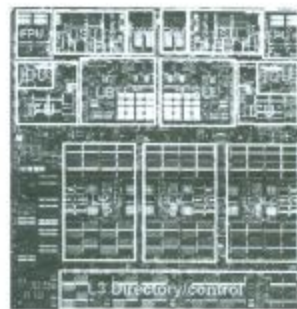
در حالتی پرسورسشن که پردازنده حافظه Cache مربوط خود را دارد و تمام اینها به حافظه اصلی به صورت اشتراکی در ارتباط هستند اما در حالتی که آنها تمام هستهها دارای یک Cache بوده و همگی آنها با یک حافظه اصلی (RAM) در ارتباط میباشند.

What is a Multicore?

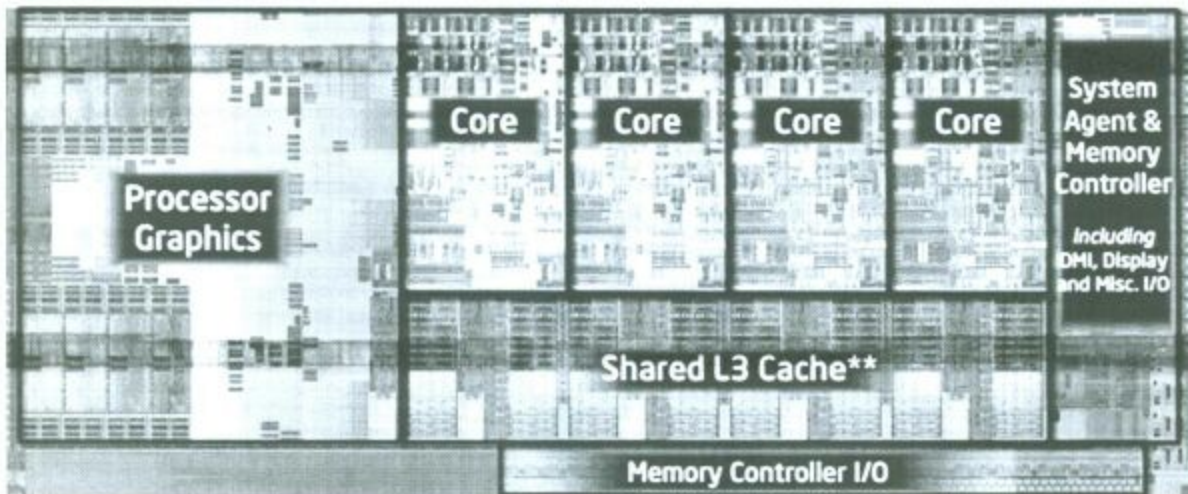
Traditional Multiprocessor



Basic Multicore IBM Power5

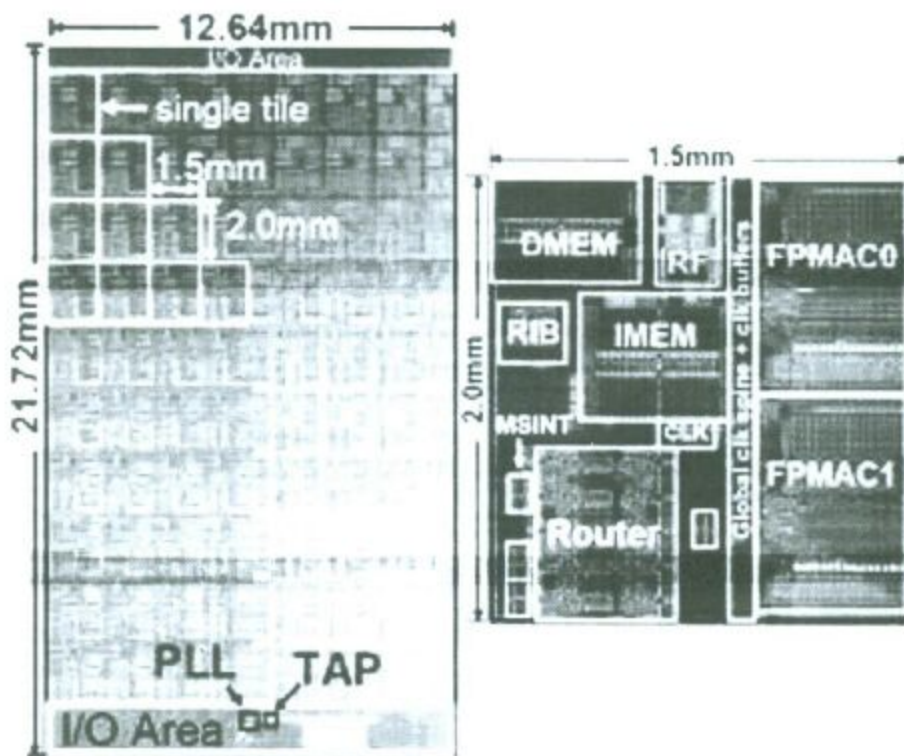


3rd Generation Intel® Core™ Processor: 22nm Process

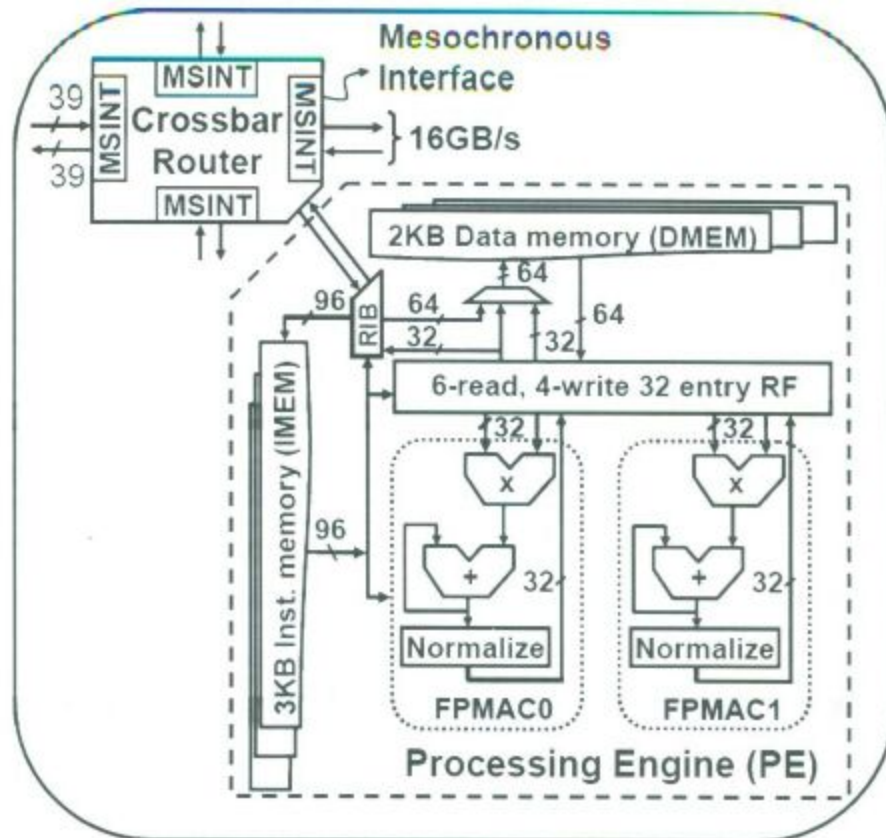


New architecture with shared cache delivering more performance and energy efficiency

Quad Core die with Intel® HD Graphics 4000 shown above
 Transistor count: 1.4Billion Die size: 160mm²
** Cache is shared across all 4 cores and processor graphics



51



Inside the SCC

24 Tiles
24 Routers
48 IA cores

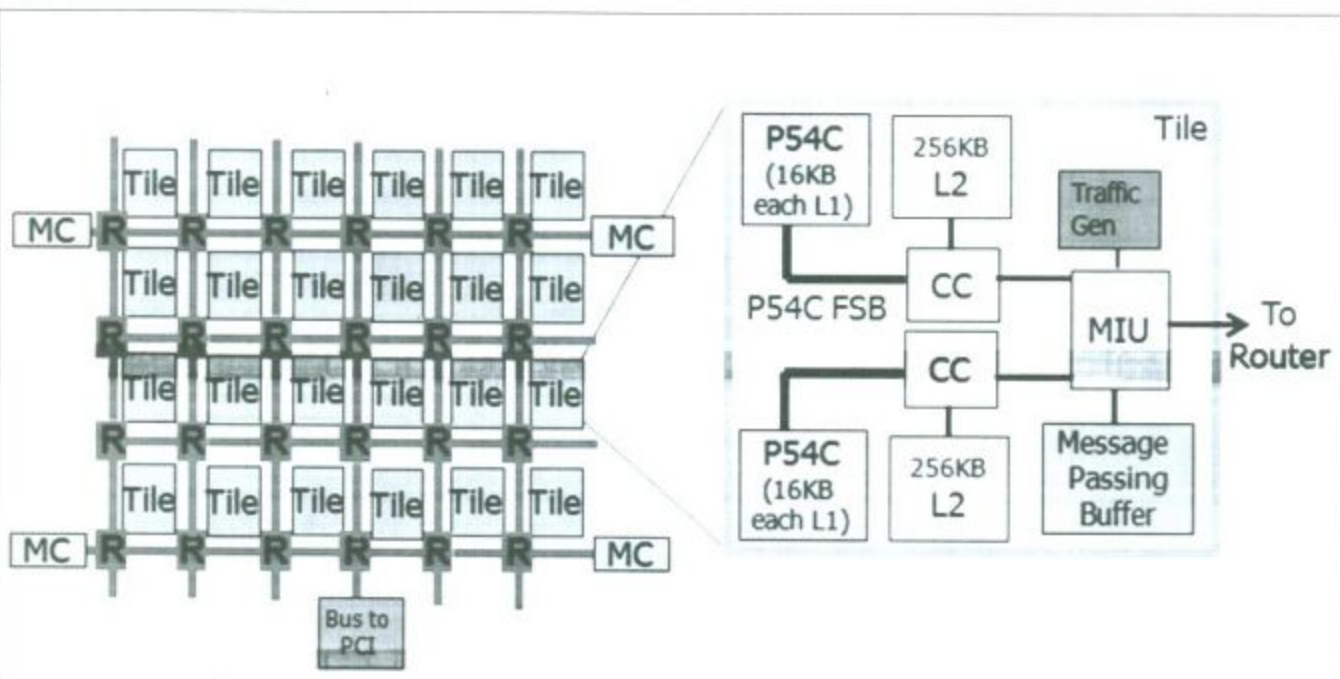
Dual-core SCDC Tile

L2 Cache	Core 1
ROUTER	ROUTER
L2 Cache	Core 2

1 TILE

- 2D mesh network with 256 GB/s bisection bandwidth
- 4 Integrated DDR3 memory controllers (64GB addressable)

dx



سیستم‌های ماتری پروسور و ماتری پردازش با هم تفاوتی ندارند اما در پیاده‌سازی کاملاً با یکدیگر متفاوت می‌باشند.

