

www.salampnu.com

سایت مرجع دانشجوی پیام نور

- ✓ نمونه سوالات پیام نور : بیش از ۱۱۰ هزار نمونه سوال همراه با پاسخنامه
- تستی و تشریحی
- ✓ کتاب ، جزوه و خلاصه دروس
- ✓ برنامه امتحانات
- ✓ منابع و لیست دروس هر ترم
- ✓ دانلود کاملاً رایگان بیش از ۱۴۰ هزار فایل مختص دانشجویان پیام نور

www.salampnu.com

معماری کامپیوتر (3 واحد درسی تئوری)

تهیه کننده: مجتبی پور محقق

این درس در مورد چیست؟

این درس در مورد ساختار و چگونگی طراحی کامپیوتر های دیجیتالی است . این مطلب معروف به ”معماری کامپیوتر“ است (که شامل معماری مجموعه دستورالعمل + سازماندهی سخت افزاری می باشد) .

چرا سازماندهی کامپیوتر را بررسی می کنیم

شاید هیچکدام از شما در شرکتهای Intel و یا AMD کار نکرده اید
اما ...

- کامپیوترهای جاسازی شده

- طراحی کامپایلر

- حتی طراحان نرم افزار

با محصولات این شرکتها در ارتباطند .

طرح کلی درس

- بازنگری کلی
- کارایی
- مجموعه دستورالعملها
- محاسبات کامپیوتر
- ماشینهای تک چرخه ای
- خط لوله ای
- سیستمهای حافظه (RAM , Caches , Virtual Memory)
- سوپر اسکالر (Superscalar/SLIW) و چند پردازنده ها
- مباحث دیگر

سرفصل 1: رؤوس مطالب در معماری و سازماندهی کامپیوتر

- یافتن توانایی ارائه اطلاعات پایه از معماری و سازماندهی کامپیوتر در جریان طراحی کامل یک کامپیوتر

- دریافتن مسئولیتهای حرفه ای و اخلاقی یک مهندس کامپیوتر (مخصوصا معمار کامپیوتر)

”معماری“ به چه معناست؟

” فن یا دانش یا ساختمان ... فن یا پرداختن به طراحی و پیاده سازی ساختارها ... ”

Webster 9th New College Dictionary •

”شامل نقشه ، طراحی ، ساخته و دکوراسیون چگونگی عملکرد“

American College Dictionary •

”معماری کامپیوتر“

- کلمه ای که توسط Fred Brooks ابداع گردید.

“معماری کامپیوتر”

“معماری کامپیوتر، یعنی کامپیوتر از دید کاربر”

- Amdhal et al, (64)

“ما بوسیله معماری، ساختار واحدهای تشکیل دهنده یک سیستم کامپیوتری را هدفمند می نمایم.”

- Stone, H. (1987)

”معماری کامپیوتر“

”معماری یک کامپیوتر عبارتست از محیط یا فضای
بین ماشین و نرم افزار“

- Andris Padgett
IBM 360/370 Architect

“معماری کامپیوتر”

ساختار: نظم و ترتیب دادن به بخشهای ثابت (نقشه)

سازماندهی: فعل و انفعال پویای این بخشها و مدیریت آنها

پیاپی سازی: طراحی کردن بلوک بخشهای دارای هدف خاص

ارزیابی کارایی: مطالعه رفتار سیستم (decorative treatment)

معماری (از دیدگاه معمار)

- پیاده سازی
- سازماندهی: منظر سطح بالا
 - سیستم حافظه
 - ساختار گذرگاه (bus)
 - طراحی داخلی CPU
- سخت افزار
 - طراحی منطقی
 - تکنولوژی بسته بندی (packaging)
- معماری مجموعه دستورالعمل

نکات مهم

- به خاطر بسیاری: نکته اینست که پیامزید چگونه معماری به مفهوم تکنولوژی موجود را ارزیابی کنید.
- شناختن روش خیلی مهم است، اما پایان کار نیست.

مراحل در سازماندهی کامپیوتر

- مفهوم ماشینهای چند سطحی
- مفاهیم ماشین مجازی

انضباط در معماری

- ساختار سخت افزار / نرم افزار
- الگوریتم ها و پیاده سازی آنها
- انتشار زبان

تصویر بزرگ

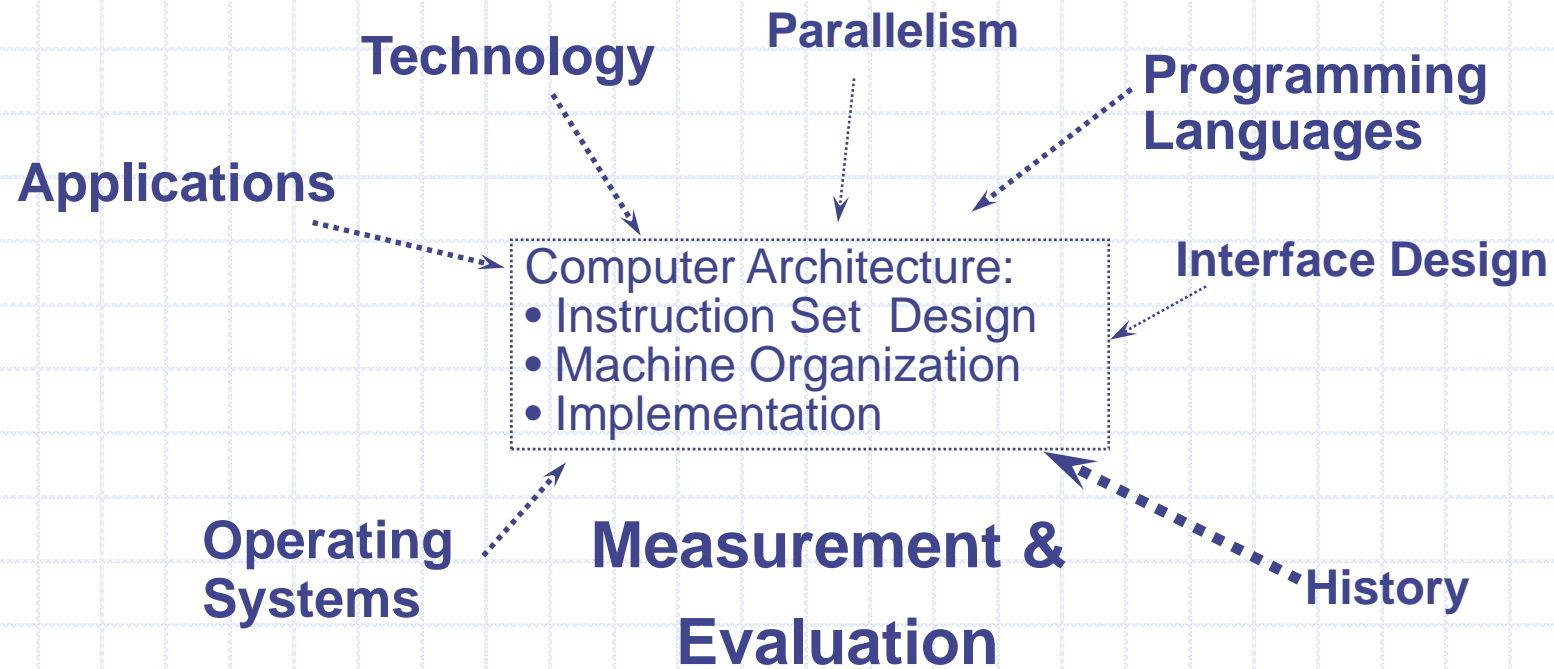
هر دوی سخت افزار و نرم افزار مرکب از لایه های سلسله مراتبی هستند، با هر لایه سطح پائینتر جزئیاتی از دید سطح بالاتر مخفی می شوند. این اصل تجرید، روشی است که طراحان سخت افزار و طراحان نرم افزار از عهده پیچیدگی سیستم های کامپیوتری بر آمدند.

یک محیط کلیدی بین لایه های انتزاعی معماری مجموعه دستورالعمل است: فضای بین سخت افزار و نرم افزار سطح پائین. این محیط مجازی توان بسیاری در پیاده سازی برای دگرگونی هزینه و کارایی در اجرای یک نرم افزار یکسان است.

John L. Hennessy

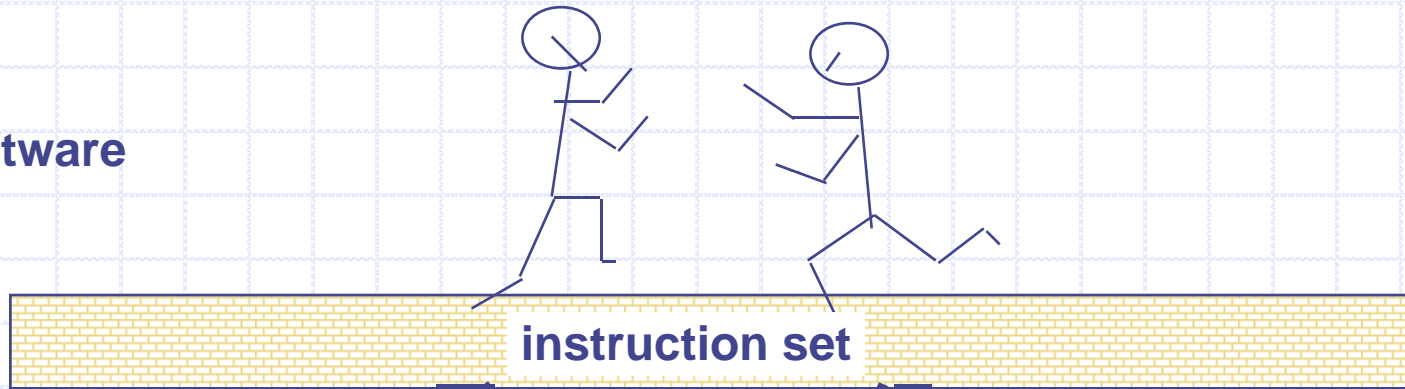
David A. Patterson

عوامل در معماری کامپیوتر



مجموعه دستورالعملها یک محیط بحرانی

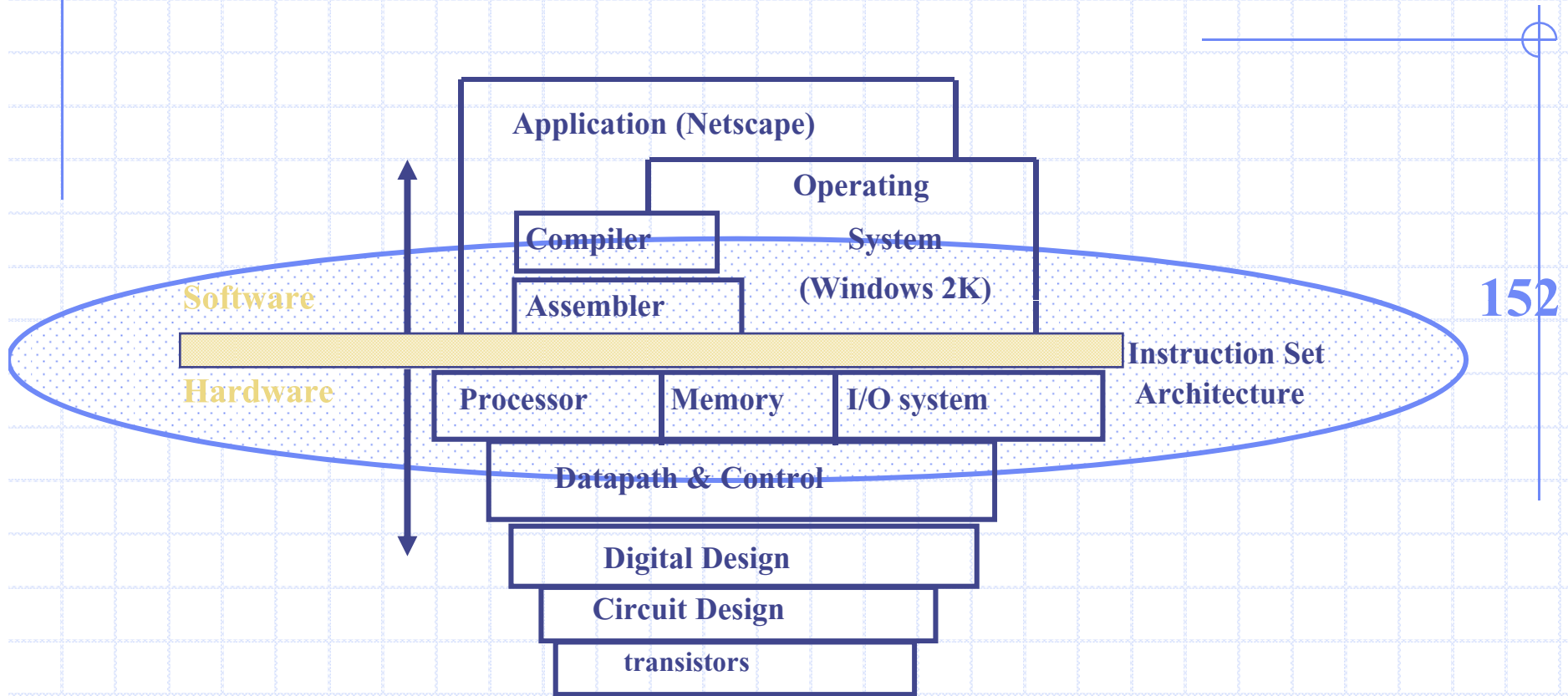
software



instruction set

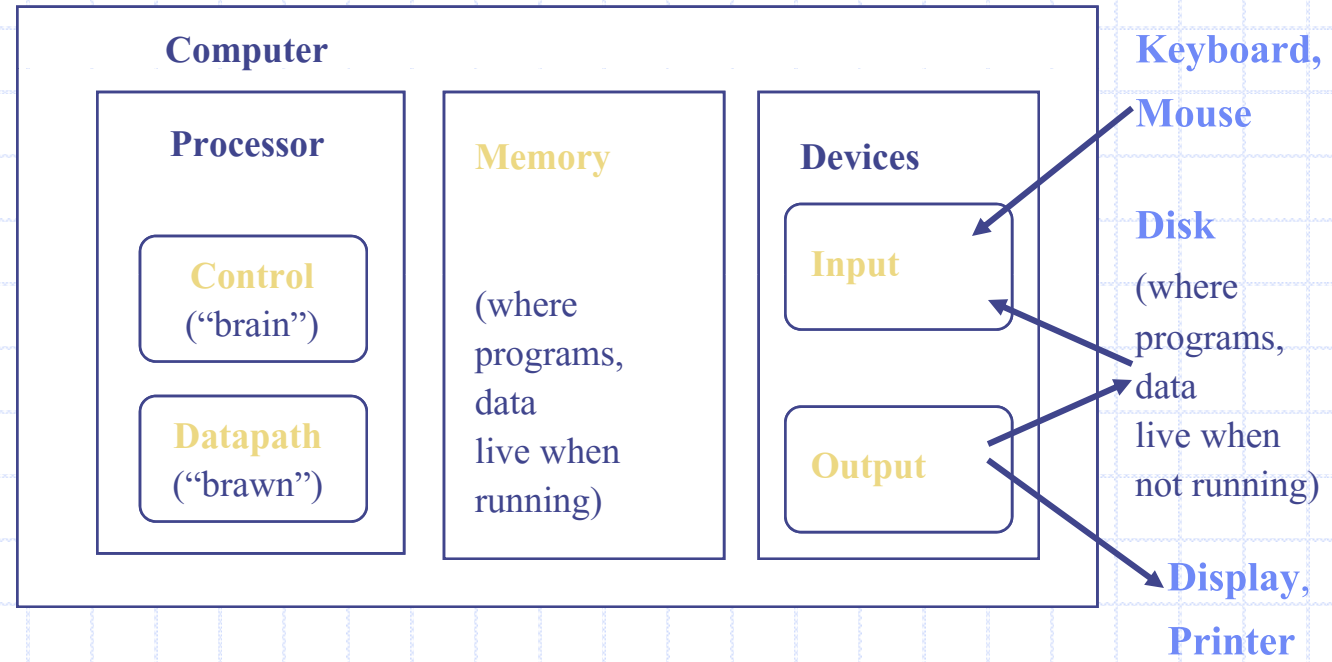
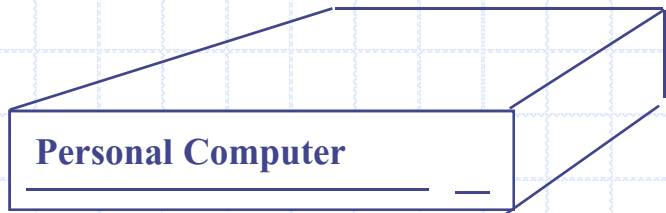
hardware

مهندسی و معماری کامپیوتر کجاست؟



*همه‌هنگی بسیاری از لایه‌های انتزاعی

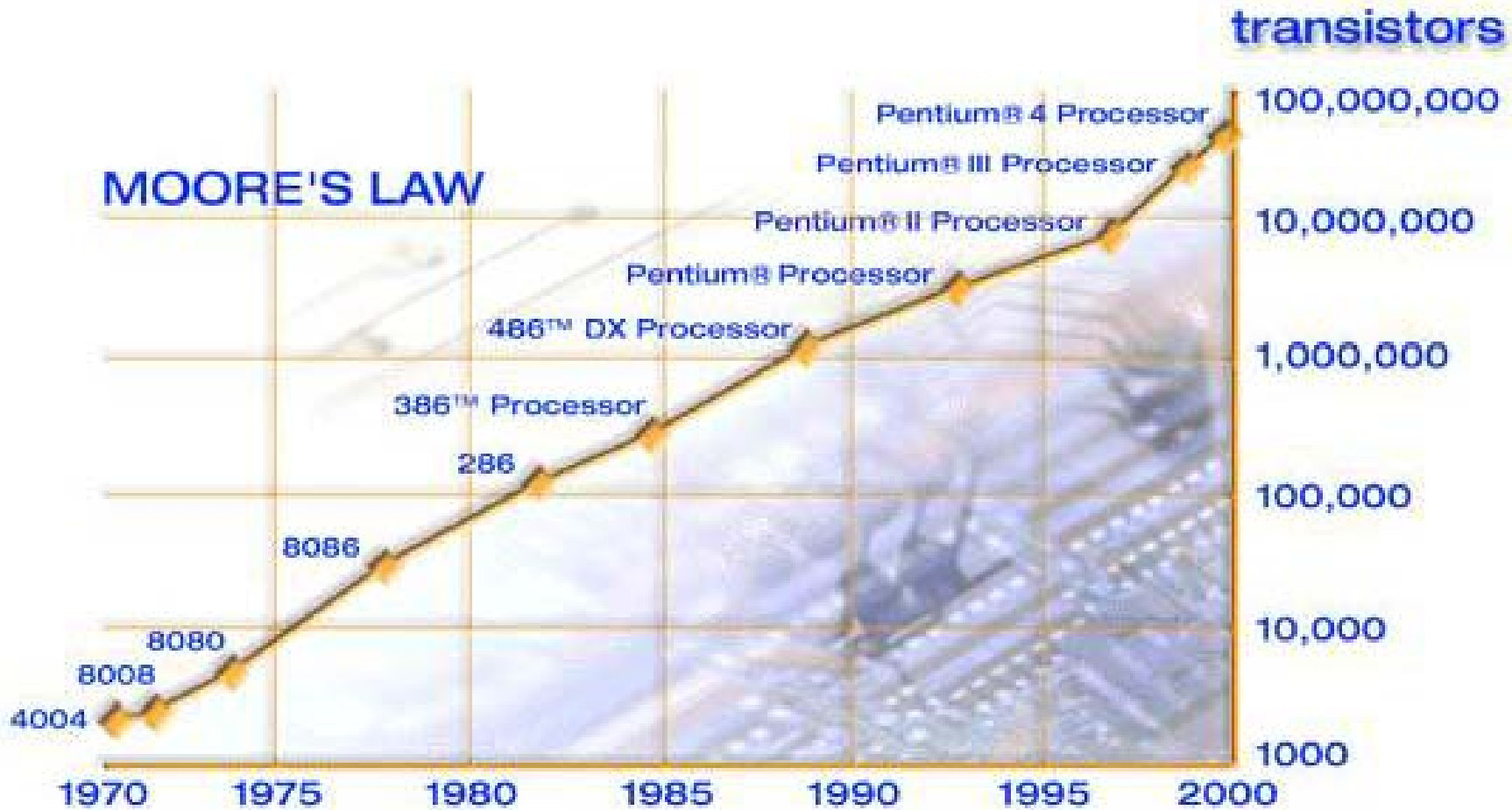
تشریح: پنج جزء ترکیب شده کامپیوتر



تکنولوژی کامپیوتر: تغییری مهیج

- پردازنده
 - هر یکسال و نیم، 2 برابر در سرعت (از سال 1985).
 - کارایی 100 برابر، در دهه گذشته
- حافظه
 - ظرفیت DRAM: 2 برابر در هر دو سال (از سال 96)
 - بهبود اندازه 64x در دهه گذشته
- دیسک
 - ظرفیت: 2 برابر در هر سال (از سال 97)
 - بهبود اندازه 250x در دهه گذشته

گرایش تکنولوژی: پیچیدگی ریزپردازنده

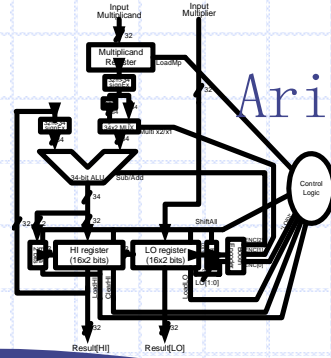
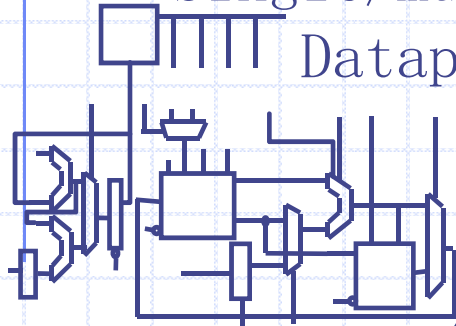


2X transistors/Chip Every 1.5 to 2.0 years

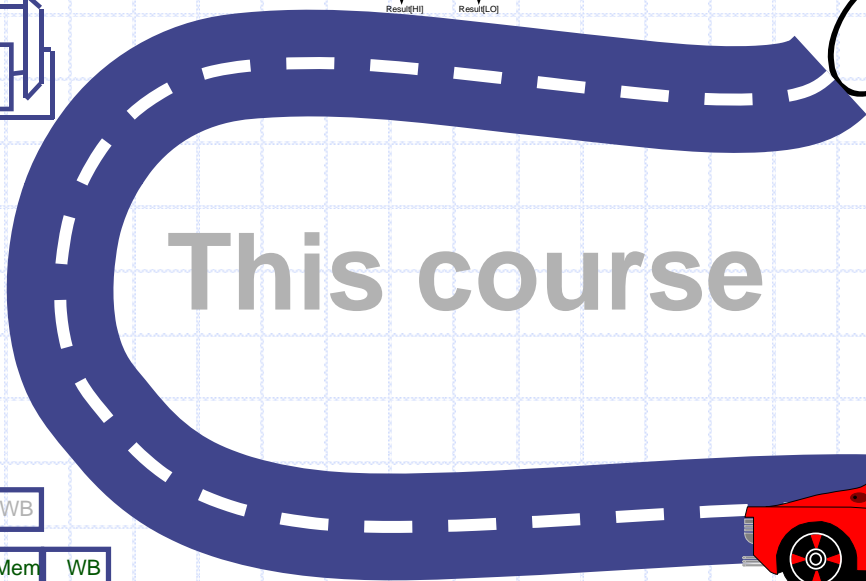
2 برابر ترانزیستور/تراشه هر 1.5 تا 2 سال "موسوم به قانون مور"

به کجا می رویم؟

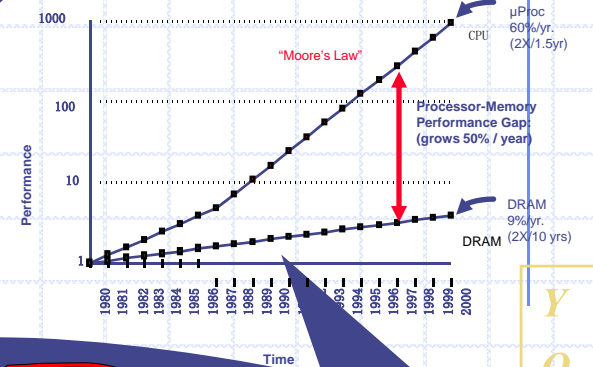
Single/multicycle
Datapaths



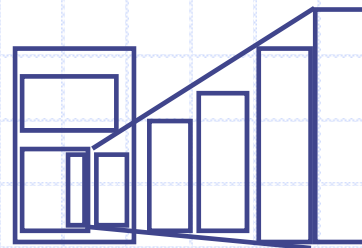
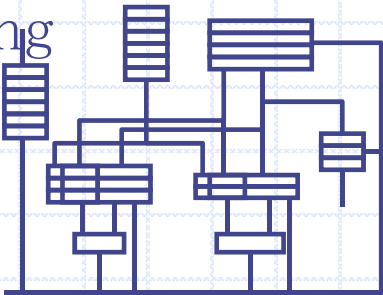
Arithmetic



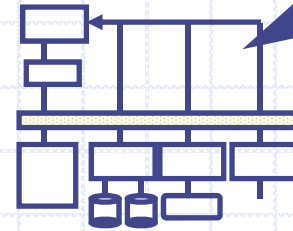
This course



Pipelining



Memory Systems



I/O

Y
O
U
R

C
P
U

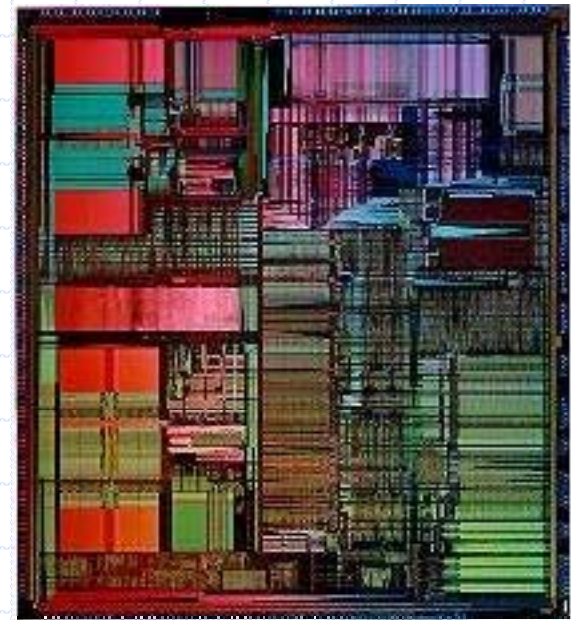
Intel 486™ DX CPU

- ◆ Design 1986 – 1989
- ◆ 25 MHz, 33 MHz
- ◆ 1.2 M transistors
- ◆ 1.0 micron
- ◆ 5 stage pipeline
- ◆ Unified 8 KByte code/data cache (write-through)
- ◆ First IA-32 processor capable of executing 1 instruction per clock cycle



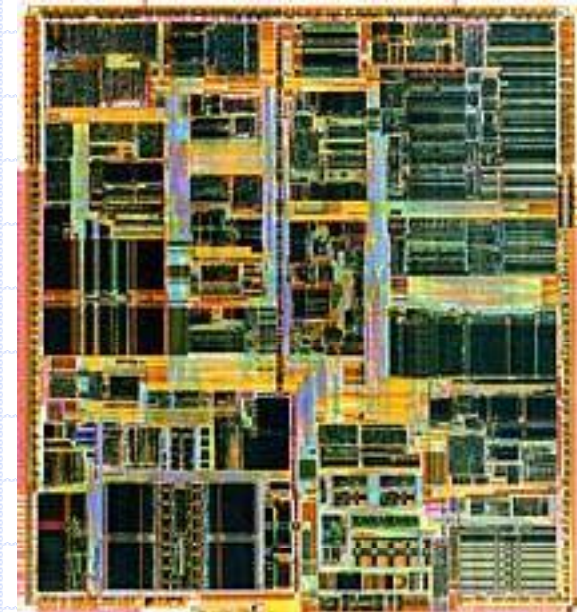
Pentium® Processor

- ◆ Design 1989 – 1993
- ◆ 60 MHz, 66 MHz
- ◆ 3.1 M transistors
- ◆ 0.8 micron
- ◆ 5 stage pipeline
- ◆ 8 KByte instruction and 8 KByte data caches (writeback)
- ◆ Branch predictor
- ◆ Pipelined floating point
- ◆ First superscalar IA-32: capable of executing 2 instructions per clock



Pentium® II Processor

- ◆ Design 1995 – 1997
- ◆ 233 MHz, 266 MHz, 300 MHz
- ◆ 7.5 M transistors
- ◆ 0.35 micron
- ◆ 16 KByte L1I, 16 KByte L1D, 512 KByte off-die L2
- ◆ First compaction of P6 microarchitecture



Pentium® III Processor (Katmai)

- ◆ Introduced: 1999
- ◆ 450 MHz, 500 MHz, 533 MHz, 600MHz
- ◆ 9.5 M transistors
- ◆ 0.25 micron
- ◆ 16 KByte L1I, 16 KByte L1D, 512 KByte off-chip L2
- ◆ Addition of SSE instructions.



SSE: Intel Streaming SIMD Extensions to the x86 ISA

Pentium® III Processor (Coppermine)

- ◆ Introduced: 1999
- ◆ 500MHz ... 1133MHz
- ◆ 28 M transistors
- ◆ 0.18 micron
- ◆ 16 KByte L1I, 16 KByte L1D, 256KByte on-chip L2
- ◆ Integrate L2 cache on chip, It topped out at 1GHz.



Pentium® IV Processor

- ◆ Introduced: 2000
- ◆ 1.3GHz ... 2GHz ... 3.4GHz
- ◆ 42M ... 55M ... 125 M transistors
- ◆ 0.18 ... 0.13 ... 0.09 micron
- ◆ Latest one: 16 KByte L1I, 16 KByte L1D, 1M on-chip L2
- ◆ Very high clock speed and SSE performance



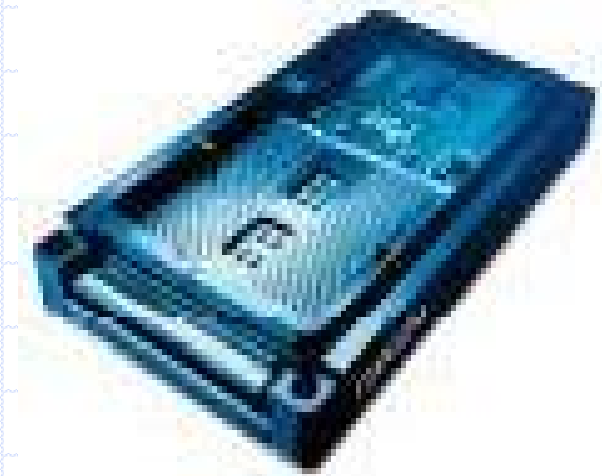
Intel® Itanium® Processor

- ◆ Design 1993 – 2000
- ◆ 733 MHz, 800 MHz
- ◆ 25 M transistors
- ◆ 0.18 micron
- ◆ 3 levels of cache
 - 16 KByte L1I, 16 KByte L1D
 - 96 KByte L2
 - 4 MByte off-die L3
- ◆ Superscalar degree 6, in-order machine
- ◆ First implementation of 64-bit Itanium architecture



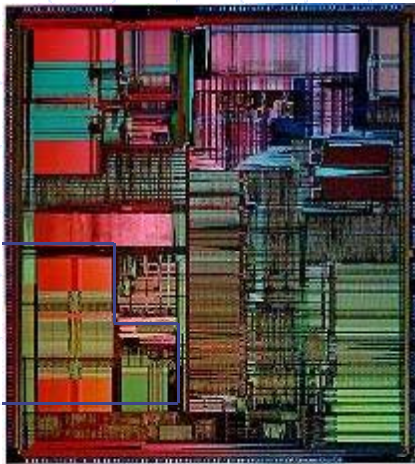
Intel® Itanium 2® Processor

- ◆ Introduced: 2002
- ◆ 1GHz
- ◆ 221 M transistors
- ◆ 0.18 micron
- ◆ 3 levels of cache
 - 32 KByte I&D L1
 - 256 KByte L2
 - integrated 1.5MByte L3
- ◆ Based on EPIC architecture
- ◆ Enhanced Machine Check Architecture (MCA) with extensive Error Correcting Code (ECC)



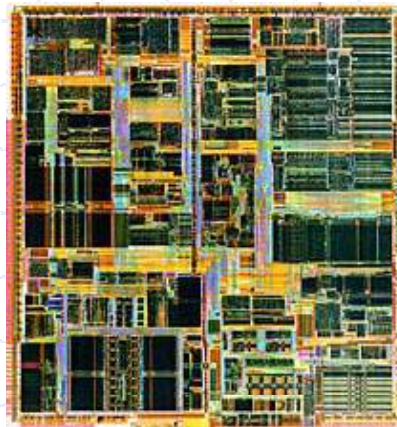
Cache Size Becoming Larger and Larger

1993: Pentium



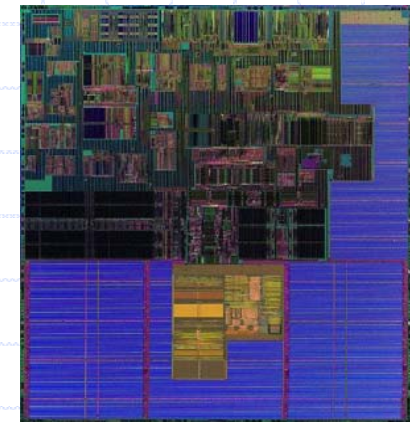
- 8 KByte I-cache and 8 KByte D-cache

1997: Pentium-II



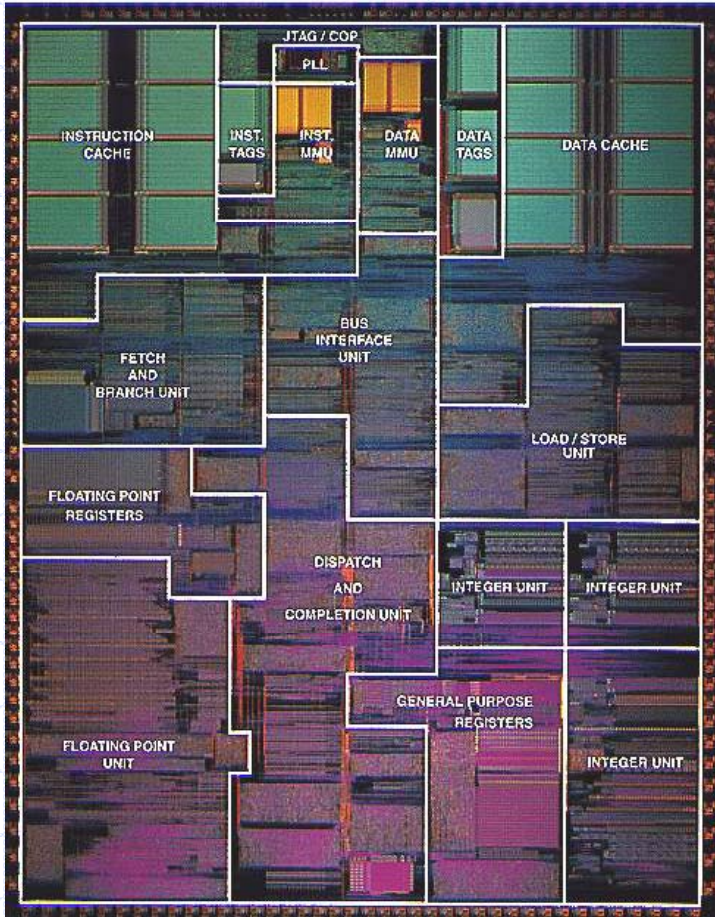
- ◆ 16 KByte L1I, 16 KByte L1D
- ◆ 512 KByte off-die L2

2002: Itanium-2

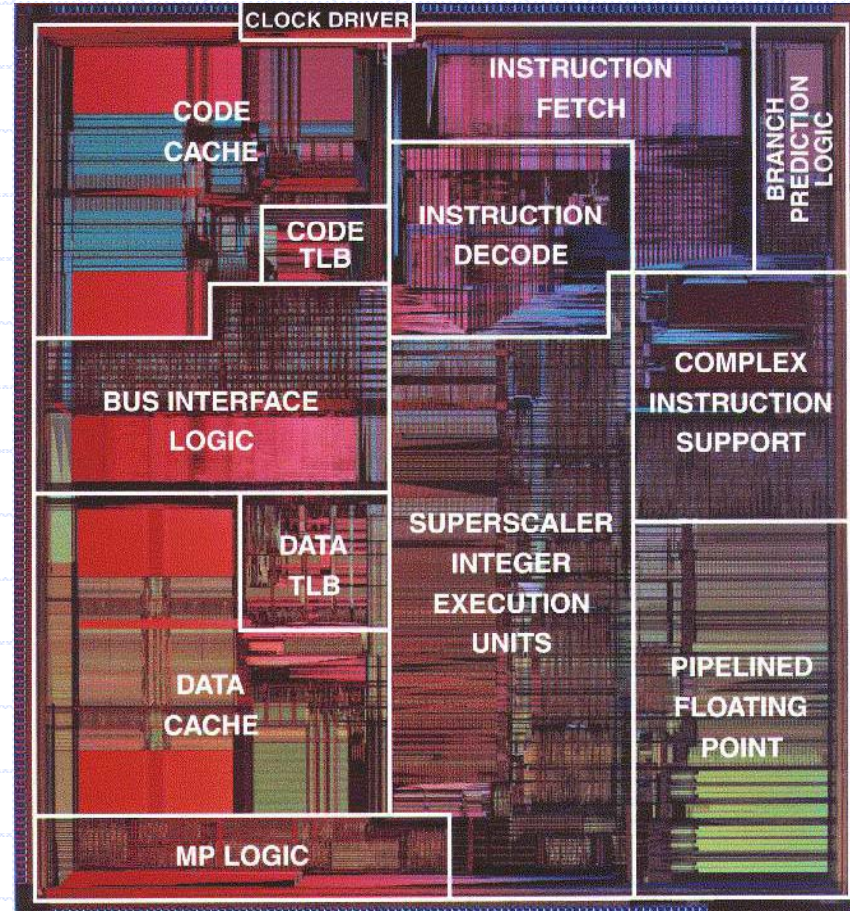


- ◆ Level 1: 16K KByte I-cache, 16 KByte D-cache
- ◆ Level 2: 256 KB
- ◆ Level 3: integrated 3 MB or 1.5 MB

Motorola's PowerPC™ 604 RISC Microprocessor



Motorola's PowerPC 604



Pentium

Intel Pentium 4 Northwood

Buffer Allocation & Register Rename

Instruction Queue (for less critical fields of the uOps)

General Instruction Address Queue & Memory Instruction Address Queue (queues register entries and latency fields of the uOps for scheduling)

Floating Point, MMX, SSE2 Renamed Register File 128 entries of 128 bit.

uOp Schedulers

FP Move Scheduler: (8x8 dependency matrix)

Parallel (Matrix) Scheduler for the two double pumped ALU's

General Floating Point and Slow Integer Scheduler: (8x8 dependency matrix)

Load / Store uOp Scheduler: (8x8 dependency matrix)

Load / Store Linear Address Collision History Table

Integer Execution Core

- (1) uOp Dispatch unit & Replay Buffer Dispatches up to 6 uOps / cycle
- (2) Integer Renamed Register File 128 entries of 32 bit + 6 status flags 12 read ports and six write ports
- (3) Databus switch & Bypasses to and from the Integer Register File.
- (4) Flags, Write Back
- (5) Double Pumped ALU 0
- (6) Double Pumped ALU 1
- (7) Load Address Generator Unit
- (8) Store Address Generator Unit
- (9) Load Buffer (48 entries)
- (10) Store Buffer (24 entries)

Execution Pipeline Start

Register Alias History Tables (2x126)
Register Alias Tables uOp Queue

Micro code Sequencer
Micro code ROM & Flash

Instruction Trace Cache

Trace Cache Fill Buffers
Distributed Tag comparators
24 bit virtual Tags

Trace Cache Access, next Address Predict

Trace Cache Branch Prediction Table (BTB), 512 entries.
Return Stacks (2x16 entries)
Trace Cache next IP's (2x)
Miscellaneous Tag Data