سوال: چرا به جای ماتری پروسورها از ماتری‌کرها استفاده می‌نماییم؟

- ۱- انتشار اطلاعات بزرگ بین پردازنده‌های ماتری پروسور باعث اتلاف توان می‌شود در صورتیکه بر روی chipها بروری وجود ندارد لذا مصرف توان نیز بسیار کم خواهد بود. و قس از ماتری پروسورها استفاده می‌کنیم که عناصری مانند المانهای منتری، op، تقسیم‌کننده جریان، نوزگیر و... قرار دارد در صورتیکه در chipها از اینها استفاده‌ای نمی‌گردد.
- ۲- تأثیر کارایی ماتری‌کرها بهتر از ماتری پروسورها می‌باشد زیرا بدلیل کاهش ارتباطات و هزینه در ماتری‌کرها کارایی مجتبر است.
- ۳- ماتری پروسورها دارای پیچیدگی بیشتری هستند و مدیریت آنها بدلیل استفاده از منابع، حافظه و هزینه کردن پردازنده دشوارتر است.
- ۴- تأخیر مربوط به سیم‌ها و ارتباطات در ماتری پروسورها بیشتر است در صورتیکه در ماتری‌کرها تأخیری وجود ندارد.
- ۵- تمایل به بهینه‌کردن اجزای دستورات مجزا زیرا در ماتری‌کرها موازی‌سازی در سطح عالی‌تری صورت می‌گیرد و جریان‌های دستورات را بهتر می‌توان به صورت موازی اجرا نمود.

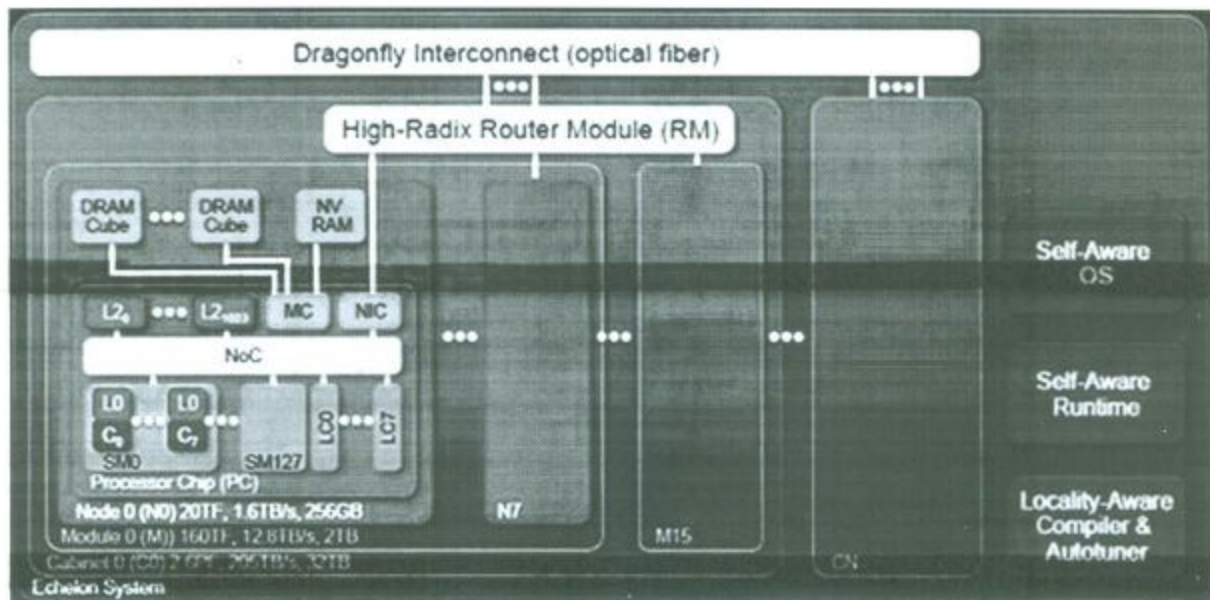


Figure 3: nVIDIA's Echelon Architecture

Why Move to Multicores

- Many issues with scaling a uncore
 - Power
 - Efficiency
 - Complexity
 - Wire Delay
 - Diminishing returns from optimizing a single instruction stream

همانگونه که مشاهده می شود سرعت انتقال اطلاعات در حالتی که حدود ۱۰ هزار برابر حالتی پررسان است. در حالتی که مقدار تقسیم های ارتباطی در درون chip بیشتر است لذا عملاً کلوگاه ارتباطی وجود ندارد و در واقع تراکم ارتباطات بیشتر بوده و به همین نسبت در یک میکرو، داده های بیشتری انتقال داده می شود.

همچنین در حالتی که مقدار Clock Rate به نسبت حالتی پررسان بیشتر است لذا در یک Clock Rate می توان اطلاعات بیشتری را منتقل نمود.

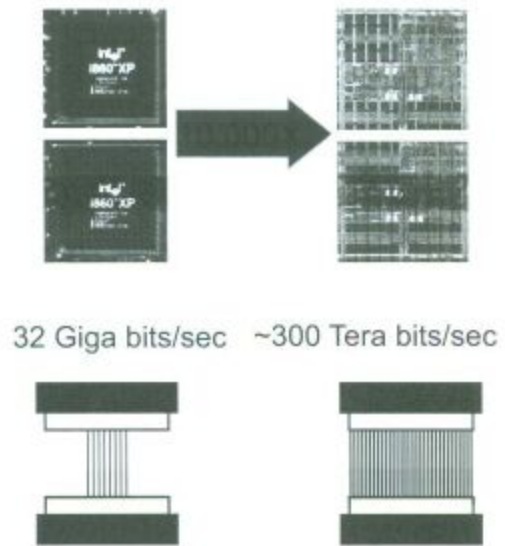
رسانایی مواردی که در حالتی که برای ارتباطات به کار می رود، دانسته بالاتری دارد. در حالتی پررسانها و متن رتیا می خواهد از پردازنده خارج شده و وارد بافر یا بورد گردد، این دورانی در استاندارد مختلف می باشد. همچنین وقت اطلاعات می خواهد از بورد واحد پردازنده گردد نیز دارای استانداردهای مختلفی می باشد. لذا در این حالت باید عمل تبدیل یا حالتی بلکی انجام شود که این کار باعث ایجاد تأخیر می گردد در صورتیکه در حالتی که چنین چیزی وجود ندارد.

سوال: با استفاده از حالتی که توانستیم سرعت انتقال داده را در پردازنده ها افزایش دهیم. حال این تغییر چه تأثیری بر روی برنامه نویسی ما دارد؟

با توجه به اینکه در حالتی که سرعت انتقال داده بین پردازنده ها بسیار زیاد می باشد لذا برای اینکه بتوان از حداکثر توان پردازنده ها استفاده نمود باید برنامه نویسی به گونه ای صورت گیرد که امکان اجرای موازی دستور را به طور کامل فراهم گردد در غیر این صورت نمی توان از حداکثر توان پردازشی پردازنده ها استفاده نمود. در حالتی که چون ارتباطات سریعتر است لذا نیاز به هنگام سازی و گسردن کردن سرعتی نیز داریم از طرف برنامه ها باید به صورت real time (رئال تایم) اجرا شوند که اینها جزو بحث های کنترلی هستند که باید رعایت شوند.

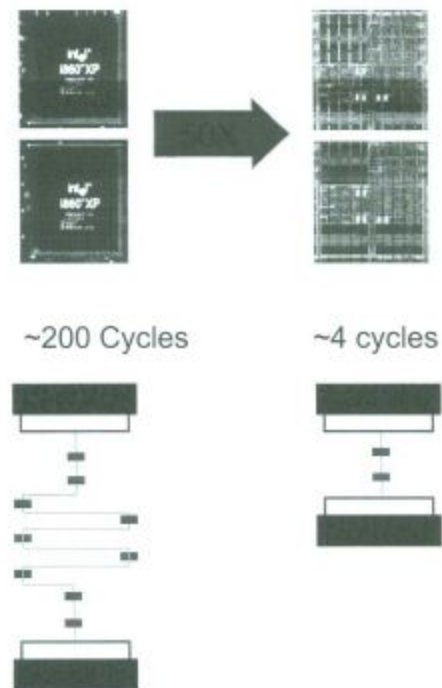
Impact of Multicores

- How much data can be communicated between two cores?
- What changed?
 - Number of Wires
 - I/O is the true bottleneck
 - On-chip wire density is very high
 - Clock rate
 - I/O is slower than on-chip
 - Multiplexing
 - No sharing of pins
- Impact on programming model?
 - Massive data exchange is possible
 - Data movement is not the bottleneck
 - Locality is not that important



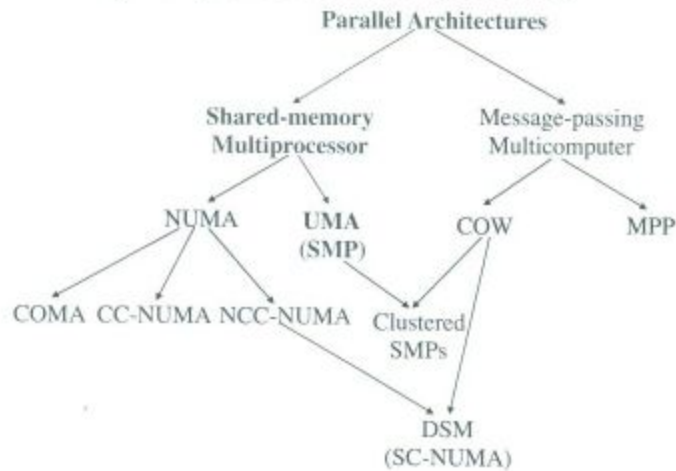
Impact of Multicores

- How long does it take for a round trip communication?
- What changed?
 - Length of wire
 - Very short wires are faster
 - Pipeline stages
 - No multiplexing
 - On-chip is much closer
- Impact on programming model?
 - Ultra-fast synchronization
 - Can run real-time apps on multiple cores

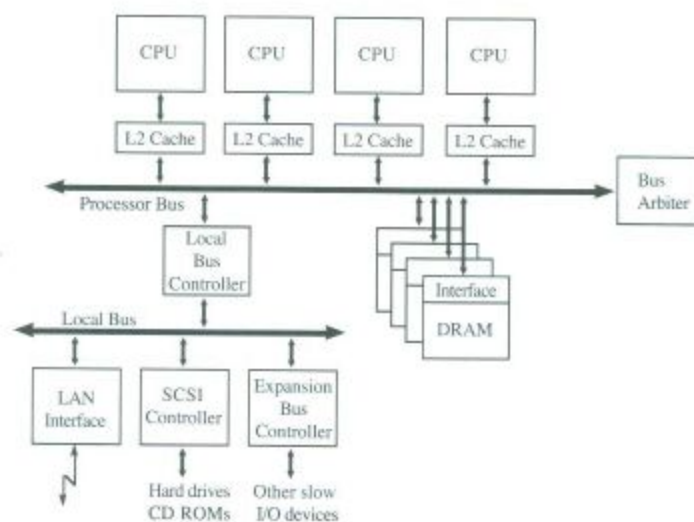


Parallel Computers

- “A parallel computer is a collection of processing elements that *cooperate* and *communicate* to solve large problems fast.”
- Taxonomy of parallel computers:



Shared-Memory Multiprocessors



Issues in Multiprocessors

- Programming
 - Need to explicitly define parallelism
- Hardware
 - Interconnection
 - Bus
 - Network
 - Synchronization
 - Test-and-set
 - Barrier synchronization
 - Fetch and add
 - Cache Coherence
 - Snoopy protocols
 - Directory-based

در حالتی پرسورها هم بصورت چند پردازنده ای و هم بصورت ماتریکس، یک سری نیازمندی ها وجود دارد در برنامه نویسی حتی امکان سفر می شود که اجرای دستورات به صورت موازی انجام شود. همچنین نیازمند یک سری تجهیزات سخت افزاری اصلی می باشیم که مثلاً شامل ارتباطات داخلی می باشد. این ارتباطات را می توان به صورت BUS یا Network ایجاد نمود که مثلاً در حالتی که ما به صورت Mesh پیاده سازی می شود در صورتیکه در حالتی پرسورها معمولاً به صورت BUS می باشد.

یک سری نیازمندی های مربوط به سنکرون کردن و همگام سازی نیاز می باشد که به چند روش می توان آن را انجام داد، همچنین وابستگی به حافظه Cache می باشد که باید آن را از بین برد که به دو طریق می توان این کار را انجام داد: ۱- استفاده از پروتکل های snoopy

۲- ارجاع با آدرس که به صورت Directory-based می باشد.

Shared-Memory Programming

- Many vendors have implemented their own proprietary versions of threads.
- A standardized C language threads programming interface, POSIX Threads or *Pthreads*
- Threaded applications offer potential performance gains and practical advantages:
 - Overlapping CPU work with I/O
 - Priority/real-time scheduling
 - Asynchronous event handling
 - Parallelization on SMPs
- Pthreads provide Over 60 routines for
 - Thread management - thread create, join, schedule, etc.
 - Mutexes - mutual exclusion
 - Conditional variables - provides communication between threads

برنامه نویسی حافظه مشترک

استفاده از زبانهای برنامه نویسی خاص مانند زبان C

(*P. threads.h*)
برنامه نویسی چندنخی

بعضی از این زبانهای برنامه نویسی چندنخی کارهایی از قبیل *overlap* کردن کارهای CPU و I/O را انجام میدهند یا بروی اجرای رستورات *real-time* اولویت ایجاد میکنند یا موازی سازی بروی SMP ها امکان پذیر میسازند .

5/9/2015 دیتا برای پردازنده ها مانند مواردی که محل می نمایند و تا زمانیکه این موارد اولی را در اختیار پردازنده دما قرار ندهیم، پردازنده هر چه قدر هم که توان پردازش بالایی داشته باشد نمی تواند کاره انجام دهد. حافظه ها مانند یک گلوگاه محل می نمایند و سرعت رسید حافظه ها به نسبت پردازنده ها خیلی کند تر صورت می گیرد.

Advanced Computer Architecture 5MD00 / 5Z033

حافظه های سلسله مراتبی

Memory Hierarchy & Caches

Henk Corporaal

www.ics.ele.tue.nl/~heco/courses/aca

h.corporaal@tue.nl

TUEindhoven

2007

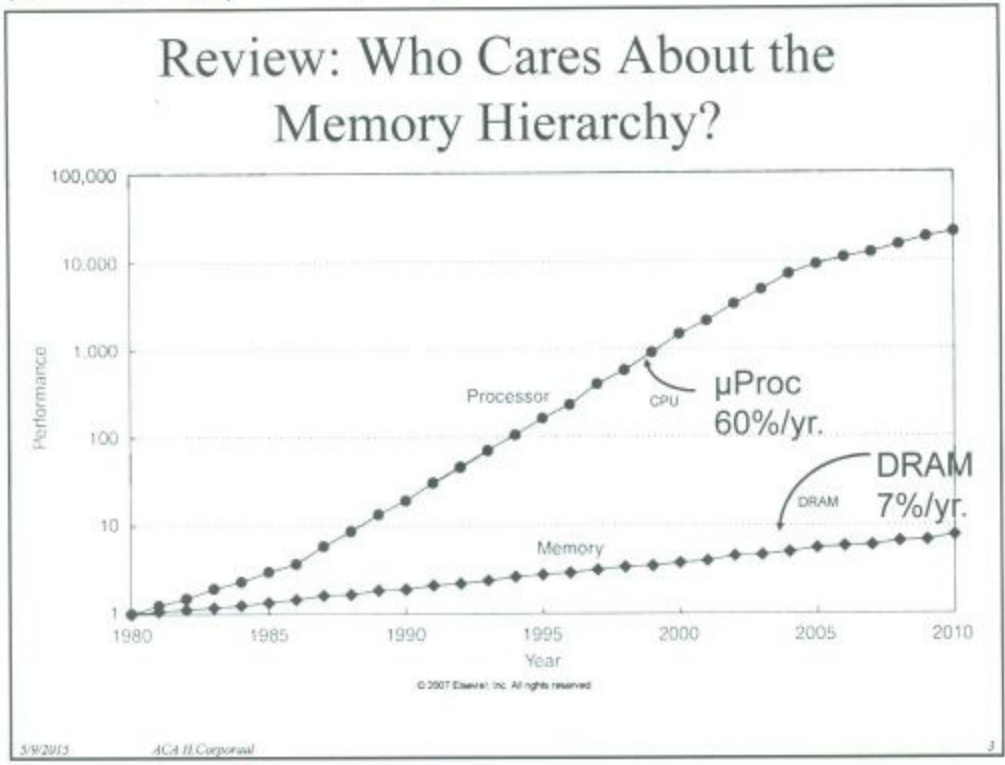
بخش از کندی حافظه ها را با استفاده از حافظه های سلسله مراتبی جبران می نمائیم.

Topics

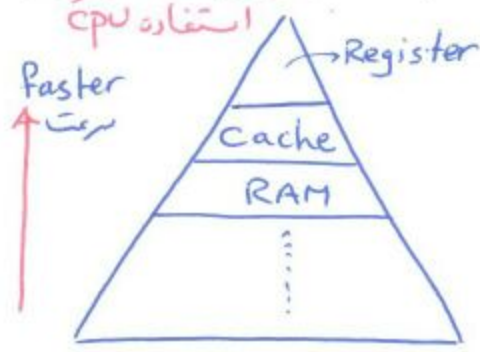
- Processor – Memory gap
- Recap of cache basics
- Basic cache optimizations
- Advanced cache optimizations
 - reduce miss penalty
 - reduce miss rate
 - reduce hit time

کارایی حافظه : نسبت ظرفیت حافظه به سرعت انتقال داده را کارایی حافظه گویند.
هر چه حافظه ای سرعت انتقال بیشتری داشته باشد و ظرفیت بیشتری را برای ما فراهم نماید عملاً دارای کارایی بالاتری خواهد بود. در عمل این دو پارامتر یعنی سرعت و کارایی در تضاد با یکدیگر هستند.
در هر ظرفیت حافظه بزرگتر باشد زمان جستجوی در حافظه طولانی تر خواهد بود. همچنین تکنولوژی ساخت حافظه ها نیز در مدت زمان انتقال داده تاثیر دارد.

هائیکه در شکل زیر مشخص شده است رشد افزایش سرعت پردازنده ها از سال ۱۹۸۰ تا ۵/۹/۲۰۱۵ سال ۲۰۱۰ به طور متوسط سالانه ۴۰٪ افزایش داشته است در صورتیکه این رشد برای سرعت حافظه ها فقط حدود ۷٪ بوده است. پس عملاً حافظه ها به صورت یک گلوگاه محمل فرغمانندند و هرچهقدر هم که سرعت پردازنده ما بالا باشد عملاً با وجود حافظه به کار آیی ۱۰۰٪ نخواهیم رسید.



حافظه ها سلسله مراتبی به لحاظ نزدیکی در استفاده CPU



ظرفیت Capacity
↓

دسترسی CPU به حافظه از طریق آدرس دهی صورت میگیرد. آدرس دهی توسط پردازنده انجام می شود. آدرس فیزیکی (واقعی) که در CPU استفاده می شود. آدرس منطقی (مجازی)

* بسته به ظرفیت حافظه به تعداد بیت جهت آدرس دهی نیاز داریم.

مثلاً $M: 2^x$ x - بیت

$M: 1 \text{ MByte} = 2^{20}$ 20 - بیت

Cache: 512KByte = 2^{19} 19 - بیت

مثلاً در این مثال مشاهده می شود که آدرس های حافظه اصلی ۲۰ بیت هستند در صورتیکه آدرس های Cache ۱۹-بیت می باشند

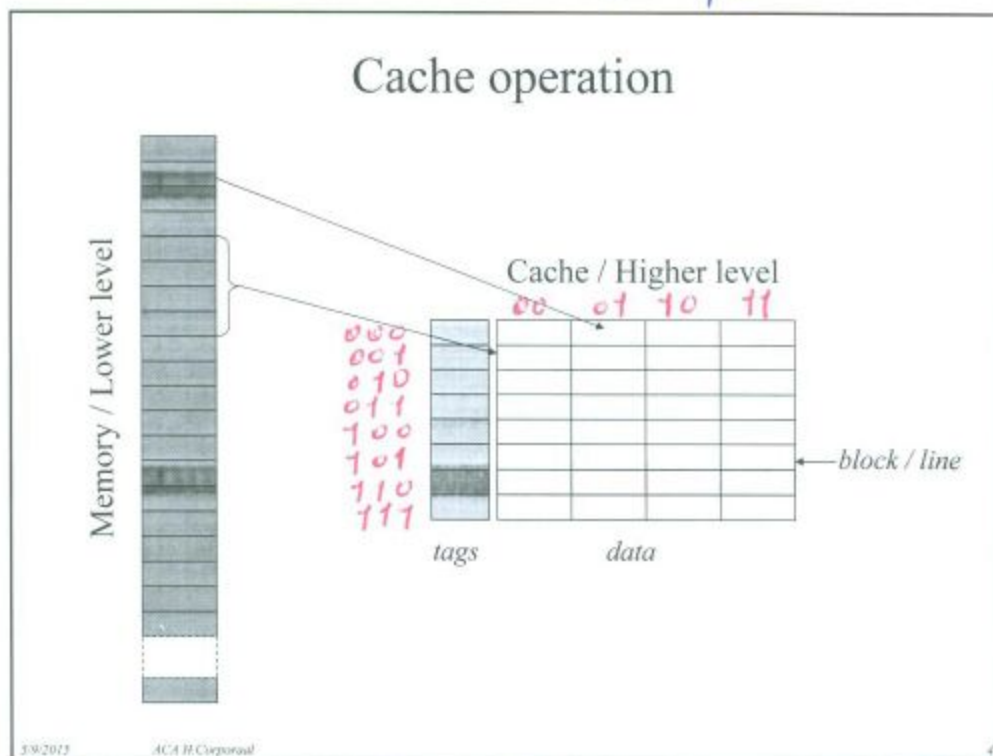
پردازنده آدرس های حافظه های اصلی را تولید می نماید منتظر ما برای دسترسی به آنها باید آنها را به Cache انتقال دهیم پس روشن باید وجود داشته باشد تا آدرس های حافظه اصلی را به آدرس های حافظه Cache تبدیل نماییم.

Direct Mapped

نیگاشت متقیم :

در نیگاشت متقیم، قسمت کم ارزش آدرس را در خانه های حافظه Cache جستجو می کنیم و قسمت با ارزش آدرس را در خانه های Tag قرار می دهیم. پس از جستجوی آدرس اگر اطلاعات مورد نظر در Cache بود که از آن استفاده می کنیم و اگر نبود باید آن را از حافظه اصلی وارد Cache کنیم.