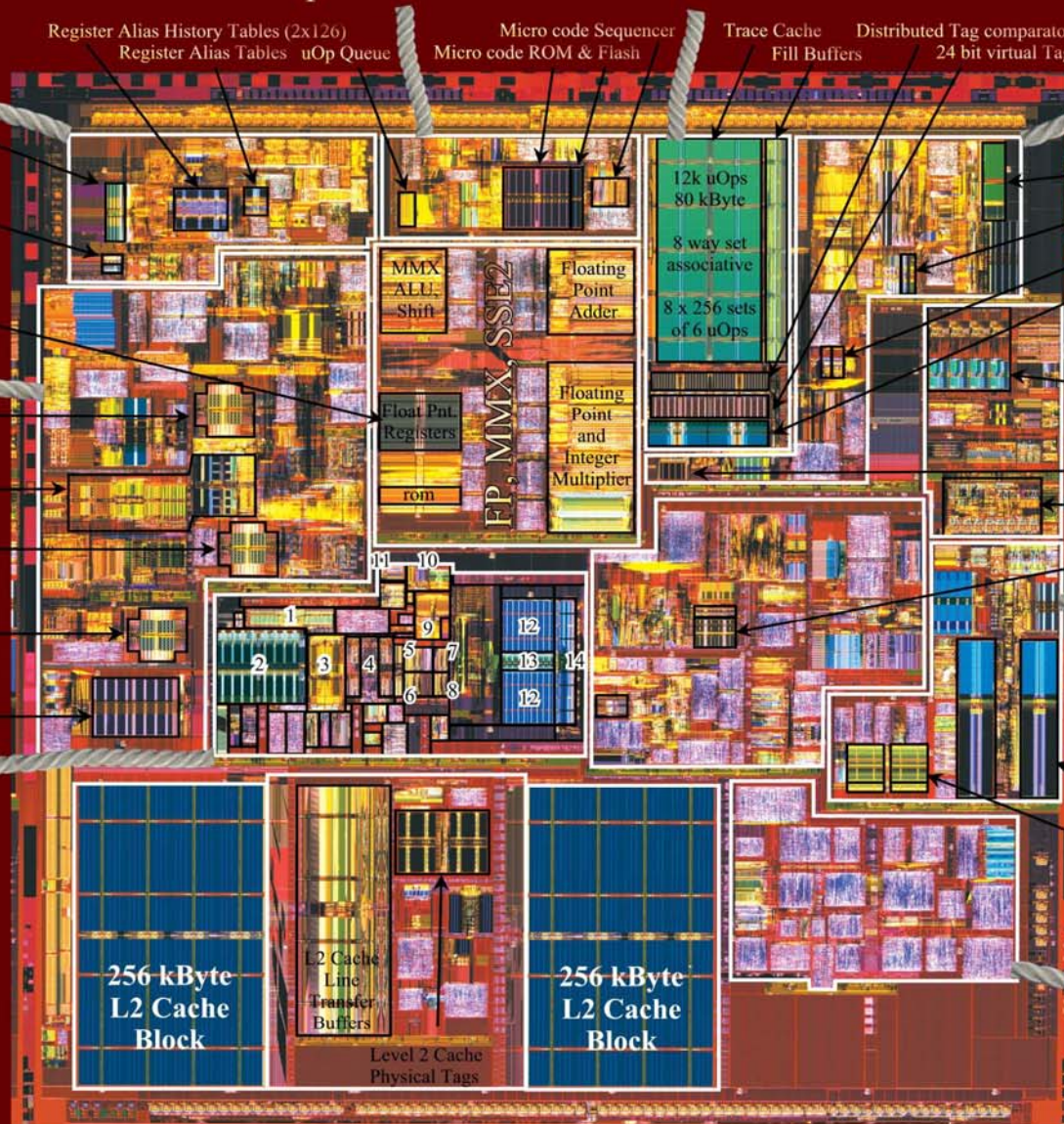
Instruction Decoder

Up to 4 decoded uOps/cycle out. (from max. one x86 instr/cycle)
Instructions with more than four are handled by Micro Sequencer
Trace Cache LRU bits
Raw Instruction Bytes in Data TLB, 64 entry fully associative, between threads dual ported (for loads and stores)

Instruction Fetch from L2 cache and Branch Prediction

Front End Branch Prediction Tables (BTB), shared, 4096 entries in total
Instruction TLB's 2x64 entry, fully associative for 4k and 4M pages. In: Virtual address [31:12] Out: Physical address [35:12] + 2 page level bits

Front Side Bus Interface, 400..800 MHz



- (11) ROB Reorder Buffer 3x42 entries
- (12) 8 kByte Level 1 Data cache
- (13) Summed Address Index decode and Way Predict
- (14) Cache Line Read / Write Transferbuffers and 256 bit wide bus to and from L2 cache

Intel Pentium 5 Prescott

Trace Cache Access, next Address Predict

Trace Cache Branch Prediction Table (BTB), 1024 entries.
Return Stacks (4 x 16 entries)
Trace Cache next IP's (4x)

Instruction Decoder

Up to 4 decoded uOps/cycle out. (from max. one x86 instr/cycle)
Instructions with more than four are handled by Micro Sequencer
Raw Instruction Bytes in Data TLB, 64 entry fully associative, between threads dual ported (for loads and stores)

Front End Branch Prediction Tables (BTB), shared, 4096 entries in total

Instruction TLB's 128 entry, fully associative for 4k and 4M pages. In: Virtual address [47:12]
Out: Physical address [39:12] + 2 page level bits

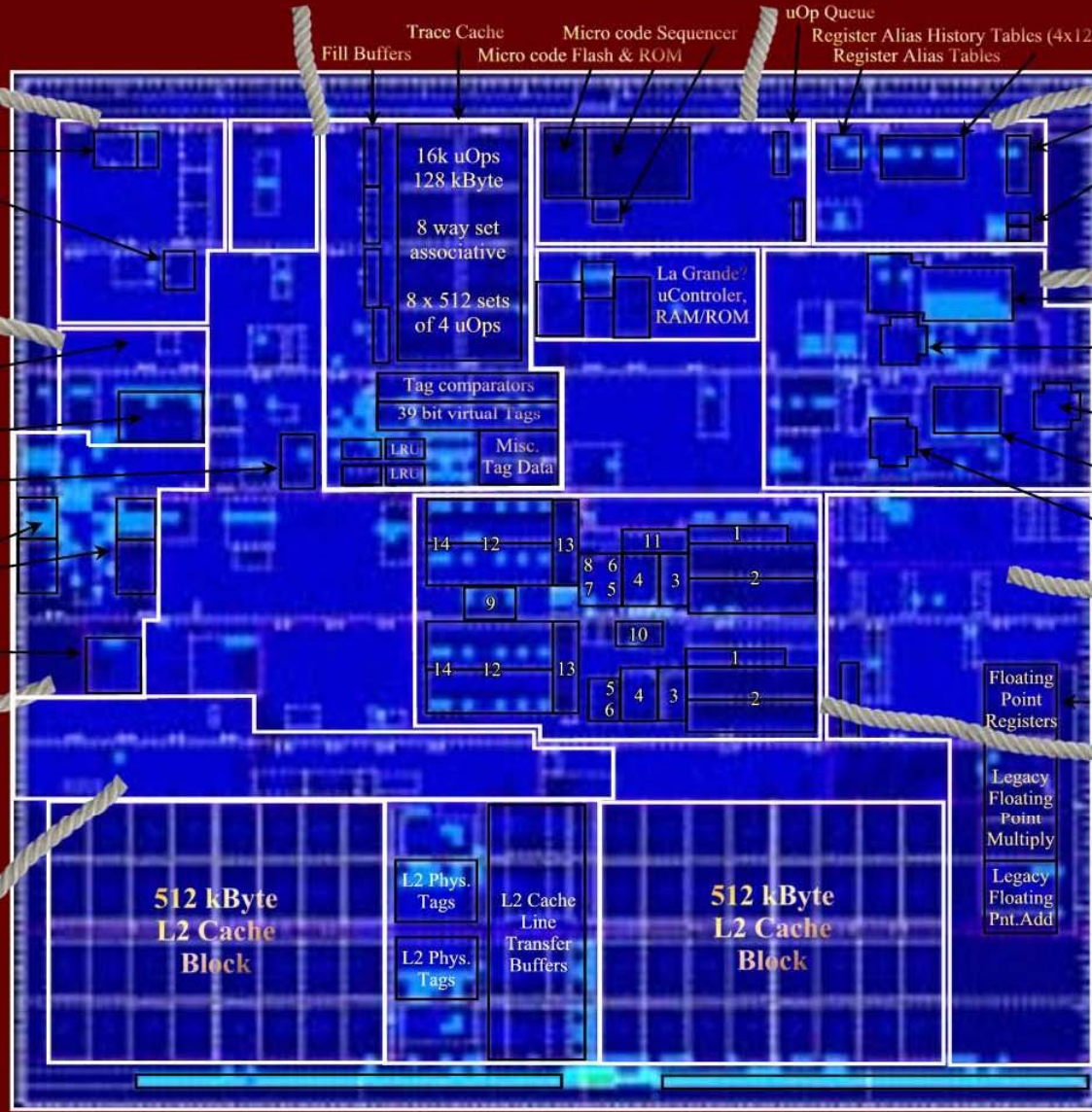
Instruction Fetch from L2 cache and Branch Prediction

Front Side Bus Interface, 533..800 MHz

Instruction Trace Cache

Execution Pipeline Start

Buffer Allocation & Register Rename



Instruction Queue (for less critical fields of the uOps)
General Instruction Address Queue & Memory Instruction Address Queue (queues register entries and latency fields of the uOps for scheduling)

uOp Schedulers

Parallel (Matrix) Scheduler for the two double pumped ALU's
General Floating Point and Slow Integer Scheduler: (8x8 dependency matrix)
FP Move Scheduler: (8x8 dependency matrix)
Load / Store Linear Address Collision History Table
Load / Store uOp Scheduler: (8x8 dependency matrix)

FP, MMX, SSE1..3

Floating Point Registers
Floating Point, MMX, SSE1..3 Renamed Register File
256 entries of 128 bit.

Integer Execution Core

- (1) uOp Dispatch unit & Replay Buffer
Dispatches up to 6 uOps / cycle
- (2) Integer Renamed Register File
256 entries of 32 bit (+ 6 status flags)
12 read ports and six write ports
- (3) Databus switch & Bypasses to and from the Integer Register File.
- (4) Flags, Write Back
- (5) Double Pumped ALU 0
- (6) Double Pumped ALU 1
- (7) Load Address Generator Unit
- (8) Store Address Generator Unit
- (9) Load Buffer (96 entries)
- (10) Store Buffer (48 entries)

- (13) Databus multiplexing
- (14) Cache Line Read / Write Transferbuffers and 256 bit wide bus to and from L2 cache
- (11) ROB Reorder Buffer 4x64 entries
- (12) 16 kByte Level 1 Data cache four way set associative. 1R/1W

یک بازبینی از جریان طراحی پردازنده

چگونه یک CPU طراحی می شود

- طراحی معماری مجموعه دستورالعمل (ISA design)
- طراحی در سطح وظایف (function-level) RTL design
- طراحی در سطح اجزاء ترکیب دهنده
- طراحی gate-level/switch-level
- طراحی در سطح مدار

روش کلاسیک طراحی مرحله ای معماری

مجموعه دستورالعمل

- انتخاب یک الگوی ساخت A
- تعریف A برای تطبیق با:
- کارایی مورد تقاضای جدید و تکنولوژی جدید
- ارزیابی (شبیه سازی معماری مجموعه دستورالعمل)
- تکرار تا کسب رضایت

کل استراتژی شبیه سازی

1. شبیه ساز در سطح دستورالعمل (ISA): این روش برای ارزیابی کارایی در سطح مجموعه دستورالعمل برای شرح بیشتر مدل سازی استفاده می شود.

2. شبیه سازی در سطح سیستم: مدلهای این شبیه ساز جزئیات مربوط به محیط سیستم شامل برخی چیزها مانند وقفه ها و مدیریت حافظه را مدل سازی می کند.

کل استراتژی شبیه سازی

(Con'd)

3. در سطح RTL: مدلهای این شبیه ساز تشریح RTL از طراحی است.

4. در سطح سوئیچ همراه با تاخیرها: بیش از همه برای شبیه سازی مولفه های طراحی به کار می رود؛ بردارهای آزمایشی از مرحله RTL تولید شده اند.

5. شبیه سازی در سطح مدار: برای مدل سازی جزئیات مسیر بحرانی به علاوه برای بازبینی مدارها در تغییرات دمایی، توان ارائه شده و غیره به کار می رود.

کارایی شبیه سازها

Simulator	Level of Accuracy	Simulation Rate
ISA	Instruction set	$> 10^6$ cycles/second
System	System level (OS instructions + interrupts)	$> 10^3$ cycles/second
RTL	Synchronous register transfer	> 10 cycles/second
Gate	Gate/switch level	> 1 cycles/second

تعداد چرخه های شبیه سازی شده در ثانیه بر روی یک کامپیوتر میزبان

کارایی دستورالعمل بر چرخه (IPC)

- ایجاد مدل کارایی که:
 - قابل انعطاف
 - پذیرای پارامتر (Parameterized) از طریق دستگیره
 - دارای دقت ساعت در مقایسه با RTL برابر با 95% باشد
 - به صورت قابل توجهی سریعتر از RTL
- مدلها مرکب از دو بخش هستند.
 - شبیه ساز مجموعه دستورالعمل - اجرا کننده محک (benchmark)
 - شبیه ساز خط لوله - "حسابدار" برای چرخه های ساعت
- سرعت های شبیه سازی
- اجرای محک ها (benchmark)، به روز رسانی ریز معماری بر طبق آن
- چرخه : کد - > شبیه ساز - > characterize - > وفق دادن (tune)

فصل دوم

معماری مجموعه دستورالعملها

- مقدمه
- یک بررسی موردی: معماری مجموعه دستورالعمل ماشین MIPS

مراحل اجرای یک دستورالعمل

واکشی دستورالعمل: برداشت دستورالعمل بعدی از حافظه

کدبرداری از دستورالعمل: بررسی دستورالعمل برای مشخص شدن اینکه:

چه عملی باید توسط دستورالعمل انجام گیرد (به عنوان مثال جمع)

چه عملوندهایی مورد نیازند، و نتایج باید کجا قرار گیرند.

واکشی عملوندها: عملوندها برداشت می شوند.

اجرا: اجرای عملیات بر روی عملوندها

بازنویسی نتیجه: نوشتن نتیجه در محل مخصوص

دستورالعمل بعدی: تعیین اینکه دستورالعمل بعدی از کجا گرفته شود.

چه چیزی در یک ISA (معماری مجموعه دستورالعمل) مشخص می شود؟

کدبرداری از دستورالعمل: اعمال و عملوندها چگونه تعیین می گردند؟

واکشی عملوندها: عملوندها ممکن است کجا باشند؟ چه تعداد؟
اجرا: چه اعمالی می تواند انجام گیرد؟ چه نوع داده و چه اندازه هایی؟

بازنویسی نتایج: نتایج کجا نوشته می شوند؟ چه تعداد؟
دستورالعمل بعدی: دستورالعمل بعدی را چگونه می توان

انتخاب نمود؟

یک ISA ساده: حافظه به حافظه

- چه عملیاتی می تواند اجرا شود؟ عملیات پایه ریاضی (برای این لحظه)
 - چه نوع داده و چه اندازه ای؟ نوع داده صحیح 32 بیتی (integer)
- (32

- عملوندها و نتایج کجا می توانند قرار گیرند؟ حافظه
- چه تعداد عملوند و نتیجه؟ 2 عملوند، 1 نتیجه
- اعمال و عملوندها چگونه مشخص می شوند؟

عمل مقصد, منبع 1, منبع 2

OP DEST, SRC1, SRC2

- چگونه می توانیم دستورالعمل بعدی را انتخاب کنیم؟ بعدی به ترتیب

مدل حافظه

حافظه را به عنوان یک آرایه بزرگ از n عدد صحیح در نظر بگیرید، که بوسیله اندیس قابل دستیابی است. (حافظه با دستیابی تصادفی موسوم به ram)

Address	Contents
0	14
1	3
2	99
⋮	⋮
⋮	⋮
$N-1$	0

به عنوان نمونه ، $M[1]$ شامل مقدار 3 است. ما می توانیم در این مکانها نوشتن و خواندن را انجام دهیم. این مکانها صرفا در دسترس ماست. تمام مکانهای "مجرد" (از قبیل متغیرها در C) باید مکانهایی را در M تعیین کنند.

ترجمه کد ساده

کد C مفروض

$$A = B + C;$$

ما می توانیم تصمیم بگیریم که متغیر A مکان 100 ، B مکان 48 و C مکان 76 را

اشغال می کند. کد بالا را به معادل کد اسمبلی آن تبدیل می کنیم:

```
ADD M[100], M[48], M[76]
```

چگونه می توانیم عبارت زیر را تبدیل کنیم

$$A = (B + C) * (D + E);$$

استفاده از یک مکان موقتی

فرض کنید ما A را در 100، B را در 48، C را در 76، D را در 20 و E را در 32 قرار می دهیم.

اکنون یک مکان بدون استفاده را انتخاب می کنیم (مثلا 84)

# A = B + C	ADD M[100], M[48], M[76]
# temp = D + E	ADD M[84], M[20], M[32]
# A = A * temp	MUL M[100], M[100], M[84]

مشکلات در معماری حافظه به حافظه

- حافظه اصلی خیلی کندتر از مدارات محاسباتی است
- این مطلب از سال 1950 تا 2003 درست است!
- خانه های زیادی برای مشخص نمودن آدرسهای حافظه گرفته می شوند.
- معمولا نتایج یک یا دو دستورالعمل بعد مورد استفاده قرار می گیرند.

به خاطر داشته باشید: بخشهای اشتراکی را سریعتر نمائید!

راه حل: نتایج میانی یا موقتی را در حافظه های سریع و نزدیک به واحد محاسبه ذخیره نمائید.

ماشینهای مبتنی بر انباشتگر (Accumulator)

یک ماشین انباشتگر، یک بافر پرسرعت واحد (مانند یک مجموعه از D latch ها یا فلیپ فلاپها، هر کدام برای یک بیت داده) را نزدیک واحد محاسبه منطق نگهداری می کند.

در ساده ترین حالت، فقط یک عملوند می تواند مشخص گردد؛ انباشتگر به صورت مجازی به مفهوم "OP Operand" می باشد یعنی:

$acc. = acc. \text{ OP operand}$

Example:

Load B into acc.

LOAD M[48]

Add C to acc. (now

ADD M[76]
has B+C)

Write acc. To A

STORE M[100]

ماشین مبتنی بر انباشتگر که $A=(B+C)*(D+E)$ را انجام می دهد

Load D into acc. LOAD M[20]

بارگذاری D در انباشتگر

Add E to acc. (now has D+E) ADD M[32]

جمع E با انباشتگر. (اکنون داریم D+E)

Write acc. To A STORE M[100]

نوشتن محتوای انباشتگر در A.

Load B into acc. LOAD M[48]

بارگذاری B در انباشتگر.

Add C to acc. (now has B+C) ADD M[76]

جمع نمودن C با انباشتگر. (اکنون داریم B+C)

Multiply A to acc. MUL M[100]

ضرب کردن A در انباشتگر.

Write (B+C) * (D+E) to A STORE M[100]

نوشتن نتیجه در A

ضعف ماشینهای مبتنی بر انباشتگر

- هنوز نیازمند ذخیره سازی بسیاری مقادیر موقتی و میانی در حافظه می باشیم
- در واقع انباشتگر فقط برای یک ترتیب از محاسبات که در آن نتیجه یکی، ورودی برای بعدی است، مفید می باشد.

ماشینهای مبتنی بر انباشتگر هنوز در کامپیوترهای اولیه معمول بودند

- یک طراحی ساده، و بنابراین محبوب، مخصوص برای
 - کامپیوترهای اولیه
 - میکروپروسسورهای اولیه (4004، 8008)
 - مدلهای Low-end (ارزان)

پیشنهادات برای ماشینهای مبتنی بر انباشتگر

- اگر منابع سخت افزاری بیشتری در دسترس است، مکانهای ذخیره سازی سریع را در کنار انباشتگر قرار دهید:
- ماشینهای مبتنی بر پشته
- ماشینهای مبتنی بر ثبات
- خاص منظوره
- عام منظوره

ماشینهای مبتنی بر پشته

ایده: یک ستون از مکانهای ذخیره سازی سریع دارای یک بالا (top) و یک پایین (bottom)

Address	Contents
top	14
2 nd from top	3
3 rd from top	99
:	:
:	:
bottom	0

یک دستورالعمل فقط از مقدار بالای پشته (top) قابلیت برداشت دارد، یا شاید دو یا سه مقدار به عنوان بالای پشته در نظر گرفته شوند.

ما می توانیم مقادیر جدید را در بالای پشته قرار دهیم ("push") یا از بالای پشته برداریم ("pop") اما فقط همین. ما نمی توانیم به مکانهایی در زیر پشته دسترسی داشته باشیم مگر اینکه هر چیز بالای آن را خارج کنیم.

معماری مجموعه دستورالعمل ماشینهای مبتنی بر پشته

اعمال اصلی شامل:

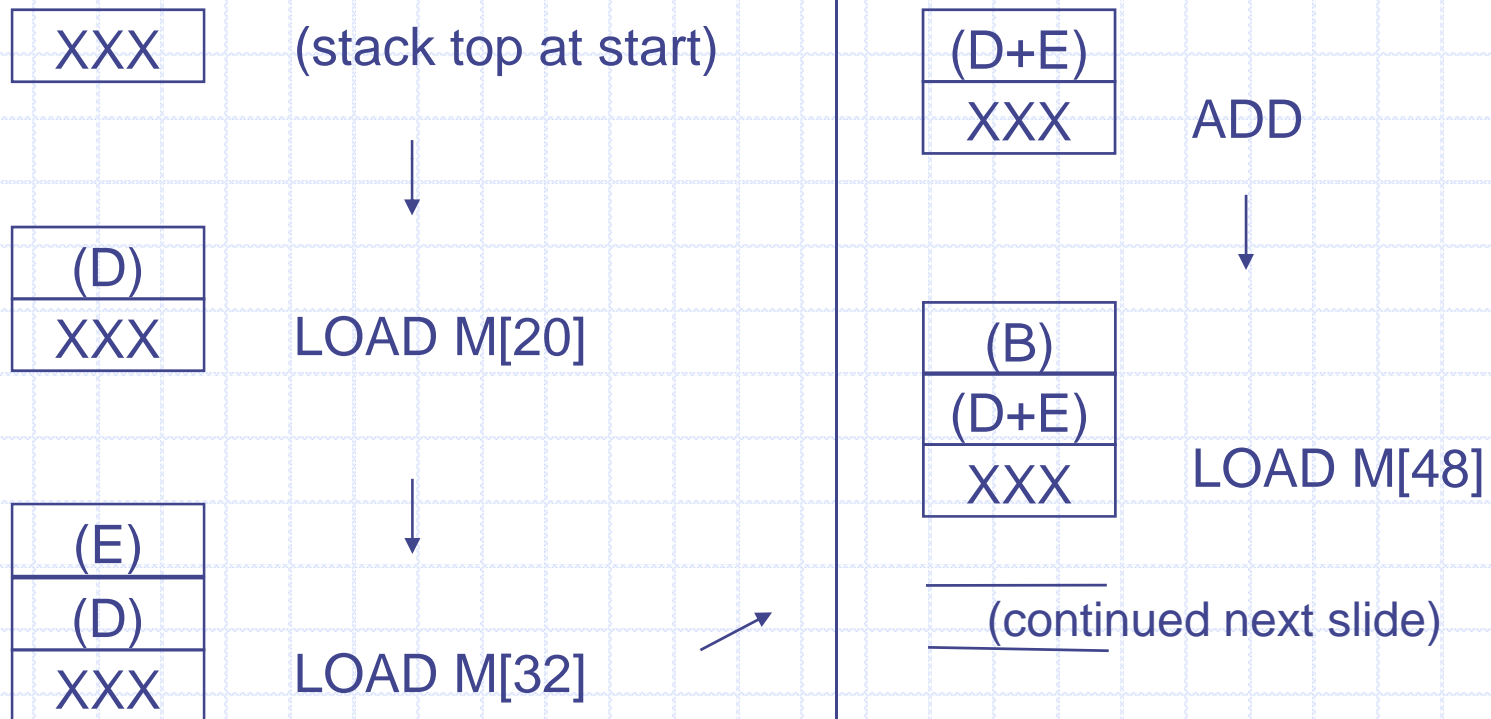
بارگذاری: برداشت مقدار از حافظه و قرار دادن آن بر روی پشته

ذخیره سازی: برداشت مقدار از پشته و ذخیره آن در حافظه

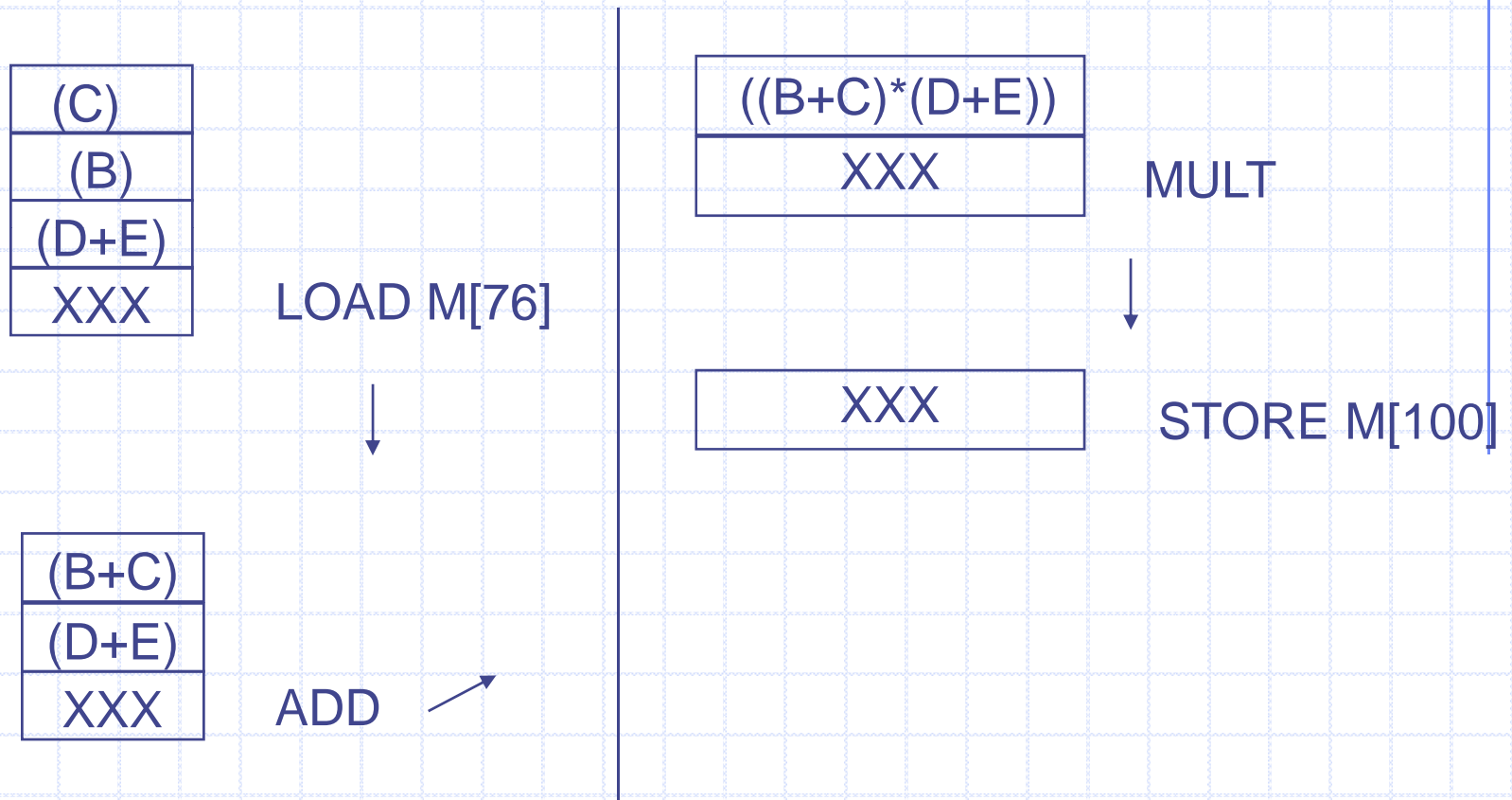
حسابی: خارج کردن یک یا دو مقدار از پشته؛ قرار دادن نتیجه روی پشته

دونسخه ای (Dup): گرفتن مقدار از بالای پشته بدون حذف آن؛ قرار دادن یک کپی جدید از آن در روی پشته (چرا این عمل کاربردی است؟)

ماشین مبتنی بر پشته که $A=(B+C)*(D+E)$ را انجام می دهد



ماشین مبتنی بر پشته



دقت کنید که اکنون پشته مشابه زمانی است که شروع نمودیم.

کاربرد ماشینهای مبتنی برپشته

- بسیاری از کامپیوترهای اولیه
- واحد ممیز شناور در 8086 (تقریبا)
- ماشین مجازی جاوا

ماشینهای مبتنی بر ثبات

- ایده: به کار بردن مکانهای ذخیره سازی زیاد ("ثباتها") نزدیک انباشتگر
- ثبات ها نام / شماره های مشخصی دارند که می توانند به جای حافظه استفاده می شوند
 - دستیابی نسبت به حافظه اصلی
 - خیلی سریعتر است

(1-2 CPU cycles vs. ~ 100 CPU cycles on PC)_

- ثبات ها نسبت به مکان های حافظه خیلی نزدیک ترند
- 32 MIPS_ تا ثبات 32 بیتی دارد
- ثبات های کمتر_ آدرس های کوچکتر_ و تعداد بیت های کمتر برای نام گذاری آنها
- منابع استفاده کمتر است و باید با دقت از آنها استفاده شود!

ثباتهای عام منظوره_خاص منظوره

- یک ثبات خاص_منظوره برای اهداف مشخص استفاده می شود و ممکن است عملیاتی را که استفاده می شود را محدود کند.
 - طراحی آسانتر سخت افزار: ثبات را در جایی که دقیقا نیاز است قرار بده
 - به منظور استفاده موثر برای کامپایلر خیلی سخت تر است.
- یک ثبات عام منظوره می تواند در بیشتر مسیرهای عملیاتی استفاده شود بنابراین مسیریابی خیلی مشکل است

ثبات های خاص منظور ه The z_80 cpu

- هفت تا ثبات 8 بیتی: (می توانند جفت شوند A,B,C,D,E,H,LBC,DE,HL).
- سه تا ثبات 16 بیتی: PC 62 (شمارنده برنامه).
- جمع, تفریق, شیفت تنها با A می تواند انجام شود (شمارنده 8 بیتی).
- افزایش و کاهش می تواند با تمام ثبات ها وجفتهای ثبات انجام شود.
- می توانند از حافظه آدرس (HL) را واکنشی کنند و در هر 8 بیت ثبات قرار دهند.
- یک واکنشی از آدرس (BC) یا (DE) تنها می تواند به A برود.
- واکنشی ها از (HL), (BC) و (IX) تعداد چرخه های متفاوتی می گیرد
- چه کسی می خواهد برای این یک کامپایلر بنویسد؟

ثبات عام منظوره ماشين هاى (GPR)

MIPS (وپر دازشگر هاى
مشابه) 32 تا ثبات عمومى
دارند

Address	Contents
\$0	0
\$1	3
\$2	99
:	:
:	:
\$31	14

(GPRs) هر 32 بيت long هستند
همه مى توانند نوشته يا
خوانده شوند به جز ثبات صفر
که هميشه صفر است ونمى
تواند تغيير کند. زمان دستيابى به
ثبات ثابت است.

ماشین GPR $A = (B + C) * (D + E)$ را انجام می دهد

#R1 = B + C M[48], M[76] \$1 ADD

#R2 = D + E M[20], M[32] \$2 ADD

#A = R1 * R2 M[100], \$1, \$2 MUL

اندازه های داده های متفاوت

چطور باید با اندازه های داده های متفاوت رفتار کنیم؟
انتخاب یک اندازه برای یک واحد ذخیره شده در یک آدرس تنها

- ذخیره کردن داده بزرگ در یک مجموعه از مکانهای همجوار حافظه

- ذخیره کردن داده کوچک در یک مکان:

use shift & mask ops

امروزه تقریباً همه ماشینها (شامل MIPS) دارای آدرس دهی بایتی
"Byte_Adressable" هستند هر مکان آدرس دهی ، در حافظه
8بیتی نگهداری می شود.

حافظه MIPS

روی یک ماشین قابل آدرس دهی بایتی از قبیل MIPS اگر ما بگوییم یک کلمه (32 بیت) در آدرس 80 ذخیره شده به این معناست که مکانهای 80 تا 83 را اشغال کرده. (کلمه بعدی از 84 شروع می شود).

به طور نرمال بارگذاری و ذخیره چندین بیت باید "تنظیم" شود. آدرس n بایتی بارگذاری یا ذخیره باید مضربی از n باشد. برای نمونه نیم کلمه تنها در آدرسهای زوج ذخیره میشود.

MIPS اجازه بارگذاری شدن و ذخیره شدن برای استفاده مخصوص دستورالعمل ها را نمی دهند اما آنها ممکن است کندتر شوند (بیشتر پردازشگرها این را برای همیشه اجازه نمی دهند!).

Byte-Order (“Endianness”)

- برای یک داده چند بایتی کدام قسمت به کدام بایت می رود؟
- اگر \$1 محتوی (F4240_H)1000000 و ما آن را در آدرس 80 ذخیره کنیم :
- در یک ماشین “Big Endian” “Big End” به آدرس 80 می رود
- در یک ماشین “Little Endian” از سوی دیگر است.

...	79	00	0F	42	40	...	
...	79	80	81	82	83	84	...

...	79	40	42	0F	00	...	
...	79	80	81	82	83	84	...

Big-Endian vs. Little-Endian

- ماشینهای MIPS, SPARC, 68000: Big Endian
- ماشینهای Little Endian: بیشتر پردازشگرهای intel, Alpha, Vax
- سازگاری مشکلات انتقال چندین بایت داده بین ماشینهای Big Endian, Little Endian

روشهای آدرس دهی

یک روش آدرس دهی ISA's این سوال را جواب می دهد
: عملوندها کجا میتوانند ذخیره شوند؟

ما دو نوع ذخیره سازی در MIPS داریم (وبیشتر ماشینهای دیگر)
: ثباتها و حافظه اصلی. ما می توانیم به هر یک از این دو یا
هر دو عملوندها برویم . یک تک عملوند می تواند با هر یک از
این دو بیاید یک ثبات یا یک مکان حافظه ، و روشهای آدرس
دهی راه های گوناگون تشخیص این مکانها را ارائه می کند.

روشهای آدرس دهی ساده

در این روشها یک مکان یا داده به طور مستقیم در یک دستورالعمل داده می شود:

Mode name	Example	Meaning	
Register	mov \$1, \$2	R2	R1
Direct (or absolute)	mov \$1, (40)	M[40]	R1
Immediate	mov \$1, #40	40	R1

روشهای آدرس دهی غیر مستقیم

در تولید یک آدرس حافظه یکی یا بیشتر ثباتها استفاده می شوند

Mode name	Example	Meaning
Reg. Indirect	mov \$1, (\$2)	M[R2] R1
Displacement	mov \$1, 40(\$2)	M[40+R2] R1
Indexed	mov \$1, \$4(\$2)	M[R4+R2] R1
Mem. Indirect	mov \$1, @(\$2)	M[M[R2]] R1

روشهای آدرس دهی پیشرفته

اجزای زیادی از آدرس دهی را در زبانهای سطح بالا پشتیبانی می کنند یا تعداد دستور العمل هارا در طول دستیابی از حافظه افزایش می دهند

Mode name	Example	Meaning	
Auto-increment	mov \$1, 4(\$2) ++	M[4+R2]	R1
Auto-decrement	mov \$1, 4(\$2) --	M{R2-4}	R1
Scaled	mov \$1, 40(\$2) [s]	M[40+R2x2s]	R1

روشهای آدرس دهی منتخب

کارهایی که انجام می شود: هر روش آدرس دهی ممکن است برای هر عملوندی در هر زمان استفاده شود

- طرح دستورات سطح بالا به طور مستقیم به دستورالعملها آسانتر است.
- برای طراحی پردازشگر سخت است بعلت اینکه باعث پیچیدگی بیشتر می شود

آدرس دهی های محدود: تنها روشهای کمی را اجازه می دهند و یا تعداد عملوندها را به روش های معین محدود میکنند

- برای کامپایلر یا برنامه نویس که پیرو قواعد مشخص هستند سخت تر است
- شاید طول کد زیاد شود

بسامد روشهای آدرس دهی

سه برنامه روی (مینی کامپیوتر) vax مشخص شده که همه انواع روشها را پشتیبانی می کنند:

Mode Name	Frequency of mode (%)		
	Min.	ave.	max.
Displacement	32	42	55
Immediate	17	33	43
Reg. Indirect	3	13	24
Scaled	0	7	16
Mem. Indirect	1	3	6
Others	0	2	3

روشهای آدرس دهی داده های تجربی

- بزرگی تغییر مکان مورد نیاز چه طور باید باشد؟
در بررسی های Spec in92 و Spec fp92 99 درصد تغییر مکانها در دامنه -215 تا $+215$ انجام می شود.
- بزرگی ثابت فوری مورد نیاز چه طور باید باشد؟
در مطالعات نشان داده شده 50 تا 60 درصد در 8 بیت و 75 تا 80 درصد در 16 بیت است.

چه طور ما نمایش دادن دستورالعمل ها را انجام می دهیم ؟

- ما به چند بیت احتیاج داریم برای اینکه بگوییم چه عملی انجام شده (مثلا جمع، تفریق، ضرب، غیره) که به این کدگذاری عملوند opcode گویند
- ما به چند بیت برای هر عملوندو نتیجه احتیاج داریم(در نمونه ما جمعا 3 تا)
- چه نوع روش آدرس دهی
- شماری از ثباتها، آدرس حافظه و یا ثابت فوری

دستور عملهای با طول متغیر

- تا کنون Vax هر روش را برای هر عملوندی اجازه داده می تواند یک دستور العمل با سه تا آدرس های 32 بیتی باشد (آدرس دهی مستقیم) 12 بایت در این دستور العمل
- اما ثباتها به بیتهای کمی برای مشخص شدن نیاز دارند بنابراین 12 بایت برای یک دستور العمل که تنها از 3 تا ثبات استفاده می کنند هدر میشود .
- باید از دستور العمل های با طول متغیر استفاده کرد. در Vax دستور العمل از 1 تا 17 بایت متغیر اند.

دستورالعمل های با طول ثابت

- اگر هر دستورالعمل با تعداد بیت‌های یکسان (ترجیحاً یک عدد زوج مانند 16 یا 32) بیشتر اجزای پردازشگر ساده تر خواهد شد. اما در هر دو صورت مقداری از فضا هدر می رود یا همه روشهای آدرس دهی را نمی تواند پشتیبانی کند!

بار گذاری اعداد صحیح کوچک

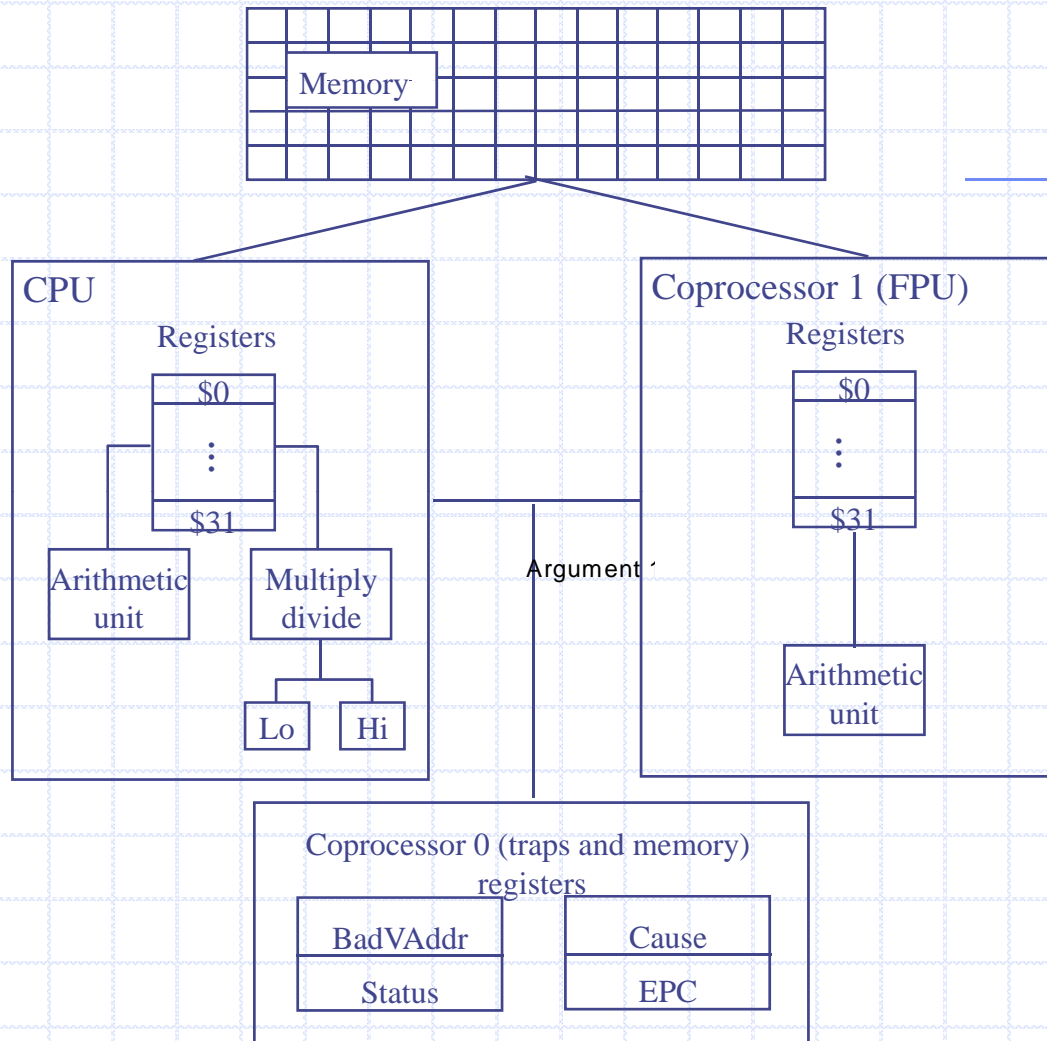
- همه ثباتها در 32 MIPS بیتی هستند.
- چه طور ما نیم کلمه یا یک بایت را در یک ثبات بار گذاری می کنیم؟
- بیتها را در حداقل 8 یا 16 بیت ثباتی بار گذاری می کنیم.
بار گذاری بدون علامت: همه بیت های بالا را صفر کنیم.
بار گذاری علامت دار: همه بیت های بالا را با توجه به علامت صفر یا یک می کنیم (نیم کلمه یا بایت MSB)

The RISC Approach

- در مجموعه دستورالعمل های کاهش یافته در دستگاه کامپیوتر همه دستورالعمل ها در اندازه های یکسانی هستند (32 بیت در MIPS)
- روشهای آدرس دهی کمی پشتیبانی می شوند (فقط آنهایی که بیشتر کاربرد دارند)
- فقط تعداد کمی از قالب های دستورالعمل (کد برداری آسانتر انجام می شود)
- دستورالعمل های محاسباتی که فقط روی ثبات کار می کنند
- داده های حافظه باید در ثبات ها قبل از پردازش شدن بارگذاری شوند که این معماری Load_Store نامیده می شود

معیارهای [COL WELL 85] RISC

- عمل تک چرخه ای
- ماشین Load_Store
- کنترل سخت افزاری
- رابطه دستور العمل های کم و روشهای آدرس دهی
- دستور العمل های با قالب ثابت
- تلاش بیشتر در زمان کامپایل



MIPS R2000 CPU and FPU

ثبات ها

32 تا ثبات با $R0=0$

ثبات های ذخیره شده: $R1, R26, R27$.