5/9/2015



اگر حافظه Cache فضای مناسب خالی را برای انتقال داده از حافظه اصلی نداشته باشد باید تعدادی از خانه های آن را خالی نموده تا فضای برای اطلاعات خواسته شده فراهم گردد.

تاری داده ها یا محسب بودن آنها با استفاده از کیت های $X-OR$ انجام می شود که در صورت مساوی بیت ها صفر و در صورت مخالف بودن آنها یک می رسد.

hit : اگر اطلاعات مورد نظر در حافظه Cache باشد اصطلاحاً گفته می شود که hit رخ داده است. یا به عبارتی اگر اطلاعات خواسته شده با اطلاعاتی که در Cache وجود دارد یکی باشد در آن صورت یک سگینال hit می آید. در این حالت تلاش CPU برای دستیابی به اطلاعات موفقیت آمیز بوده است.

$\text{hit Rate یا hit Ratio} < 1$
نسبت بظهور داده موفق نرخ بظهور داده موفق

miss : زمانی که اطلاعات مورد نظر در حافظه Cache نباشد، miss رخ داده است.

$$\text{miss Rate یا miss Ratio} = 1 - \text{Hit Rate}$$

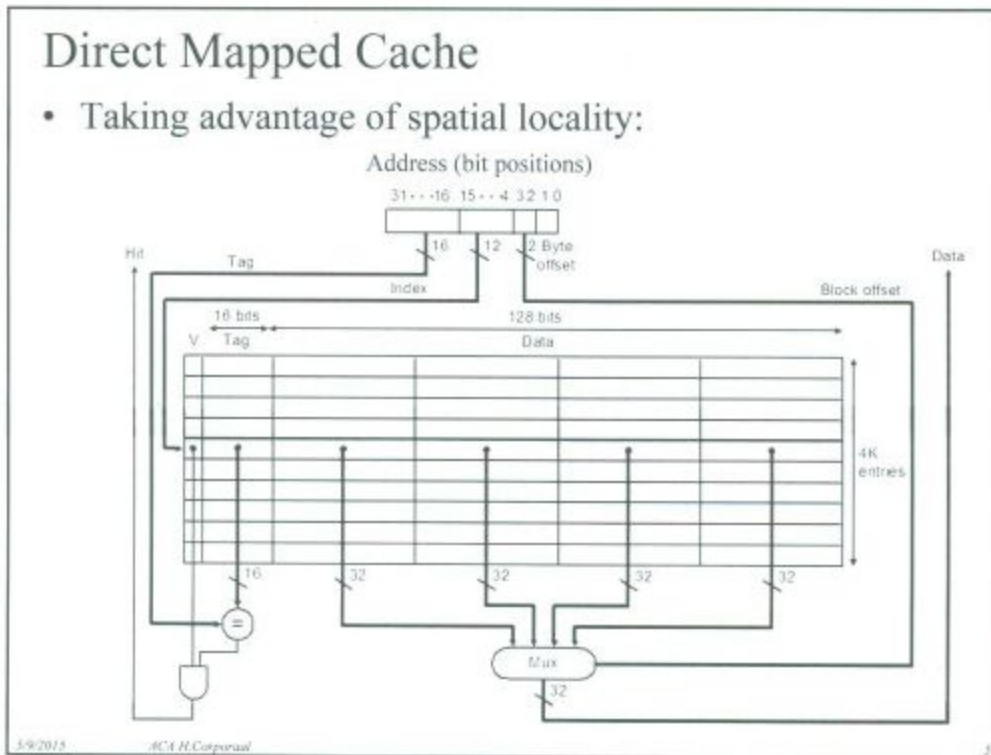
Set - Associative Cache

استفاده از حافظه های انجینی :

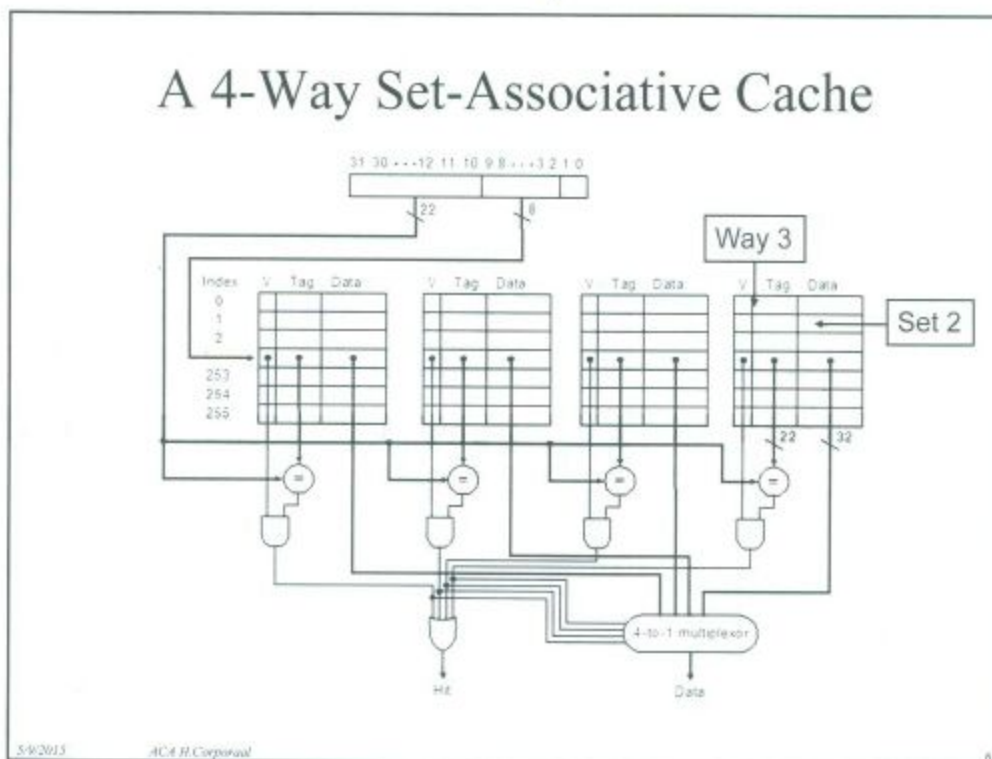
در این حالت استفاده از حافظه ها بر اساس بلوک بندی است و ما متدین است که از چند Cache موازی هم استفاده شده است. در حافظه های انجینی نیز از یک سگینال جهت مشخص نمودن شماره بلوک حافظه Cache جهت آدرس دهی استفاده می شود.

آدرس دهی انجمنی سریعتر از آدرس دهی بردش نکاست مستقیم است
 سمت افزار حافظه های انجمنی پیچیده تر است

5/9/2015



کلاً Cache ها یا به صورت نکاست مستقیم یا نکاست انجمنی آدرس دهی می شوند



* سرعت نکاست مستقیم نسبت به حافظه های انجمنی بالاتر است
 در حافظه های با ظرفیت Cache کم، نکاست مستقیم بهتر است در صورتیکه در حافظه های Cache با ظرفیت بالاتر، روش حافظه های انجمنی مناسب تر بوده و باعث افزایش کارایی می گردد
 توان مصرف در حافظه های انجمنی بیشتر است، حافظه های انجمنی اطلاعات را سریعتر جستجو می کنند

تمرکز موصف زبانی (حذف اطلاعات از Cache)
 در زمان حذف، اطلاعاتی را که قدیمی تر هستند و زمان طولانیتری از آنها استفاده نشده است، حذف می‌شوند.
 - تمرکز موصف مکانی (امضانه کردن اطلاعات به Cache)

6 basic cache optimizations (App. C)

- Reduces miss rate
 - Larger block size
 - Bigger cache
 - Higher associativity
 - reduces conflict rate
- Reduce miss penalty
 - Multi-level caches
 - Give priority to read misses over write misses
- Reduce hit time
 - Avoid address translation during indexing of the cache

5/9/2015

ACI H Corporation

7

چگونه می‌توان کارایی یک حافظه Cache را بهبود کرد؟

۱- کاهش miss rate یا درخواست‌های ناموفق

- بزرگتر کردن اندازه بلوک‌های حافظه که در این حالت نرخ miss به همان میزان کاهش می‌یابد زیرا در این حالت احتمال آنکه به خانه‌های متعلق شده نیاز باشد، زیاد است. البته باید توجه نمود که اندازه این بلوک‌های حافظه را نمی‌توان خیلی بزرگ انتخاب نمود یعنی مثلاً نمی‌توان اندازه بلوک را برابر حافظه Cache یا نصف آن در نظر گرفت و این کار را معمولاً به روش‌های آماری انجام می‌دهند.
- بزرگتر کردن حافظه Cache زیرا مسلم است هر چه اندازه Cache بزرگتر باشد اطلاعات بیشتری را می‌تواند در آن قرار دهد و به همان نسبت مقدار miss rate نیز کاهش خواهد یافت.
- استفاده از حافظه‌های انحصاری بزرگتر

۲- کاهش جریمه یا هزینه مربوط به درخواست‌های ناموفق

- استفاده از حافظه‌های Cache چندسطحی یا لول‌های مختلف
- ایجاد اولویت بر خواندن اطلاعات در زمان نوشتن اطلاعات

۳- کاهش زمان درخواست‌های موفق یا hit time

- یعنی جستجو در Cache سریع‌تر انجام شود که برای این منظور از روش‌های استفاده می‌شود که آدرس‌دهی سریع‌تری را ایجاد نماید.

11 Advanced Cache Optimizations (&5.2)

- | | |
|--|---|
| <ul style="list-style-type: none"> • Reducing hit time 1. Small and simple caches 2. Way prediction 3. Trace caches | <ul style="list-style-type: none"> • Reducing Miss Penalty 7. Critical word first 8. Merging write buffers |
| <ul style="list-style-type: none"> • Increasing cache bandwidth 4. Pipelined caches 5. Multibanked caches 6. Nonblocking caches | <ul style="list-style-type: none"> • Reducing Miss Rate 9. Compiler optimizations |
| | <ul style="list-style-type: none"> • Reducing miss penalty or miss rate via parallelism 10. Hardware prefetching 11. Compiler prefetching |

5/9/2015

ACA H. Corporal

8

کاهش زمان Hit time

- 1- استفاده از Cache های کوچک و ساده البته باید توجه نمود که استفاده کردن از cache های کوچکتر باعث آدرس دهی سریعتر شده ولی miss ها را نیز افزایش خواهد داد . در هر صورت همزمان نمی توان کمه شرایط را ایده آل نمود .
- 2- راههای پیش بینی کردن اطلاعات یعنی بتوان پیش بینی نمود که بلائی را که به cache منتقل می نمایم در دستورات بعدی هم مورد استفاده قرار گیرد .
- 3- Trace کردن Cache که باعث بالا رفتن سرعت جستجوی مورد البته تو مصرفی را هم افزایش می دهد .

افزایش پهنای باند Cache

- 4- ورود و خروج اطلاعات را به صورت pipeline انجام دهیم .
- 5- استفاده از Cache های چند بانگی یعنی چندین فیلد Tag داشته باشند .
- 6- از Cache های بزرگ نشده استفاده شود زیرا در استفاده از Cache به صورت بلاک درست است که ممکن است اطلاعات منتقل شده مورد استفاده قرار گیرد ولی این احتمال نیز وجود دارد که اطلاعات اضافی منتقل شده مورد استفاده قرار نگیرد .

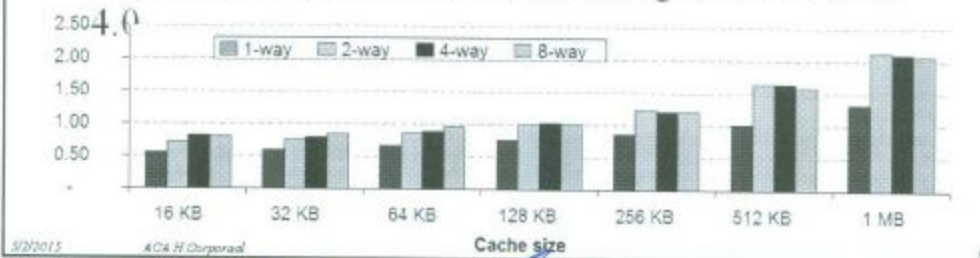
کاهش جریدهای مربوط به درخواست های ناموفق

- 7- ابتدا داده های بحرانی یا Critical را به Cache منتقل نمایم
- 8- استفاده از بافرین Cache و حافظه اصلی زیرا با فریب صورت کب حافظه میانی عمل کرده و نیاز به جستجو هم ندارد .

کاهش miss rate
 ۹- با بهینه سازی کامپایلر می توان دستورالعمل را به صورت موازی اجرا نمود و miss rate را کاهش داد
 کاهش هزینه های مربوطه در خواست های ناموفق با استفاده از موازی سازی
 ۱۰- پیش واکش سخت افزاری مثلاً با استفاده از پردازنده های آرایه ای می توان برای دو عدد، چهار محل جمع، تفریق، منب و تقسیم را انجام داد پس یکی از آنها را که مورد نیاز است استفاده نمود.

1. Fast Hit via Small and Simple Caches

- Index tag memory and then compare takes time
- ⇒ Small cache is faster
 - Also L2 cache small enough to fit on chip with the processor avoids time penalty of going off chip
- Simple ⇒ direct mapping
 - Can overlap tag check with data transmission since no choice
- Access time estimate for 90 nm using CACTI model



۱۱- با استفاده از کامپایلر می توان پیش واکش نرم افزاری داشته باشیم

انجام نگاشت مستقیم،
 انجام جستجو و مقایسه
 اطلاعات را سریعتر
 می نماید.

2. Fast Hit via Way Prediction

- Make set-associative caches faster
- Keep extra bits in cache to predict the "way," or block within the set, of next cache access.
 - Multiplexor is set early to select desired block, only 1 tag comparison performed
 - Miss ⇒ 1st check other blocks for matches in next clock cycle
- Accuracy ≈ 85%
- Drawback: CPU pipeline is hard if hit takes 1 or 2 cycles
 - Used for instruction caches vs. L1 data caches
 - Also used on MIPS R10K for off-chip L2 unified cache, way-prediction table on-chip

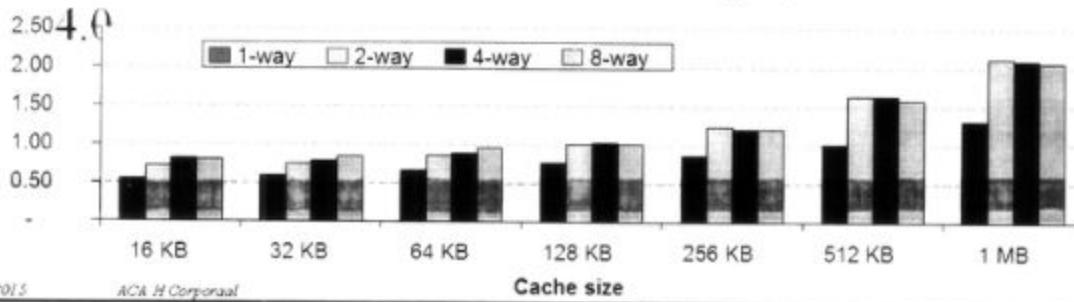


هرچه اندازه Cache کوچکتر باشد تعداد tag های که باید محتوای آنها با قیمت پردازش آدرس مقایسه شود، کمتر خواهد بود و لذا سریعتر این کار انجام می شود و نیز سریعتر مشخص می شود که اطلاعات مورد نظر در Cache قرار دارند یا خیر
 استفاده از حافظه های انحصاری و راههای پیش بینی می تواند باعث سریعتر شدن زمان Hit گردد.

کاهش miss rate با بهینه سازی حمایتی تر دستورات را به صورت موازی اجرا نمود و miss rate را کاهش داد
 5/20/2015
 کاهش جریمه های مربوط به درخواست های ناموفق با استفاده از موازی سازی
 ۱۰ - پیش واکش نسبت افزایشی مثلاً با استفاده از پردازنده های آرایه ای می توان برای دو عدد، چهار عمل جمع، تفریق، ضرب و تقسیم را انجام داد و سپس یکی از آنها را که مورد نیاز است استفاده نمود.

1. Fast Hit via Small and Simple Caches

- Index tag memory and then compare takes time
- ⇒ Small cache is faster
 - Also L2 cache small enough to fit on chip with the processor avoids time penalty of going off chip
- Simple ⇒ direct mapping
 - Can overlap tag check with data transmission since no choice
- Access time estimate for 90 nm using CACTI model



۱۱ - با استفاده از کامپایلر می توان پیش واکش نرم افزار را دسته بندی کرد.

* انجام نگاشت متتیم، انجام جستجو و مقایسه اطلاعات را سریعتر می نماید.

به طور کلی برای آنکه بتوان کارایی یک Cache را به یونیت خود می توان کارهای زیر را انجام داد:

- ۱- کاهش زمان Hit که زمان لازم برای جستجوی داده ای که در حافظه Cache قرار دارد می باشد که اگر بتوانیم این زمان را کاهش دهیم در واقع کارایی Cache خود را می توانیم بالا ببریم زیرا یکی از دلایل افزایش کارایی، سرعت دسترسی به اطلاعات است.
- ۲- افزایش پهنای باند حافظه Cache یعنی بتواند داده های مفیدی را در حافظه Cache جای داد.
- ۳- کاهش جریمه مربوط به درخواست های ناموفق که باعث پرداخت هزینه می شود یا مجبوریم زمانی را تحمل نمائیم تا داده از حافظه اصلی به حافظه Cache منتقل پیدا کرده و سپس در اختیار پردازنده قرار گیرد. لذا اگر بتوانیم این زمان را کاهش دهیم در واقع عملکرد Cache خود را افزایش داده ایم.
- ۴- کاهش نرخ miss یعنی داده هایی را که به احتمال زیاد پردازنده برای اجرای دستور به آنها نیاز دارد را از قبل در حافظه Cache آماده کنیم.
- ۵- جریمه مربوط به درخواست های ناموفق یا نرخ miss را با استفاده از موازی سازی کاهش دهیم.

2. Fast Hit via Way Prediction

- Make set-associative caches faster
- Keep extra bits in cache to predict the "way," or block within the set, of next cache access.
 - Multiplexor is set early to select desired block, only 1 tag comparison performed
 - Miss \Rightarrow 1st check other blocks for matches in next clock cycle
- Accuracy $\approx 85\%$
- Drawback: CPU pipeline is hard if hit takes 1 or 2 cycles
 - Used for instruction caches vs. L1 data caches
 - Also used on MIPS R10K for off-chip L2 unified cache, way-prediction table on-chip