• کاربرد مخصوص:

اشاره گر به ناحیه سراسری: $R28$

اشاره گر پشته: $R29$

اشاره گر قاب: $R30$

آدرس بازگشت: $R31$

قراردادهای ثابت استاندارد

- 32 تا ثابت صحیح در MIPS عام منظوره هستند_ هر کدام می توانند به عنوان یک عملوند یا نتیجه یک محاسبه عملیاتی استفاده شوند
- اما ساخت تکه های متفاوت از نرم افزار که باهم کار می کنند آسانتر است اگر قراردادهای معینی که دنبال می شود درباره آن ثباتهایی باشد که برای آن اهداف استفاده می شود
- این قراردادها معمولاً به وسیله فروشنده پیشنهاد می شوند و به وسیله کامپایلرها پشتیبانی می شوند

قراردادهای ثابتی در MIPS

Names	Regs	Purpose
\$zero	0	Constant 0
-	1	(Reserved for assembler)
\$v0-\$v1	2-3	Return values/expression eval
\$a0-\$a3	4-7	Args to functions
\$t0-\$t9	8-15, 24-25	Temporaries (NOT SAVED)
\$s0-\$s7	16-23	Saved values
-	26-27	(Reserved for OS kernel)
\$gp	28	Global pointer
\$sp	29	Stack pointer
\$fp	30	Frame pointer
\$ra	31	Return address

MIPS registers and usage convention

Register name	Number	Usage
zero	0	Constant 0
at	1	Reserved for assembler
v0	2	Expression evaluation and results of a function
v1	3	Expression evaluation and results of a function
a0	4	Argument 1
a1	5	Argument 2
a2	6	Argument 3
a3	7	Argument 4
t0	8	Temporary (not preserved across call)
t1	9	Temporary (not preserved across call)
t2	10	Temporary (not preserved across call)
t3	11	Temporary (not preserved across call)
t4	12	Temporary (not preserved across call)
t5	13	Temporary (not preserved across call)
t6	14	Temporary (not preserved across call)
t7	15	Temporary (not preserved across call)
s0	16	Saved temporary (preserved across call)
s1	17	Saved temporary (preserved across call)
s2	18	Saved temporary (preserved across call)
s3	19	Saved temporary (preserved across call)
s4	20	Saved temporary (preserved across call)
s5	21	Saved temporary (preserved across call)
s6	22	Saved temporary (preserved across call)
s7	23	Saved temporary (preserved across call)
t8	24	Temporary (not preserved across call)
t9	25	Temporary (not preserved across call)
k0	26	Reserved for OS kernel
k1	27	Reserved for OS kernel
gp	28	Pointer to global area
sp	29	Stack pointer
fp	30	Frame pointer
ra	31	Return address (used by function call)

عملیات MIPS

- بارگذاری/ذخیره
- عملیات ALU
- انشعاب ها/پرش ها

قالب های دستورالعمل MIPS

R-Format



I-Format



J-Format



فیلدها در دستورالعملهای MIPS

- op: مشخصات عمل: بیان اینکه کدام قالب استفاده شود
- rs: اولین ثبات منبع
- rt: دومین ثبات منبع
- rd: ثبات مقصد
- shamt: مقدار شیفت
- funct: جزئیات اضافی opcode
- address: ثابت فوری و تغییر مکان یافته و یا انشعاب

ماشین نمایش دستورالعملهای MIPS

فیلدهای MIPS نام گذاری می شوند که بحث در مورد آنها آسانتر شود:

op	rs	rt	rd	shamt	funct
6 bits	5 bits	5 bits	5 bits	5 bits	6 bits

ینجا منظور هر نامی از فیلدهای دستورالعمل های MIPS است:

- OP: عملگر هر دستورالعمل
- RS: اولین عملوند ثبات منبع
- RT: دومین عملوند ثبات منبع
- RD: عملوند ثبات مقصد آن نتایج عملیات را فراهم می کند.
- SHAMT: مقدار شیفت
- FUNCT: تابع: این فیلد عملیات مختلف در فیلد OP را انتخاب می کند.

عملوندهای ALU

R_1, R_2, R_3 ADD

effect: $R_1 = R_2 + R_3$

MIPS مثال (شکل اسمبلر در)

```
ADD $t1, $s1,  
$s2
```

نمایش دهنده: 0 17 18 8 0 32

نمایش ماشینی

نمایش دهنده

0	17	18	8	0	32
6 bits	5 bits	5 bits	5 bits	5 bits	6 bits

نمایش باینری

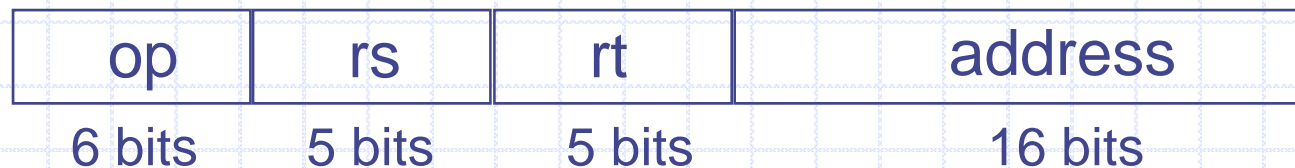
000000	10001	10010	01000	00000	100000
6 bits	5 bits	5 bits	5 bits	5 bits	6 bits

ضرب و تقسیم صحیح در MIPS

- ضرب 2 تا عدد 32 بیتی تا حدود 64 بیت می تواند باشد
- 2MIPS تا مخصوص دارد: LO و HI
- نتایج ضرب: بیت های پایین تر به LO می روند و بیت های بالاتر به HI می روند
- نتایج تقسیم: خارج قسمت به LO می رود و باقی مانده به HI می رود
- استفاده کردن از عملیات اضافی (از قبیل mflo) برای جابجایی به بالا و پایین ثبات های منظوره

دستور عملهای انتقال داده ها

I-type (base + 16 bit offsets)



مثال:

`/w t0, 8 ($s3) --- # Temporary reg t0 gets A[8]`

یادداشت:

همچنین `rs` ثبات پایه است (`$s3` در این نمونه _ همچنین ثبات شاخص نامیده می شود) و `rt` (در این نمونه `$t0` نتایج را ذخیره می کند (به عنوان ثبات مقصد)

An Example

MIPS انجام می دهد $A = (B + C) * (D + E)$

Assembly

1w \$8, 48(\$0)

1w \$9, 76(\$0)

add \$8, \$8, \$9

1w \$9, 20(\$0)

1w \$10, 32(\$0)

add \$9, \$9, \$10

add \$8, \$8, \$9

sw \$8, 100(\$0)

op	rs	rt	rd	sh.	Ft.
35	0	8		48	
35	0	9		76	
0	8	9	8	0	32
35	0	9		20	
35	0	10		32	
0	9	10	9	0	32
0	8	9	8	0	32
43	0	8		100	

انشعاب ها

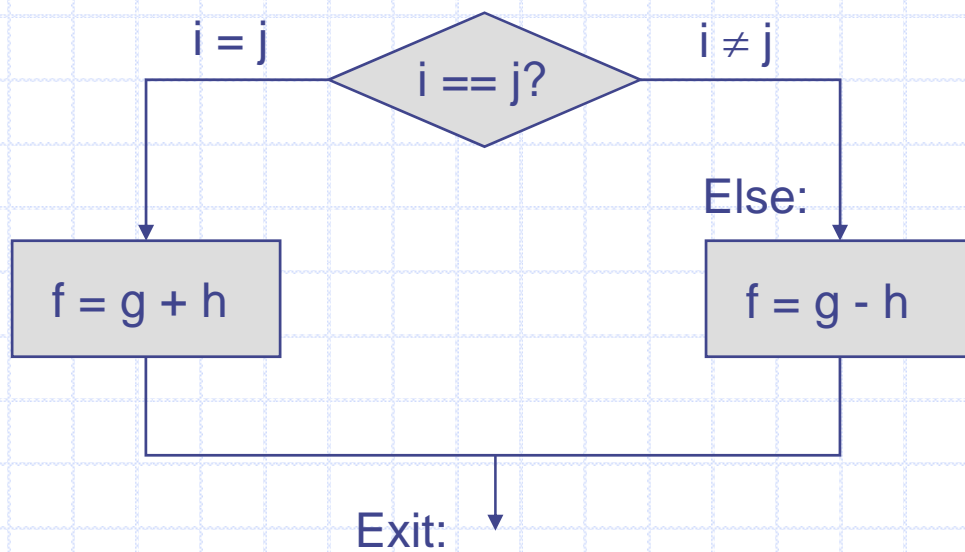
- در بیشتر پردازشگرها "شمارنده برنامه" (PC) آدرس دستور العمل بعدی را نگه می دارد: واکنشی از $M[(PC)]$
- به طور عادی بعد از اینکه یک دستور العمل تمام شد CPU n تا به PC اضافه می کند n تعداد بایت ها در دستور العمل است
- انشعاب ها به یک برنامه اجازه می دهند واکنشی کنند از مکانهای متفاوت
- انشعابها استفاده میشوند برای به کار بردن همه روند کنترلی فرمانهای زبان های سطح بالا از قبیل if-then-else, for, switch, etc .

طبقه بندی انشعاب ها

- دو نوع اساسی از پرشها:
- غیر شرطی همیشه به آدرس مشخص شده جهش می کند
- شرطی: اگر شرط درست باشد به آدرس مشخص شده می رود.
به عبارت دیگر با دستور العمل بعدی ادامه می دهد.
- آدرس های مقصد با روش مشابهی می تواند مشخص شود به عنوان عملوندهای دیگر (جمع آوری ثباتها، ثابت های فوری، و مکان های حافظه) بستگی دارد که چه چیزی در ISA پشتیبانی شود

مثال همگردانی انشعاب

- کامپایل زیر را دنبال کنید



```
if ( i == j)
    f = g + h;
else
    f = g - h;
```

MIPS `if Then Else`

Assume f,g,h,i,j in R8-R12 (respectively)

<code># Branch if i<>j</code>	<code>bne \$11, \$12, Else</code>
<code> # f = g + h;</code>	<code>add \$8, \$9, \$10</code>
<code> # Jump to Exit</code>	<code> j Exit</code>
<code># f = g - h;</code>	<code>sub \$8, \$9, \$10 Else:</code>
<code># Code after if</code>	<code> ... Exit:</code>

بررسی انشعاب ها

- بیشتر انشعابهای شرطی برای آدرس های ثابت و کوتاه روی می دهند
- اما اغلب این طور استفاده نمی شوند
- استفاده ای برای افزایش یا کاهش پیشوندی ندارند
- پس مطابق با فلسفه سادگی MIPS, RISC تعداد کمی انشعاب اساسی دارد

انواع انشعاب در MIPS

- انشعاب شرطی: `beq/bne reg1, reg2, addr`
 - If $reg1 \neq reg2$, jump to $PC + addr$ (PC-relative)
- پرش ثابتی: `jr reg`
 - واکشی آدرس از ثبات مشخص شده و پرش به آن
- پرش غیر شرطی: `z addr`
- همیشه به آدرس پرش می کند (استفاده "pseudodirect" آدرس دهی)

دستورالعمل های انشعاب

- انشعاب های شرطی

- beq R1, R2, L1 # if R1 = R2 go to L1

- bne R1, R2, L1 # if R1 \neq R2 go to L1

- این ها دستورالعملهای R_type هستند

- انشعاب های غیر شرطی

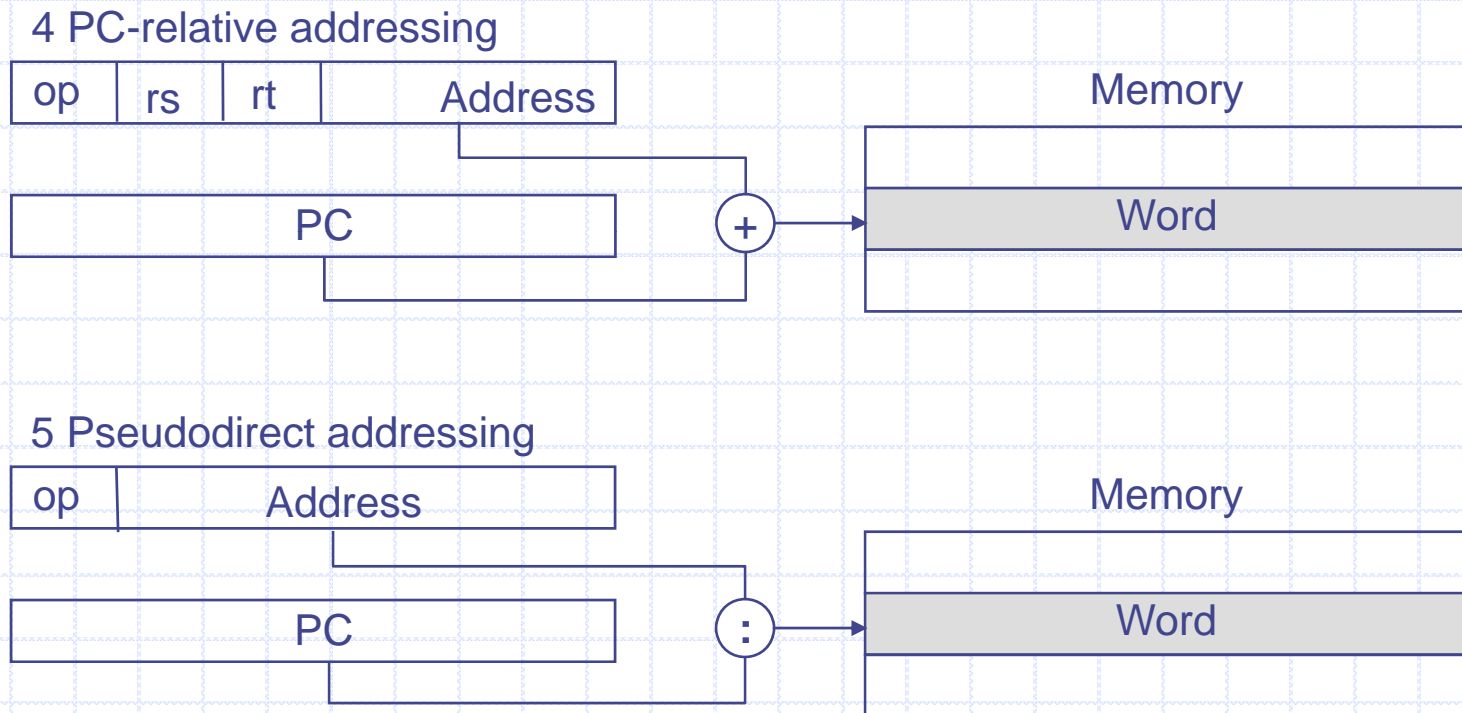
JR R8 # Jump based on register 8

Test if < 0

slt R1, R16, R17 # R1 gets 1 if R16 < R17
(slt: set-less-than)

bne R1, 0, less # branch to less if R1 \neq 0

مولد اهداف انشعاب در MIPS



مثال تلفیق سوئیچ

• کامپایل را دنبال کنید:

```
switch (k) {  
    case 0: f = f +  
1;    break;  
    case 1: f = f -  
2;    break;  
    case 3: f = -f;  
        break;  
}
```

Note the gap (case 2);

بدنه سوئیچ در MIPS

1 to r8 (f) immed. \$8, \$8, 1 add addi L0:
jump to Exit (break) j Exit
\$8, \$8, 2 subtract imm. 2 from r8 subi L1:
Another break j Exit
f = 0 - f sub \$8, \$0, \$8 L3:
Another break j Exit

ساختن جدول مراجعه در حافظه

1000	address of L0
1004	address of L1
1008	address of Exit
1012	address of L3

سوئیچ کامپایل شده برای MIPS

(Assume k in r13)

```
# set r14 if r13 < 0
# Go to Exit if k < 0
# set r14 if k < 4
# Go to Exit if k ≥ 4
# r14 = 2*k
# r14 = 4*k
# Base of table at 1000
# Jump to the address
```

```
slti $14, $13, 0
bne $14, $0, Exit
slti $14, $13, 4
beq $14, $0, Exit
add $14, $13, $13
add $14, $14, $14
lw $14, 1000($14)
jr $14
```

کامپایل کردن دستورات کنترلی دیگر

- حلقه ها:
- For,while: تست کردن قبل از بدنه حلقه. پرش کن به بعد از بدنه حلقه اگر شرط نادرست است.
- Do: تست کردن شرط در انتهای بدنه حلقه. پرش کن به ابتدا اگر درست است
- سوئیچ: ("case" فراخوانده می شود در بعضی زبانهای دیگر)
- ساختن یک جدول آدرس
- استفاده jr (or equiv. In non-MIPS processor)
- مطمئن باشید از چک کردن برای پیش فرض یا case های بدون استفاده

پشتیبانی دستورالعملهای فراخونی پردازش

Jump and link

آدرس پردازشگر
jal
یادداشت:

آدرس بازگشت در R31 ذخیره می شود

Return

jr R31

ذخیره کردن آدرس بازگشت روی پشته

R29 به عنوان اشاره گر پشته استفاده می شود
پارامتر زود گذر:

R₄ ~ R₇ استفاده می شوند برای اینها

روش آدرس دهی MIPS های دیگر

- عملوندهای ثابت یا فوری

LW R24, AddrConstant4(0)

Addi R3, 5 (I type)

ثابتها 16 بیتی هستند

Lui R8 255 Load_upper_immediate

j_type J 10000 10000 می رود به مکان

عملوندهای MIPS

Name	Example	Comments
32 registers	\$s0-\$s7, \$t0-\$t9, \$zero, \$a0-\$a3, \$v0-\$v1, \$gp, \$fp, \$sp, \$ra	Fast locations for data. In MIPS, data must be in registers to perform arithmetic. MIPS register \$zero always equals 0. \$gp (28) is the global pointer, \$sp(29) is the stack pointer, \$fp (30) is the frame pointer, and \$ra (31) is the return address.
2 ³⁰ memory words	Memory [0], Memory [4],..., Memory[42949672920]	Accessed only by data transfer instructions. MIPS uses byte addresses, so sequential words differ by 4. Memory holds data structures, such as arrays, and spilled register, such as those saved on procedure calls.

MIPS زبان اسمبلی

Category	Instruction	Example	Meaning	Comments
Arithmetic	add	add \$s1, \$s2, \$s3	$\$s1 = \$s2 + \$s3$	Three operands; data in registers
	subtract	sub \$s1, \$s2, \$s3	$\$s1 = \$s2 - \$s3$	Three operands; data in registers
Data transfer	load word	lw \$s1, 100 (\$s2)	$\&s1 = \text{Memory} [\$s2 + 100]$	Data from memory to register
	store word	sw \$s1, 100 (\$s2)	$\text{Memory} [\$s2 + 100] = \$s1$	Data from register to memory
Conditional branch	branch on equal	beq \$s1, \$s2, L	if ($\$s1 == \$s2$) go to L	Equal test and branch
	branch on not equal	bne \$s1, \$s2, L	if ($\$s1 \neq \$s2$) go to L	Not equal test and branch
	set on less than	slt \$s1, \$s2, \$s3	if ($\$s2 < \$s3$) $\$s1 = 1$; else $\$s1 = 0$	Compare less than: for beq, bne
Unconditional jump	jump	j 2500	go to 10000	jump to target address
	jump register	jr \$ra	go to \$ra	For switch, procedure return
	jump and link	jal 2500	$\$ra = PC + 4$; go to 1000	For procedure call

MIPS زبان ماشین

Name	Format	Example						Comments
add	R	0	18	19	17	0	32	add \$s1, \$s2, \$s3
sub	R	0	18	19	17	0	34	sub \$s1, \$s2, \$s3
lw	I	35	18	17	100			lw \$s1, 100 (\$s2)
sw	I	43	18	17	100			sw \$s1, 100 (\$s2)
beq	I	4	17	18	25			beq \$s1, \$s2, 100
bne	I	5	17	18	25			bne \$s1, \$s2, 100
slt	R	0	18	19	17	0	42	slt \$s1, \$s2, \$s3
j	J	2	2500					j 10000 (see section 3.8)
jr	R	0	31	0	0	0	8	jr \$ra
jal	J	3	2500					jal 10000 (see section 3.8)
field size		6 bits	5 bits	5 bits	5 bits	5 bits	6 bits	All MIPS instructions 32 bits
R-format	R	op	rs	rt	rd	shamt	funct	Arithmetic instruction format
I-format	I	op	rs	rt	address			Data transfer, branch format

فراخوانی تابع در MIPS

- فراخوانی تابع یک ساختار ضروری زبان های برنامه نویسی است.
برنامه یک تابع را برای انجام چند وظیفه فرامی خواند.
موقعی که تابع انجام میشود CPU دقیقاً از بعد از جایی که برنامه فراخوانی شده ادامه می دهد

فراخوانی تابع در MIPS

استفاده از دستور العمل Jal (“jump and link”)
Jal addr فقط J addr است به جز موارد زیر
“آدرس بازگشت” $R31 + 4 + \text{pc}$
این آدرس دستور العمل بعدی بعد از jal است
از $\$31$ برای بازگشت استفاده می کنیم

مثال فراخوانی

Callee

F: 1w \$6, 0(\$4)
1w \$7, 0(\$5)
sw \$6, 0(\$5)
sw \$7, 0(\$4)
jr \$31

Caller

add \$4, \$0, 1000
add \$5, \$0, 1200
add \$1, \$0, 1
sw \$1, 0(\$4)
add \$1, \$1, \$
sw \$1, 0 (\$5)
jal F
sub \$1, \$1, \$2

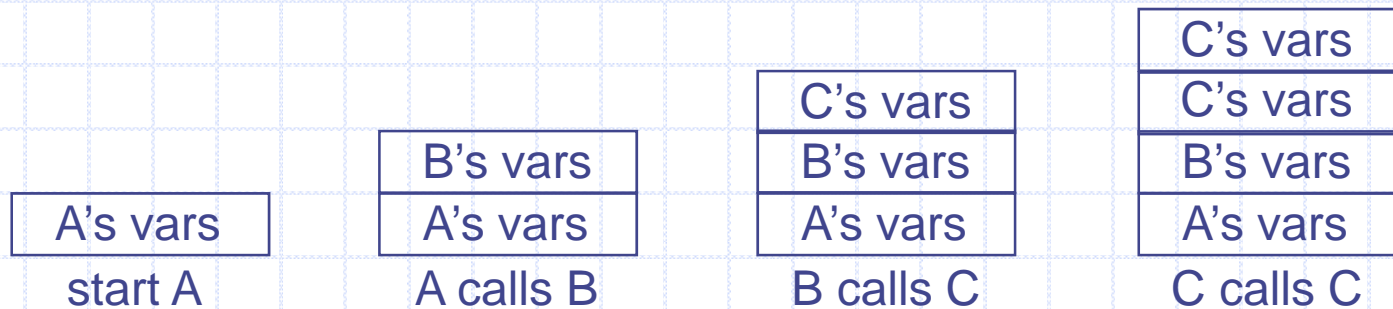
What does F do?

مشکلات فراخوانی تابع ها

این مثال درست کار می کند اما چه طور؟
تابع F تابع دیگر را فراخوانی می کند
فراخوانی کننده یک نکته مهم در ثباتهای R6 و R7 دارد؟
تابع فراخوانی شده خودش را فرامی خواند؟
هر نسخه از یک تابع باید کپی متغیر هایش را داشته
باشد این ها در یک پشته به عنوان یک ستون از قاب ها
مرتب شده اند.

مثال پشته

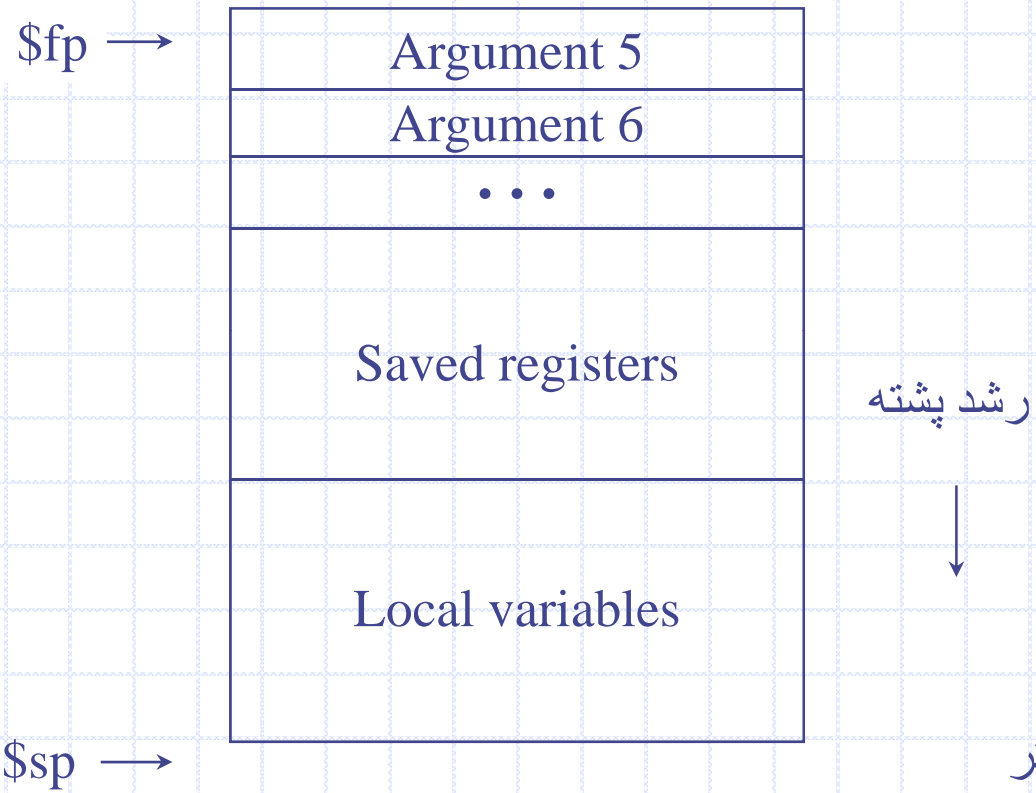
فرض کنید تابع A, B را فراخوانی کند و C, B را فراخوانی کند. تابع C خودش را فراخوانی می کند.



فراخوانی پردازش

- مکان پارامترها
- کنترل انتقال
- ذخیره سازی Acquire
- انجام دادن کار
- مکان بازگشت نتایج
- بازگشت کنترل بر گشتی به فراخوان

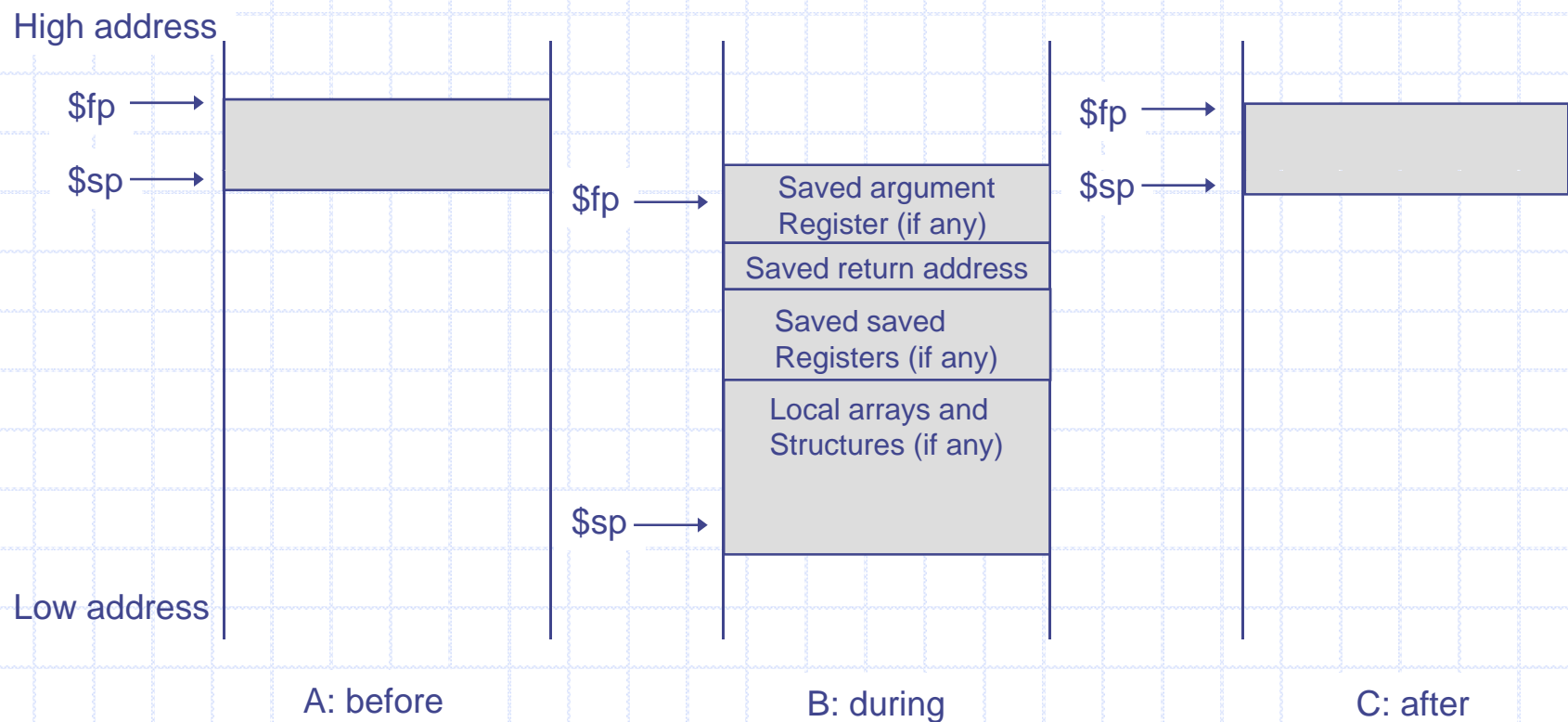
آدرسهای حافظه بالاتر



آدرسهای حافظه پایینتر

آرایش قاب پشته

جزئیات پشته



(From Patterson and Hennessy, p. 139; COPYRIGHT 1988 MORGAN KAUFMANN PUBLISHERS, INC. ALL RIGHTS RESERVED)

قراردادهای ثابت استاندارد

- 32 تا ثابت صحیح در MIPS عام منظوره هستند_ هر کدام می توانند به عنوان یک عملوند یا نتیجه یک محاسبه عملیاتی استفاده شوند
- اما ساخت تکه های متفاوت از نرم افزار که باهم کار می کنند آسانتر است اگر قراردادهای معینی که دنبال می شود درباره آن ثباتهایی باشد که برای آن اهداف استفاده می شود
- این قراردادها معمولاً به وسیله فروشنده پیشنهاد می شوند و به وسیله کامپایلرها پشتیبانی می شوند

پارامترهای میانی

- ذخیره کردن فراخوانی کننده: پردازنده فراخوان (caller) مسئول است برای ذخیره کردن و دوباره ذخیره کردن هر کدام از ثباتهایی که باید در طی فراخوانی نگهداری شوند. پردازنده فراخوانی شده (callee) می تواند هر ثباتی را بدون محدودیت تغییر دهد.
 - ذخیره کردن فراخوانی شده: پردازنده ای که فراخوانی شده مسئول است برای ذخیره کردن و دوباره ذخیره کردن هر کدام از ثباتهایی که ممکن است استفاده شوند.
- فراخوان ثباتها را بدون هیچ نگرانی درباره دوباره ذخیره کردن آنها بعد از یک فراخوانی استفاده می کند.

ثباتهای ذخیره سازی

- اگر شما یک تابع را فراخوانی کنید هر آنچه شما در $\$S7$ تا $\$S0$ دارید در آنجا تضمین شده زمانی که تابع برمی گردد به شما
- اما ثباتهای $\$t9$ تا $\$t0$ fair game هستند به وسیله تابع دوباره استفاده می شود
- چه گزینه‌هایی هستند؟
- هیچ ذخیره کردن؟
- همه چیز ذخیره کردن؟

اعلام یک تابع در MIPS

قرار بدهید آرگومانهای تابع را در a_0 تا a_3
اگر یک چیز مهم در t_9 باشد اکنون آنها را انتقال می دهد!
jal به تابع

ذخیره کردن s_0 تا s_7 و r اگر لازم باشد شما کاری انجام دهید
و مقادیر بازگشت را (اگر داشته باشد) در v_0 تا v_1
دوباره ذخیره میکنیم s_0 تا s_7 و r اگر لزومی داشته باشد
ادامه دهید تا جایی که تمام شود

سوال: چه کسی آن را به عنوان یک طرح می سازد؟
چه اتفاقی می افتد اگر ثباتهای ما تمام شود؟
واقعا چه چیزی عمل ذخیره کردن را انجام می دهد؟

فصل سوم

نمایش اعداد و محاسباتی کامپیوتر

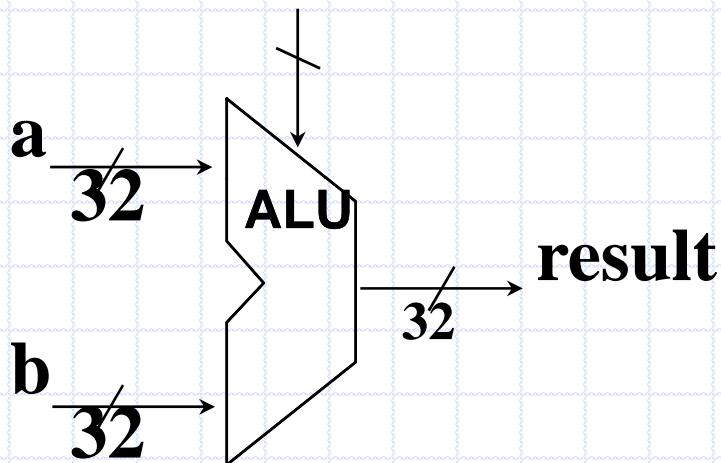
چکیده

- ◆ مقدمه
- ◆ نمایش عددی و مکمل دو
- ◆ طراحی واحد محاسبه، منطق
- ◆ ضرب و تقسیم صحیح
- ◆ اعداد اعشاری

مرور

◆ بازبینی

Decoded opcode ■ مراحل اجرای یک دستور العمل



◆ فیلهای یک دستور العمل

■ کارایی؟

■ مبانی عملوندها؟

◆ آنچه پیش روی ماست:

■ نمایش عددی

■ الگوریتم های محاسباتی

■ پیاده سازی سخت افزاری

■ دستور العملها

نمایش عددی

◆ آنچه پیچیده به نظر می رسد:

- اعداد منفی را چگونه نمایش دهیم
- دامنه (Range) اعداد چیست؟
- اگر یک عدد خارج از محدوده بود؟
- در مورد اعداد منطقی و غیر منطقی چطور؟
- سخت افزار چگونه این اعداد را جمع، تفریق، ضرب و تقسیم میکند؟

بیت ها و اعداد

◆ بیتها فقط بیت هستند: هیچ معنای ذاتی ندارند.
■ قراردادهای رابطه بین بیتها و اعداد را تعیین می کنند.

■ سخت افزار چگونه تشخیص می دهد که چه قراردادی در حال استفاده است؟

Add \$s1, \$s2, \$s3

Addu \$s1, \$s2, \$s3

علامت \$s1، \$s2 و \$s3 چیست؟

نمایش عددی

- ◆ صحیح بدون علامت
- ◆ صحیح با علامت
- ◆ BCD دهدهی کد شده به باینری
- ◆ ممیز ثابت
- ◆ ممیز شناور
- ◆ انواع دیگر داده:
 - کاراکترها (اسکی و یونیکد)
 - پیکسل ها (گرافیک ها)
 - گروهی از بیتها

صحيح بدون علامت

◆ چرا صحيح بدون علامت؟

■ دستيابی به حافظه PC, SP, RA,

■ In C, unsigned int

◆ چگونه نمایش می دهیم؟

■ انگشتهای خود را بشمارید.

■ مبنای اولیه سیستم اعداد

◆ دهدهی: 10 نماد مختلف: 0 1 2 3 4 5 6 7 8 9

$$\sum_{i=0}^{n-1} (d_i * 10^i) = \text{عدد}$$

صحيح بدون علامت (مبنای اولیه سیستم اعداد)

◆ در حالت کلی، در یک سیستم عددی مبنای K

$$\sum_{i=0}^{n-1} (d_i * k^i) = \text{عدد}$$

- ◆ چه تعداد d_i متفاوت؟
- ◆ بزرگترین عدد کدامست؟
- ◆ کوچکترین عدد کدامست؟

◆ مبنای خاص:

- مبنای 2 (Binary)
- مبنای 8 (octal)
- مبنای 16 (hexadecimal)
- ◆ آنها چه تعداد نماد مختلف دارند؟

صحيح بدون علامت (مبنای اولیه سیستم اعداد)

◆ چگونه می توان مبنای دو را به مبنای 8 و 16 تبدیل کرد؟
◆ مثال:

$$010100 = (0 * 2^2 + 1 * 2^1 + 0 * 2^0) * 2^3 + (1 * 2^2 + 0 * 2^1 + 0 * 2^0) \\ = 2 * 8^1 + 4 * 8^0 = 24_8$$

Hexadecimal: 1 9 4 8 B 6
Binary: 000110010100100010110110
Octal: 0 6 2 4 4 2 6 6

صحيح بدون علامت (BCD: binary coded decimal)

نمایش: 127_{10} ◆

- چه تعداد بیت برای BCD مورد نیاز است؟
برای دودویی چطور؟
- فضای ذخیره سازی بکار گرفته شده؟

ضمیمه: ◆

- چگونه می توان carry را مشخص نمود؟
“If $a_i + b_i > 10$ ” vs. “If $a_i + b_i > 2$ ”

کدامیک برای کامپیوتر کارا تر است؟ ◆

BCD	Number
0000	0
0001	1
0010	2
0011	3
0100	4
0101	5
0110	6
0111	7
1000	8
1001	9
1010-1111	No use

صحيح با علامت (علامت مقدار)

12: -(sign) 12(absolute value) ◆

■ يك بيت علامت مجزا

■ يك مقدار

◆ براي سخت افزار

■ بيت علامت را كجا قرار دهيم؟ چپ يا راست؟

■ علامت حاصلجمع را چگونه تعيين مي كنيم؟ (مثلا يك مرحله اضافي)

■ علامت صفر چيست؟ مثبت يا منفي؟

مکمل 2

اعداد و نمایش آنها یکی نیست.

- نمایش عدد، یک قرارداد تعریف شده برای بیان آن است.
- یک عدد می تواند دارای نمایشهای متفاوتی باشد مثلا BCD مبنای 16، مکمل دو و غیره.
- مکمل دو: یک نوع از نمایش که بصورت زیر تعریف شده است:

عدد	نمایش
$b_{n-1} * (-2^{n-1}) + \sum_{i=0}^{n-2} (b_i * 2^i)$	$b_{n-1} b_{n-2} \dots b_1 b_0$

مکمل دو

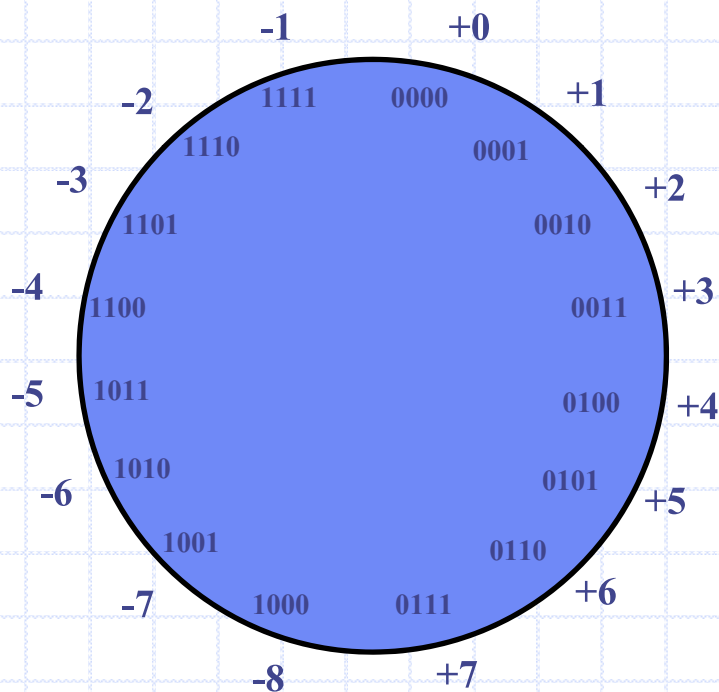
$$b_{n-1} * (-2^{n-1}) + \sum_{i=0}^{n-2} (b_i * 2^i)$$

$b_{n-1} b_{n-2} \dots b_1 b_0$

0	0000	-8	1000
1	0001	-7	1001
2	0010	-6	1010
3	0011	-5	1011
4	0100	-4	1100
5	0101	-3	1101
6	0110	-2	1110
7	0111	-1	1111

صحيح با علامت (مکمل دو)

(برای 4 بیت)



مرجع: Katz: Contemporary
Logic Design, p243

چرا 1111 و 0000 همسایه اند؟
مفهوم بیت سوم چیست؟
علامت صفر چیست؟
کوچکترین عدد کدامست؟
بزرگترین عدد کدامست؟

مکمل دو

نمایش یک عدد منفی $-x$

■ از نمایش عدد x تا نمایش $-x$

عدد	نمایش
x	$b_{n-1} b_{n-2} \dots b_1 b_0$
$-x$	$b'_{n-1} b'_{n-2} \dots b'_1 b'_0$

هر b_i را معکوس کرده، سپس با 1 جمع کنید.