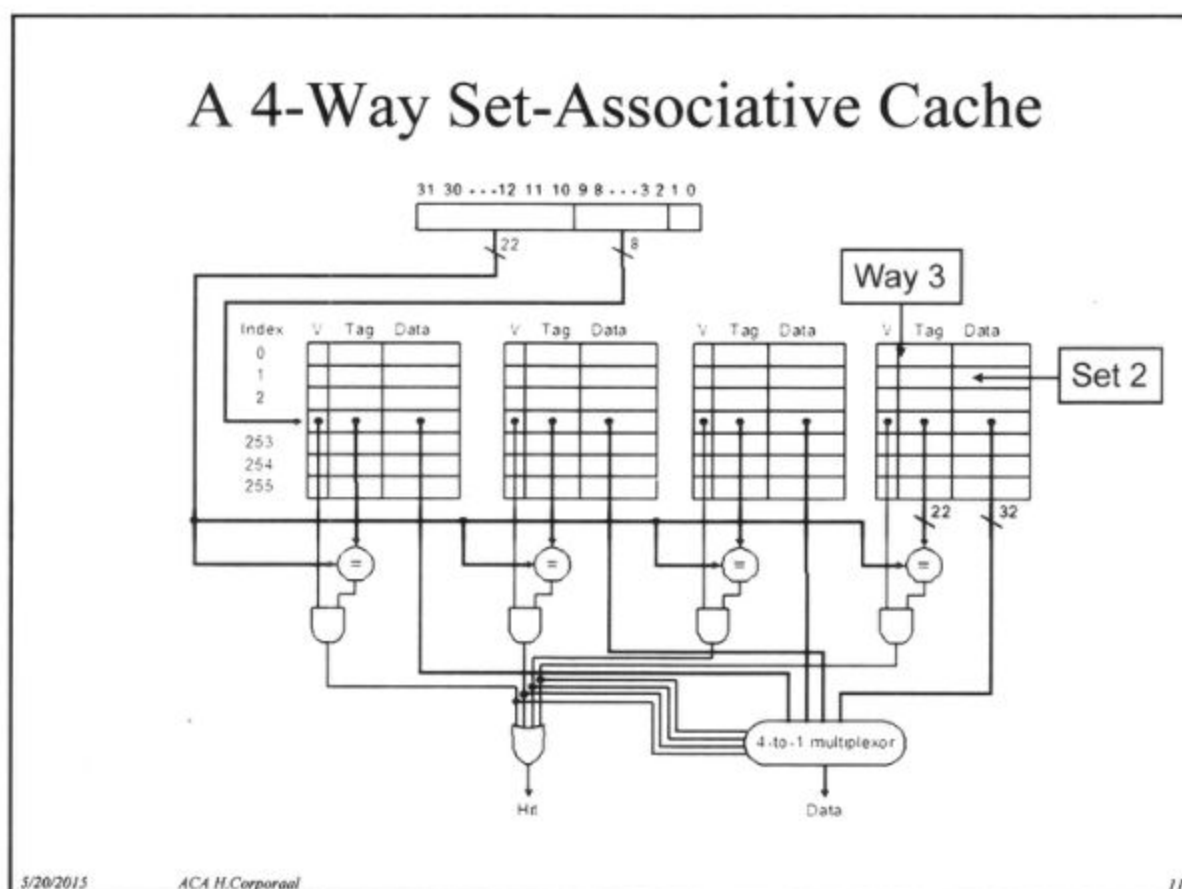
5/20/2015

ACA H. Corporaal

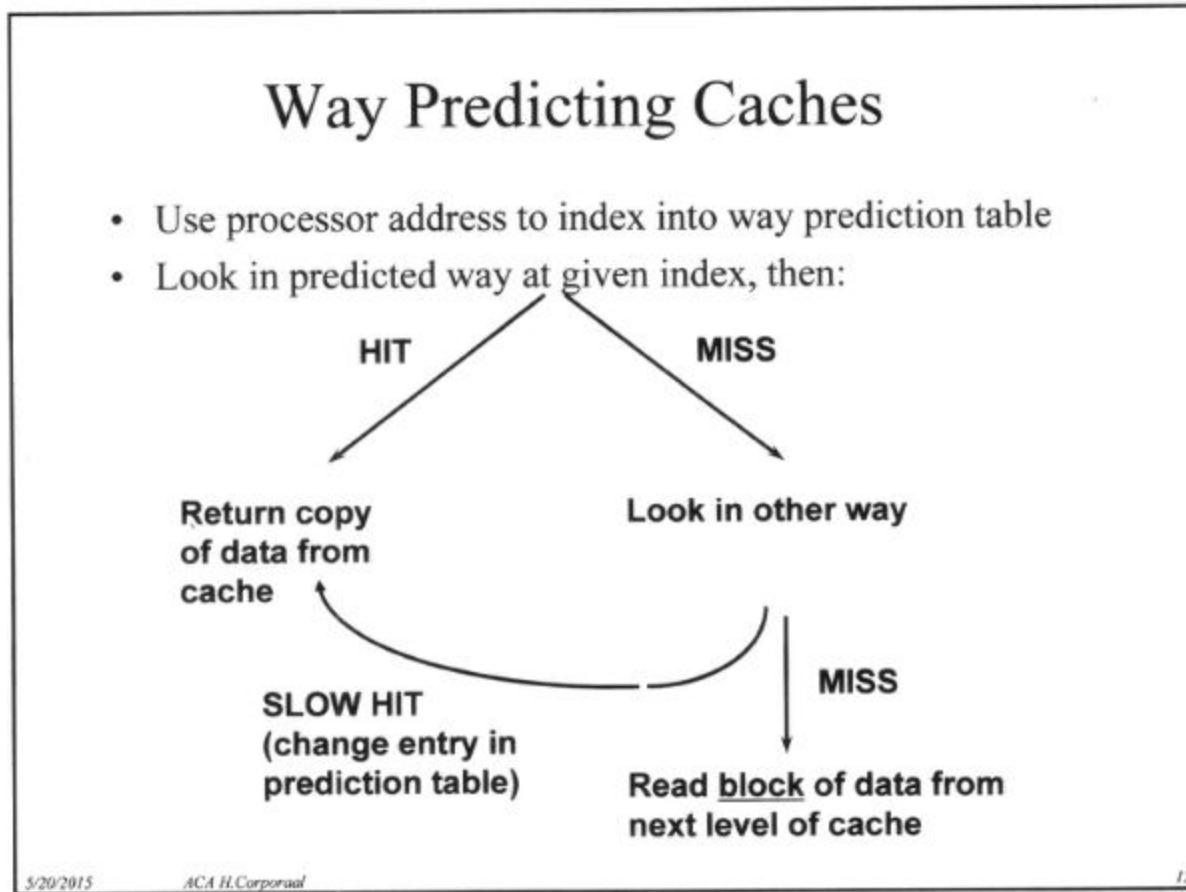
10

هرچه اندازه Cache کوچکتر باشد تعداد Tag های که باید محتوای آنها با قسمت پردازش آدرس مقایسه شود، کمتر خواهد بود و لذا سرعت این کار انجام می شود و نیز سرعت مشخص می شود که اطلاعات مورد نظر در Cache قرار دارند یا خیر

استفاده از حافظه های انجمنی و راههای پیش بینی می تواند باعث سرعت شدن زمان Hit گردد



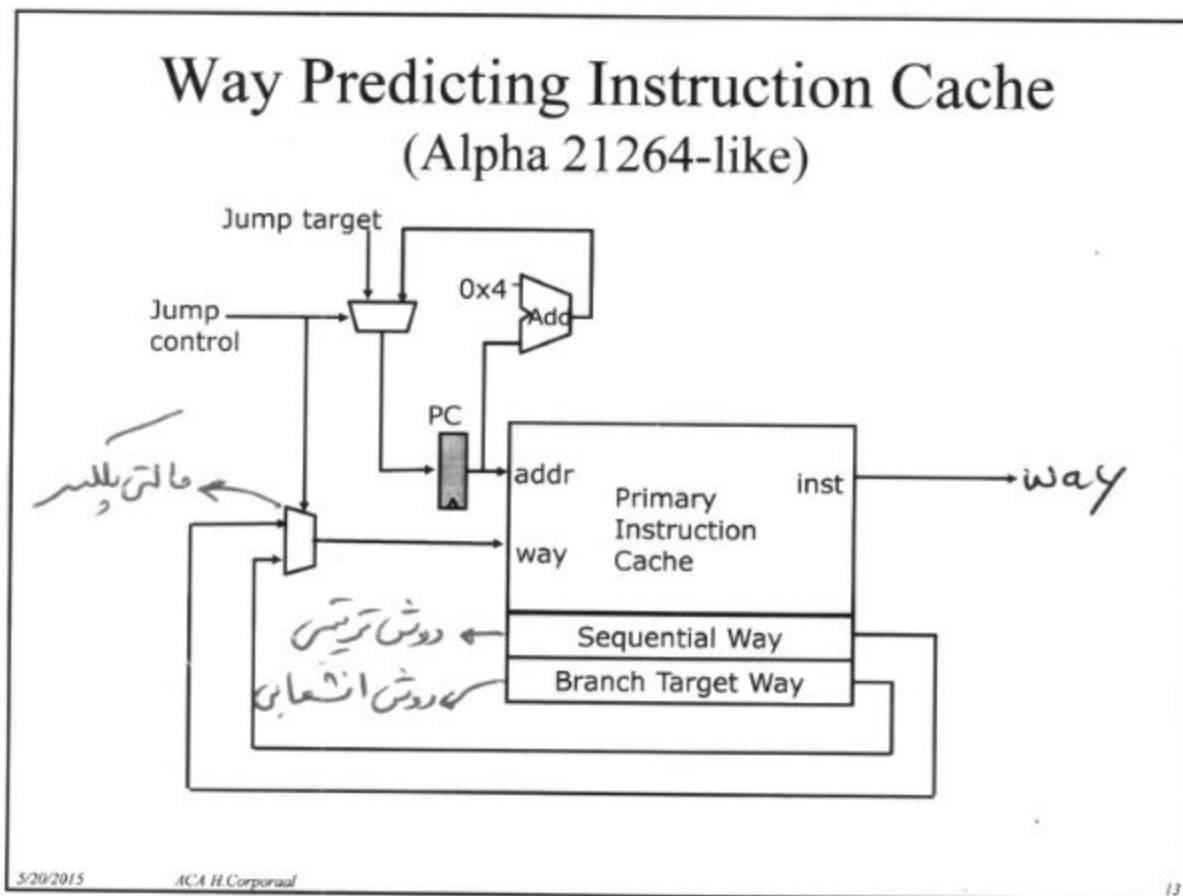
همانگونه که گفتیم سداگر بخواهیم Hit های خود را افزایش دهیم می توانیم از روش های پیش بینی
 نیز روش های آدرس دهی استفاده کنیم که برای این منظور می توان از Cache های انجمن استفاده
 نمود که هر یک از آنها به صورت یک Cache نگاشت مستقیم عمل کند که جستجوی داده در آنها
 به صورت موازی انجام شده و در نتیجه سرعت Hit یا جستجوی داده در حافظه افزایش می یابد



وقتی به دنبال داده‌ای در Cache می‌گردیم دو حالت اتفاق می‌افتد:

۱ - داده مورد نظر در Cache قرار دارد و به عبارتی Hit رخ می‌دهد که در این حالت یک کپی از داده‌ها تهیه شده و به پردازنده ارسال می‌گردد.

۲ - Miss رخ می‌دهد که در این حالت باید برای یافتن داده مورد نظر، در سطح‌های بعدی حافظه جستجو کنیم مثلاً اگر اطلاعات مورد نظر در Cache نبود باید به دنبال آن در RAM بگردیم و اگر آنجا نبود باید در حافظه جانبی جستجو کنیم و به همین ترتیب در لول‌های بعدی حافظه جستجو کنیم.

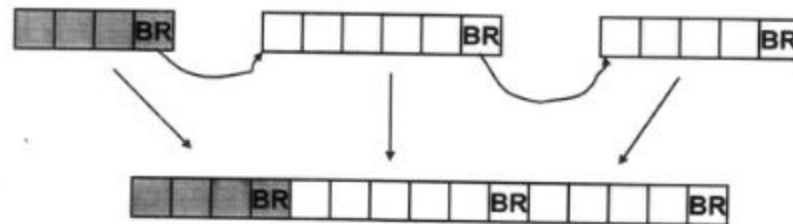


جستجو در حافظه به دور روش ترتیبی و انشعابی انجام می‌شود و در سطح حافظه باید به این
 از این دور روش عمل جستجو را انجام دهیم. در بلوک ریگرام فوق روش برای پیش بینی حافظه‌های
Cache مربوط به دستورالعمل‌ها نشان داده شده است. توجه شود که حافظه‌های مربوط به دارن
 و دستورالعمل‌ها با هم متفاوت می‌باشند و برای درک از آنها حافظه‌های سلسله مراتبی جداگانه

داریم

3. Fast (Inst. Cache) Hit via Trace Cache

Key Idea: Pack multiple non-contiguous basic blocks into one contiguous trace cache line



- Single fetch brings in multiple basic blocks
- Trace cache indexed by start address and next n branch predictions

5/20/2015

ACA H. Corporaal

14

روش برای افزایش سرعت Hit یا کاهش زمان Hit ، استفاده از Trace کردن است .
 Trace کردن یعنی جستجو کردن موردی یعنی جستجو را به صورت پیش از روی بعضی خانه های
 حافظه انجام دهیم . در روش معمولی برای یافتن داده در حافظه ، جستجو در حافظه به صورت
 خانه به خانه صورت می گیرد اما در روش Trace کردن با آدرس که به هر قسمت حافظه مربوط به یک
 دستور العمل داده می شود می توان از روی بعضی از خانه های حافظه پیش نموده و لذا دستور العمل
 مورد نظر را سریعتر پیدا نمود زیرا در اکثر اوقات اطلاعات مربوط به یک دستور العمل بدلیل حجم زیاد به
 طور پشت سر هم قرار ندارند و هر قسمتی آن ممکن است در یک قسمت از حافظه ذخیره شود . در این حالت
 در انتهای هر قسمت ، آدرس شروع قسمتی بعدی را قرار می دهیم . اگر بخواهیم از Cache بدین صورت
 استفاده نماییم می توانیم از اشاره گر استفاده کنیم بدین صورت که اگر اطلاعاتی را که نیاز داریم در یک
 قسمت نبود ، در آن صورت باید برویم مارا به قسمتی بعدی راهنمایی کند . در این حالت لازم نیست
 تمام بلاک های حافظه برای یافتن داده مورد نظر جستجو شود یعنی اگر آدرس قسمتی بعدی در بلاک
 بعدی نباشد از روی آن پیش نموده و جستجو را در بلاک مورد اشاره انجام می دهیم . لذا استفاده از
 Trace کردن باعث افزایش سرعت جستجو می گردد .