■ اثبات درستی

■ **میانبر:** از با ارزشترین بیت تا کم ارزشترین 1 هر بیت را معکوس کنید و کم ارزشترین یک را 1 نگاه دارید.

Example: $-(0100) = 1011 + 1 = 1100$

Example: $100_{10} = 01100100 \rightarrow \underline{10011100} = -100_{10}$

گنجاندن اعداد کوچکتر در بیت‌های بیشتر

◆ چرا لازم است؟

- مقایسه یک integer با یک long integer : تغییر نوع
- بارگذاری یک بایت در یک کلمه
- برای اضافه کردن یک بخش فوری به یک عدد 32 بیتی

◆ واحد محاسبه و منطق MIPS فقط با مقادیر موجود در رجیسترهای 32 بیتی کار می‌کند.

- چگونه می‌توان با اندازه‌های کوچکتر کار کرد؟
- در مورد اندازه‌های بزرگتر چطور؟

گنجاندن اعداد کوچکتر در بیت‌های بیشتر

◆ Example(2-bit): $10 = 1 * (-2^1) + 0 * 2^0 = -2_{10}$

نمایش 2- مبنای 10 در 4 بیت کلمه: آیا می‌توانید آنرا انجام دهید؟ چه چیزی کشف می‌شود؟

◆ تکرار بیت علامت در بیت‌های دیگر (گسترش علامت)

$$0010 \rightarrow 0000\ 0010$$

$$1010 \rightarrow 1111\ 1010$$

($1 \times (2^{-7}) + 111\ 1010 = ?$) نکته! اثبات؟

جمع و تفریق

جمع : درست مانند دوران شیرین دبستان

$$\begin{array}{r} 0110 \\ + 0001 \\ \hline \end{array}$$

تفریق : $a-b=a+(-b)$

■ اعمال مکمل دو ساده است

$$0111-0110=0111$$

$$\begin{array}{r} \text{---} \\ + 1010 \\ \hline 10001 \end{array}$$

■ بنابراین $0111-0110=0001$

■ آیا شما مخصوصاً بیت‌های علامت را دستکاری نموده اید؟

جمع و تفریق

❖ صبر کنید! یک بیت اضافی گم شده است!

$$0111-0110= 0111$$

چرا 1 می تواند حذف شود؟ $+1010$

$$10001$$

❖ زیرا

اما رفتار می کند در جمع مانند

$$\begin{array}{r} 0*2^3 \\ + 1*(-2^3) \\ \hline 1*2^3 \\ \hline 0 \end{array}$$

$$\begin{array}{r} 0*2^3 \\ + 1*2^3 \\ \hline 1*2^3 \\ \hline 10 \end{array}$$

■ از دیدگاه 1 بیت : $1-1=1+1=0$

■ بنابراین تا زمانی که شما carry بوجود آمده بوسیله $1+1$ (1 سبز) را در نظر

نگرفته اید، درست انجام داده اید.

جمع و تفریق

سئوالات:

- آیا شما یک جمع کننده و یک تفریق کننده دارید؟
- در مورد اعداد بدون علامت چطور؟ (یک جمع کننده دیگر؟)
- مزایای استفاده از مکمل 2
- تفریق می تواند از منطق مشترکی با جمع استفاده کند
- بیت علامت می تواند مانند یک بیت معمولی عدد رفتار کند
- اینها نکات زیرکانه ای هستند.

سرریزی

◆ جمع:

$$\begin{array}{r} 0111 \\ +1000 \\ \hline \end{array} \quad \begin{array}{r} 1111 \\ + 0110 \\ \hline \end{array}$$

علامت ورودی چیست؟

علامت خروجی چیست؟

◆ آیا امکان وقوع سرریزی در جمع یک عدد مثبت و یک عدد منفی وجود دارد؟

◆ امکان وقوع سرریز با صفر وجود دارد؟
■ اعمال $A+B$ و $A-B$ را ملاحظه کنید.

سرریزی

جمع 

$$\begin{array}{r} 1111 \\ + 1000 \\ \hline 10111 \end{array} \quad \begin{array}{r} 0111 \\ + 0110 \\ \hline 1101 \end{array}$$

علامت ورودی چیست؟

علامت خروجی چیست؟

کشف سرریزی

سرریزی زمانی اتفاق می افتد که:

- جمع دو عدد منفی یک عدد مثبت بدهد.
- جمع دو عدد مثبت یک عدد منفی بدهد.
- یا یک عدد منفی را از مثبت تفریق کنیم و نتیجه منفی بگیریم
- یا یک عدد مثبت را از منفی تفریق کنیم و نتیجه مثبت بگیریم

در مورد جمع و تفریق اعداد بدون علامت چطور؟

کار CPU در برابر سرریزی چیست؟

◆ نادیده بگیرد؟

- هیچ وقت نمی خواهد سرریزی آشکار شود
- Addu, addiu, subu (MIPS: عدم تولید سرریز)
- تولید یک تله که بر اساس آن برنامه نویس بتواند با آن کار کند.
- یک استثنا (وقفه) رخ دهد
- ◆ کنترل به آدرس از پیش تعریف شده استثنا پرش کند
- ◆ ذخیره آدرس مکان وقفه برای امکان بازگشت
- ◆ جمع ، تفریق

دستور العمل ها

◆ سنجش

■ اعداد بدون علامت

sltu: set on less than unsigned

sltiu: set on less than immediate unsigned

■ اعداد با علامت

slt: set on less than

slti: set on less than immediate

◆ مثال: مقادير \$s0 و \$s1 چیست؟

\$s0 = 1111 1111 1111 1111 1111 1111 1111 1111

\$s1 = 0000 0000 0000 0000 0000 0000 0000 0001

(1) slt \$t0, \$s0, \$s1 (2) sltu \$t1, \$s0, \$s1

مثال

◆ مقادیر s_0 و s_1 چیست؟

$s_0 = 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111$

$s_1 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0001$

پاسخ:

-- $1 \text{slt } t_0, s_0, s_1$ است t_0 اگر هر دو با علامت باشند

-- $0 \text{sltu } t_1, s_0, s_1$ است t_1 اگر هر دو بدون علامت باشند

دستور العمل ها

بارگذاري/ذخيره سازي ◆

lb: load byte ■

lbu: load byte unsigned ■

lb \$s1, 100(\$s2) Example: ◆

مقادير \$s0 چيست

When memory[\$s2+100] = 0000 0000? ◆

When memory[\$s2+100] = 0000 0001? ◆

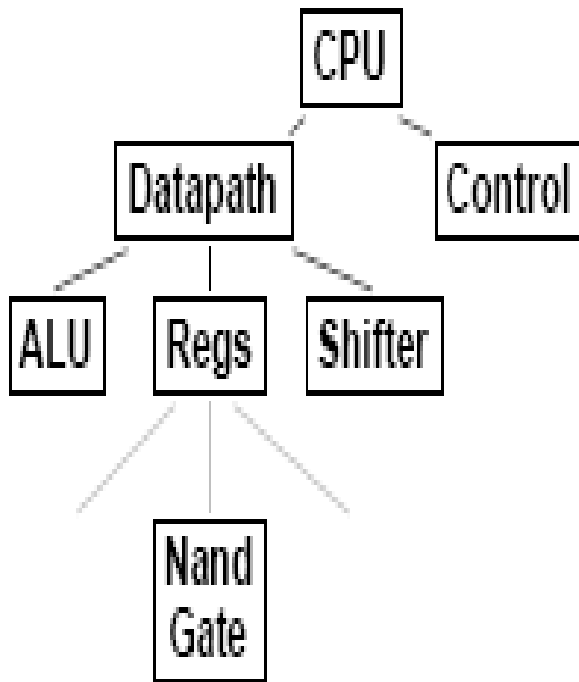
When memory[\$s2+100] = 1111 1111? ◆

lbu \$s1, 100(\$s2) Example: ◆

مقادير \$s1 چيست

when memory[\$s2+100] = 1111 1111? ◆

طراحی پردازش



♦ طراحی

■ اجزا

■ آنها چگونه در کنار هم قرار می گیرند

♦ تجزیه بالا به پایین بخشهای مختلط

♦ ترکیب پایین به بالای بخشهای ساده

مسئله: طراحی واحد محاسبه و منطق

عملیات 

add, addu, sub, subu, addi, addiu ■

◆ جمع کننده تفریق کننده مکمل 2 با آشکار سازی سرریز

and, or, andi, ori ■

◆ اعمال بیتی

مجموع اعمال = 10

طراحی: روش تقسیم و غلبه

◆ شکستن مسئله به بخشهای ساده تر

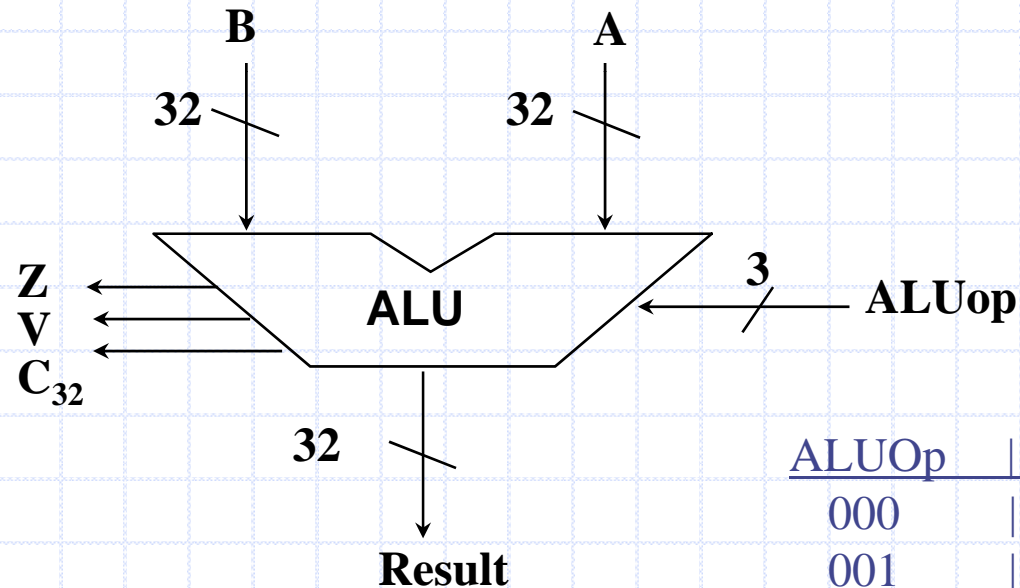
- کار با اجزا
- قرار دادن قطعات با یکدیگر
- بازبینی راهکارها در نهایت

◆ مثال: جدا کردن دستورات فوری از بقیه دستورات

- پردازش فوریها قبل از ALU
- ◆ اکنون ورودیهای ALU یک شکل هستند.
- 6 عمل غیر فوری باقی مانده است
- ◆ نیازمند 3 بیت برای تعیین مد ALU است

محیط ALU (ALU Interface)

♦ ما یک ALU 32 بیتی طراحی خواهیم کرد با محیط زیر:



Z = 1, if Result=0

V = 1, if Overflow

C₃₂ = 1, if Carry-Out

ALUOp	Function
000	AND
001	OR
010	ADD
110	SUBTRACT
111	SET-ON-LESS-THAN

طراحی: کاهش مسئله به یک مسئله ساده تر

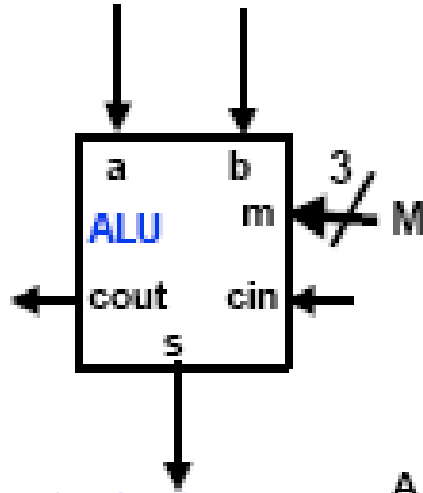
◆ برای این مسئله، ما 32 ALU بیتی را به بخش ساده تر یک بیتی کاهش می دهیم.

■ مسائل ترکیبی بزرگ را به یک مسئله ترکیبی کوچکتر تغییر می دهیم

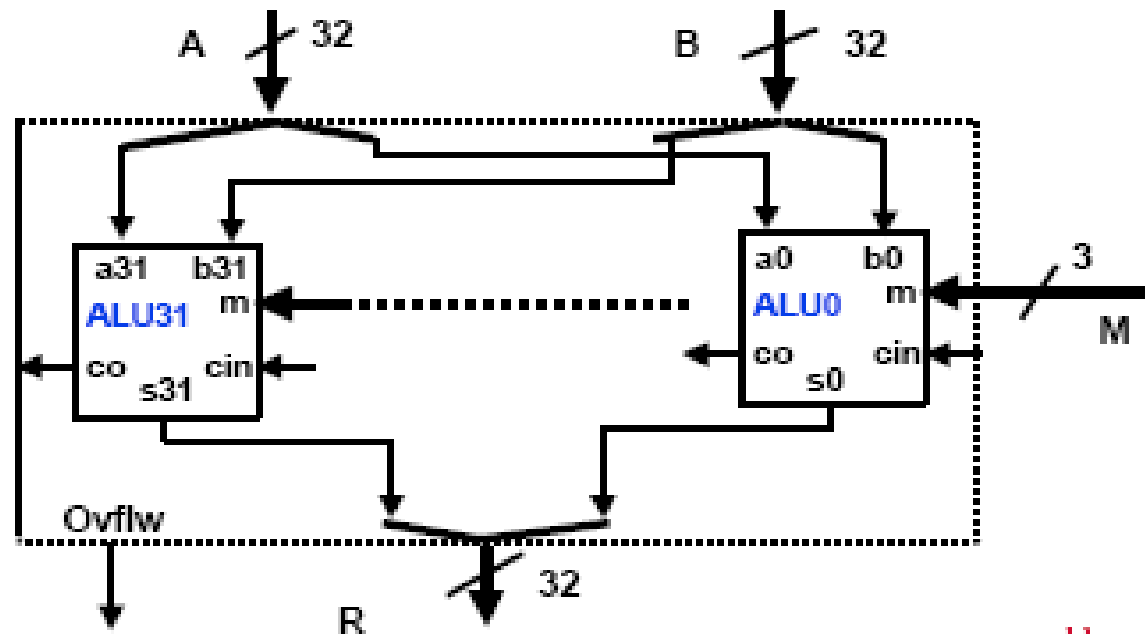
◆ بخشها را برای حل مسئله بزرگ در کنار هم قرار می دهیم.

طراحی با بلاک دیاگرام سطح پایین تر

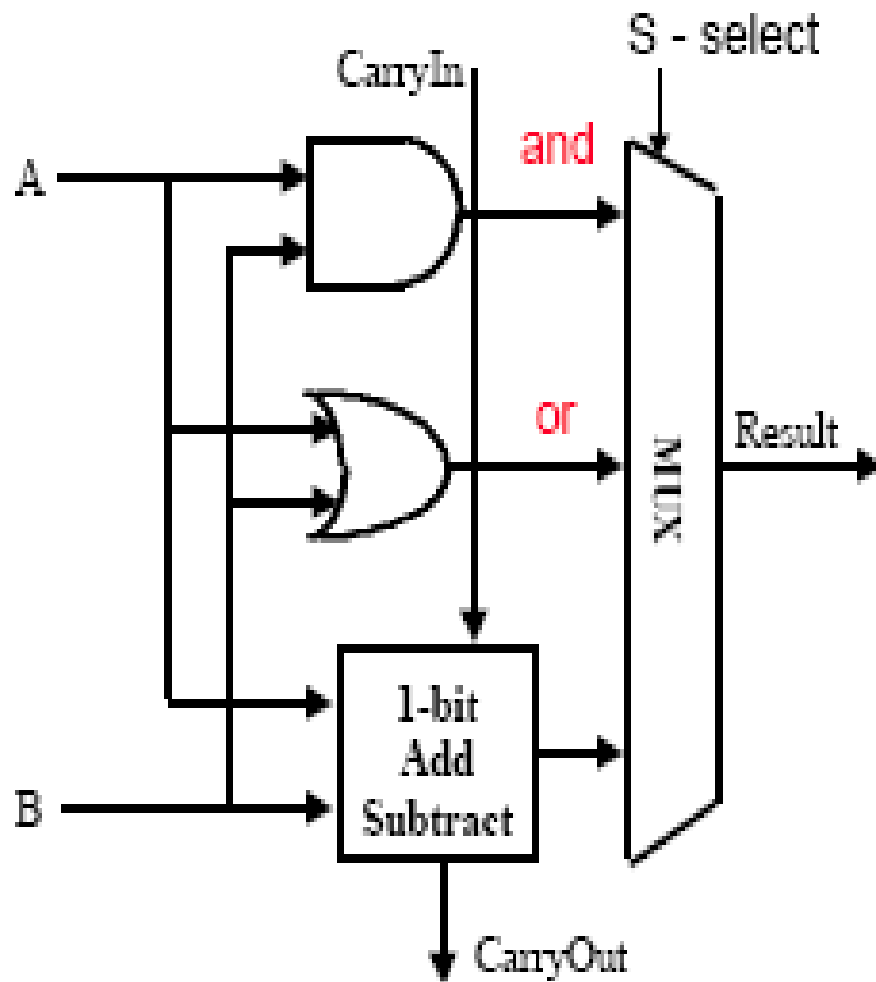
بلاک ALU
یک بیتی



برای یک ALU 32 بیتی
32 بار تکرار می کنیم



بلاک ALU یک بیتی



◆ تقسیم به بلاکهای جدا و مستقل

■ منطق

■ محاسبه

◆ تکمیل هر بلاک در این مرحله برای

تخلیص بیشتر

■ تکمیل بلاکهای منطقی

■ تکمیل انتخاب عمل

■ تجزیه بلاکهای محاسباتی به بخشهای

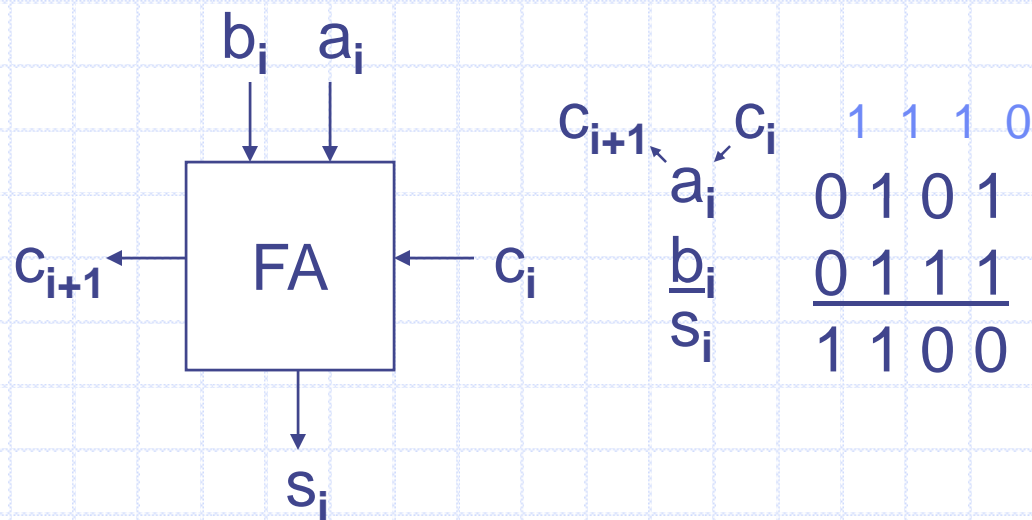
ساده تر

تمام جمع کننده Full Adder

◆ تمام جمع کننده، یک بلاک ساخته بنیادی در ALU است.

◆ یک تمام جمع کننده، جمع یک بیت را انجام می دهد.

$$a_i + b_i + c_i = 2c_{i+1} + s_i$$



جدول درستی تمام جمع کننده

◆ S_i ، "1" است اگر یک، تک عدد از ورودیها یک باشد.

◆ C_{i+1} ، "1" است اگر 2 عدد (یا بیشتر از) ورودیها 1 باشند.

a_i	b_i	c_i	C_{i+1}	S_i
0	0	0	0	0
0	0	1	0	1
0	1	0	0	1
0	1	1	1	0
1	0	0	0	1
1	0	1	1	0
1	1	0	1	0
1	1	1	1	1

$$s_i = \bar{a}_i \bar{b}_i c_i + \bar{a}_i b_i \bar{c}_i + a_i \bar{b}_i \bar{c}_i + a_i b_i c_i$$

$$s_i = a_i \otimes b_i \otimes c_i$$

$$c_{i+1} = \bar{a}_i b_i c_i + a_i \bar{b}_i c_i + a_i b_i \bar{c}_i + a_i b_i c_i$$

$$c_{i+1} = a_i b_i + a_i c_i + b_i c_i$$

$$c_{i+1} = a_i b_i + c_i (a_i + b_i)$$

$$c_{i+1} = a_i b_i + c_i (a_i \otimes b_i)$$

طراحی تمام جمع کننده

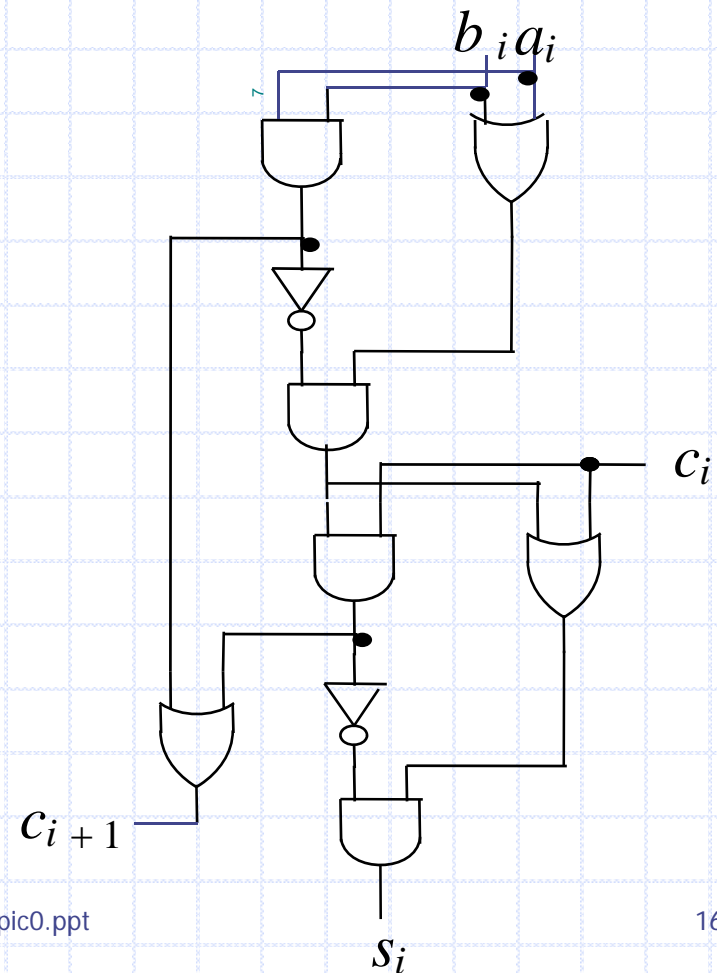
یک پیاده سازی ممکن از یک تمام جمع کننده که 9 gate بکار برده

است

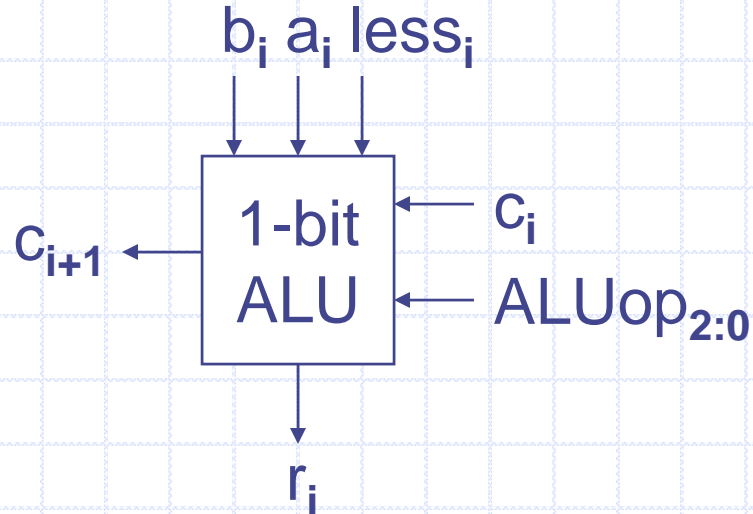
$$s_i = a_i \otimes b_i \otimes c_i$$

$$c_{i+1} = a_i b_i + c_i(a_i \otimes b_i)$$

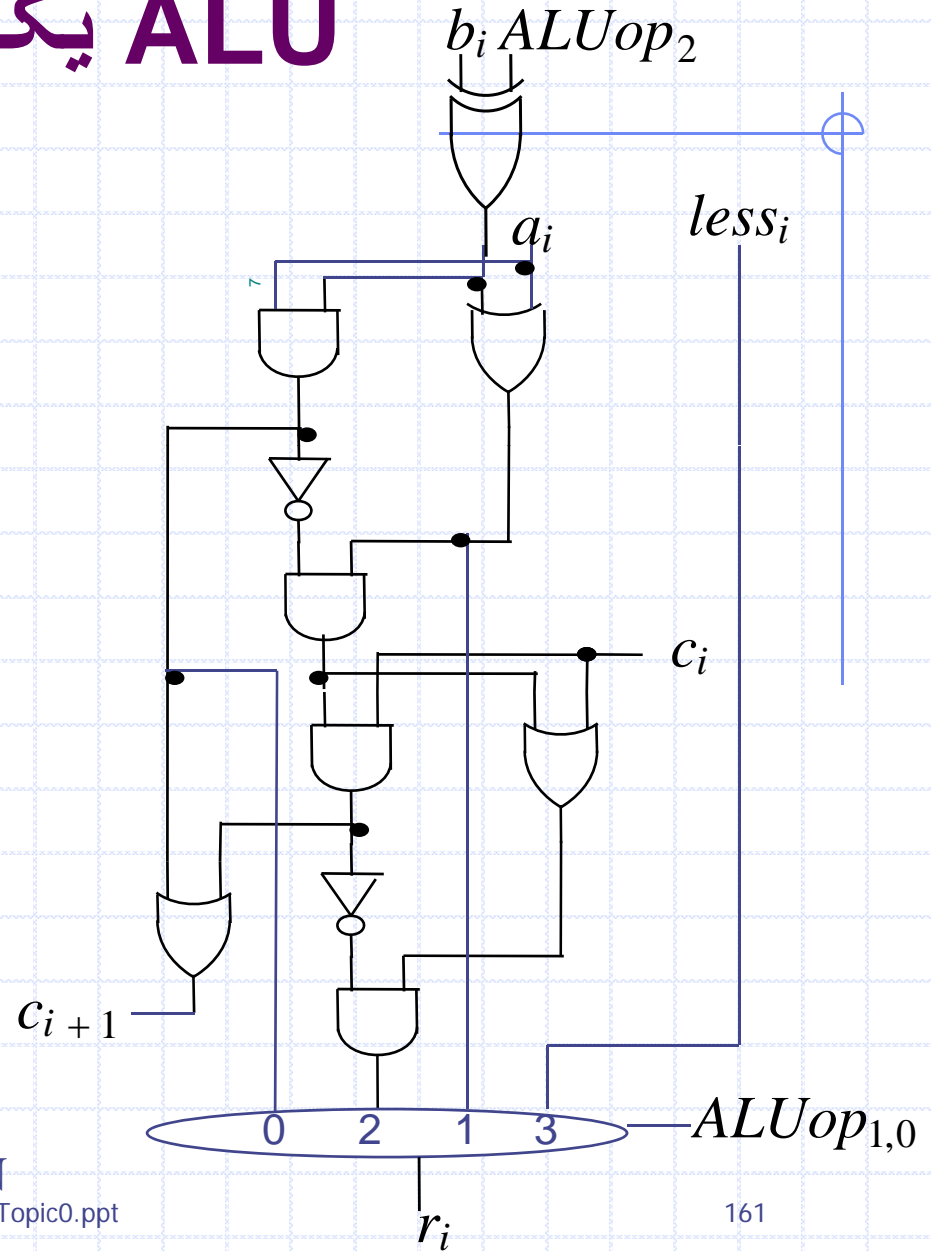
$$a_i \otimes b_i = \overline{(a_i + b_i)} a_i b_i$$



ALU یک بیتی



ALUOp	Function
000	AND
001	OR
010	ADD
110	SUBTRACT
111	SET-ON-LESS-THAN



ALU یک بیتی برای MSB

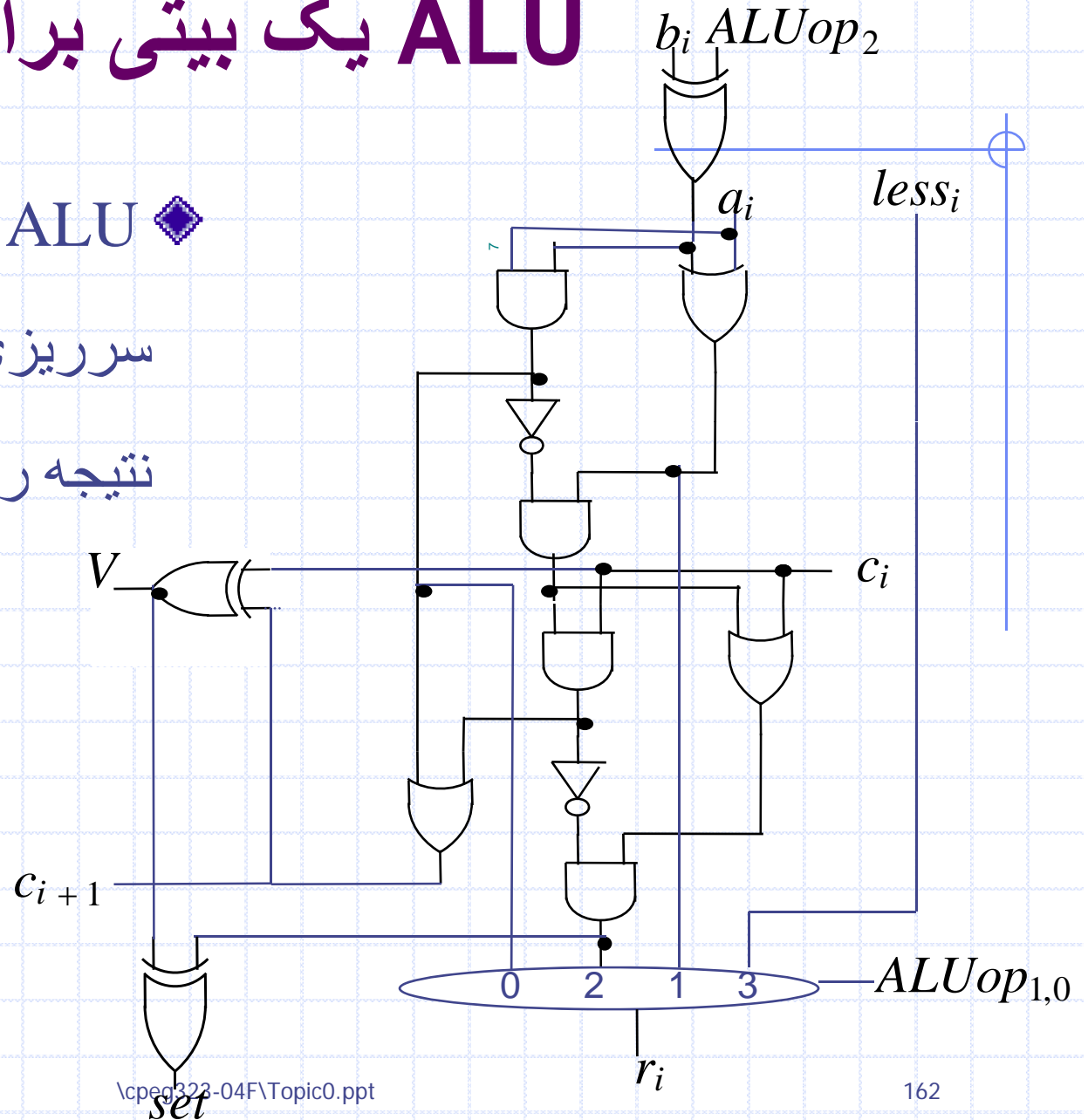
ALU برای MSB باید

سرریزی را آشکار و علامت

نتیجه را نیز نشان دهد.

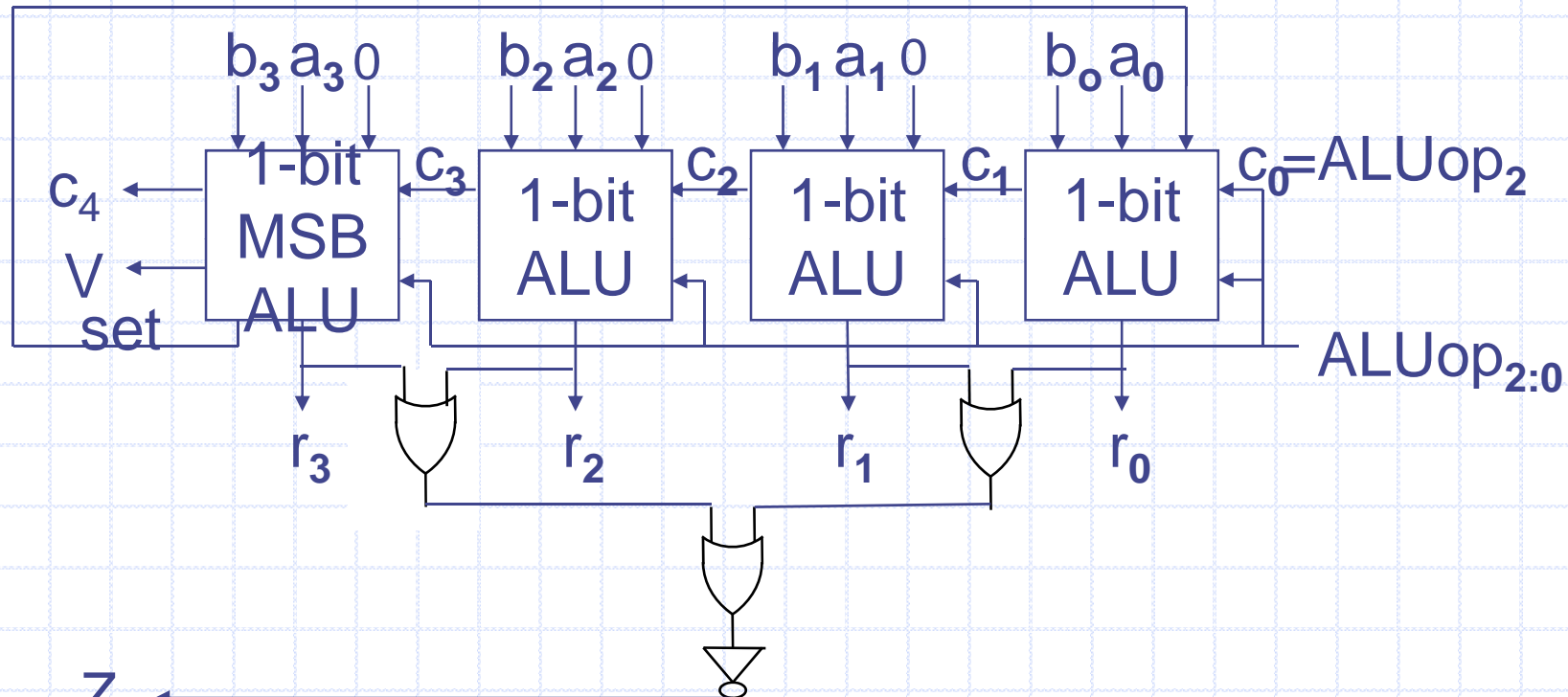
$$V = C_n \otimes C_{n-1}$$

$$set = (A < B)$$



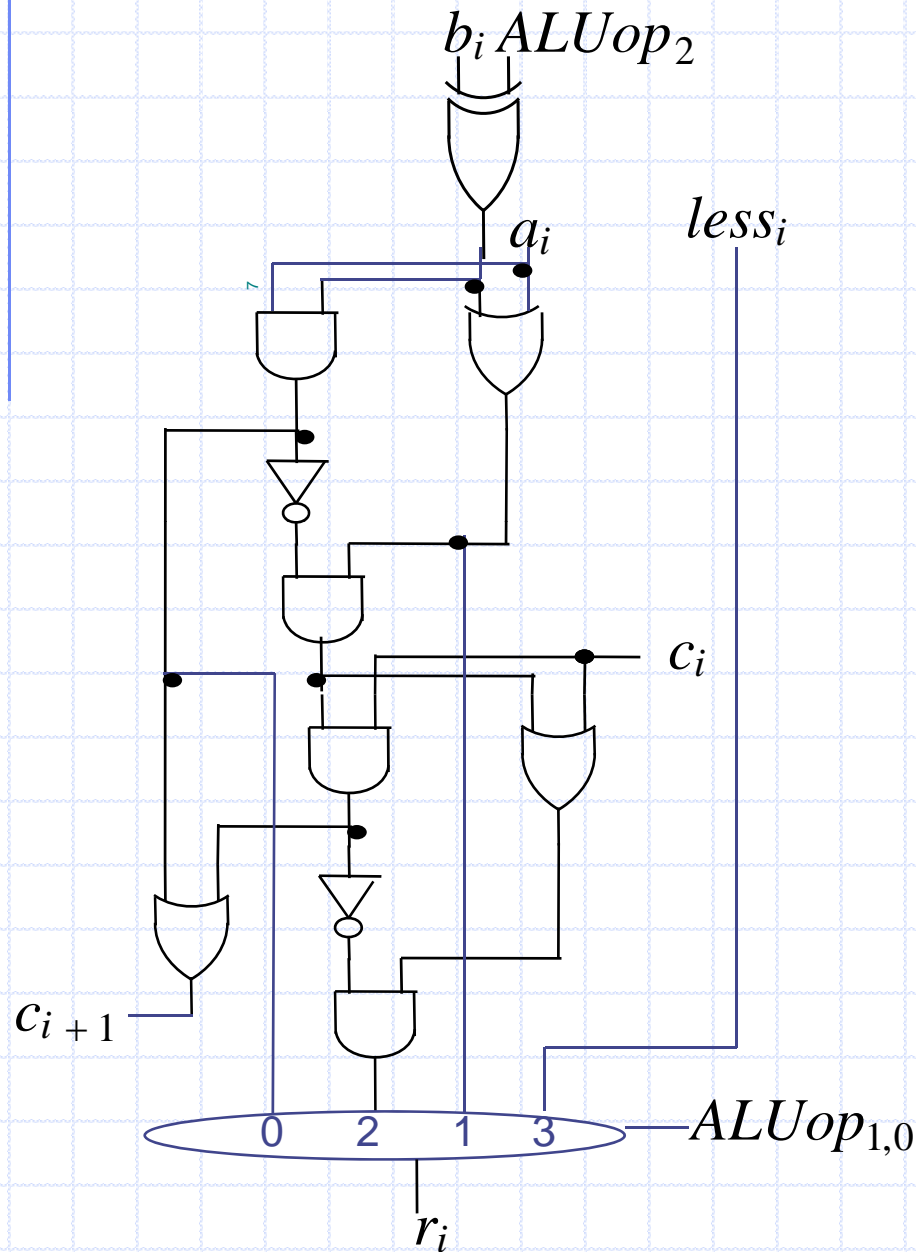
واحدهای محاسبه منطق (ALU) بزرگتر

سه تا واحد محاسبه و منطق یک بیتی و یک واحد محاسبه و منطق یک بیتی
 MSB و 4 تا گیت NOR میتوانند باهم الحاق شوند و تشکیل یک واحد
 محاسبه و منطق (ALU) 4 بیتی را دهند.