3. Fast Hit times via Trace Cache

- Trace cache in Pentium 4
 1. Dynamic instr. traces cached (in level 1 cache)
 2. Cache the micro-ops vs. x86 instructions
 - Decode/translate from x86 to micro-ops on trace cache miss
- + ⇒ better utilize long blocks (don't exit in middle of block, don't enter at label in middle of block)
- ⇒ complicated address mapping since addresses no longer aligned to power-of-2 multiples of word size
- ⇒ instructions may appear multiple times in multiple dynamic traces due to different branch outcomes

5/20/2015

ACA H. Corporaal

15

استفاده از Trace کردن اولین بار در Pentium 4 مورد استفاده قرار گرفت .

نکات مثبت در Trace کردن :

هرچه بلوک های داده بزرگتر باشد بهره وری بیشتر خواهد بود . توجه شود که در روش Trace کردن وقتی در یک بلوک موفول جستجو هستیم امکان اینکه از وسط بلوک عمل پرش را انجام داد امکان ندارد و لازم است جستجوی بلوک تا انتها انجام شود یعنی یک بلوک را کامل جستجو نموده و سپس بر اساس اطلاعات و آدرس که از این بلوک بدست می آید ، به بلوک دیگر پرش می نماید .

نکات منفی در Trace کردن :

— هرچه نگاشت آدرس دهی پیچیده تر باشد ، Hit ها طولانی تر خواهد بود زیرا هرچه تعداد بخش های آدرس یا تعداد مؤلفه های آدرس کمتر باشد ، Hit سریعتری خواهیم داشت زیرا در آدرس دهی پیچیده نیاز به تعداد مالتی پلکس های بیشتری داریم که این امر باعث تأخیر بیشتر شده و زمان جستجوی طولانی تر خواهد بود .

— اگر Cache مورد استفاده برای ذخیره دستورالعمل باشد ، یک دستورالعمل مشخص می تواند در زمانهای متفاوت در بخش های مختلف برنامه مورد استفاده قرار گیرد . در روش Trace کردن

چون مجبور به استفاده از روش ترتیبی هستیم و چون معیار Trace کردن محضاً در دستورالعملها است

1 همیشه دارای معیار ثابت و مشخصی نمی باشد لذا برای دستورالعمل ها روش Trace کردن روش مناسبی برای کاهش Hit-time نمی باشد در صورتیکه برای Cache های که مربوط به دیتا هستند مناسب می باشد .

4: Increasing Cache Bandwidth by Pipelining

- Pipeline cache access to maintain bandwidth, but higher latency
- Nr. of Instruction cache access pipeline stages:
 - 1: Pentium
 - 2: Pentium Pro through Pentium III
 - 4: Pentium 4
- ⇒ greater penalty on mispredicted branches
- ⇒ more clock cycles between the issue of the load and the use of the data

5/20/2015

ACA H. Corporaal

16

یکی دیگر از روش‌های افزایش بهره‌وری Cache‌های پیشرفته، افزایش پهنای باند با استفاده از pipeline می‌باشد. همانطور که گفته شد می‌توان حافظه‌های چندسطحی را که به صورت پشت سرهم قرار می‌دهیم مانند ماشین‌هایی که به صورت pipeline عمل می‌کنند، در نظر بگیریم.

استفاده از روش pipeline کردن در Cache‌هایی که برای حافظه‌های دیتایم می‌باشد، روش مناسبی نیست و برعکس این روش برای Cache‌های مربوط به دستورالعمل‌ها مناسب می‌باشد.

برای هر عملی که بتوانیم آن را با آبی تعریف نمود. مثلاً برای پردازنده، کلاسی عبارت است از مقدار انجام دستور العمل در واحد زمان

$$\text{cpu time} = \left(\begin{array}{l} \text{مدت زمان اجرای دستور} \\ \text{توسط پردازنده} \end{array} + \begin{array}{l} \text{مدت زمان واکنش دستور} \\ \text{از حافظه} \end{array} \right)$$

یا

$$\text{cpu time} = \left(\begin{array}{l} \text{مقدار کلاک های مورد} \\ \text{نیاز برای اجرای دستور} \end{array} + \begin{array}{l} \text{مقدار کلاک های مورد} \\ \text{نیاز برای واکنش دستور} \\ \text{از حافظه} \end{array} \right) * \begin{array}{l} \text{دسترز نام کلاک} \\ \text{clock} \\ \text{cycle} \end{array}$$

خواندن از حافظه
Read

نوشتن در حافظه
Write

لازم است در یک کامپیوتر مقدار کلاک های مورد نیاز برای خواندن و نوشتن با هم برابر نباشند و مقدار آنها کاملاً متفاوت می باشد. لذا در مقدار کلاک های مورد نیاز برای دسترسی به حافظه، باید هم کلاک های مربوط به خواندن و هم کلاک های مربوط به نوشتن را در نظر بگیریم که در واقع مجموع آنها در نظر گرفته می شود.

$$\text{مقدار کلاک های مورد نیاز برای دسترسی به حافظه} = (\text{Read clocks} + \text{write clocks})$$

کارایی یک حافظه باید در اجرای یک برنامه مورد بررسی قرار گیرد نه یک دستور لذا اگر دو حافظه مقدار کلاک خواندن و نوشتن یکسان داشته باشند، هر یک از آنها که دارای حجم حافظه بیشتری است دارای کارایی بالاتری خواهد بود زیرا مقدار Miss های که رخ بدهد کمتر است زیرا مقدار دستورات بیشتری از برنامه که مورد نیاز است در حافظه قرار دارد.

برنامه $P_i \Rightarrow n$ دستور
مقدار دستورات ارجاع در حافظه m
تعداد دستورات ارجاع در حافظه خوانده n
از حافظه

$$\text{Read clock} = \frac{m_r}{n} \times (\text{Read hit Rate clock}) \times (\text{Read Miss Rate Penalty})$$

مقدار دستورات ارجاع در حافظه خوانده \leftarrow $\frac{m_r}{n}$
تعداد کلاک های مورد نیاز زمانی که Miss رخ بدهد \leftarrow (Read Miss Rate Penalty)

$$\text{write clock} = \frac{m_w}{n} \times (\text{write hit Rate clock}) \times (\text{write Miss Rate Penalty}) + \text{write Buffer clock}$$

۷۴

باید توجه نمود که بین خواندن و نوشتن در حافظه تفاوتی وجود دارد و آن این است که وقتی پردازنده ای نیاز به خواندن از حافظه دارد، آدرس آن تولید شده و جستجویی گردد و با توجه به کلاک های hit و یا در صورت ایجاد miss، یک مقدار کلاکی ایجاد می شود و در نهایت داده مورد نظر پیدا شده و توسط پردازنده استفاده می گردد که در این حالت پردازنده منتظر خواندن می ماند تا در ادامه، اجرای دستورات خودش را کامل کند یعنی در زمان خواندن تا پردازنده چیزی دریافت نکند کاری انجام نمی دهد پس منتظر می ماند تا عمل خواندن انجام شود در صورتیکه در زمان نوشتن، پردازنده منتظر نمی ماند بلکه اطلاعاتی که قرار است در حافظه نوشته شود را در درون بافر قرار می دهد لذا مدت زمانی برای نوشتن داده ها در بافر نیاز می باشد و چون این عمل یک بار نیاز است لذا تأثیر

آن در رابطه به صورت جمع قرار داده می شود.

AMT : Access Memory Time میانگین مدت زمان دسترسی به حافظه

در حافظه ای که دارای یک مدت زمان hit و یک نرخ miss می باشد، میانگین مدت زمان دسترسی به حافظه (AMT) به صورت زیر بدست می آید. در این حالت صحبتی از کارایی در میان نمی باشد.

$$AMT = \underbrace{hit\ time}_{\substack{\text{مدت زمان دسترسی} \\ \text{موفق به حافظه}}} + \left(\underbrace{(miss\ rate)}_{\substack{\text{زمان دسترسی به حافظه} \\ \text{یا وقوع miss}}} \times \right)$$

توجه شود برای هر miss حتماً یک hit رخ می دهد یعنی ابتدا باید حافظه جستجو شود تا مشخص گردد که داده مورد نظر در حافظه وجود ندارد و سپس بر روی مربوطه جستجوی آن در حافظه های سطح های پایین تر بررسی گردد.

حافظه های Cache حافظه های محدودی هستند و همین محدودیت باعث ایجاد دورخیز hit و miss می شود. زمانیکه miss رخ می دهد یعنی داده های مورد نظر ما در حافظه قرار ندارند و باید به یک روش این اطلاعات از حافظه های سطح های بعدی به حافظه Cache منتقل شود. از طرفی در این انتقال ممکن است ظرفیت حافظه Cache تکمیل باشد پس مجبوریم ابتدا داده ای را از حافظه Cache حذف نموده و آن را به حافظه اصلی انتقال دهیم و به جای آن داده مورد نظر را از حافظه اصلی به حافظه Cache منتقل کنیم.

اینکه کداسی از داده‌ها از حافظه Cache حذف شود تا فضای خالی برای انتقال داده مورد نظر از حافظه اصلی به حافظه Cache فراهم شود در واقع تکنیکی است که برای جای‌گذاری یا Replacement استفاده می‌کنیم که این امر بسیار روی کارایی Cache تأثیر دارد زیرا اگر داده‌ای را که بعداً مورد نیاز است از حافظه Cache بدهیم جای کافی خارج نمی‌مانیم عملاً باعث می‌شود که در زمانهای بعد miss رخ دهد.

اگر بتوانیم اطلاعاتی را که جای‌های می‌کنیم در جای‌های قرار دهیم که در زمانهای بعد تداخلی با داده‌های دیگر نداشته باشند توان معنی بود که میزان miss Rate را کاهش داده ایم

جای‌گذاری

Cache Replacement Policies

Prof. Mikko H. Lipasti
University of Wisconsin-Madison

ECE/CS 752 Spring 2012

علاوه بر Replacement چند پارامتر دیگر نیز وجود دارد که بر روی کارایی Cache تأثیر دارند.

Cache Design: Four Key Issues

- These are: یعنی اینکه داده‌ها جدیدی را که وارد حافظه می‌کنیم در کجا قرار دهیم
 - Placement - جای‌دهی
 - Identification - شناسایی یا جستجو در حافظه
 - Where can a block of memory go?
 - How do I find a block of memory?
 - Replacement
 - Write Policy - سیاست در نوشتن در حافظه
 - How do I make space for new blocks?
 - How do I propagate changes?
- Consider these for caches
 - Usually SRAM
- Also apply to main memory, disks

© 2005 Mikko Lipasti 2

پارامترهای تأثیرگذار بر کارایی حافظه Cache عبارتند از :

— placement یا جای دهن یعنی اینکه داده‌های جدیدی را که وارد حافظه Cache می‌کنیم در کجای حافظه قرار دهیم.

— Identification شناسایی یا جستجو در حافظه یعنی جستجوی در حافظه به چه صورت

... باشد. آیا به صورت ترتیبی یا بلیک بندی شده باشد و یا به صورت موازی مانند حافظه‌های انحصاری صورت گیرد که این امر بسیار زیاد بر روی سرعت و بهره‌وری Cache تأثیر دارد.