31 عدد واحد محاسبه و منطق یک بیتی و یک واحد محاسبه و منطق یک
 بیتی MSB و 32 تا گیت NOR می توانند باهم الحاق شوند و تشکیل یک
 واحد محاسبه و منطق (ALU) 32 بیتی را می دهند

محاسبه تعداد گیتها



فرض کنید

- مالتی پالکسر 4 ورودی = 5 گیت
- گیت XOR = 3 گیت
- گیت AND/OR = 1 گیت
- وارونه ساز = 5/1 گیت

چند تا گیت نیاز است به وسیله

- یک واحد محاسبه منطق 1 بیتی؟
- یک واحد محاسبه منطق 4 بیتی؟
- یک واحد محاسبه منطق 32 بیتی؟
- یک واحد محاسبه منطق n بیتی؟

محاسبه تعداد گیتها

◆ فرض کنید

- مالتی پالکسر 4 ورودی = 5 گیت
- گیت XOR = 3 گیت
- گیت AND/OR = 1 گیت
- وارونه ساز = 5/1 گیت

◆ چند تا گیت نیاز است به وسیله

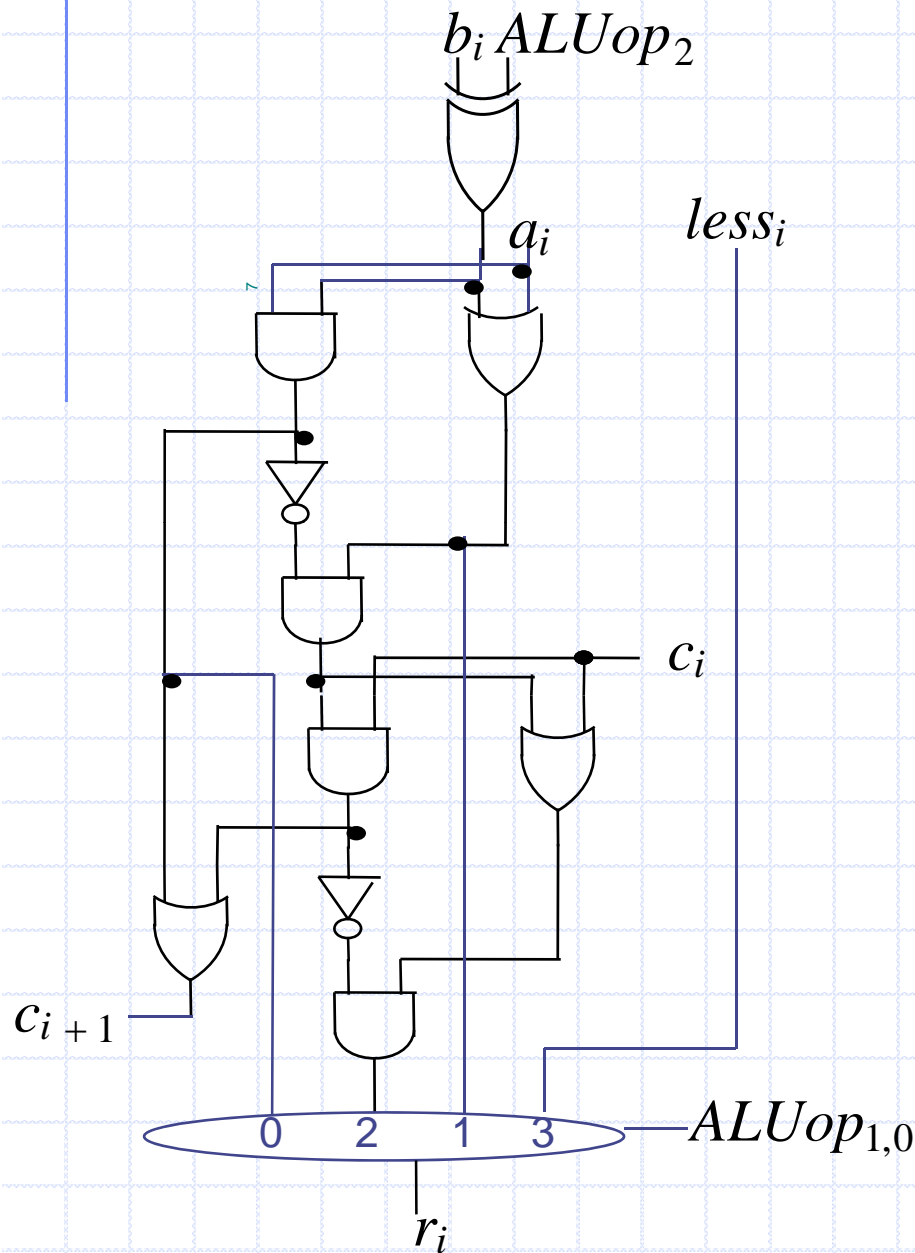
- یک واحد محاسبه منطق 1 بیتی؟ 16
- یک واحد محاسبه منطق 4 بیتی؟ $4 * 16$
- یک واحد محاسبه منطق 32 بیتی؟ $32 * 16$
- یک واحد محاسبه منطق n بیتی؟ $n * 16$

در یک واحد محاسبه و منطق n

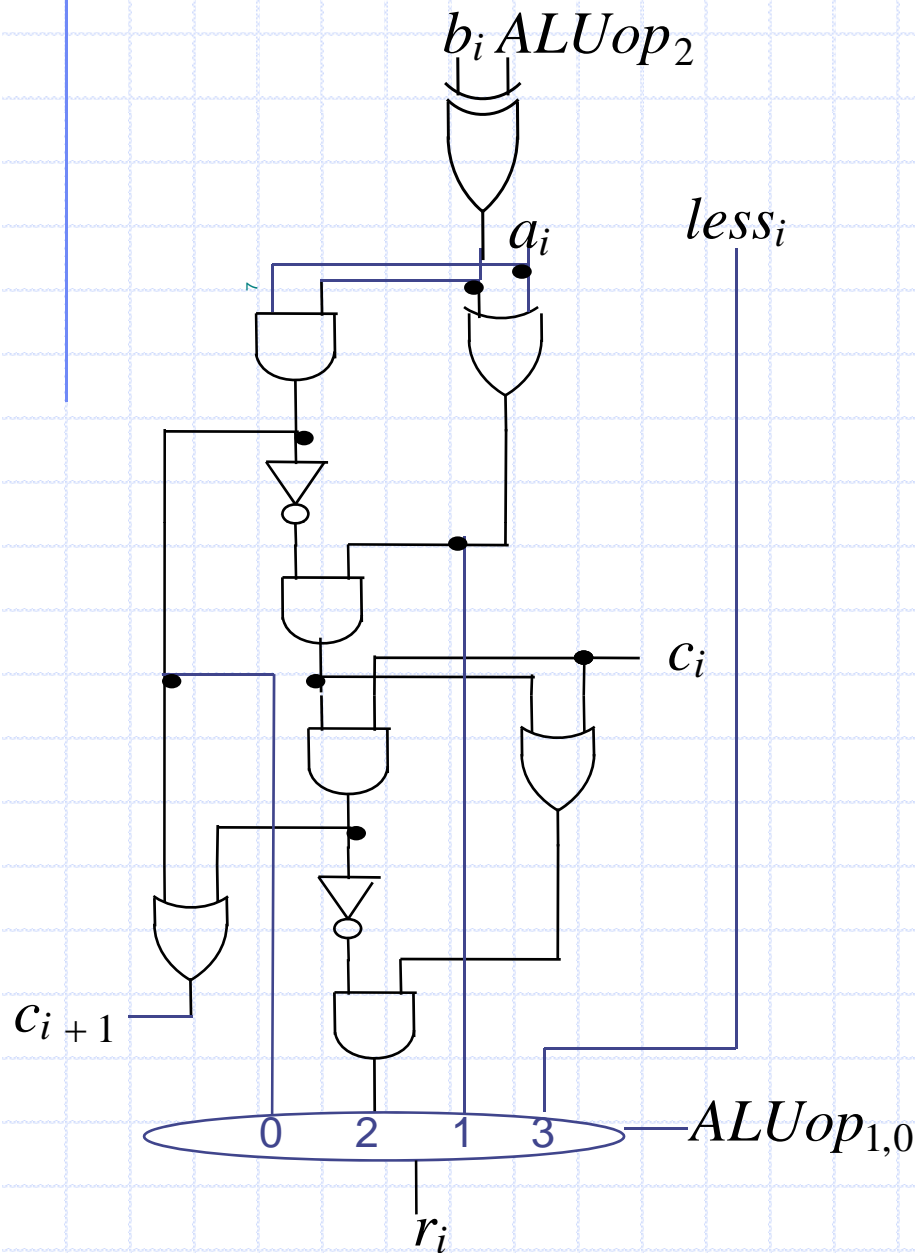
بیتی (n-1) عدد گیت OR با 2

ورودی و یک وارونه ساز و یک

گیت XOR



تاخیر گیت ها



فرض کنید تاخیر های

- مالتی پالکسر 4 ورودی $2t =$
- گیت XOR $2t =$
- گیت AND/OR $1t =$
- وارونه ساز $1t =$

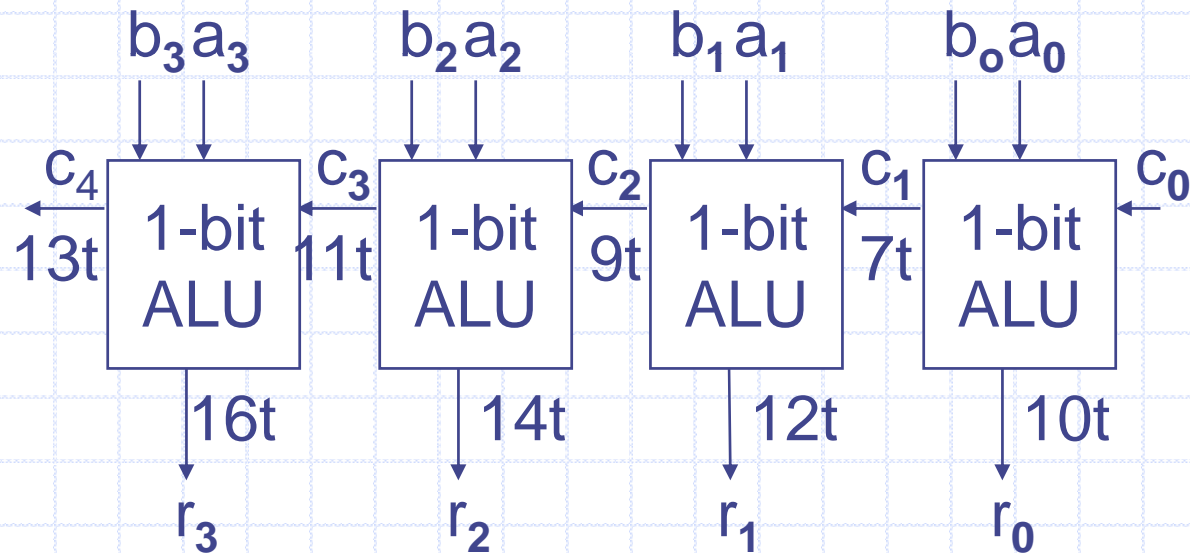
چقدر است تاخیر

- یک واحد محاسبه منطق 1 بیته $10t =$ ؟
- یک واحد محاسبه منطق 4 بیته؟
- یک واحد محاسبه منطق 32 بیته؟
- یک واحد محاسبه منطق n بیته؟

برای محاسبه Z تاخیر اضافه نیاز است

جمع کننده کری موج دار

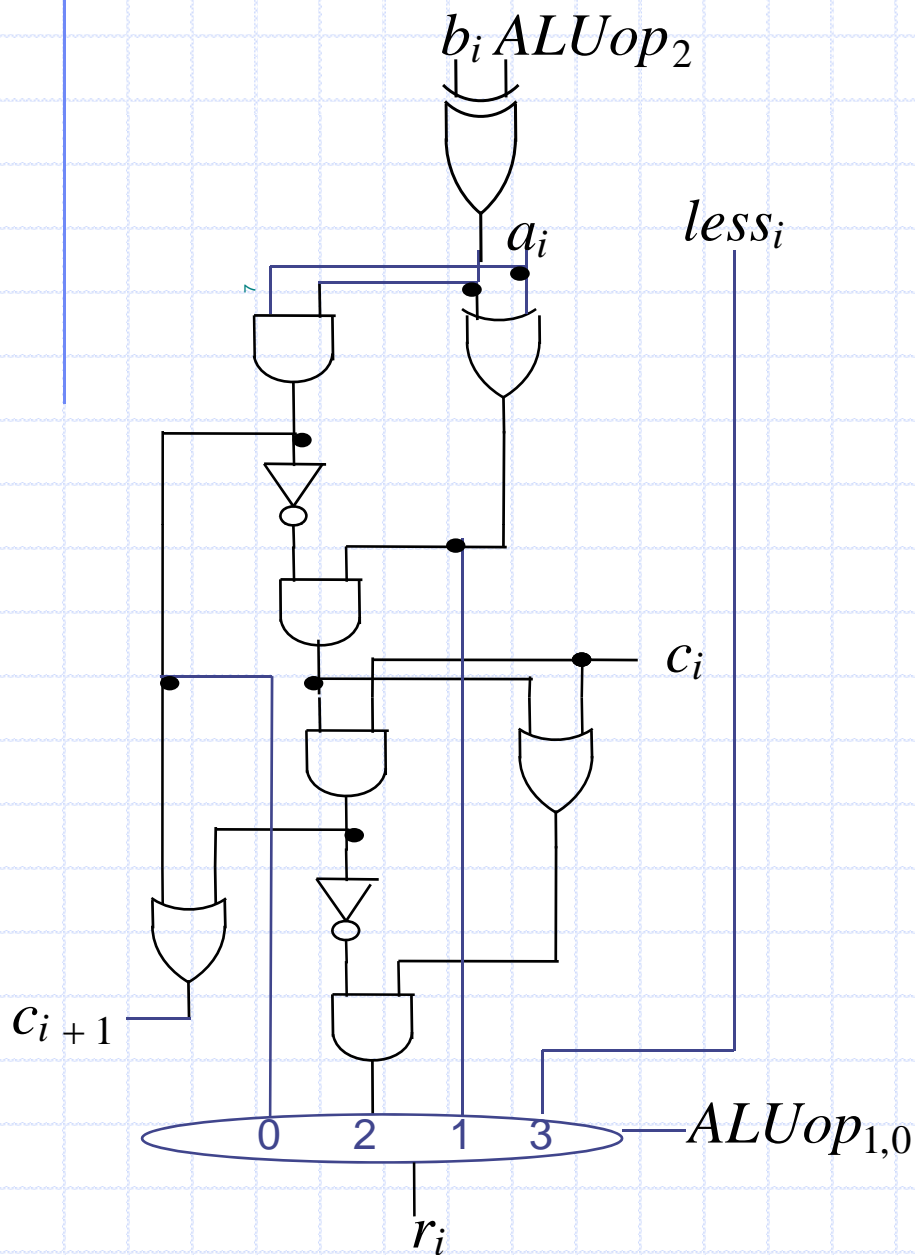
◆ در طراحی قبلی کری "موج دار" شده "از یک واحد محاسبه و منطق یک بیٹی برای بعدی تولید می شود



◆ این ها موجب می شوند که تا اندازه ای طراحی کند شود

◆ Z در لحظه $19t$ آماده است

تاخیر گیت



فرض کنید تاخیر های

- مالتی پالکسر 4 ورودی $2t = 2t$
- گیت XOR $2t = 2t$
- گیت AND/OR $1t = 1t$
- وارونه ساز $1t = 1t$

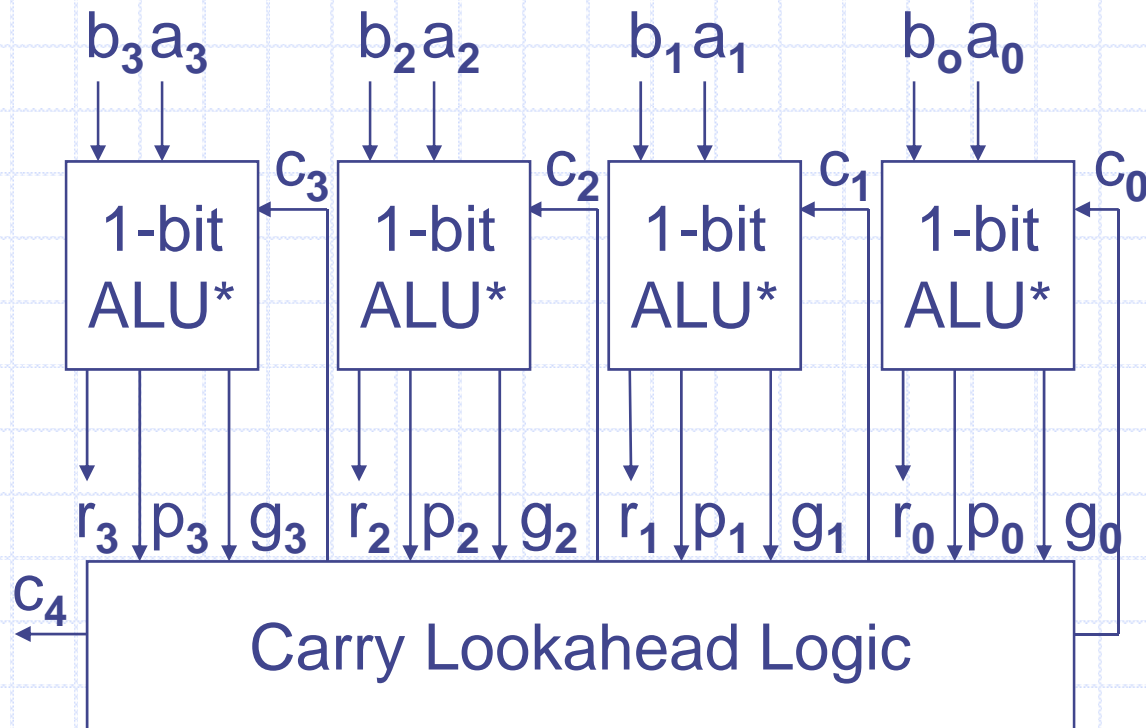
چقدر است تاخیر

- یک واحد محاسبه منطق 1 بیتی $10t = 10t$
- یک واحد محاسبه منطق 4 بیتی $16t = 16t$
- یک واحد محاسبه منطق 32 بیتی $72t = (2 \cdot 32 + 8)t = 72t$
- یک واحد محاسبه منطق n بیتی $(2n + 8)t = (2n + 8)t$

$\lceil \log_2(n) \rceil$ سطح OR دو ورودی و 1 وارونه ساز برای محاسبه Z نیاز داریم

Carry Lookahead Adder (CLA)

◆ در یک CLA کوری ها به صورت همزمان از یک carry lookahead logic (CLL) محاسبه می شوند



معادله منطقی کری

◆ معادله منطقی کری هست:

$$C_{i+1} = a_i b_i + (a_i + b_i) C_i$$

◆ ما یک سیگنال مشابه تعریف می کنیم propagate

$$p_i = a_i + b_i$$

◆ و یک سیگنال تولید میشود generate

$$g_i = a_i b_i$$

◆ این امکان را ایجاد می کند که معادله منطقی کری دوباره به صورت زیر نوشته شود

$$C_{i+1} = g_i + p_i C_i$$

Carry Lookahead Logic

◆ برای یک carry lookahead adder 4بیتی کری ها به

صورت زیر محاسبه می شوند

$$c_1 = g_0 + p_0c_0$$

$$c_2 = g_1 + p_1c_1 = g_1 + p_1(g_0 + p_0c_0)$$

$$= g_1 + p_1g_0 + p_1p_0c_0$$

$$c_3 = g_2 + p_2c_2 = g_2 + p_2(g_1 + p_1g_0 + p_1p_0c_0)$$

$$= g_2 + p_2g_1 + p_2p_1g_0 + p_2p_1p_0c_0$$

$$c_4 = g_3 + p_3c_3 = g_3 + p_3(g_2 + p_2g_1 + p_2p_1g_0$$

$$+ p_2p_1p_0c_0)$$

$$= g_3 + p_3g_2 + p_3p_2g_1 + p_3p_2p_1g_0 +$$

$$p_3p_2p_1p_0c_0$$

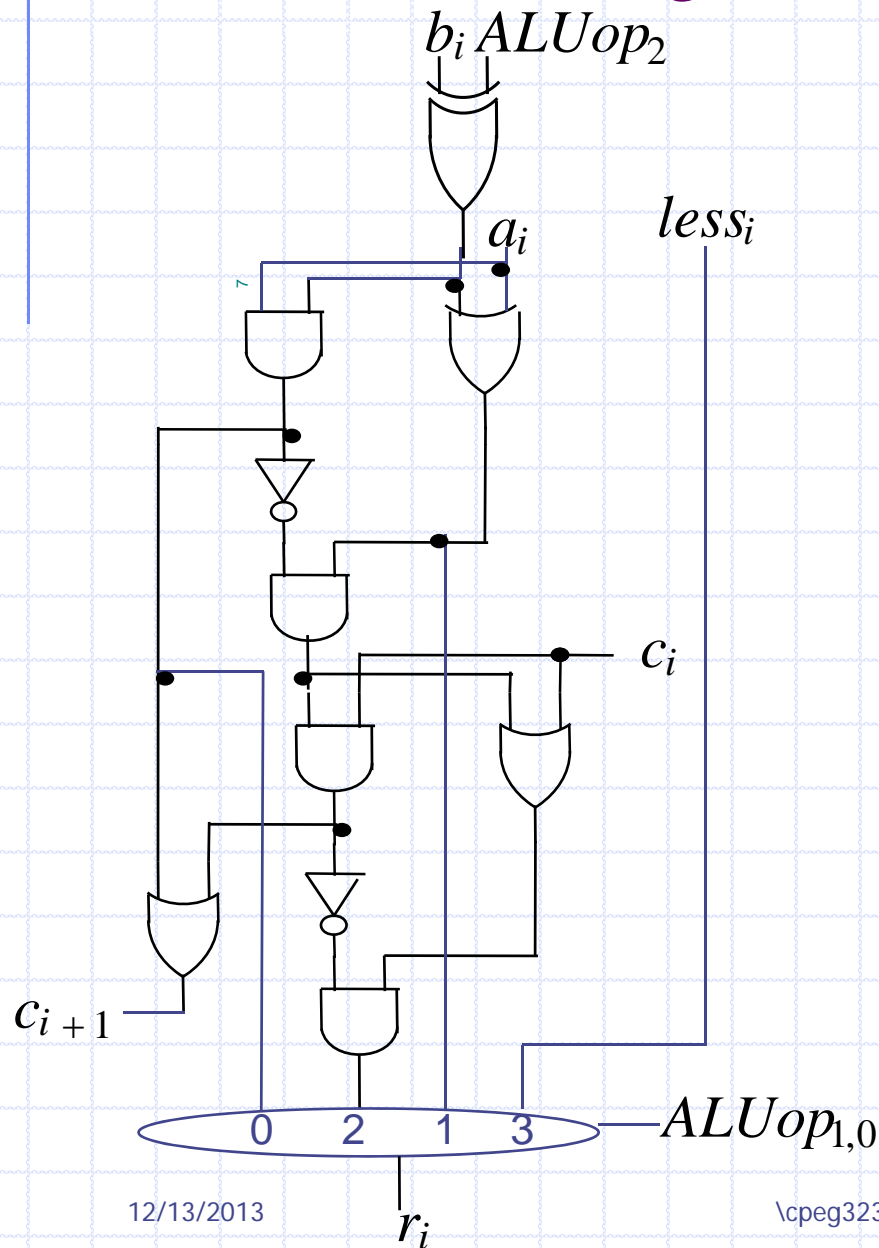
◆ چند تا گیت CLL 4بیتی نیاز داریم اگر که گیت ها بتوانند گنجایش

ورودی نا محدود داشته باشند؟

◆ اگر هر سطح منطقی تنها تاخیر t داشته باشد CLL دارای $2t$ تاخیر

است در نتیجه این تمرین ممکن است واقع بینانه نباشد.

تعدیل کردن واحد محاسبه و منطق یک بیتی



◆ چه طور ما می توانیم یک واحد محاسبه و منطق یک بیتی را تعدیل کنیم اگر که در CLA مورد استفاده باشد؟

◆ چند تا گیت تعدیل شده برای واحد محاسبه و منطق یک بیتی نیاز است؟

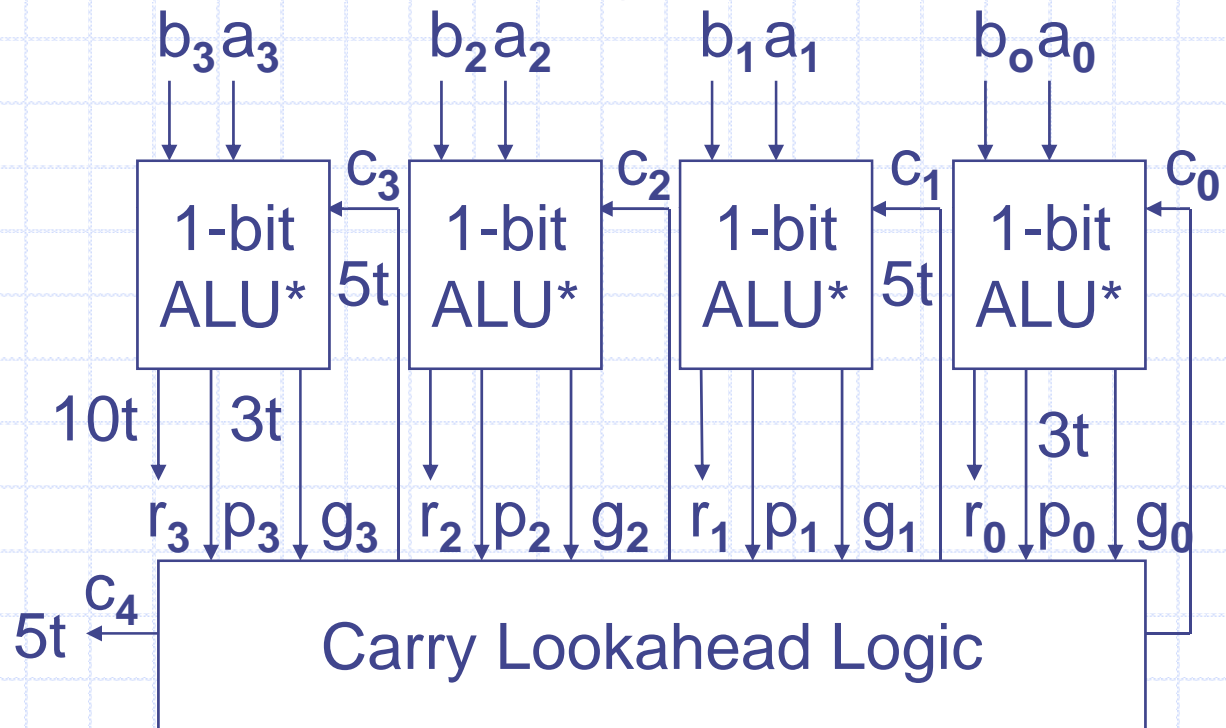
◆ چند تا گیت برای CLA 4 بیتی نیاز است؟

◆ چقدر تاخیر داریم تا زمانی که P_i و G_i آماده شوند؟

تعیین زمان CLA 4بیتی

در یک Carry lookahead adder کاری ها به صورت

همزمان محاسبه می شوند از Carry lookahead logic

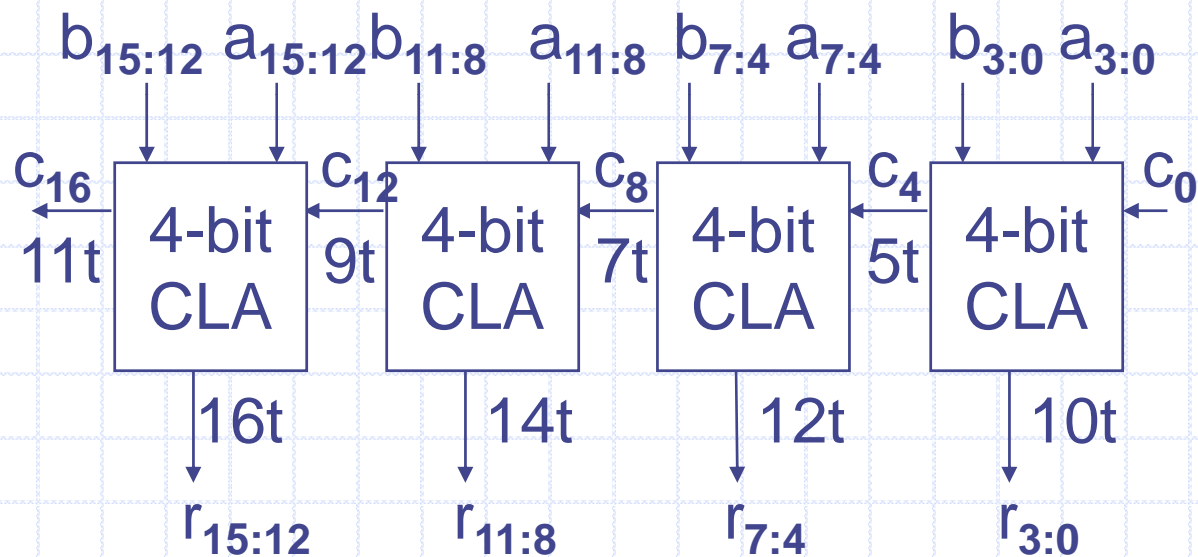


این طراحی به $15 \cdot 4 + 14 = 74$ گیت احتیاج دارد بدون محاسبه

ویا z

واحد محاسبه و منطق 16 بیتی. ورژن 1

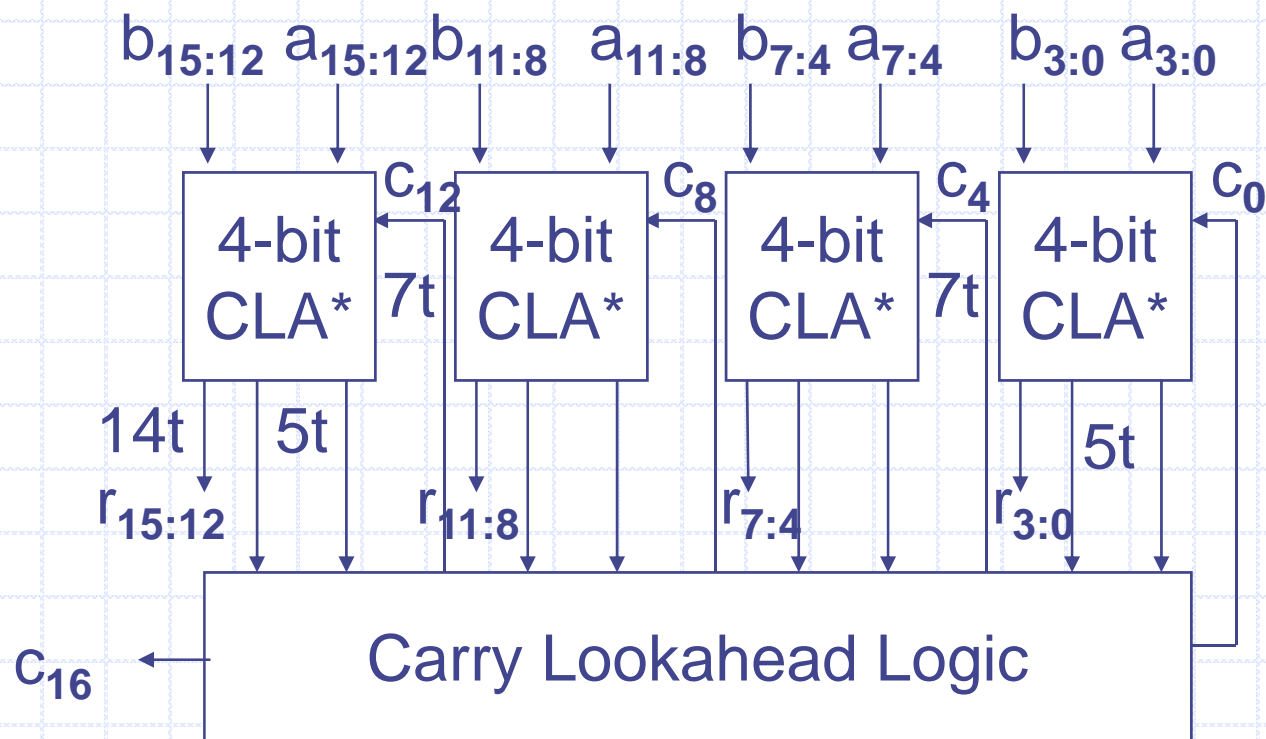
◆ یک ALU 16 بیتی ساخته شده از الحاق کردن 4 تا CLA 4 بیتی و اجازه انتقال کری "موج دار" بین "بلاک های" 4 بیتی



◆ این طراحی نیاز دارد به $74 \times 4 = 296$ گیت بدون محاسبه V یا Z

واحد محاسبه و منطق 16 بیتی . ورژن 2

- ◆ پیشنهاد دیگر استفاده از دومین سطح Carry lookahead logic است.
- ◆ این پیشنهاد سریعتر است اما به گیت های بیشتری نیاز دارد
 $16 \times 15 + 5 \times 14 = 310$ تا گیت.



CLA 4بیتی

◆ CLA 4بیتی (BLOCK CLA) شبیه اولین CLA 4بیتی است به جز محاسبات CLL یک "تولید بلاک" و "منتشر کردن بلاک" به جای خروجی کری

◆ بدین ترتیب محاسبه

$$c4 = g3 + p3g2 + p3p2g1 + p3p2p1g0 + p3p2p1p0c0$$

◆ جایگزین شده به وسیله

$$P3:0 = p3p2p1p0$$

$$G3:0 = g3 + p3g2 + p3p2g1 + p3p2p1g0$$

$$\text{Note: } c4 = G3:0 + P3:0c0$$

◆ این پیشنهاد محدود می کند حداکثر Fan را به 4 و CLL هنوز نیاز به 14 گیت دارد

نتایج

◆ یک n بیتی می تواند به وسیله الحاق کردن n تا ALU یک بیتی طراحی شود. CLL می تواند برای افزایش سرعت محاسبات استفاده شود.

◆ انواع گوناگون گزینش های طراحی برای طراحی ALU وجود دارد

◆ بهترین طراحی بستگی دارد به ناحیه تاخیر و توان نیازمندی ها که بر اساس تنوع تکنولوژی اصولی است

ضرب و تقسیم صحیح

ضرب صحيح بدون علامت

◆ مثال مداد و كاغذ

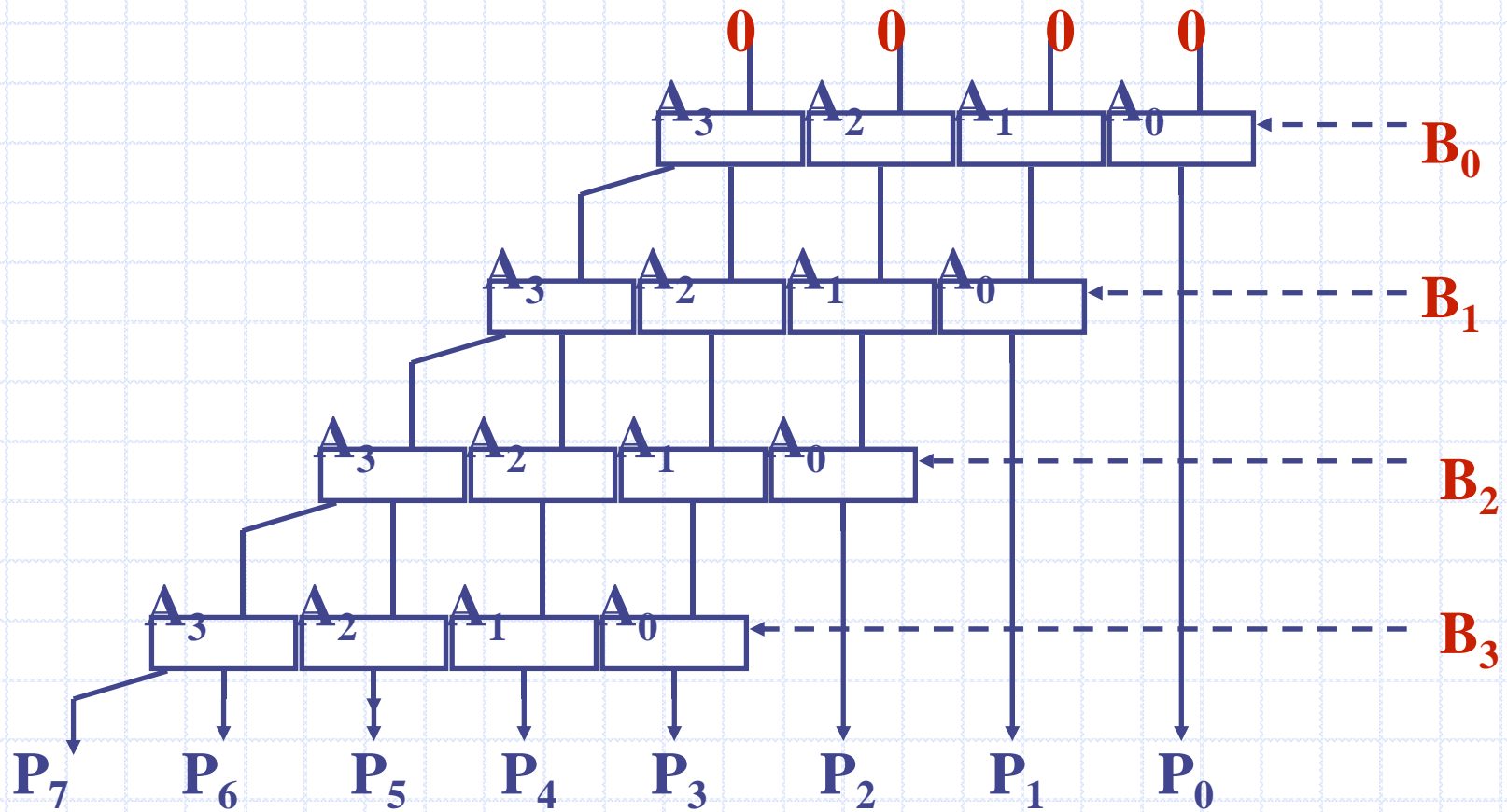
$$\begin{array}{r} \text{مضروب} \\ \text{مضروب فيه} \end{array} \quad \begin{array}{r} 1000 \\ * 1001 \\ \hline 1000 \\ 0000 \\ 0000 \\ 1000 \\ \hline 01001000 \end{array}$$

حاصل ضرب

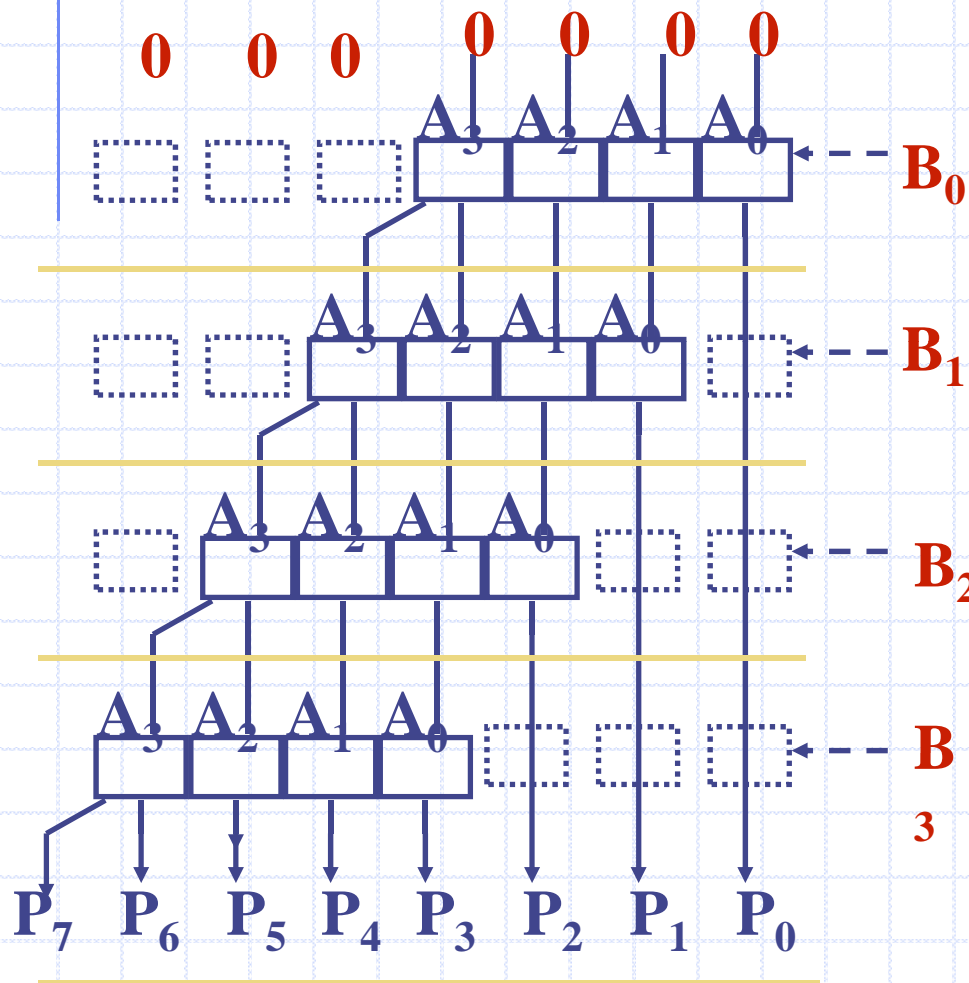
مشاهدات

◆ m بیت n * بیت $m+n$ حاصل ضرب

◆ انباشته های i مرحله $A * 2^i$ if $B_i == 1$



این چگونه کار می کند؟



◆ در هر مرحله 2 بار به چپ شیفت میدهد.

◆ استفاده از بیت بعدی B برای اینکه تعیین کند که آیا با مضروب شیفت داده شده جمع شود یا نه.

◆ جمع کردن $2n$ بیت حاصل ناتمام در هر مرحله

محاسبات ساده ریاضی

$A * B$

اگر B تا N بیت داشته باشد

$$B = \sum_{i=0}^{n-1} (2^i * B_i) \quad Product_n = \sum_{i=0}^{n-1} (2^i * B_i * A)$$

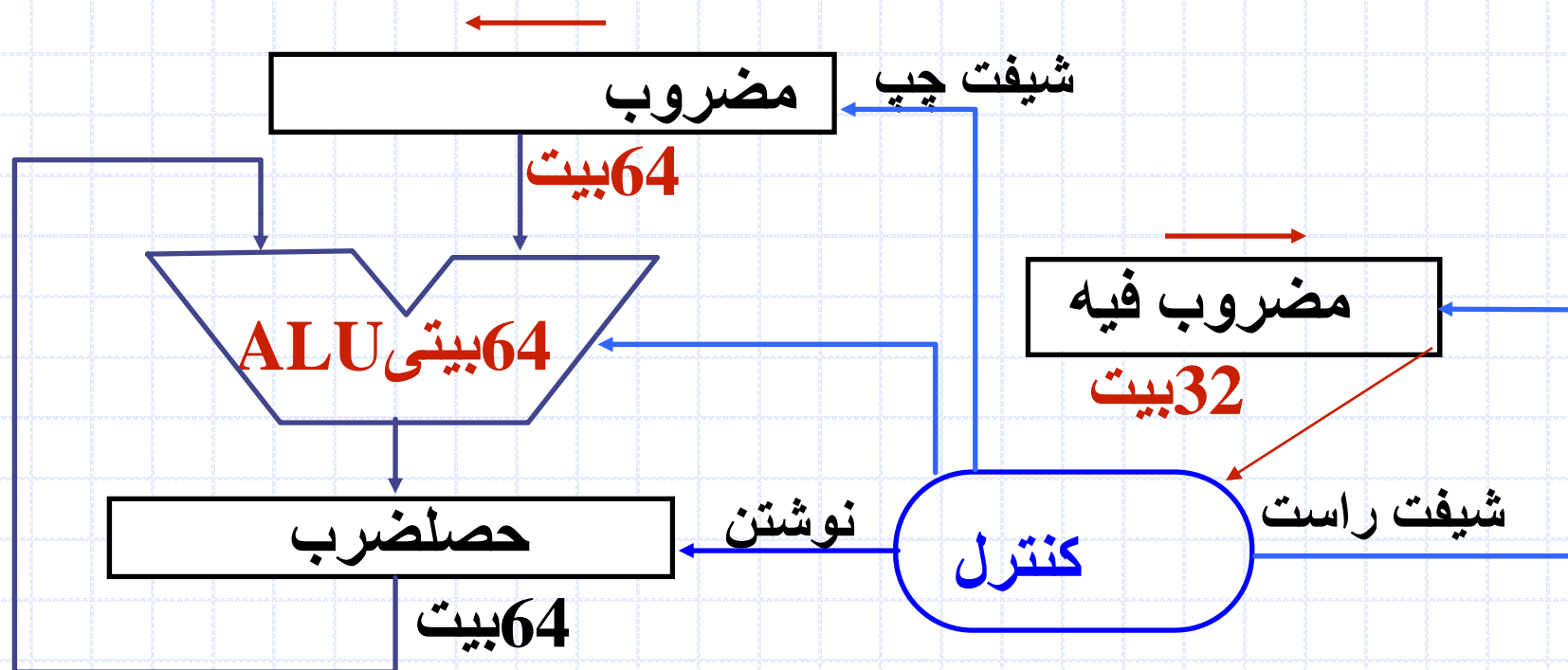
سپس زمانی که B $N+1$ بیت دارد

$$Product_{n+1} = \sum_{i=0}^n (2^i * B_i * A) = 2^n * B_n * A + Product_n$$

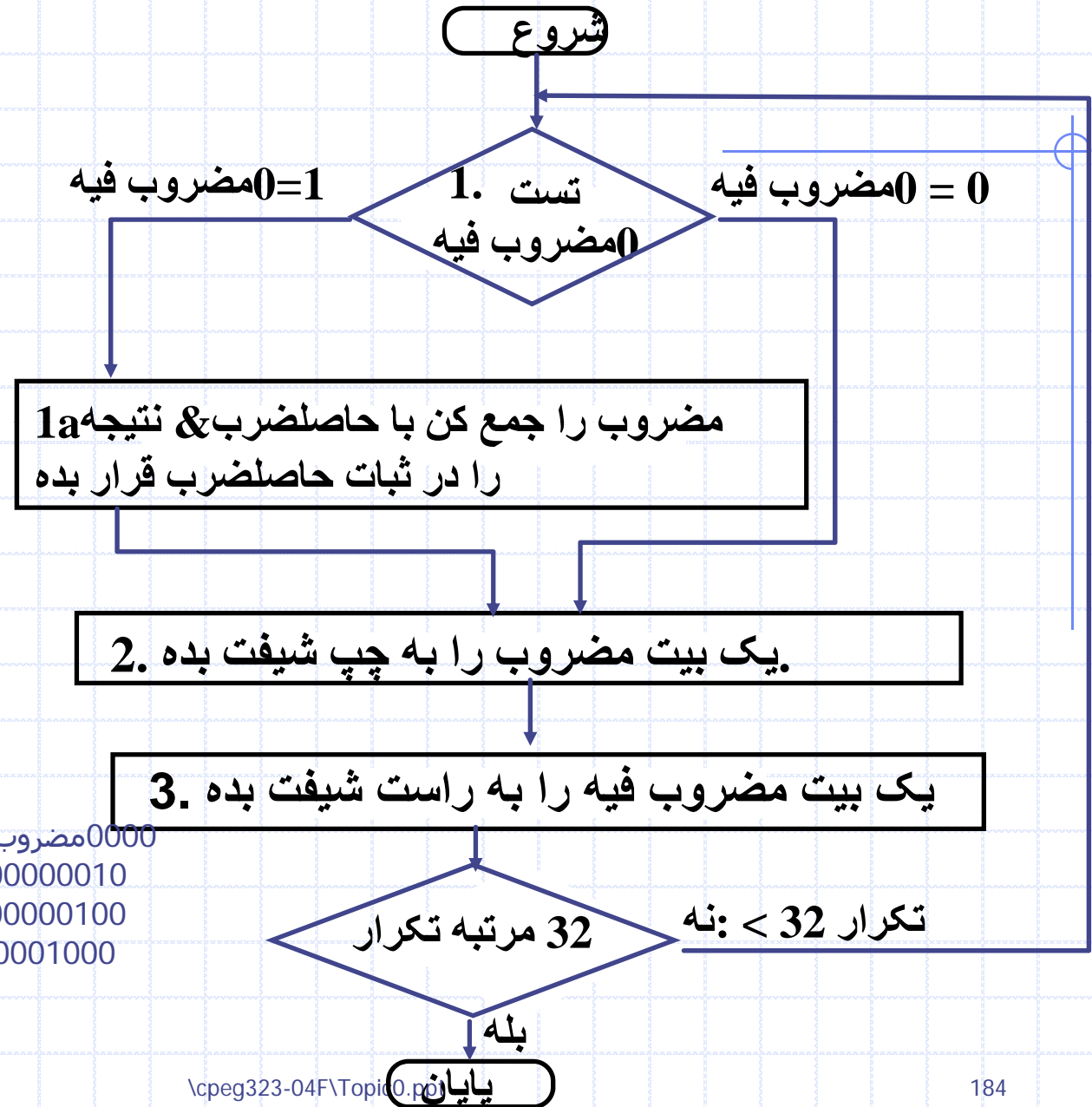
◆ A را یک مرتبه به چپ شیفت دادن: $2^i * A$.

◆ این واضح است که ضرب مشتمل بر شیفت و جمع تکراری است.

سخت افزار ضرب ورژن 1



الگوریتم ضرب ورژن 1



مضروب فیه	حاصلضرب
0011	0000
0001	00000010
0000	00000110
	00000110

مضروب
00000010
00000100
00001000

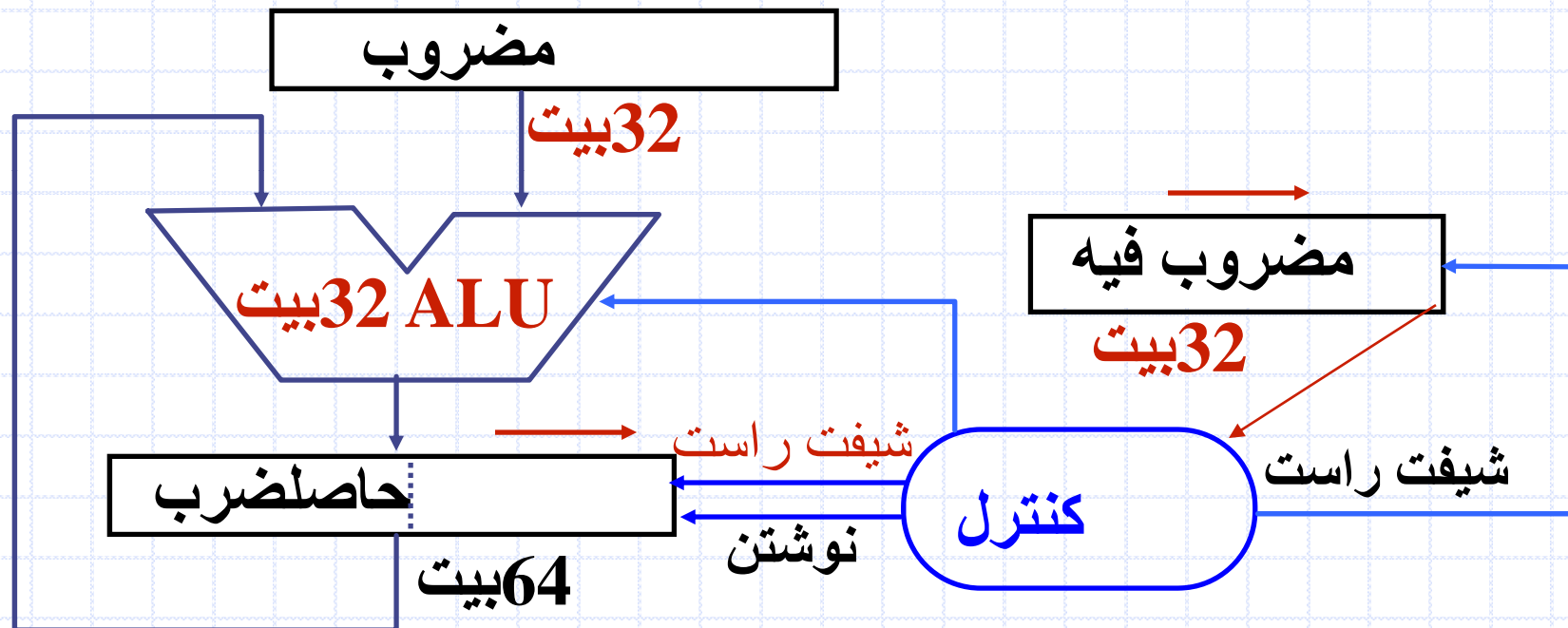
نظارت بر روی ورژن 1 ضرب

◆ نصف بیتها در مضروب همیشه صفر است ≤ 64 بیت جمع کننده هدر رفته است

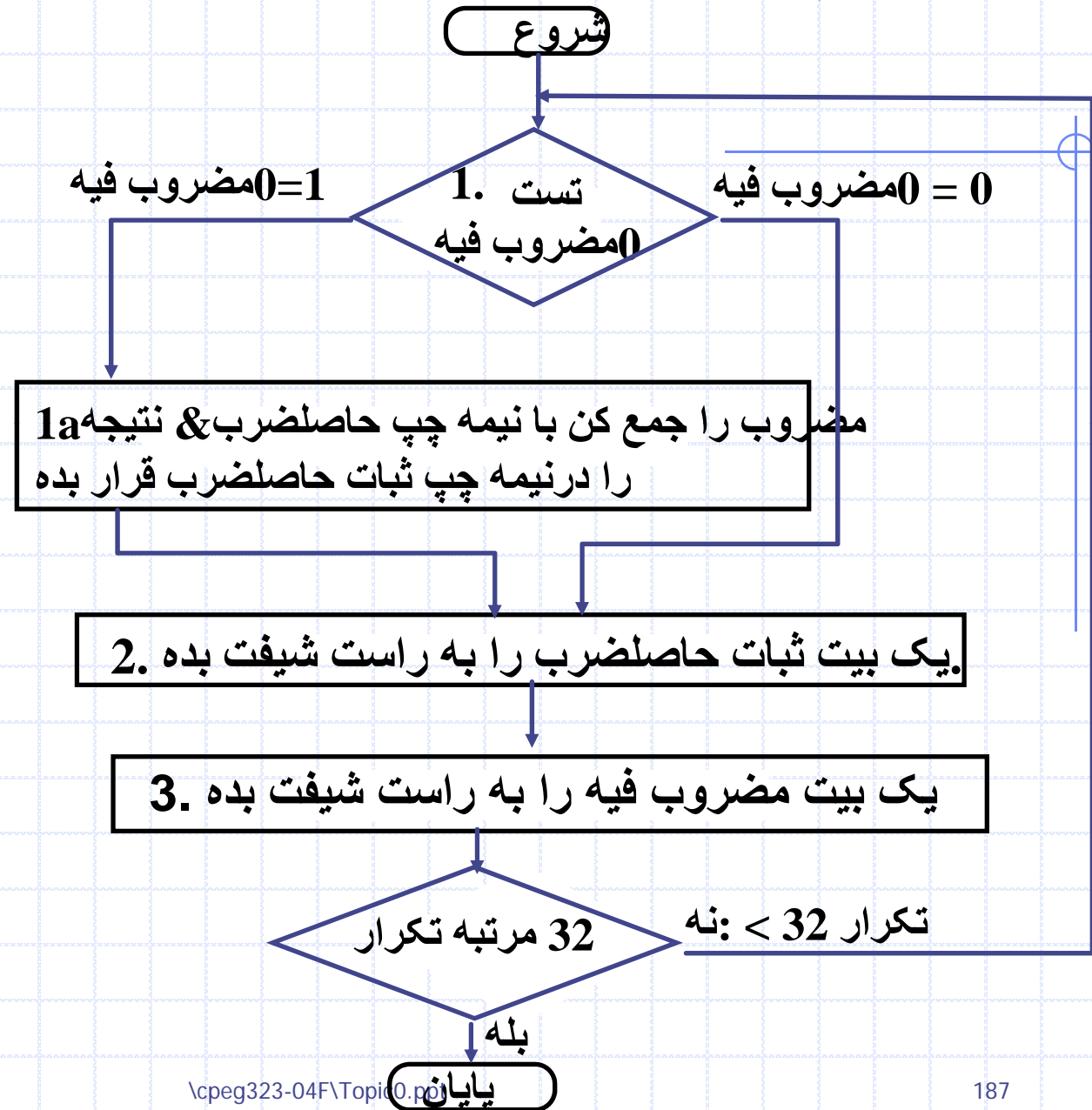
◆ مقدار دادن در سمت چپ مضروب به عنوان شیفت دادن \leq کم ارزش ترین بیت حاصل ضرب هرگز تغییر نمی کند

◆ به جای اینکه مضروب را به چپ شیفت دهیم. حاصل را به راست شیفت می دهیم؟

سخت افزار ضرب ورژن 2



الگوریتم ضرب ورژن 2

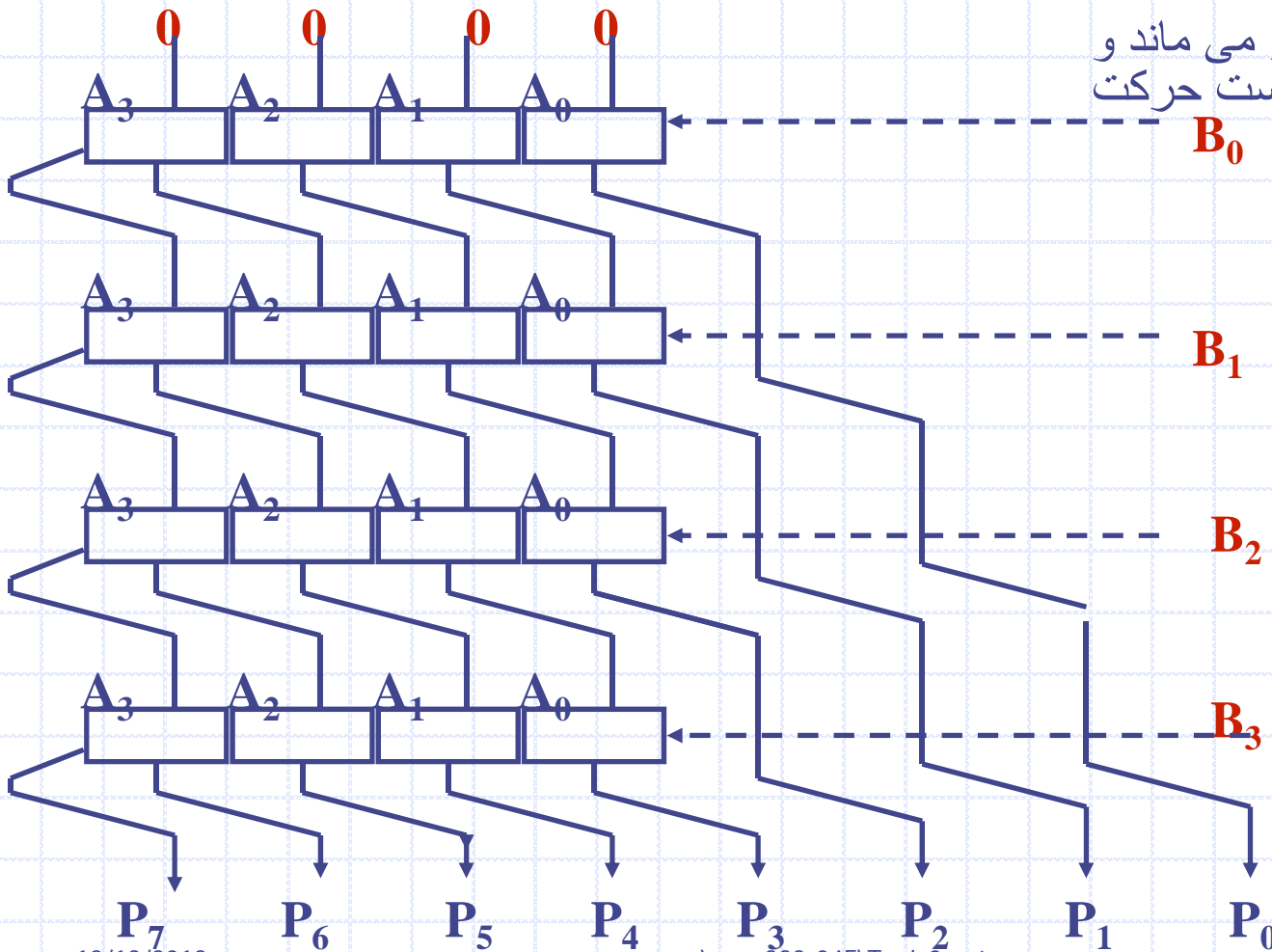


مضروب مضروب فيه حاصل ضرب

	0000 0000	0011	0010
1:	0010 0000	0011	0010
2:	0001 0000	0011	0010
3:	0001 0000	0001	0010
1:	0011 0000	0001	0010
2:	0001 1000	0001	0010
3:	0001 1000	0000	0010
1:	0001 1000	0000	0010
2:	0000 1100	0000	0010
3:	0000 1100	0000	0010
1:	0000 1100	0000	0010
2:	0000 0110	0000	0010
3:	0000 0110	0000	0010
	0000 0110	0000	0010

چه چیزی روی می دهد؟

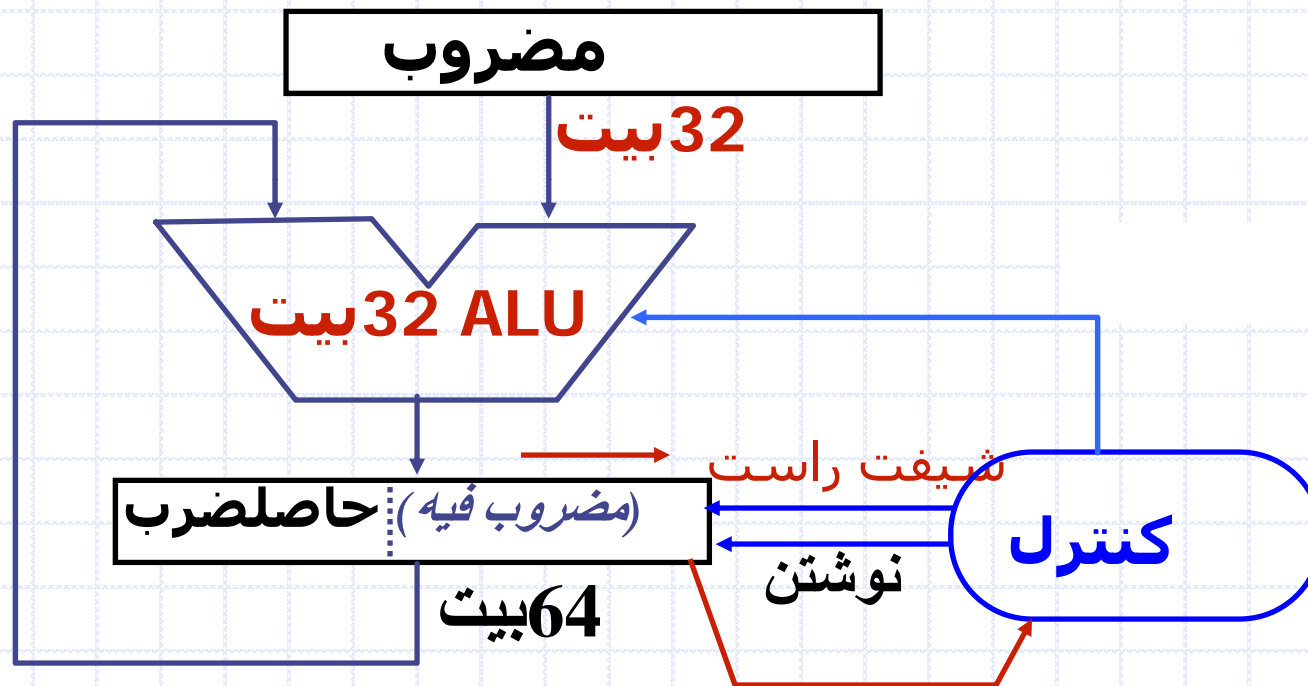
مضروب فیه هنوز می ماند و حاصلضرب به راست حرکت می کند



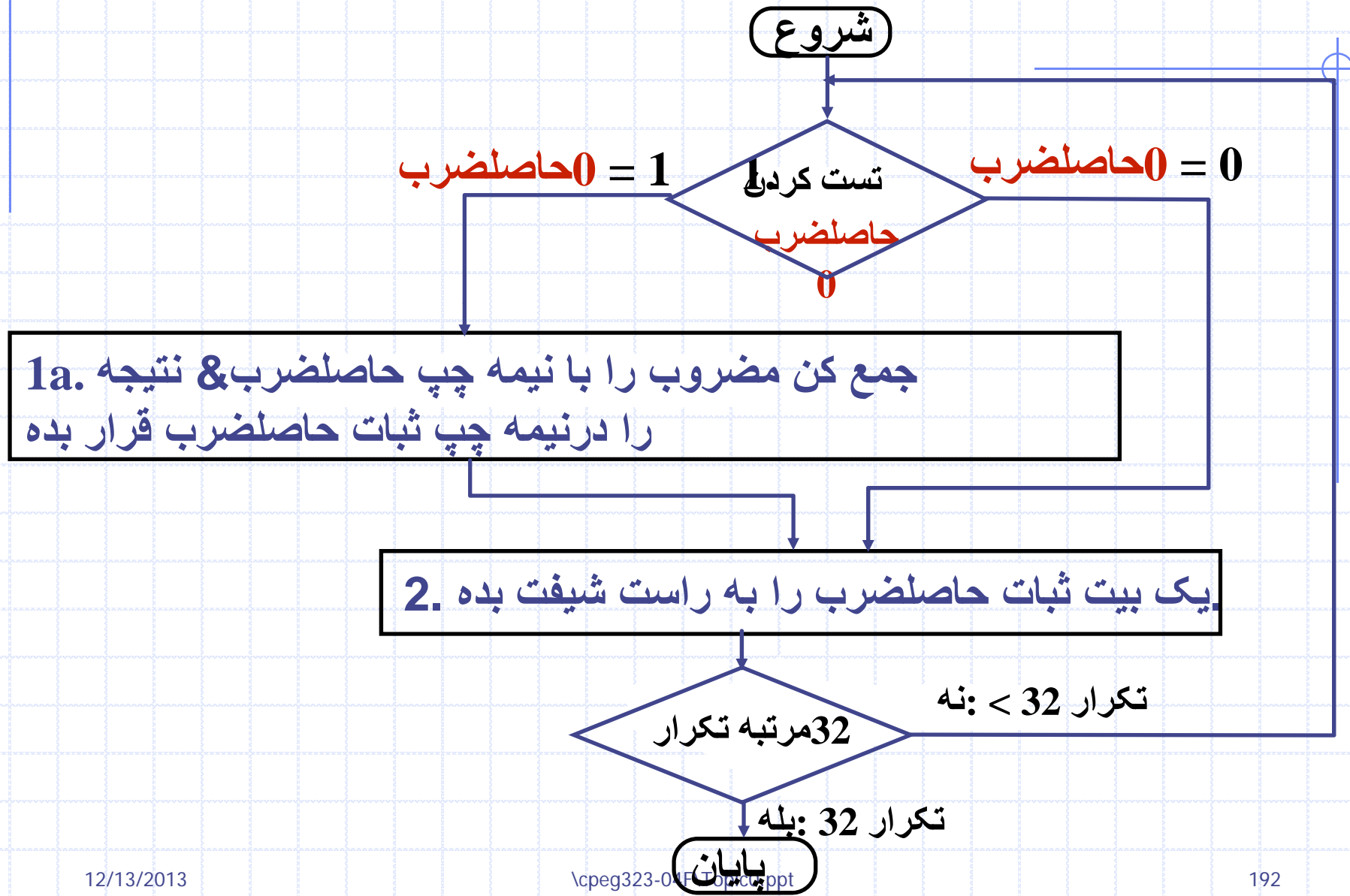
مشاهدات ورژن 2 ضرب

◆ فضای هدر رفته ثبات حاصلضرب که دقیقاً هم اندازه مضرب فیه است $=$ به هم پیوستن ثبات مضروب فیه و ثبات حاصلضرب

سخت افزار ضرب ورژن 2



الگوریتم ضرب ورژن 3

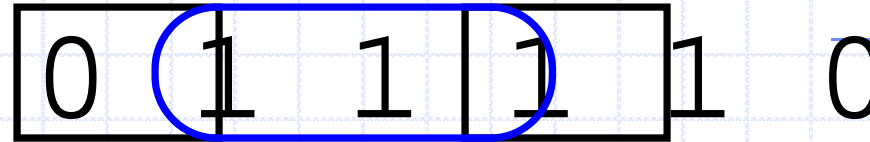


مشاهدات ورژن 3 ضرب

- ◆ دو مرحله در هر بیت است برای اینکه مضروب فیه & حاصلضرب به هم می پیوندد.
- ◆ ثبات های Lo,Hi MIPS هستند نیمه چپ و راست حاصلضرب
- ◆ MIPS دستور العمل ضرب بدون علامت را به ما می دهد
- ◆ درمورد ضرب علامت دار چی؟
- آسانترین راه حل زمانی است که هر دو مثبت باشند و به یاد داشته باشید داده تکمیل شده اجرا می شود(بیت علامت را دور بریز اجرا برای 31 مرحله)
- ◆ کاربرد تعریف متمم 2
- نیاز به گسترش علامت داده های بخصوص و طبقه بندی آنها در انتها
- ◆ الگوریتم بوت راه فوق العاده برای ضرب اعداد علامت دار است و از سخت افزار یکسان مانند قبلی استفاده میکند و چرخه ها را ذخیره می کند
- می تواند تعدیل شود به **to handle multiple bits at a time**
-

الگوریتم بوت

ابتدای اجرا وسط اجرا پایان اجرا



OP	مثال	توضیح	بیت جاری (بیت راست)	بیت راست
sub	0001111000	شروع اجرا از 1s	1	0
none	0001111000	وسط اجرا از 1s	1	1
add	0001111000	پایان اجرای 1s	0	1
none	0001111000	وسط اجرا از 1s	0	0

♦ ابتکاری برای سرعت (موقعی که شیفتم سریعتر از add بود). یک رشته از

1s ها را در مضروب فیه با نخستین تفریق جایگزین کنید

$$\begin{array}{r}
 -1 \\
 + 10000 \\
 \hline
 01111
 \end{array}$$

♦ ما ابتدا یک 1 می بینیم و سپس در جمع بعدی //

مثال بوت (2*7)

عملگر	مضروب	حاصل	بعدي؟
0. initial value	0010	0000 0111 0	10 -> sub
1a. $P = P - m$	1110	<u>+1110</u> 1110 0111 0	shift P (sign ext)
1b.	0010	1111 0011 1	11 -> nop, shift
2.	0010	1111 1001 1	11 -> nop, shift
3.	0010	1111 1100 1	01 -> add
4a.	0010	<u>+0010</u> 0001 1100 1	shift
4b.	0010	0000 1110 0	done

مثال بوت (3-2*)

عملگر	مضروب	حاصل	? بعدی
0. initial value	0010	0000 1101 0	10 -> sub
1a. $P = P - m$	1110	+1110 1110 1101 0	shift P (sign ext)
1b.	0010	1111 0110 1	01 -> add
2a.		+ 0010 0001 0110 1	shift P
2b.	0010	0000 1011 0	10 -> sub
3a.		+1110 0010 1110	
1011 0	shift P		
3b.	0010	1111 0101 1	11 -> nop
4a		1111 0101 1	shift
4b.	0010	1111 1010 1	done

ضرب علامت دار

◆ آسانترین راه حل این است که

■ هر دو مثبت باشد

■ به یاد داشته باشید که حاصل تکمیل شده اجرا می شود

■ علامت حاصلضرب را محاسبه می کند عملوندها را به اعداد

مثبت تبدیل می کند ، بیت علامت را دور می ریزد ، برای 31

مرتبه اجرا می شود ، سپس نتیجه ثابت می شود.

الگوریتم های سریعتر یا عمل ضرب

◆ الگوریتم بوت

- ضرب اعداد علامت دار مانند قبل داز سخت افزار یکسان استفاده می کند و چرخه ها را ذخیره می کند
 - بیت های مضرب را در یک لحظه می تواند دستکاری کند
- ## ◆ استفاده کردن از یک آرایه آدرس
- ## ◆ مشاهدات این گونه نشان می دهد که :چه اضافه شود یا نشود مضروب را یک واحد شیفت می دهیم_ این تصمیم می تواند کامل به صورت موازی انجام شود.

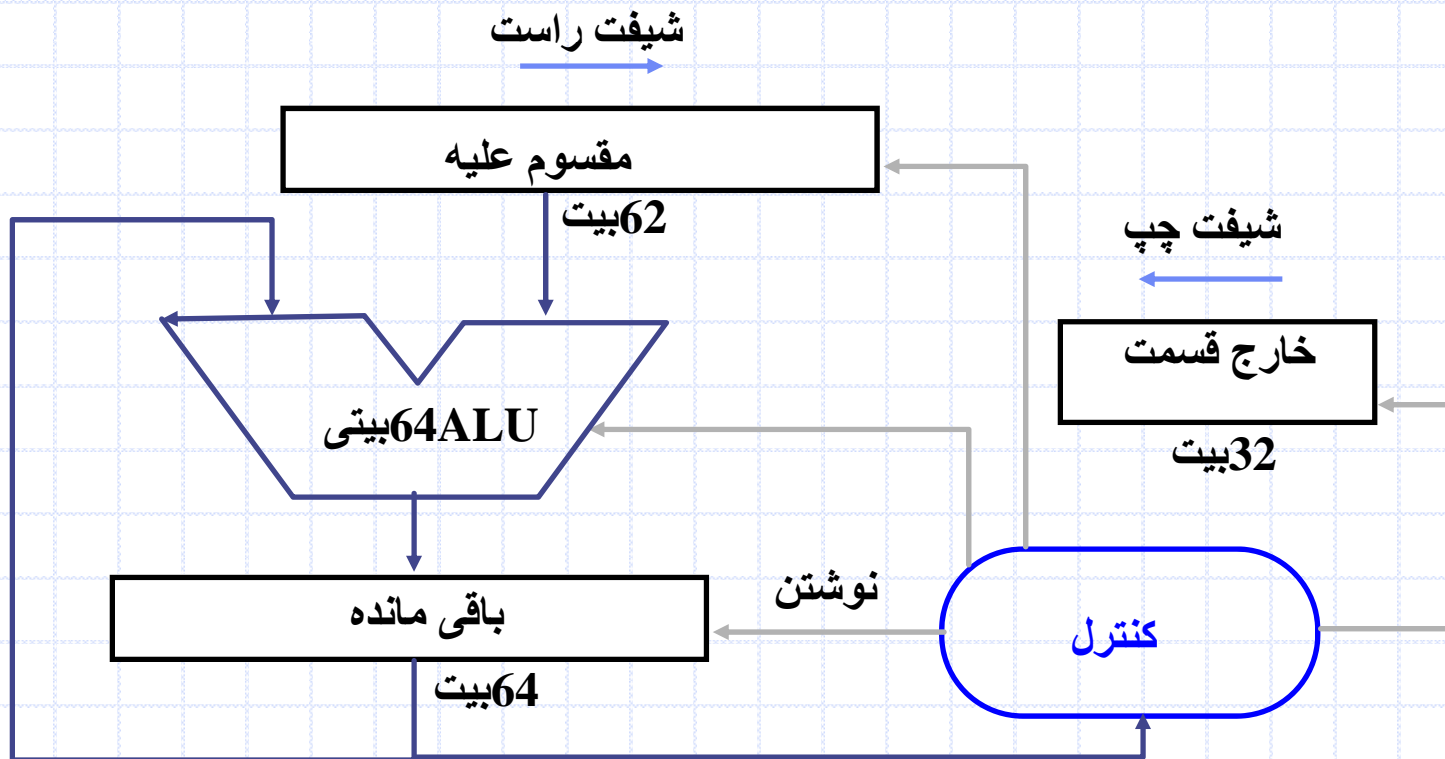
تقسیم: مداد و کاغذ

$$\begin{array}{r} \text{Quotient} \quad 1001 \\ \text{Divisor } 1000 \overline{) 1001010} \\ \underline{-1000} \\ 10 \\ 101 \\ \underline{1010} \\ \underline{-1000} \\ \text{Remainder } 10 \end{array}$$

$$\text{Dividend} = \text{Quotient} * \text{Divisor} + \text{Remainder}$$

سخت افزار تقسیم ورژن 1

◆ 64 بیت ثبات مقسوم علیه، واحد محاسبه و منطق 64 بیتی،
ثبات باقی مانده، ثبات خارج قسمت 32 بیتی



الگوریتم تقسیم ورژن 1

N+1 مرحله برای n بیت خرج قسمت و باقی مانده

شروع: مقسوم را در باقی مانده قرار دهید

کم کنید ثبات مقسوم علیه را از ثبات باقی مانده 1. نتیجه را در ثبات باقی مانده قرار دهید

تست کردن باقی مانده ≥ 0 باقی مانده < 0

2a. شیفت به چپ
ثبات خارج قسمت و 1
کردن آخرین بیت جدید
سمت راست

2ب. بازیابی مقدار اولیه با جمع ثبات مقسوم علیه و قرار
دادن مجموع در ثبات باقی مانده. شیفت ثبات خارج قسمت
به چپ و صفر کردن کم ارزش ترین بیت جدید

3- شیفت به راست ثبات مقسوم علیه به اندازه یک بیت

آیا سی و سومین تکرار مقسوم علیه است؟
بله
خیر

باقی مانده 0000 0111
خرج قسمت 0000 0010 0000

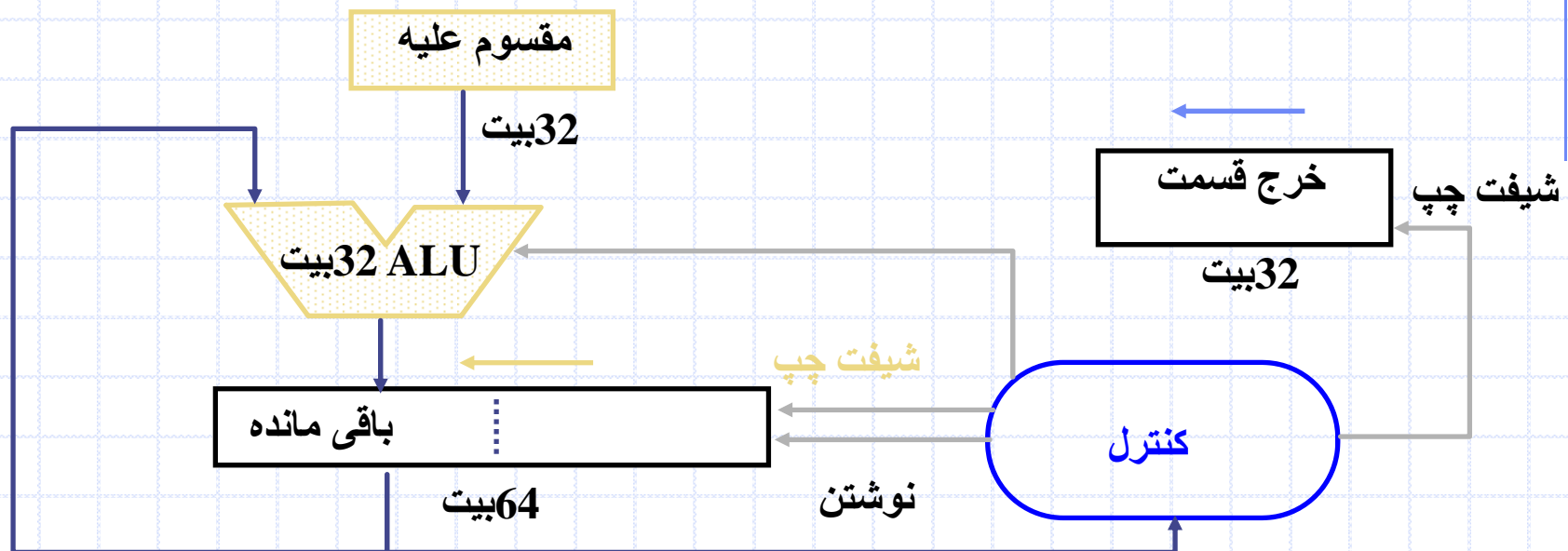
پایان

مشاهدات تقسیم ورژن 1

- ◆ نصف بیت ها در مقسوم علیه همیشه صفر است
- ◆ \leq نصف 64 بیت جمع کننده هدر رفته است
- ◆ \leq نصف مقسوم علیه هدر رفته است
- ◆ به جای اینکه مقسوم علیه را شیفِت به راست دهیم ، باقی مانده را شیفِت به چپ می دهیم ؟
- ◆ مرحله اول نمی تواند تولید یک خارج قسمت کند(به عبارت دیگر خیلی بزرگ است)
- ◆ \leq به وسیله اجرای فرمان ها با شیفِت اول و سپس تفریق می توان یک تکرار ذخیره کرد.

سخت افزار تقسیم ورژن 2

◆ 32 بیت ثابت مقسوم علیه ، 32 بیت ALU ، 64 بیت ثابت
باقی مانده ، 32 بیت ثابت باقی مانده



الگوریتم تقسیم ورژن 2

شروع: مقسوم علیه را در باقی مانده قرار دهید

شیفت بده ثبات باقی مانده را به چپ

2- ثبات مقسوم را از نیمه چپ
ثبات باقی مانده کم کن & قرار بده نتیجه را در نیمه چپ
ثبات باقی مانده

آزمایش باقی مانده ≥ 0 باقی مانده < 0

الف- شیفت به چپ ثبات
خارج قسمت و یک
کردن آخرین بیت جدید
سمت راست

ب- بازیابی مقدار اولیه با جمع ثبات مقسوم علیه و قرار دادن مجموع
در نیمه سمت چپ ثبات باقیمانده همچنین شیفت ثبات خارج قسمت به
چپ و صفر کردن کم ارزش ترین بیت جدید

خیر

آیا می‌توانیم دوباره تکرار کنیم؟

بله

مقسوم خارج قسمت باقی مانده

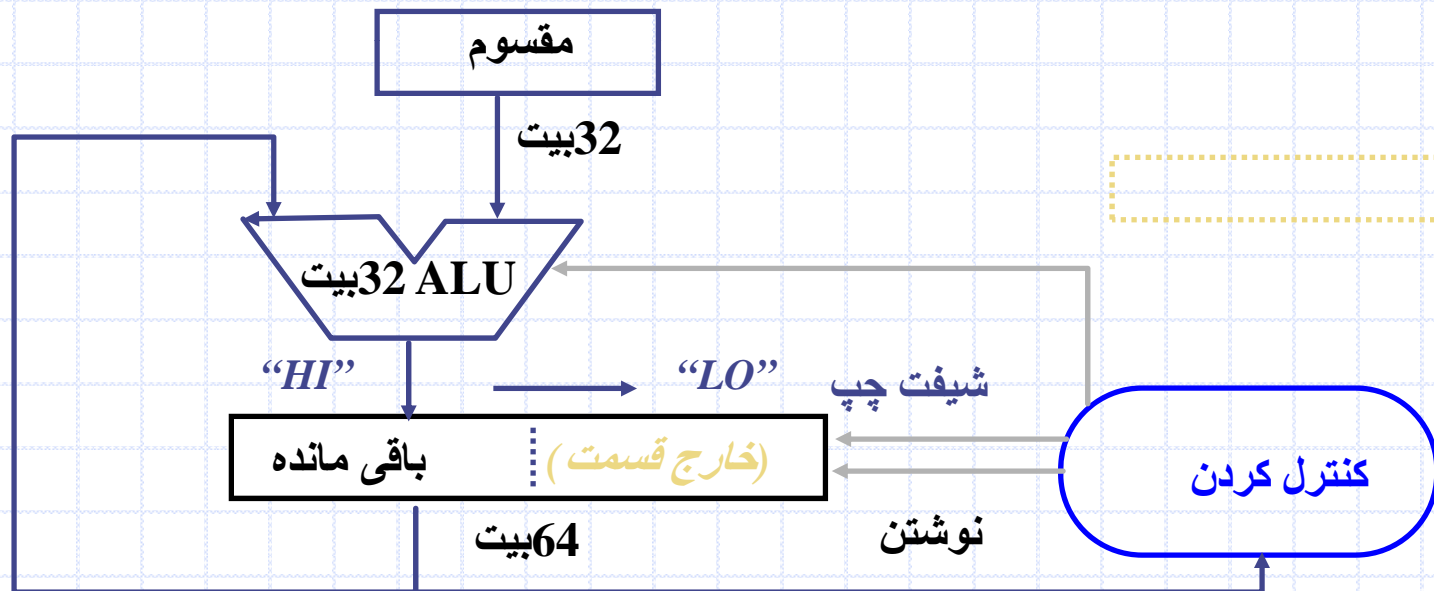
0000 0111 0000 0010

مشاهدات تقسیم ورژن 2

- ◆ حذف کردن ثبات خارج قسمت به وسیله بهم پیوستن با باقیمانده به عنوان شیفیت چپ
- ◆ آغاز شیفیت دادن باقی مانده به چپ مانند قبل
- ◆ پس از آن حلقه فقط محتوی 2 مرحله است چون شیفیت دادن ثبات باقی مانده هر دو شیفیت ها باقی مانده در نیمه چپ و خارج قسمت در نیمه راست.
- ◆ در نتیجه به هم پیوستن 2 ثبات به یکدیگر و فرمان های جدید عملگرها در این حلقه باقی مانده ای است که در یک لحظه بارها به چپ شیفیت خواهد یافت
- ◆ بنابراین در مرحله تصحیح پایانی باید شیفیت به عقب بدهیم که تنها باقی مانده در نیمه چپ ثبات باشد

سخت افزار تقسیم ورژن 2

◆ 32 بیت ثبات مقسوم ، 32 بیت ALU ، 64 بیت ثبات باقی مانده ، (صفر بیت خارج قسمت).