— Replacement یا جای‌گزینی یعنی اینکه وقتی حافظه Cache پر است کدام یک از خانه‌های حافظه را خالی کنیم که بعداً به آن نیاز نداشته باشیم.

— Write Policy یعنی سیاست در نوشتن

حافظه Cache یک حافظه کاملاً موقت است و زمانی که جریان الکتریکی قطع شود،

اطلاعات حافظه Cache هم از بین می‌رود پس بنابراین زمانی که یک کپی از اطلاعات را

درون Cache قرار می‌دهیم که برای اجرای دستورات، این اطلاعات باید تغییر پیدا نماید، حال

به سیاست در نوشتن را به کار ببریم که این تغییرات بر روی داده اصلی نیز اعمال شود که این امر

مربوط به write Policy می‌باشد.

در نیات ما ، کامپایلر و مدیریت برنامه تصمیم می گیرند که جای دهن در کجا انجام شود .
 در حافظه Cache ، جای دهن به صورت یک روش ثابت و از قبل تعیین شده می باشد مثلاً
 به صورت نگاشت مستقیم یا استفاده از حافظه های انجینی و یا تمام انجینی انجام می گیرد .
 در حافظه اصلی ، عمل جای دهن توسط سیستم عامل مشخص می گردد .

Placement

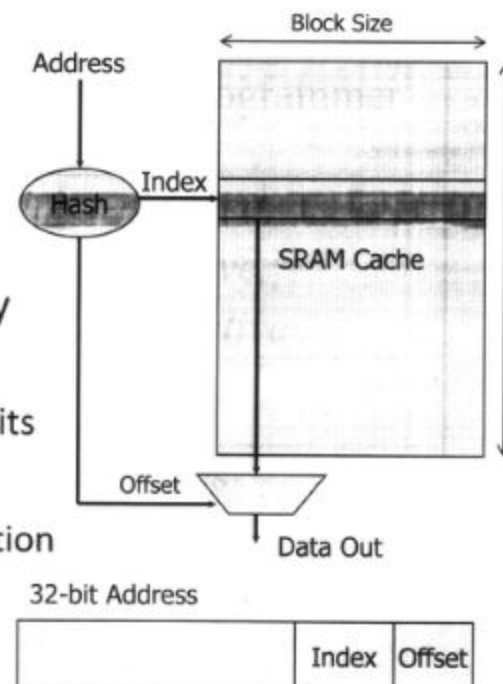
Memory Type	Placement	Comments
Registers	Anywhere; Int, FP, SPR	Compiler/programmer manages
Cache (SRAM)	Fixed in H/W	<i>Direct-mapped, set-associative, fully-associative</i>
DRAM	Anywhere	O/S manages
Disk	Anywhere	O/S manages

© 2005 Mikko Lipasti

3

Placement

- Address Range
 - Exceeds cache capacity
- Map address to finite capacity
 - Called a *hash*
 - Usually just masks high-order bits
- *Direct-mapped*
 - Block can only exist in one location
 - Hash collisions cause problems

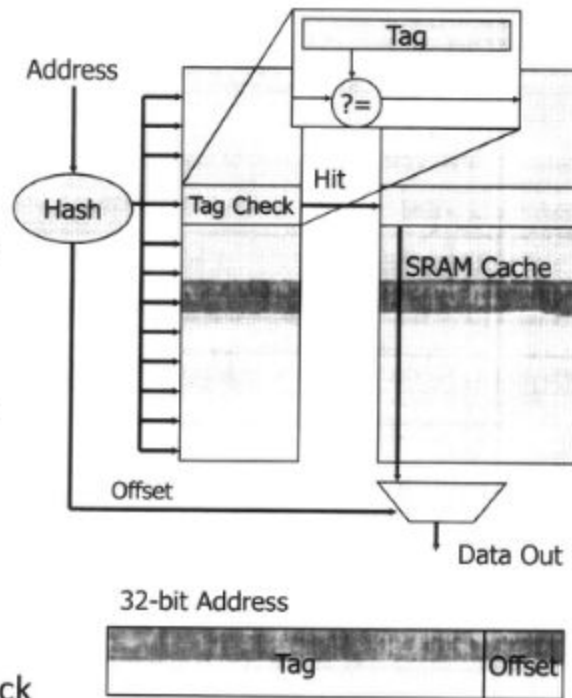


© 2005 Mikko Lipasti

4

Identification

- *Fully-associative*
 - Block can exist anywhere
 - No more hash collisions
- *Identification*
 - How do I know I have the right block?
 - Called a *tag check*
 - Must store address tags
 - Compare against address
- **Expensive!**
 - Tag & comparator per block

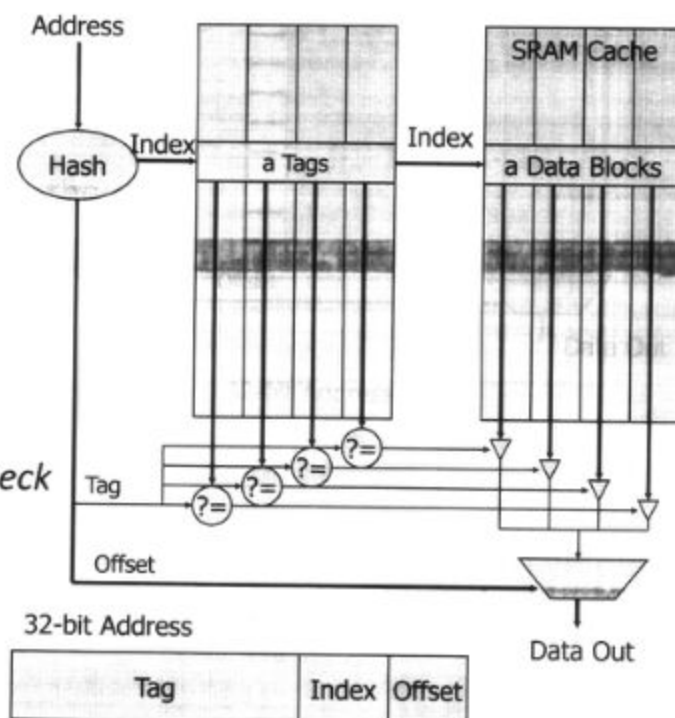


© 2005 Mikko Lipasti

5

Placement

- *Set-associative*
 - Block can be in *a* locations
 - Hash collisions:
 - *a* still OK
- *Identification*
 - Still perform *tag check*
 - However, only *a* in parallel



© 2005 Mikko Lipasti

6

سایز حافظه Cache محدود است. وقتی به یک حافظه Cache برسی شود و نیاز است داده‌ای را در آن قرار دهیم، چه استراتژی را برای آن به کار ببریم؟ کدام داده را باید خالی یا از Cache خارج نموده و داده جدید را به جای آن قرار دهیم. استراتژی Replacement بسیار زیاد بر روی hit و miss اثرگذار است.

Replacement

- Cache has finite size
 - What do we do when it is full?
- Analogy: desktop full?
 - Move books to bookshelf to make room
 - Bookshelf full? Move least-used to library
 - Etc.
- Same idea:
 - Move blocks to next level of cache

© 2005 Mikko Lipasti

7

در زمانیکه Cache خالی است و پردازنده شروع به کار می‌کند، در ابتدا برای هر داده‌ای miss رخ می‌دهد یعنی برای اولین ارجاع حتماً miss رخ می‌دهد. این امر به هیچ عنوان قابل چشم پوشی و حذف نمی‌باشد.

Cache Miss Rates: 3 C's [Hill]

- Compulsory miss or Cold miss
 - First-ever reference to a given block of memory
 - Measure: number of misses in an infinite cache model
- Capacity
 - Working set exceeds cache capacity
 - Useful blocks (with future references) displaced
 - Good replacement policy is crucial!
 - Measure: additional misses in a fully-associative cache
- Conflict
 - Placement restrictions (not fully-associative) cause useful blocks to be displaced
 - Think of as *capacity within set*
 - Good replacement policy is crucial!
 - Measure: additional misses in cache of interest

8

استفاده از ظرفیت‌های بزرگ حافظه Cache اگر چه در ابتدا باعث افزایش miss rate می‌شود اما در درازمدت اگر از روش‌های مناسب جای ردهی استفاده شود باعث کاهش نرخ miss می‌شود.

Replacement

- How do we choose *victim*?
 - Verbs: *Victimize, evict, replace, cast out*
- Many policies are possible
 - FIFO (first-in-first-out)
 - LRU (least recently used), pseudo-LRU
 - LFU (least frequently used)
 - NMRU (not most recently used)
 - NRU
 - Pseudo-random (yes, really!)
 - Optimal
 - Etc

© 2005 Mikko Lipasti

9

Optimal Replacement Policy?

[Belady, IBM Systems Journal, 1966]

- Evict block with longest reuse distance
 - i.e. next reference to block is farthest in future
 - Requires knowledge of the future!
- Can't build it, but can model it with trace
 - Process trace in reverse
 - [Sugumar&Abraham] describe how to do this in one pass over the trace with some lookahead (Cheetah simulator)
- Useful, since it reveals **opportunity**
 - (X,A,B,C,D,X): LRU 4-way SA \$, 2nd X will miss

© 2005 Mikko Lipasti

10

مخواهم ببینیم جای گذاری به چه روشی انجام شود زیرا الگوریتم های مختلفی برای جای گذاری وجود دارد یعنی در زمانی که حافظه Cache پر است ، کدام اطلاعات را حذف نمائیم که اطلاعات جدید جایگزین آن شود .

- FIFO ساده ترین روش آن است که هر کدام که اول آمده است زودترین خارج شود . این روش برای استفاده از دستورالعمل ها بسیار مناسب می باشد زیرا دستورات به صورت ترتیبی و خط به خط و پشت سر هم اجرا می شوند در صورتیکه برای رتبه بندی چنین نیست و متغیرها به ترتیبی که تعریف می شوند ، استفاده نمی گردند و ممکن است یک متغیر در چندین جای برنامه استفاده شود و یا فقط در یک خط از برنامه استفاده گردد .

- LRU در این روش داده ای که در مدت زمان دورتری استفاده شده و یا مدت زمان زیادی گذشته و از آن استفاده نشده است جهت جای گذاری انتخاب می شود . این روش یک استراتژی است که بیان می کند داده ای که خیلی وقت پیش استفاده شده و مدت زیادی است که استفاده نشده است به احتمال زیاد در آینده نزدیک هم استفاده نمی شود .

- LFU روش دیگری که وجود دارد آن است که داده ای که اخیراً استفاده شده ممکن است تا مدت زیادی از آن استفاده نشود

- NMRU این روش به صورت الگوریتمی عمل می کند و چون استفاده از الگوریتم باعث صرف زمان از سوی پردازنده است لذا کارایی پردازنده را کاهش می دهد .

- NRU تقریباً شبیه NMRU است .

- Pseudo-Random در این روش بکلی از داده ها را به صورت رندوم انتخاب می نمائیم که به صورت 50-50 می باشد . البته توابع رندوم باید عادلانه بوده و دارای توزیع یکسانی باشند .

- optimal در این روش می توان اولویت ^ن دس نمود که کدام یک از داده ها خارج تا داده جدید جایگزین آن شود .

* نکته ای که باید توجه نمود آن است که همه این روش ها همیشه و در همه جا مناسب نیستند .