الگوریتم تقسیم ورژن 3

شروع: مقسوم علیه را در باقی مانده قرار دهید

شيفت بده ثبات باقی مانده را به چپ

2- ثبات مقسوم را از نیمه چپ
ثبات باقی مانده کم کن & قرار بده نتیجه را در نیمه چپ
ثبات باقی مانده

آزمایش باقی مانده ≥ 0 باقی مانده < 0

الف- شيفت به چپ ثبات باقی مانده و یک کردن آخرین بیت جدید سمت راست

ب- بازیابی مقدار اولیه با جمع ثبات مقسوم علیه و قرار دادن مجموع در نیمه سمت چپ ثبات باقیمانده همچنین شيفت ثبات باقی مانده به چپ و صفر کردن کم ارزش ترین بیت جدید

آیا می و دومین تکرار است؟

خیر

بله

باقی مانده	مقسوم
0000 0111	0010

پایان

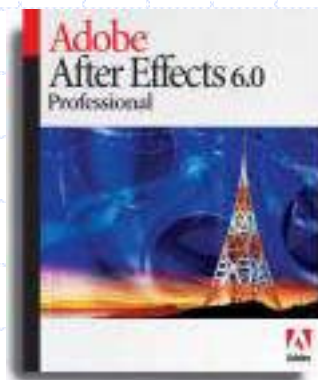
فصل چهارم

سنجش کارایی

- سنجش: چه، چرا، چگونه.
- معادله کارایی
- قانون امدال
- چگونه انرژی کارایی را محدود می کند.

اغلب چه کسی CPUها را ارتقا می دهد

یک حرفه ای که چرخه های CPU را به پول تبدیل می کند، و محدود به چرخه است



**Artist tool:
animation, video
special effects.**



برای خرید یک ماشین جدید چگونه تصمیم می گیرید؟

اندازه گیری "زمان اجرای" After Effects در "ظرفیت کاری" تولید خروجی تصویر
.....
////////////////////////////////////
////////////////////////////////////

***Night flight
Movie
Goes
Here***

"Night flight"
نقشه شهر و ابرها
در آسمان بوسیله
تصاویر فرموله
محاسبه می شوند

CPU شدید
I/O ناچیز

کارایی و هزینه از نظر مشتری/CPU

◆ از لحاظ خرید

■ معین کردن یک مجموعه از ماشینها که دارای

◆ بهترین کارایی

◆ کمترین هزینه

◆ بهترین کارایی / هزینه

◆ از لحاظ طراحی

■ مواجهه با انتخاب طراحی است که دارای:

◆ بهترین بهبود کارایی

◆ کمترین هزینه

◆ بهترین کارایی / هزینه

◆ نیاز هر دو

■ پایه ای برای مقایسه

■ معیاری برای ارزیابی

◆ هدف ما فهماندن مفهوم هزینه و کارایی از انتخاب های معماری است.

دو تصور از “کارایی”

Plane	DC to Paris	Speed	Passengers	Throughput (pmp)
Boeing 747	6.5 hours	610 mph	470	286,700
BAD/Sud Concodre	3 hours	1350 mph	132	178,200

کدامیک کارایی بالاتری دارد؟

° زمان انجام کار (زمان اجرا)

— زمان اجرا ، زمان پاسخ ، زمان برگشت

° کارها در روز، ساعت ، هفته ، ثانیه ، نانو ثانیه ، ... (کارایی)

— کارایی ، توان عملیاتی ، پهنای باند

اغلب زمان پاسخ و توان عملیاتی در عکس هم قرار دارند، چرا؟

تعاریف کارایی

◆ کارایی میزان چیزهایی در ثانیه است

■ بزرگتر بهتر است

◆ اگر به صورت عمده، زمان پاسخ برای ما مهم باشد:

$$\text{performance}(x) = \frac{1}{\text{execution_time}(x)} \quad \blacksquare$$

◆ "X" n بار سریعتر از "Y" است به معنی اینست که :

$$\text{execution_time}(Y) = \frac{\text{performance}(X)}{n} = \frac{\text{execution_time}(X)}{\text{performance}(Y)}$$

◆ چه زمانی توان عملیاتی مهمتر از زمان اجراست؟

◆ چه زمانی زمان اجرا از توان عملیاتی مهمتر است؟

مثال هایی از کارایی

- زمانهای Concorde در مقابل Boeing 747 ؟
 - Concord is $1350 \text{ mph} / 610 \text{ mph} = 2.2$ بار سریعتر
 $= 6.5 \text{ hours} / 3 \text{ hours}$
- توان عملیاتی Concorde در مقابل Boeing 747 ؟
 - Concord is $178,200 \text{ pmph} / 286,700 \text{ pmph} = 0.62$ “times faster”
 - Boeing is $286,700 \text{ pmph} / 178,200 \text{ pmph} = 1.6$ “times faster”
- در مورد توان عملیاتی، Boeing 1.6 بار (“60%”) سریعتر است.
- در مورد زمان پرواز، Concorde 2.2 بار (“120%”) سریعتر است.
- زمانیکه بحث کارایی پردازنده مطرح است، عمده تمرکز ما بر روی زمان اجرا برای “یک کار” است. چرا؟

درک کارایی

◆ موارد زیر احتمالاً، چگونه بر زمان پاسخ و توان عملیاتی اثر می گذارد؟

- افزایش سرعت ساعت در یک پردازشگر معین.
- افزایش تعداد کارها در یک سیستم (مثلاً داشتن یک سرویس واحد کامپیوتری دارای چند کاربر).
- افزایش تعداد پردازنده ها در یک سیستم که از چند پردازنده استفاده می کند (مانند یک شبکه از ماشینهای خود پرداز (ATM))

◆ اگر یک پنتیوم 3 یک برنامه را در 8 ثانیه اجرا کند و یک PowerPc برنامه ای مشابه را در 10 ثانیه اجرا کند، محصول پنتیوم چند بار سریعتر است؟

$$n = 10 / 8 = 1.25 \text{ times faster (or 25\% faster)}$$

◆ با این وجود چه چیزی باعث می شود که کسی PowerPC را برای خرید انتخاب کند؟

تعریف زمان

◆ زمان به چند طریق قابل تعریف است، وابسته به اینکه چگونه مورد سنجش قرار دهیم:

■ **زمان پاسخ:** مجموع زمان صرف شده برای تکمیل یک کار، شامل زمان صرف شده برای اجرا در CPU دستیابی به دیسک و حافظه، انتظار برای I/O و دیگر پردازشها، و سربار سیستم عامل.

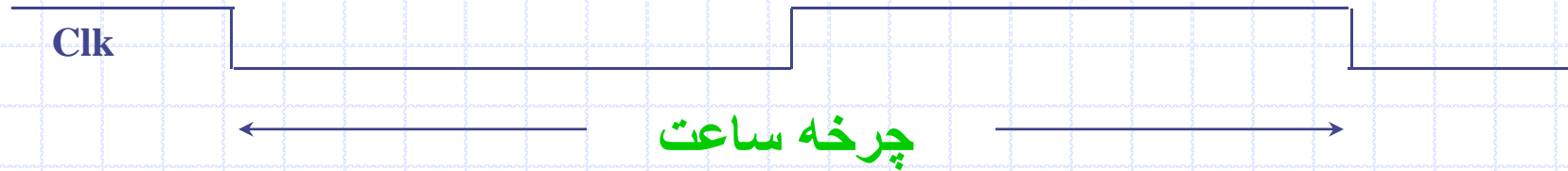
■ **زمان اجرای CPU:** مجموع زمانی که یک CPU برای یک کار صرف می کند (به غیر از زمان I/O یا اجرای برنامه های دیگر). این زمان به طور ساده معروف به زمان CPU است.

■ **زمان کاربر CPU:** مجموع زمانهایی که CPU در برنامه صرف می کند.

■ **زمان اجرای سیستمی CPU:** مجموع زمانی که سیستم عامل صرف اجرای کارها در برنامه می کند.

◆ برای مثال: یک برنامه می تواند دارای **زمان سیستمی CPU** معادل 22 ثانیه باشد، یک زمان کاربر CPU برابر 90 ثانیه، یک زمان اجرایی CPU، معادل 112 ثانیه، و یک **زمان پاسخ** برابر 162 ثانیه باشد.

ساعت کامپیوتر



◆ یک ساعت کامپیوتر در یک نرخ ثابت انجام می گیرد و تعیین می کند چه زمانی رخ داده‌ها در سخت افزار قرار گیرند.

° زمان چرخه ساعت مقدار زمان برای سپری شدن یک چرخه ساعت است (مثلا 5 نانوثانیه).

° نرخ ساعت، عکس زمان چرخه ساعت است.

° مثلا، اگر یک کامپیوتر دارای زمان چرخه ساعت معادل 5 نانوثانیه باشد، نرخ ساعت برابر است با:

$$\frac{1}{5 \times 10^{-9} \text{ sec}} = 200 \text{ MHz}$$

معیارهای پایه ارزیابی

◆ مقایسه معیارهای ماشینها

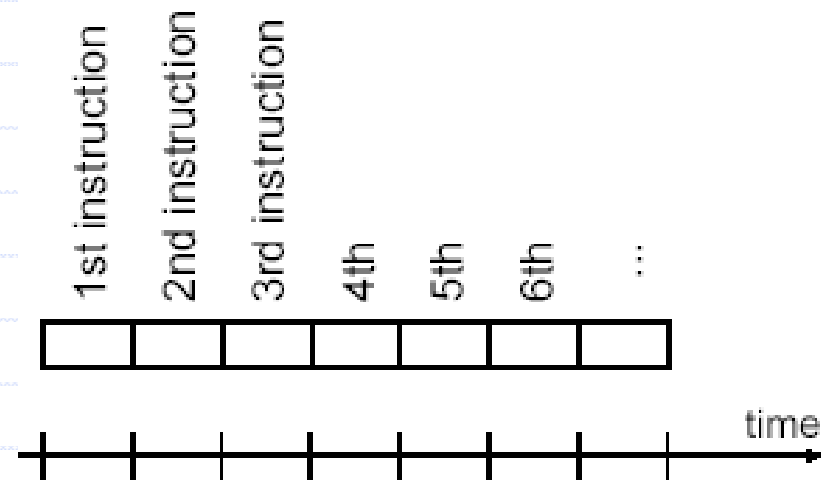
- زمان اجرا
- توان عملیاتی
- زمان CPU
- MIPS : مخفف : میلیونها دستورالعمل در ثانیه
- MFLOPS : میلیونها عمل ممیز شناور در هر ثانیه

◆ مقایسه مجموعه برنامه های استفاده شده در ماشینها

- مفهوم محاسباتی
- مفهوم محاسباتی سنگین
- Benchmarks

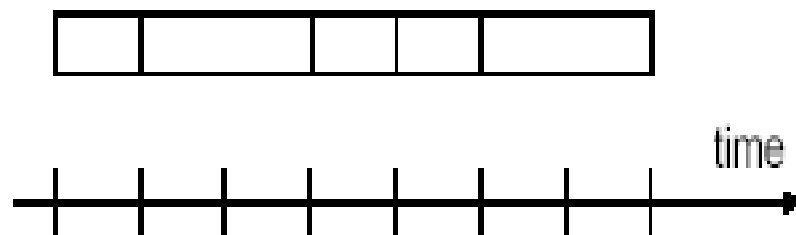
چه تعداد چرخه برای یک برنامه نیاز است؟

ممکن است تصور کرده باشید که تعداد چرخه ها = تعداد دستورالعملها باشد



این تصور اشتباه است، دستورالعملهای متفاوت، در ماشینهای متفاوت زمانهای متفاوتی را صرف می کنند.

تعداد چرخه های متفاوت برای دستورالعمل های متفاوت



- تقسیم زمان خیلی بیشتری از جمع صرف می کند.
- اعمال ممیز شناور، زمان بیشتری نسبت به نوع صحیح آن مصرف می کند.
- دسترسی به حافظه زمان بسیار بیشتری از دسترسی به ثباتها صرف می کند.

اکنون که ما با مفهوم چرخه آشنا شدیم

◆ یک برنامه معین نیاز دارد به
تعدادی از دستورالعملها
تعدادی از چرخه های ساعت
تعدادی از ثانیه ها

◆ ما یک لغت نامه برای شرح این کمیت ها داریم:
زمان چرخه ساعت (ثانیه در هر چرخه)
نرخ ساعت (سیکل در ثانیه)
CPI (تعداد چرخه ها در هر دستورالعمل)

یک برنامه ممیز شناور ممکن است دارای CPI بالاتری باشد

محاسبه زمان CPU

زمان اجرای یک برنامه معین می تواند مانند زیر محاسبه شود

زمان چرخه ساعت \times چرخه های ساعت CPU = زمان CPU

از آنجا که زمان چرخه ساعت و نرخ ساعت عکس یکدیگرند

نرخ ساعت/چرخه های ساعت CPU = زمان CPU

تعداد چرخه های ساعت CPU می تواند مانند زیر محاسبه گردد

تعداد چرخه های ساعت CPU = (دستورالعملها/برنامه) \times (چرخه های ساعت/دستورالعمل)

= تعداد دستورالعملها \times CPI

بنابراین

زمان CPU = تعداد دستورالعملها \times CPI \times زمان چرخه ساعت

زمان CPU = تعداد دستورالعملها / CPI \times نرخ ساعت

واحدهای این، اینگونه هستند

$$= \frac{\text{ثانیه}}{\text{چرخه ثانیه}} \times \frac{\text{چرخه های ساعت}}{\text{دستورالعمل}} \times \frac{\text{دستورالعملها}}{\text{برنامه}} \text{ ثانیه}$$

مثالی از محاسبه زمان CPU

◆ اگر یک کامپیوتر دارای نرخ ساعت 50MHz باشد، چقدر طول می کشد تا یک برنامه دارای 1000 دستورالعمل اجرا شود؟ CPI برای این برنامه 3.5 است.

◆ استفاده از معادله

CPU زمان = تعداد دستورالعمل / CPI x نرخ ساعت

پس

$$\text{CPU زمان} = 1000 \times 3.5 / (50 \times 10^6)$$

◆ اگر نرخ ساعت CPU ی یک کامپیوتر از 200MHz به 250 MHz افزایش یابد و بقیه

فاکتورها یکسان باقی بمانند، افزایش سرعت کامپیوتر چقدر است؟

$$\frac{\text{قدیمی CPU زمان}}{\text{جدید CPU زمان}} = \frac{\text{نرخ جدید ساعت } 250 \text{ MHz}}{\text{نرخ قدیمی ساعت } 200 \text{ MHz}} = 1.25$$

◆ تصور ساده انگارانه ما چه بود؟

فاکتورهای موثر در کارایی CPU

چه فاکتورهایی از موارد زیر تاثیر می گیرند؟

	clock rate	CPI	instr. count
			Program
			Compiler
			ISA
			Organization
			Technology

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

فاکتورهای موثر در کارایی CPU

چه فاکتورهایی از موارد زیر تاثیر می گیرند؟

	clock rate	CPI	instr. count	Program
				*
		Compiler	*	*
		ISA	*	*
Organization			*	*
Technology				*

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

محاسبه CPI

◆ CPI، میانگین تعداد چرخه ها در هر دستورالعمل است
◆ اگر برای هر نوع دستورالعمل، لازم است ما فرکانس ها و تعداد چرخه ها را برای اجرای آن بدانیم، ما می توانیم CPI را به صورت زیر محاسبه کنیم:

$$CPI = \sum CPI_i \times F_i$$

◆ برای مثال:

Op	F_i	CPI_i	$CPI_i \times F_i$	% Time
ALU	50%	1	.5	23%
Load	20%	5	1.0	45%
Store	10%	3	.3	14%
Branch	20%	2	.4	18%
Total	100%		2.2	100%

مسائل با میانگین حسابی

◆ تقاضاهایی که احتمال یکسانی برای اجرا ندارند. 2 تا ماشین که زمان بندی شدن برای 2 آزمون کارایی:

	Machine A	Machine B
Program1	2 seconds(%20)	6 seconds
Program2	12 seconds(%80)	10 seconds

میانگین زمان اجرا = $(12+2)/2 = 7$ ثانیه

میانگین زمان اجرا = $(10+6)/2 = 8$ ثانیه

میانگین زمان اجرای وزنی A = $2*0.2 + 12*0.8 = 10$ ثانیه

میانگین زمان اجرای وزنی B = $6*0.2 + 10*0.8 = 9.2$ ثانیه

خلاصه ای از کارایی

- ◆ 2 سنجش اصلی کارایی
 - ◆ زمان اجرا: زمان انجام کار
 - ◆ توان عملیاتی: تعداد کارهایی هستند که در یک واحد زمانی کامل می شوند
 - ◆ کارایی و زمان اجرا متقابل اند افزایش کارایی کاهش زمان اجرا را در پی دارد
 - ◆ زمان اجرایی که به یک برنامه اختصاص داده می شود می تواند به صورت زیر محاسبه شود
 - زمان CPU = تعداد دستورالعمل ها * CPI * زمان هر چرخه
 - زمان CPU = تعداد دستورالعمل ها * CPI / آهنگ پالس ساعت
- این عوامل متاثر است از تکنولوژی کامپایلر ، معماری مجموعه دستورالعمل ها ، سازمان دهی ماشین ، و تکنولوژی زیر بنایی . زمانی که کارایی افزایش می یابد دقت کنید که چه رخدادهایی به طور متناوب اتفاق می افتد => نمونه متداول سریعی ساخته می شود.
- به خاطر داشته باشید: بخشهای اشتراکی را سریعتر نمائید!

آزمون های کارایی کامپیوتر

- ◆ آزمون کارایی ، یک برنامه یا مجموعه ای از برنامه هایی است که برای دستیابی کامپیوتر به کارایی استفاده می شود
- ◆ آزمون های کارایی به ما اجازه می دهند سنجش کارایی بر پایه زمان های اجرایی را می دهد
- ◆ آزمون های کارایی باید نمایش دهنده نوع درخواست های در حال اجرا در یک کامپیوتر باشند
- ◆ آزمون های کارایی نباید به عوامل دیگری در کامپیوتر وابسته باشند ، مثل موس ، کیبورد ، و...
- ◆ آزمون های کارایی می توانند در گونه های پیچیده و سودمندشان بسیار گوناگون باشند.

انواع آزمون های کارایی

مزایا

- نمایشگر

- قابل حمل
- کاملا کاربردی
- پیشرفت های مفید در واقعیت

- آسان برای اجرا در اوایل طراحی چرخه

- قابلیت به اوج رسانیدن همانندی و پتانسیل گلوگاه ها

12/13/2013

هدف اصلی بار کاری

درخواست کامل آزمون های کارایی
(مثل آزمون های کارایی spec)

Small "Kernel"
Benchmarks

آزمون های کارایی میکرو

\\cpeg323-04F\Topic0.ppt

معایب

- خیلی مشخص
- غیر قابل حمل
- برای اجرا یا سنجش سخت است
- برای همانندسازی سخت است

- نشان دهنده کم

- سیستم حافظه قابل اندازه گیری نیست

- به اوج رسیدن "می تواند یک راه طولانی از درخواست کارایی باشد

231

SPEC

(System Performance Evaluation Cooperative)

موسسه ارزیابی کارایی سیستم

- ◆ آزمون های کارایی SPEC به شدت برای گزارش کارایی و کارایی PC به کار می رود
- ◆ اولین دوره SPEC CPU89
 - شامل 10 برنامه روی اعداد single
- ◆ دومین دوره SPEC CPU92
 - Spec CINT92 (6 تا برنامه صحیح) و Spec CFP92 (14 تا برنامه ممیز شناور)
 - فلگ های کامپایلر می توانند قرار بگیرند به طور متفاوت برای برنامه های متفاوت
- ◆ سومین دوره SPEC CPU95
 - قرار گیری جدید برنامه ها: SPEC CINT 95 (8 تا برنامه صحیح) و SPEC CFP 95 (10 تا برنامه ممیز شناور)
 - تنها فلگ کامپایلری که قرار می گیرد برای همه برنامه ها .

SPEC

◆ چهارمین دوره SPEC CPU2000

- قرار گیری جدید برنامه ها: SPEC CINT2000 (12 تا برنامه صحیح) و SPEC CFP2000 (14 تا برنامه ممیز شناور)
- تنها فلگ کامپایلری که قرار می گیرد برای همه برنامه ها .
- گزارشات مستند نسبت ارزیابی برای SPEC هستند
- زمان CPU ماشین مورد ارزیابی / زمان CPU ماشین مرجع

دیگر آزمون های کارایی SPEC

◆ :JVM98

● اندازه گیری کارایی ماشین های مجازی

◆ :SFS97

◆ اندازه گیری قوانین فایل سرویس دهند شبکه (NFS)

◆ :Web99

◆ اندازه گیری کارایی در خواست های شبکه گسترده جهانی

◆ :HPC96

◆ اندازه گیری کارایی گسترده در خواست های صنعتی

◆ APC, MEDIA, OPC

◆ اندازه گیری کارایی درخواست های گرافیکی

◆ برای اطلاعات بیشتر درباره آزمون های کارایی SPEC

به <http://www.spec.org> مراجعه کنید

مثال هایی از آزمون های کارایی SPEC95

◆ نسبت های SPEC نشان داده می شوند در پردازش
گره های Pentium و Pentium pro (+)

Clock Rate	Pentium SPECint	Pentium+ SPECint	Pentium SPECfp	Pentium+ SPECfp
100 MHz	3.2	N/A	2.6	N/A
150 MHz	4.4	6.0	3.0	5.1
200 MHz	5.5	8.0	3.8	6.8

از این اطلاعات چه چیزی را می توان فرا گرفت؟

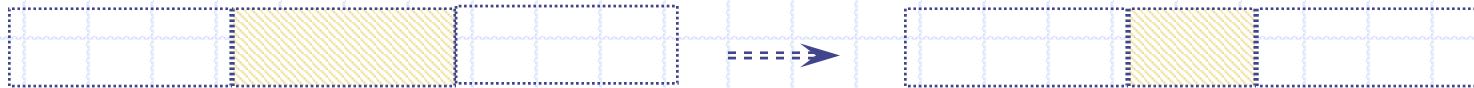
عوامل ضعیف ارزیابی کارایی

- اندازه گیری کارایی کامپیوتر شامل MIPS و MFLOPS می شود
- MIPS: میلیون ها دستورالعمل در هر ثانیه
_MIPS = (زمان اجرا * 10^6) / تعداد دستورالعمل ها
_ برای مثال یک برنامه که 3 میلیون دستورالعمل را در 2 ثانیه محاسبه می کند درجه MIPS آن 1.5 است
_ مزایا: آسان است برای فهمیدن و اندازه گیری
_ معایب: ممکن است کارایی واقعی را منتقل نکند در حالی که دستورالعمل های ساده تر را بهتر اجرا می کند.
- MFLOPS: میلیون ها عملگر ممیز شناور در هر ثانیه
_MFLOPS = (زمان اجرا * 10^6) / عملگرهای ممیز شناور
_ برای مثال برنامه S که 4 میلیون دستورالعمل را در 5 ثانیه محاسبه می کند درجه MFLOPS آن 0.8 است.
_ مزایا: آسان است برای فهمیدن و اندازه گیری
_ معایب: مانند MIPS تنها اندازه گیری های ممیز شناور

بالا بردن میزان تسریع به صورت زیر تعریف می شود

زمان اجرای قدیم کارایی جدید

$$\text{میزان تسریع} = \frac{\text{زمان اجرای قدیم}}{\text{زمان اجرای جدید}} = \text{کارایی جدید}$$



◆ این عامل انکسار و انحراف را برای میزان تسریع چنین در نظر می گیریم

$$\text{میزان تسریع افزوده شده} = \frac{\text{Fraction}_{\text{enhanced}} + (1 - \text{Fraction}_{\text{enhanced}}) \times \text{زمان اجرای قدیم}}{\text{زمان اجرای جدید}}$$

$$\text{میزان تسریع} = \frac{\text{زمان اجرای قدیم}}{\text{زمان اجرای جدید}} = \frac{1}{\text{Fraction}_{\text{enhanced}} + (1 - \text{Fraction}_{\text{enhanced}})}$$

میزان تسریع افزوده شده

مثالی از قانون آمدال

◆ دستورالعمل‌های ممیز شناور برای 2 برابر شدن سرعت اجرا بهبود داده شده اند، اما تنها 10% زمان برای این دستورالعمل‌های اصلی استفاده می‌شود. ماشین جدید چقدر سریعتر است؟

$$\text{میزان تسریع} = \frac{\text{زمان اجرای قدیم}}{\text{زمان اجرای جدید}} = \frac{1}{(1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}}}$$

$$\text{میزان تسریع} = \frac{1}{(1 - 0.1) + 0.1/2} = 1.053$$

ماشین جدید 1.053 سرعت دارد، یا 5.3% سریعتر است

اگر سرعت دستورالعمل‌های ممیز شناور 100 برابر سریعتر شود ماشین جدید چقدر سریعتر می‌شود؟

$$\text{میزان تسریع} = \frac{1}{(1 - 0.1) + 0.1/100} = 1.109$$

تخمین بهبود کارایی

◆ فرض کنید یک پردازشگر به طور جاری به 10 ثانیه برای اجرای یک برنامه نیاز دارد و کارایی یک پردازشگر هر سال 50% بهبود می یابد.

◆ به وسیله چه عاملی کارایی پردازشگر در 5 سال بهبود پیدا می کند؟

$$(1+0.5)^5=7.59$$

بعد از 5 سال یک پردازشگر همان برنامه را در چه زمانی اجرا می کند؟
ثانیه $10/7.59=1.32$ = زمان اجرای جدید

چه فرضیاتی در مسئله بالا پدید می آیند؟

مثال کارایی

کامپیوتر های M1 و M2 اجزای مجموعه دستورالعمل های یکسانی دارند
M1 آهنگ پالس ساعت 50MHz دارد و M2 آهنگ پالس ساعت 75MHz دارد. M1 دارای CPI 2.8 است و
M2 دارای CPI 3.2 برای برنامه دارد.
چقدر زمان M2 سریعتر است نسبت به M1 در این برنامه؟

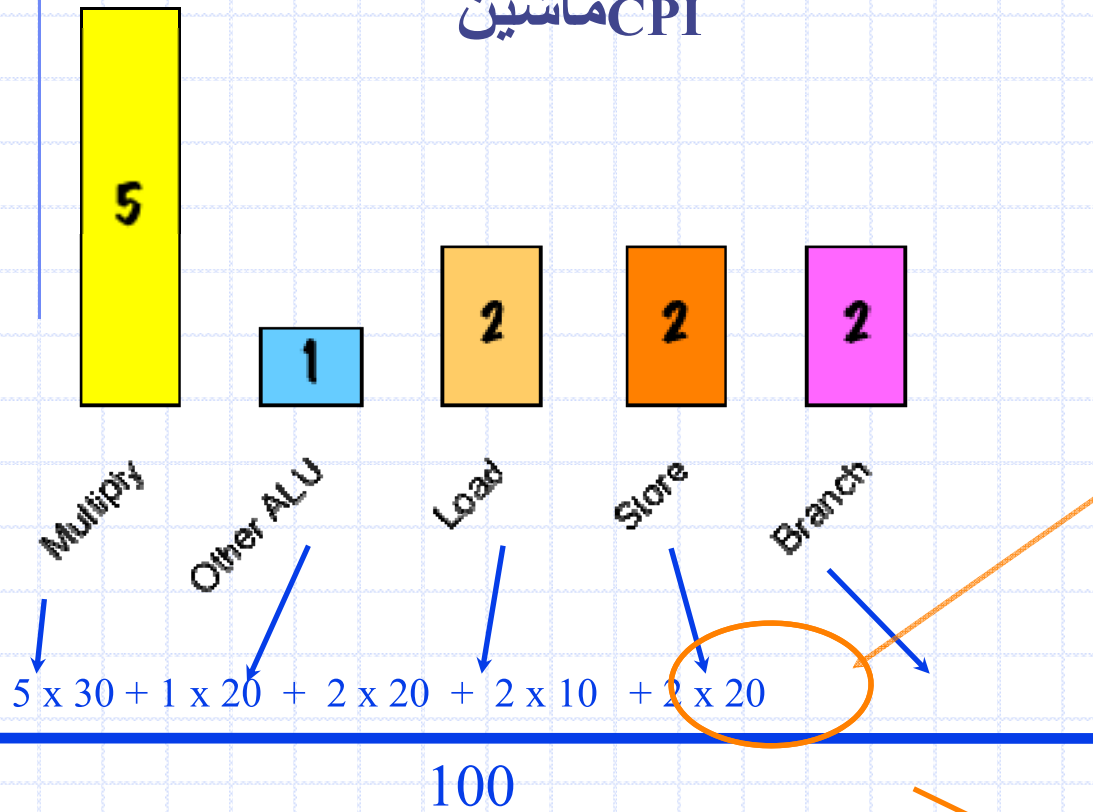
$$\frac{\text{زمان اجرای M1}}{\text{زمان M2}} = \frac{\text{M1 آهنگ پالس ساعت} / \text{CPI}_{M1} \times \text{تعداد دستورالعملها}}{\text{M2 آهنگ پالس ساعت} / \text{CPI}_{M2} \times \text{تعداد دستورالعملها}}$$
$$\frac{2.8/50}{3.2/75} = 1.31$$

آهنگ پالس ساعت M1 چقدر باید باشد تا زمان اجرای یکسانی برای آنها داشته باشد؟

CPI به عنوان یک ابزار تحلیلگر برای راهنمای طراحی

CPI ماشین

میکس دستورالعمل های برنامه



$$5 \times 30 + 1 \times 20 + 2 \times 20 + 2 \times 10 + 2 \times 20$$

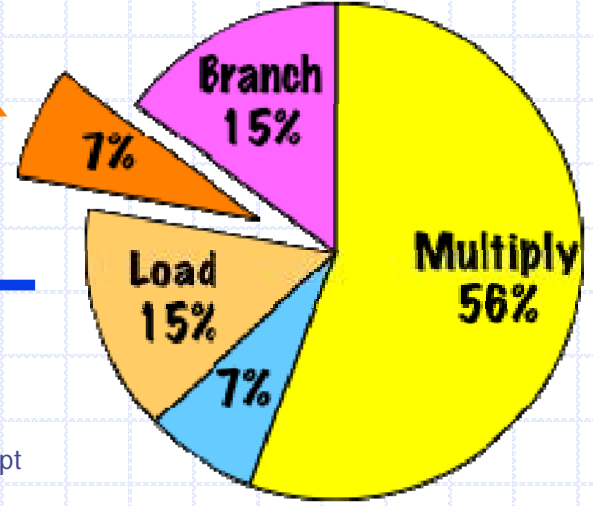
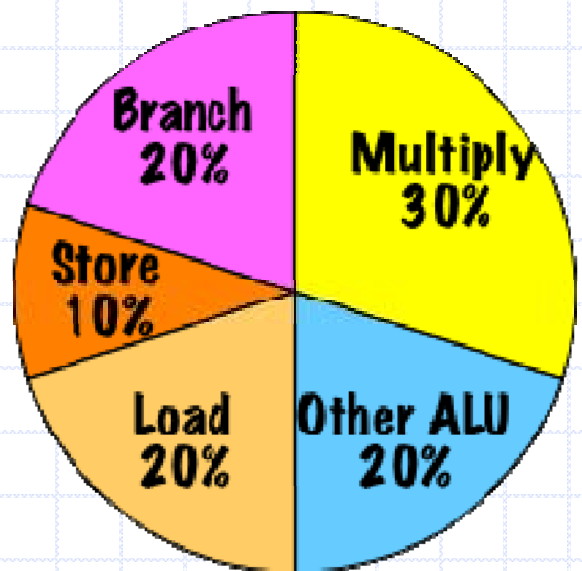
100

= 2.7 cycles/instruction

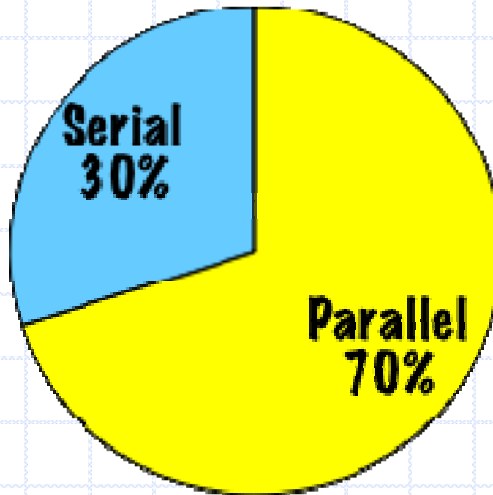
$$5 \times 56 + 1 \times 7 + 2 \times 15 + 2 \times 7 + 2 \times 15$$

100

= 3.61 cycles/instruction
12/13/2013



دستور العمل PEER: قانون آمدال



برنامه 30% از زمانش را به کدهایی اختصاص می دهد که نمی توانند دوباره برای اجرا به صورت موازی کد شوند

میزان تسریع محاسبه می شود برای

$$N = 2, 3, 4, 5, \text{ and } \infty.$$

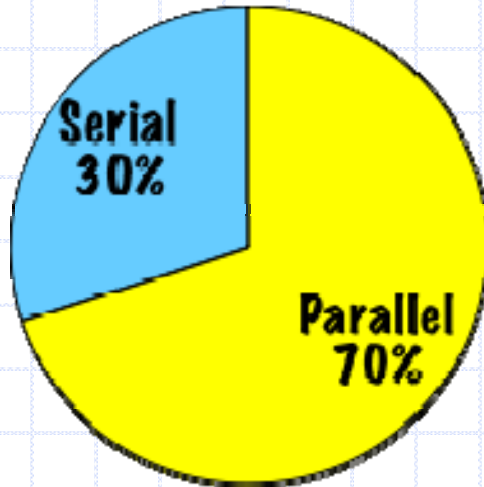
CPUs	2	3	4	5	∞
میزان تسریع					

دستور العمل PEER: قانون آمدال

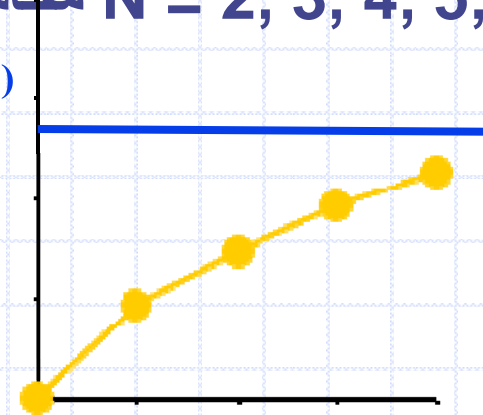
برنامه 30% از زمانش را به صورت سریال استفاده می کند

میزان تسریع را برای

$N = 2, 3, 4, 5, \infty$ محاسبه می کند



$S(\infty)$



$$S = \frac{1}{0.30 + (0.70 / N)}$$

CPUs	2	3	4	5	∞
میزان تسریع	1.54	1.85	2.1	2.3	3.3

یک ژول انرژی: یک جریان 1 آمپری است که از وسط یک مقاومت
1 اهمی عبور می کند و گرمایی معادل 1 ژول تولید می کند

و همچنین یک وات برای یک ثانیه

1 وات: عبور 1 آمپر جریان از یک مقاومت 1 اهمی

انرژی و کارایی

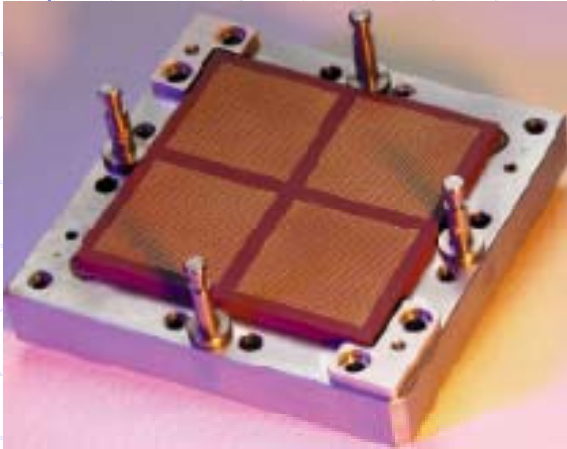
کالری $0.24 =$ یک ژول

یک کالری دمای یک گرم آب را یک درجه افزایش می دهد

یک مورد ناخواسته: کامپیوترها انرژی الکتریکی را به گرما تبدیل
که به عنوان فرآورده های جانبی محاسبه می شود. می کنند.

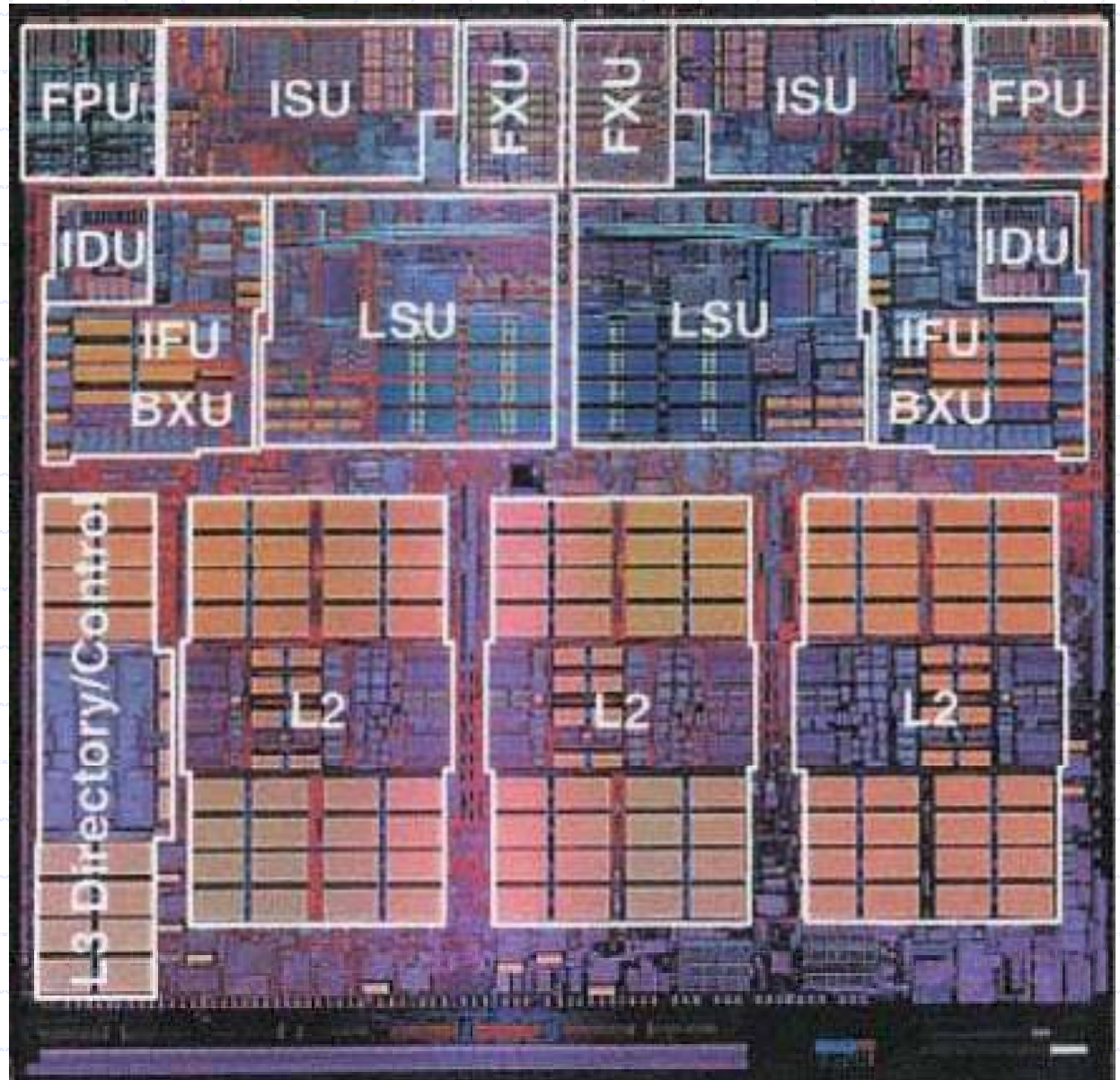
هوا یا آب می توانند مانع ذوب شدن چیپ در اثر انتقال حرارت شود

IBM Power 4: How does die heat up?

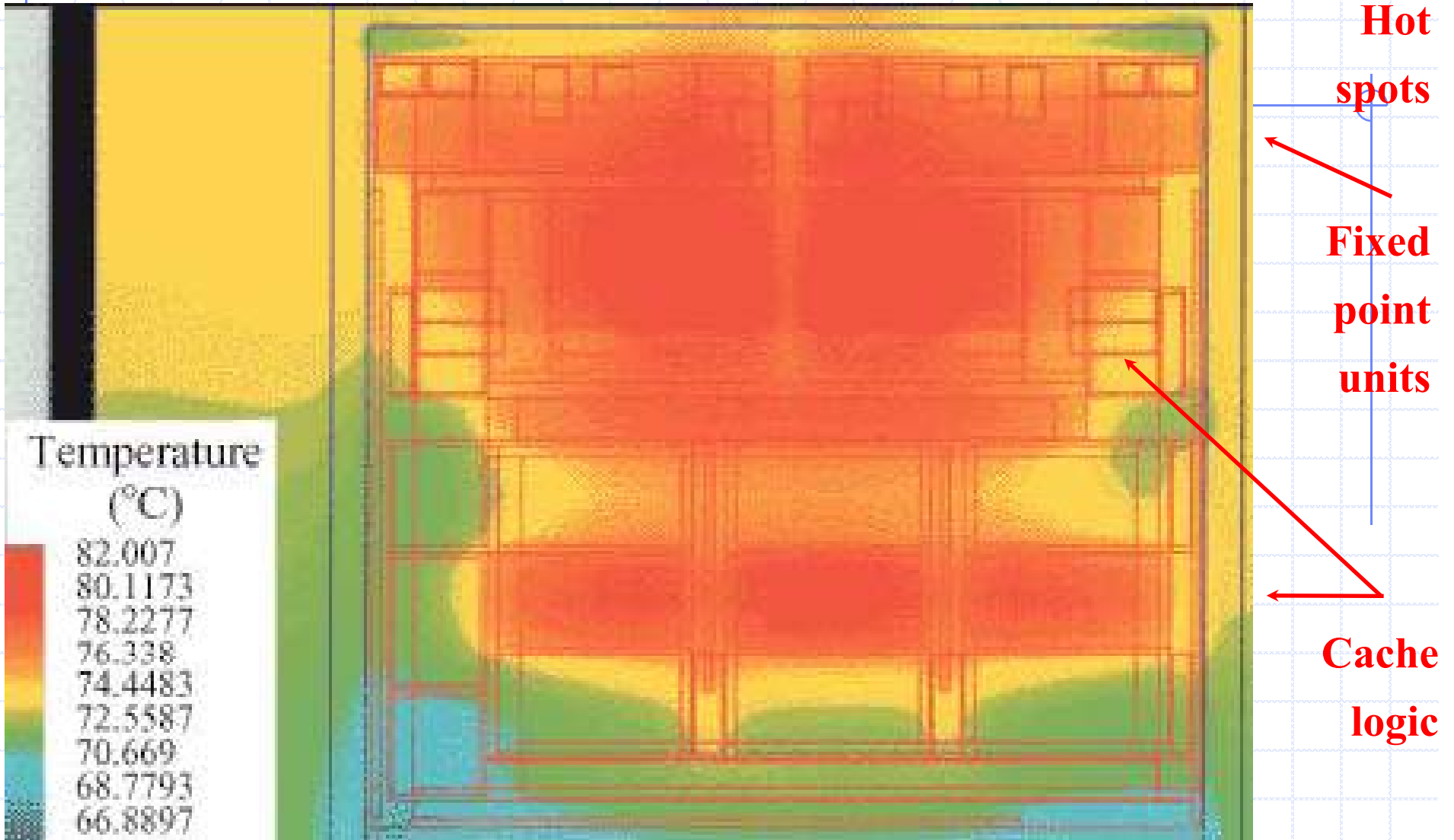


↑
**4 dies on a
multi-chip
module**

**2 CPUs
per die** →

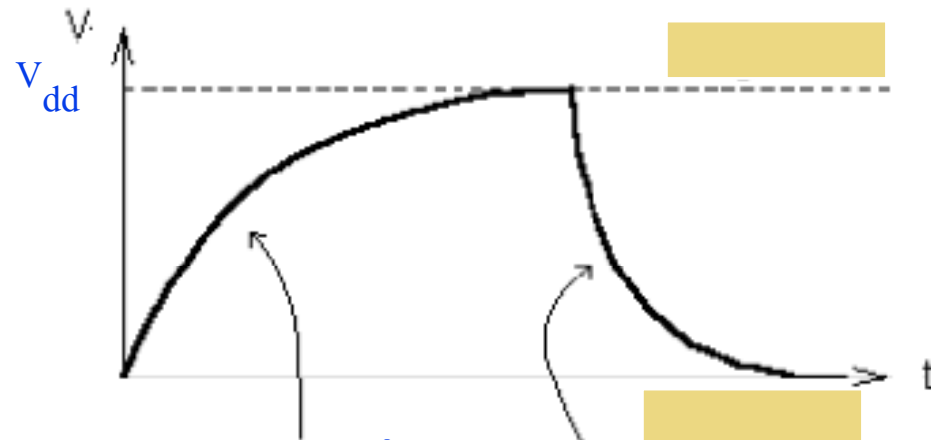
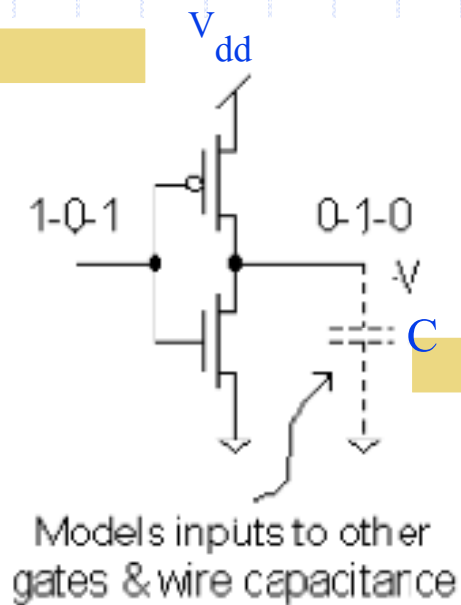


IBM Power 4: Dissipating 115 Watts



اتصال انرژی: فیزیک های بنیادی

انتقال منطقی اتلاف انرژی



$$E_{0 \rightarrow 1} = \frac{1}{2} C V_{dd}^2 \quad E_{1 \rightarrow 0} = \frac{1}{2} C V_{dd}^2$$

نتیجه قطعی: مستقل از تکنولوژی

How can we limit switching energy?

State-of-the-art CPUs (90 nm):

Switching energy is 70% of total energy.

Remainder: at 90nm, “switches” are “dimmers”!

خلاصه ای از ارزیابی کارایی

◆ آزمون های کارایی خوب از قبیل آزمون های کارایی SPEC می تواند روش خاصی را برای ارزیابی و مقایسه کارایی کامپیوتر تولید کند

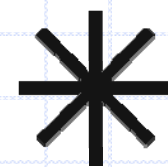
◆ MIPS و MFLOPS برای استفاده آسان هستند اما برای نشان دادن کارایی نا صحیح هستند.

◆ قانون آمدال یک روش موثر را برای توضیح میزان تسریع وابسته به عامل انحراف ایجاد می کند.

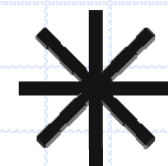
به خاطر داشته باشید: بخشهای اشتراکی را سریعتر نمائید!

نتیجه گیری

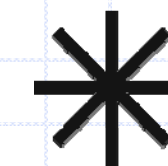
مشتری ها: قیمت را می سنجند
معمارها : طرح را می سنجند



ابزار آلات: CPI و معادله کارایی



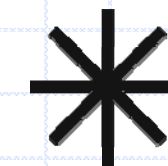
قانون آمدال: محاسبه می کند



انرژی:

$$E_{0 \rightarrow} = \frac{1}{2} C V_{dd}^2$$

$$E_{1 \rightarrow} = \frac{1}{2} C V_{dd}^2$$



فصل پنجم

روشهای پیاده سازی معماری پردازشگر

چکیده

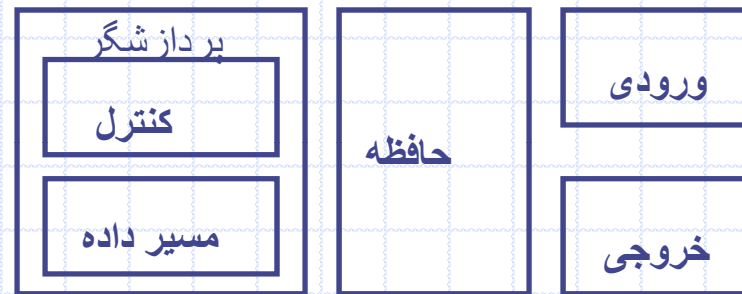
- مقدمه
- بخشهای سازنده مسیر داده
- پیاده سازی یک مسیر داده ساده (تک چرخه ای)
- پیاده سازی چند چرخه ای

مقدمه

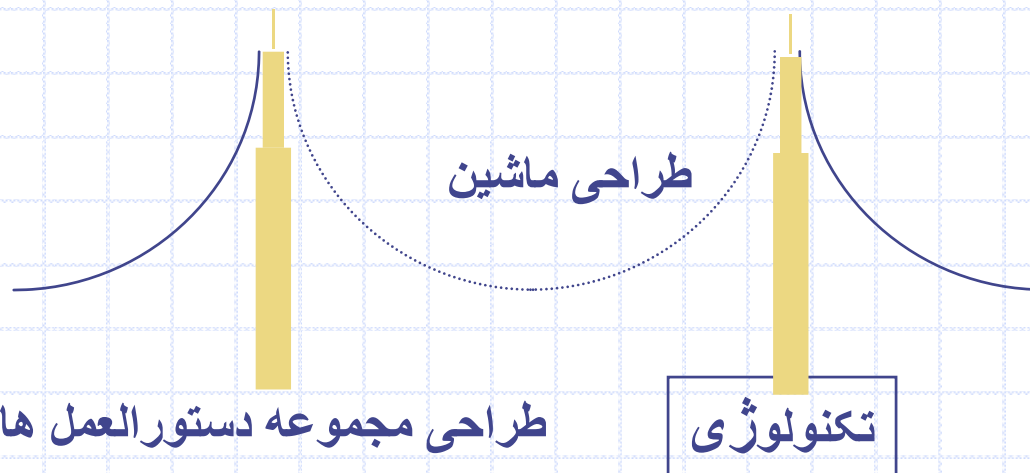
- روند پایه اجرای دستورالعملها
- برخی قراردادهای
- مسیر داده 32 بیتی
- استراتژی ساعت زنی (حساس به لبه)
- ما به زیر مجموعه ای از MIPS تمرکز می کنیم:
- دستورالعملهای دستیابی به حافظه: lw ، sw
- اعمال ALU: add, sub, and, or
- دستورات انشعاب در صورت تساوی (beq) و دستورات پرش (j)

یک تصویر بزرگ: ما الان کجا هستیم؟

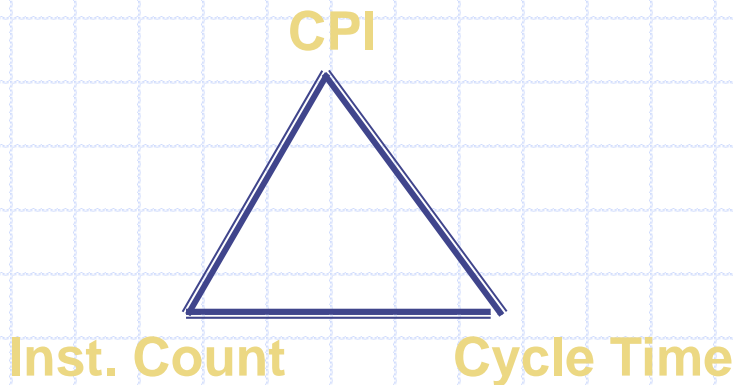
◆ پنج نمونه از بخشهای سازنده کامپیوتر



◆ موضوع روز: طراحی یک پردازشگر تک چرخه ای



یک شمای اصلی: چشم انداز کارایی



◆ کارایی یک ماشین تعیین می شود به وسیله:

■ تعداد دستور العمل ها (Inst. Count)

■ زمان هر چرخه (Cycle Time)

■ چرخه های هر دستور العمل (CPI)

◆ طراحی پردازشگر (مسیر داده و کنترل) تعیین خواهد شد:

■ زمان هر چرخه

■ چرخه های هر دستور العمل

◆ پردازشگر تک چرخه ای یک کلاک دارد برای هر دستور العمل

■ مزایا: طراحی ساده و CPI کم.

■ معایب: زمان چرخه طولانی است و این باعث می شود که به وسیله کندترین

دستور العمل محدود شود

چه طور یک پردازشگر طراحی می شود: مرحله به مرحله

- ◆ تجزیه و تحلیل مجموعه دستورالعمل = <پیش نیازهای مسیر داده
- مفهوم هر دستورالعمل داده میشود با انتقال دادن ثبات ها

$$R[rd] \leftarrow R[rs] + R[rt];$$

- مسیر داده باید شامل محیط ذخیره سازی برای ثبات های ISA شود
- مسیر داده باید انتقال دادن هر ثبات را پشتیبانی کند
- ◆ انتخاب کردن مجموعه ای از بخشهایی از مسیر داده و بنیاد نهادن روشهای ساعت زنی
- ◆ طراحی مسیر داده با در نظر گرفتن پیش نیاز ها

Analyze implementation of each instruction to determine setting of control points that effects the register transfer

◆ طراحی کنترل منطقی

یادآوری: قالب های دستورالعمل های MIPS

همه دستورالعمل های MIPS طولشان 32 بیت است. 3 تا قالب دستورالعمل هستند.



دارای فیلدهای متفاوتی هستند:

- OP: عملکرد هر دستورالعمل
- Rs, rt, rd: منبع و مقصد مخصوص ثبات ها
- shamt: اندازه شیفت
- Funct: انتخاب کردن عملیات متفاوت در فیلد "op"
- address / immediate: آدرس یا مقدار فوری
- target address: آدرس هدف دستورالعمل پرش

ثبات منطقی انتقال (RTL)

RTL مفهوم دستور العمل را می دهد

همه دستور العمل ها با واکنشی دستور العمل آغاز می شوند

$op \mid rs \mid rt \mid rd \mid shamt \mid funct = MEM[PC]$

$op \mid rs \mid rt \mid Imm16 = MEM[PC]$

inst Register Transfers

addu $R[rd] \leftarrow R[rs] + R[rt];$ $PC \leftarrow PC + 4$

subu $R[rd] \leftarrow R[rs] - R[rt];$ $PC \leftarrow PC + 4$

ori $R[rt] \leftarrow R[rs] + zero_ext(imm16);$ $PC \leftarrow PC + 4$

load $R[rt] \leftarrow MEM[R[rs] + sign_ext(imm16)];$ $PC \leftarrow PC + 4$

store $MEM[R[rs] + sign_ext(imm16)] \leftarrow R[rt];$ $PC \leftarrow PC + 4$

beq $if (R[rs] == R[rt]) then PC \leftarrow PC + 4 + sign_ext(imm16) \parallel 00$

مرحله 1: تعیین پیش نیازهای مجموعه دستورالعمل

حافظه

- دستورالعمل و داده ها

- ثبات ها (32 x 32)

- خواندن rs

- خواندن rt

- نوشتن rd یا rt

- PC شمارنده برنامه

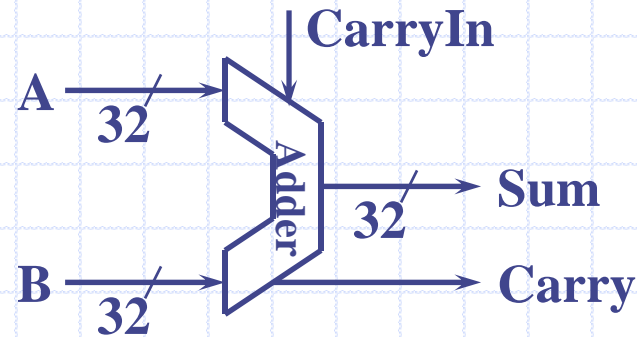
- گسترش دهنده (گسترش علامت یا گسترش صفر)

- جمع و تفریق یا مقدار فوری گسترده شده

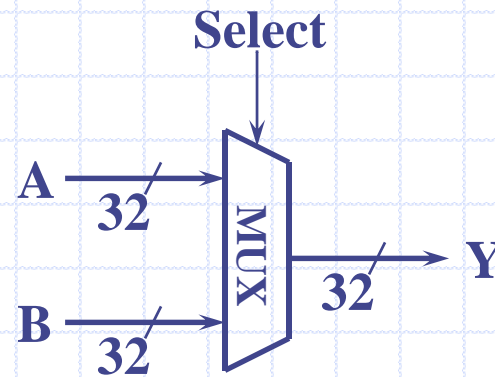
- Add 4 or shifted extended immediate to PC

مرحله 2: اجزاء مسیر داده

◆ جمع کننده

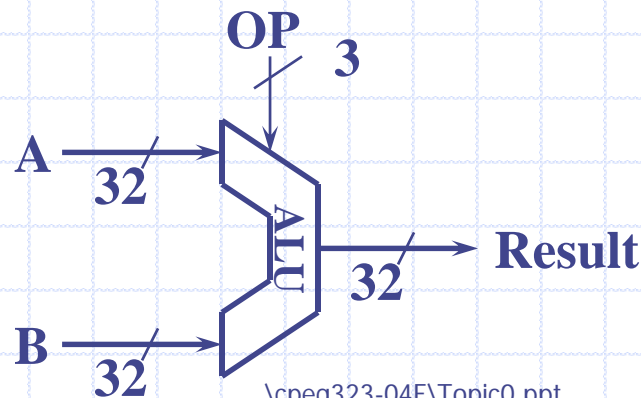


◆ تسهیم کننده



Combinational Logic:
Does not use a clock

◆ واحد محاسبه و منطق



مرحله 3

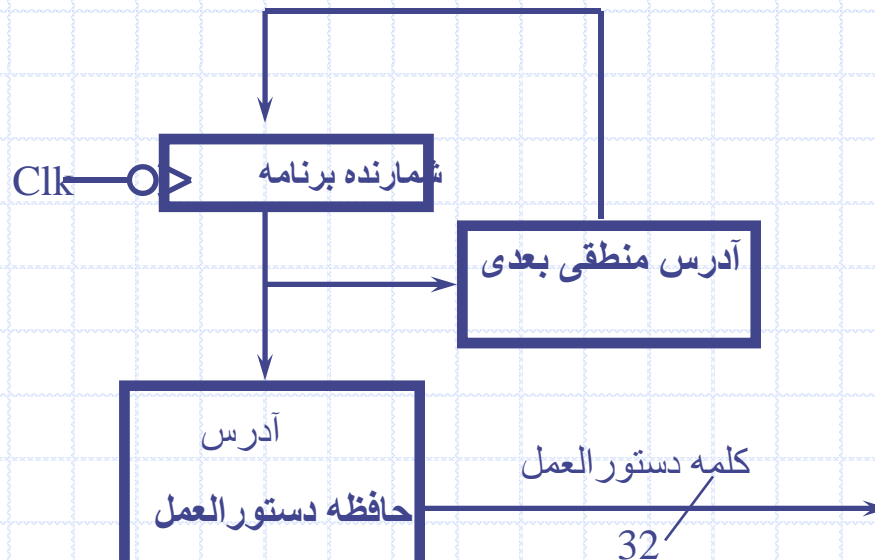
پیش نیازهای ثبات انتقال - < طراحی مسیر داده

- واکنشی دستورالعمل
- کدبرداری دستورالعمل ها و خواندن عملوندها
- محاسبه کردن عمل
- بازنویسی نتیجه

الف 3: مقدمه واحد واكشى دستور العمل

عمليات متداول RTL

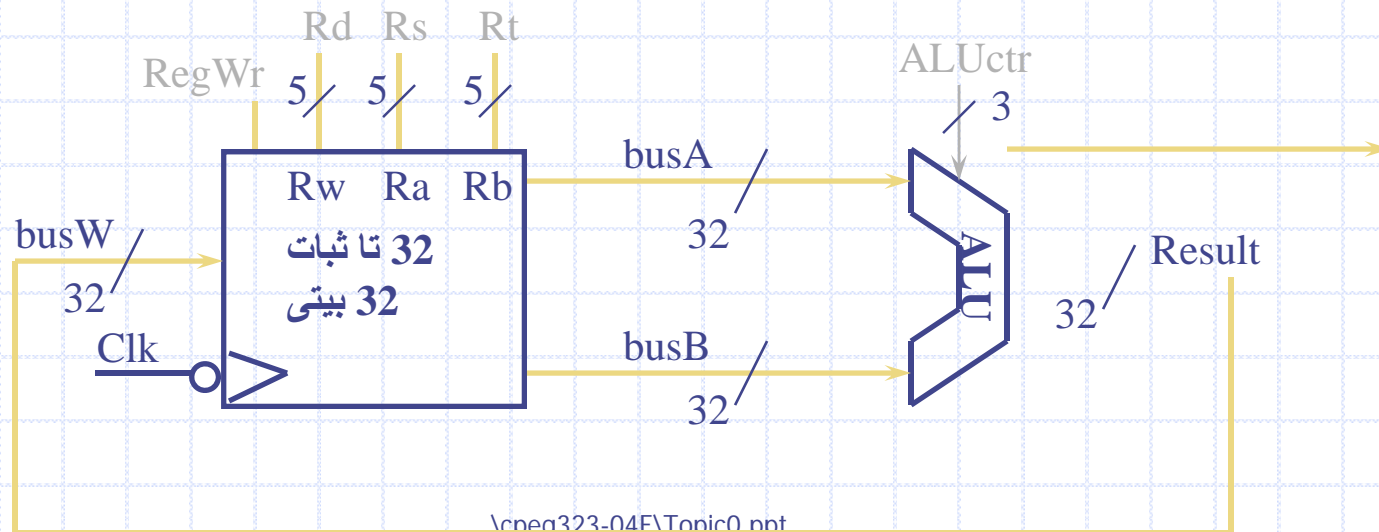
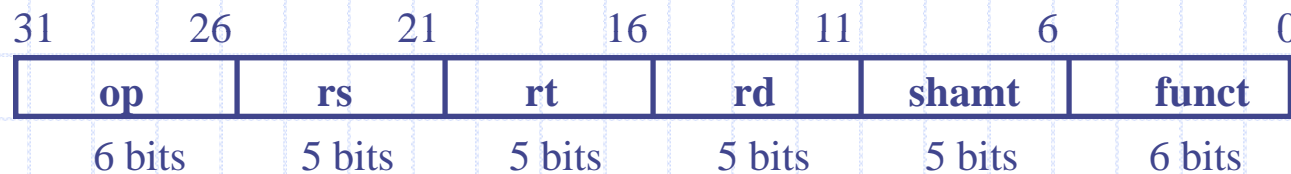
- واكشى دستور العمل: $mem[PC]$
- جديد كرون شمارنده برنامه:
- ◆ كد پشت سرهم: $PC <- PC + 4$
- ◆ انشعاب و پرش: $PC <- \text{"something else"}$



ب3: جمع و تفریق

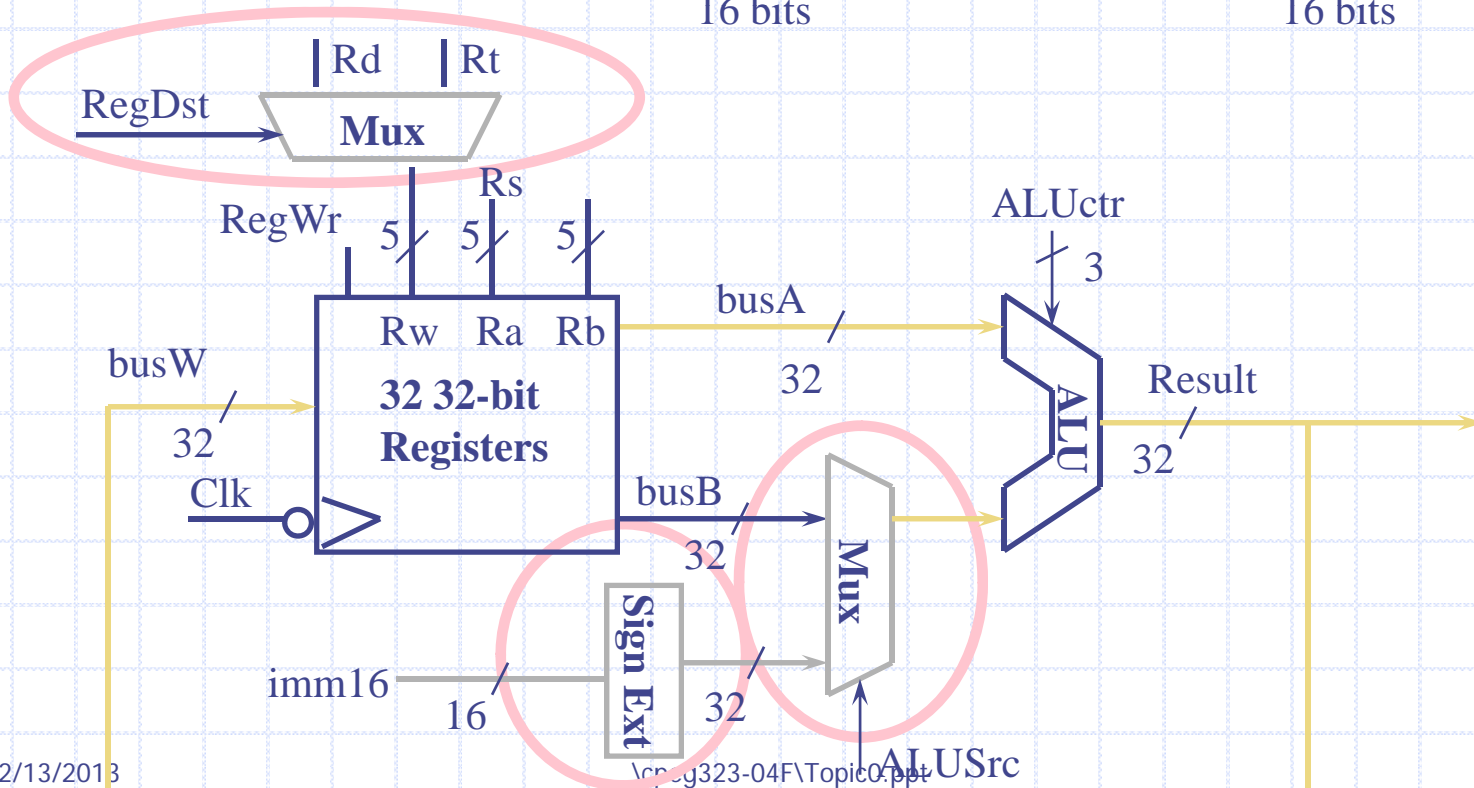
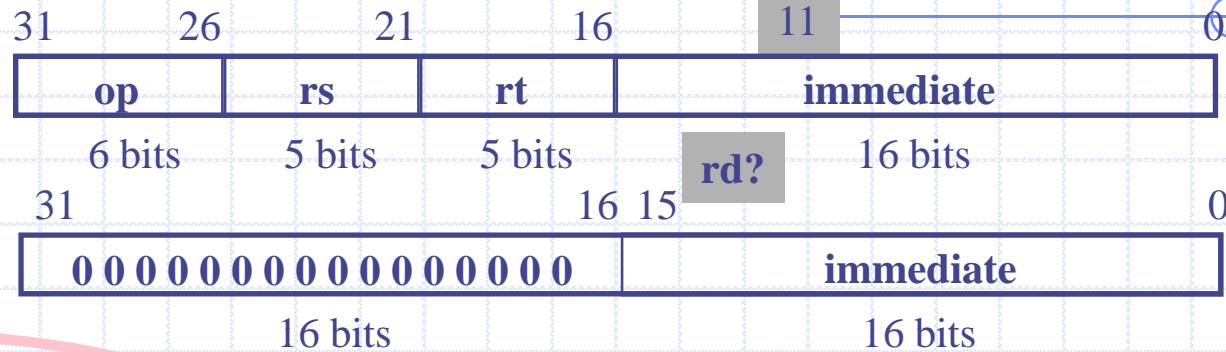
مثال: $R[rd] \leftarrow R[rs] \text{ op } R[rt]$: addu rd, rs, rt

- Ra, Rb, Rw می‌آیند از فیلدهای دستور العمل های rs, rt, rd و
- $ALUctr$ and $RegWr$: بعد از کد برداری دستور العمل کنترل منطقی می کنند



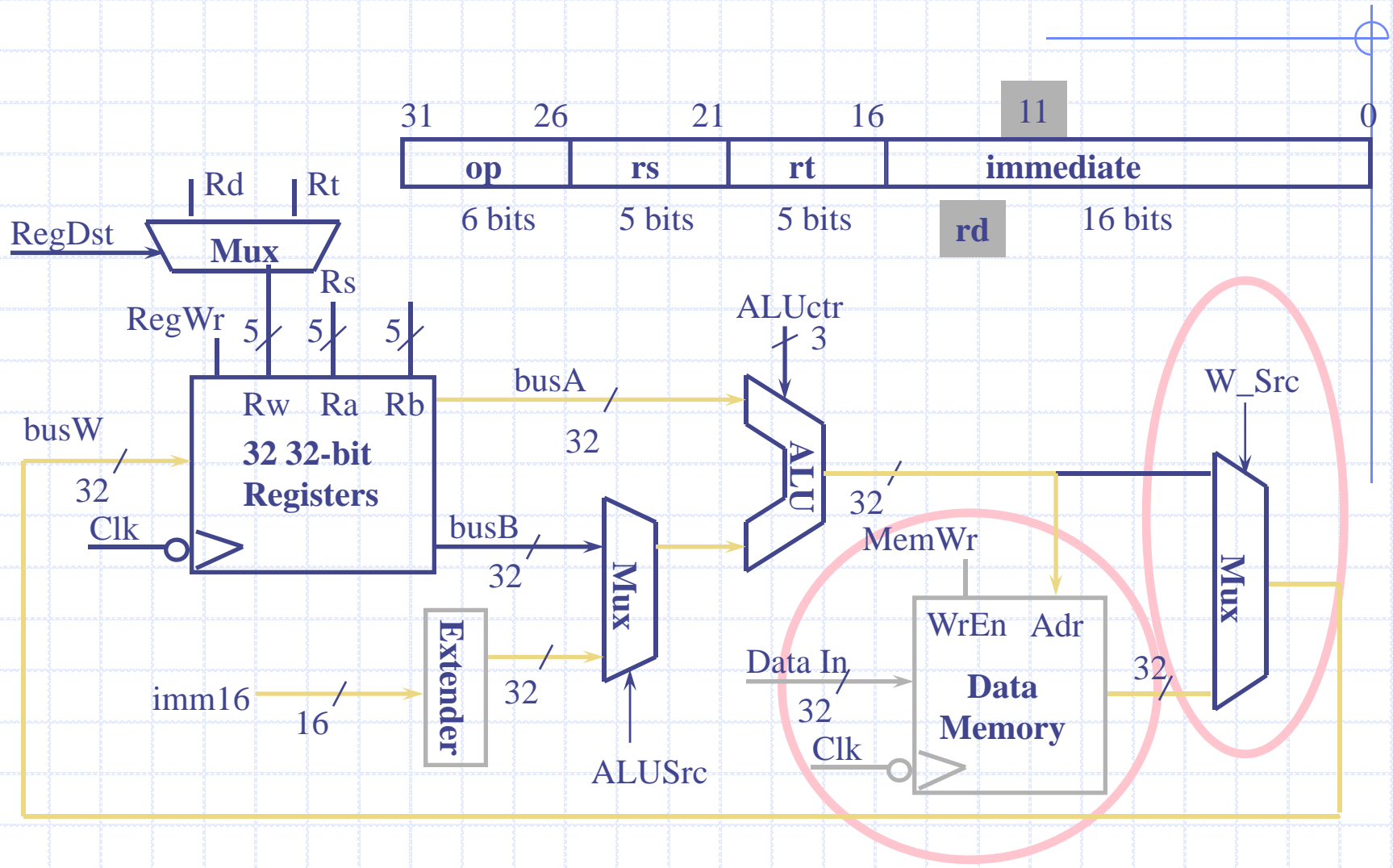
پ3: عمل منطقی بلافاصله

مثال $R[rt] \leftarrow R[rs] \text{ op } \text{imm16}$: ori rt, rs, imm16



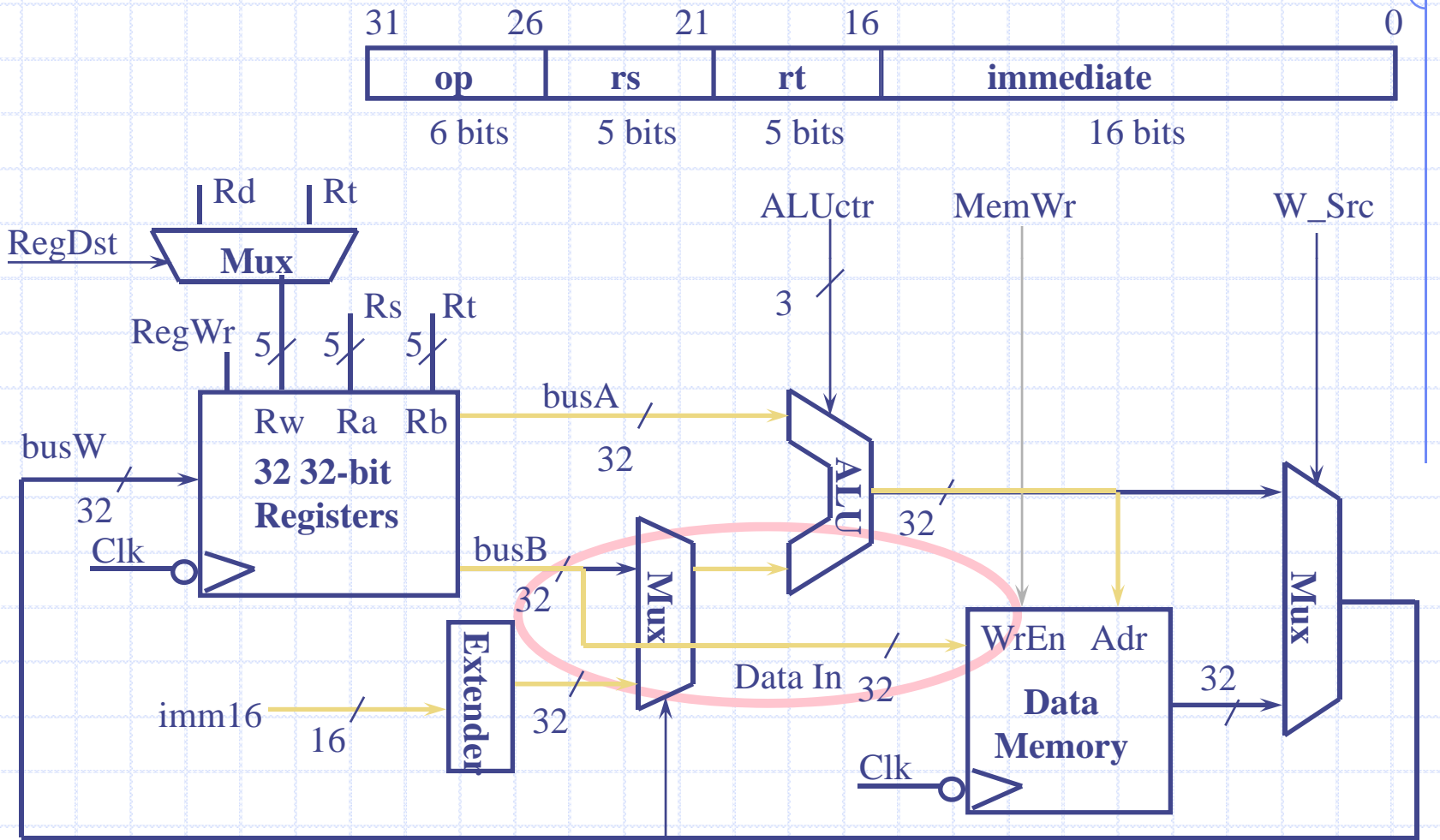
ت3: عمل بار کردن

: lw rt, rs, imm16 [imm16] $R[rt] \leftarrow Mem[R[rs] +$

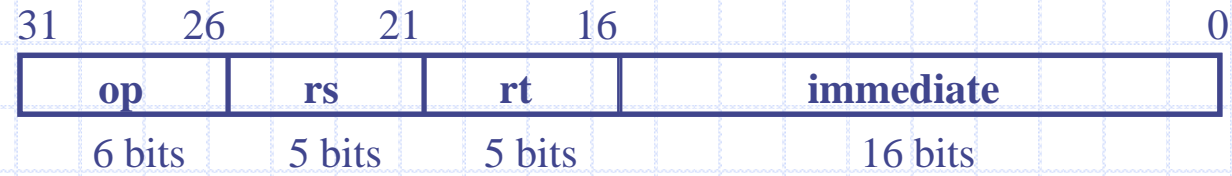


ت3: عمل ذخیره سازی

Mem[R[rs] +] گسترش علامت [imm16] <- R[rt] مثال sw rt, rs, imm16 :



ج3: دستور العمل انشعاب



rs, rt, imm16 beq

واکشی دستور العمل از حافظه

mem[PC]

محاسبه شرط انشعاب

Equal <- R[rs] == R[rt]

محاسبه آدرس دستور العمل بعدی

if (COND eq 0)

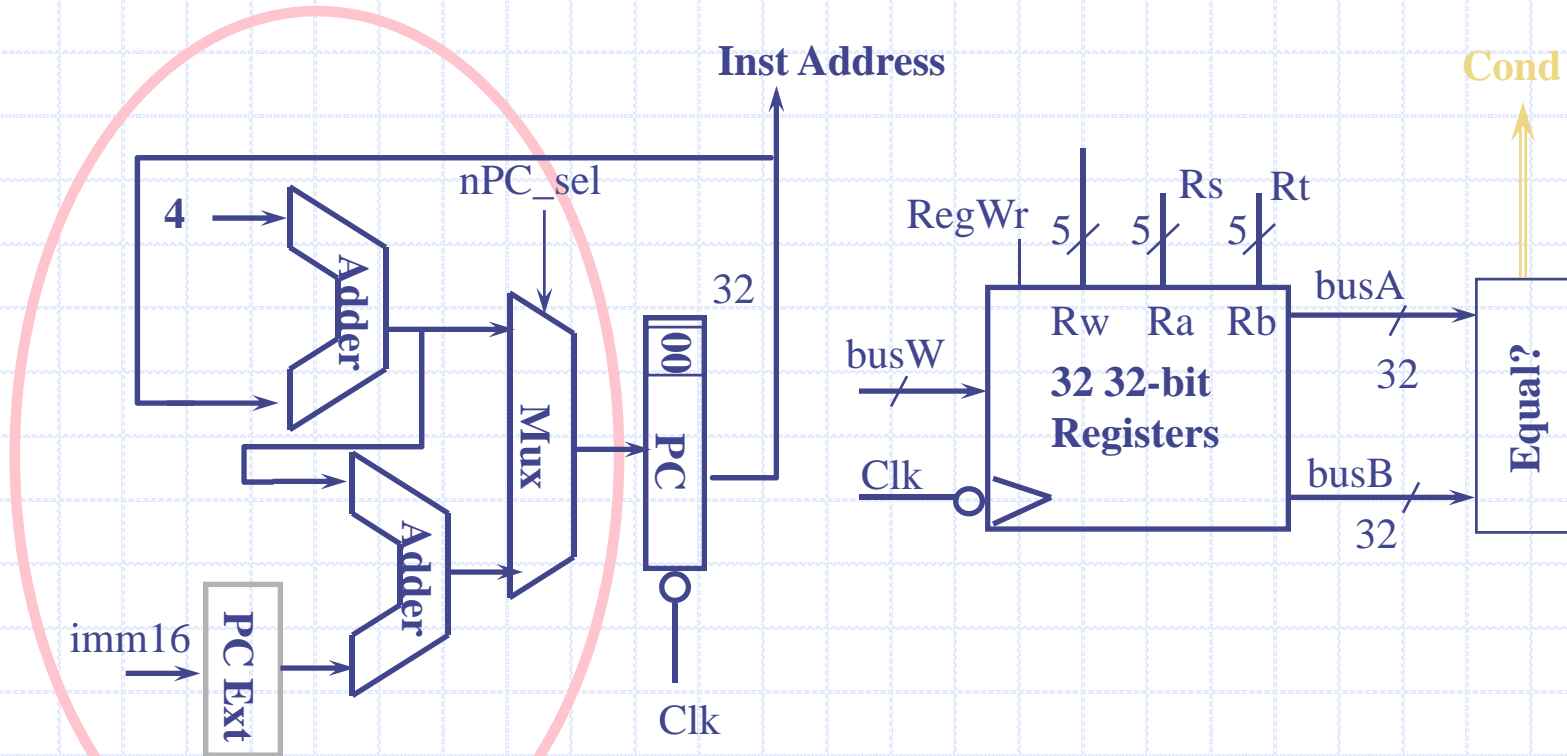
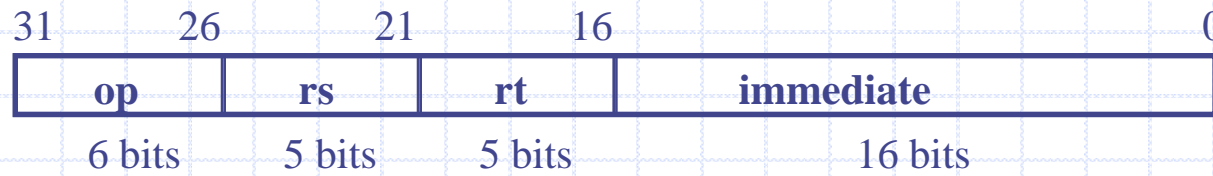
PC <- PC + 4 + (گسترش علامت) (imm16) x 4

else

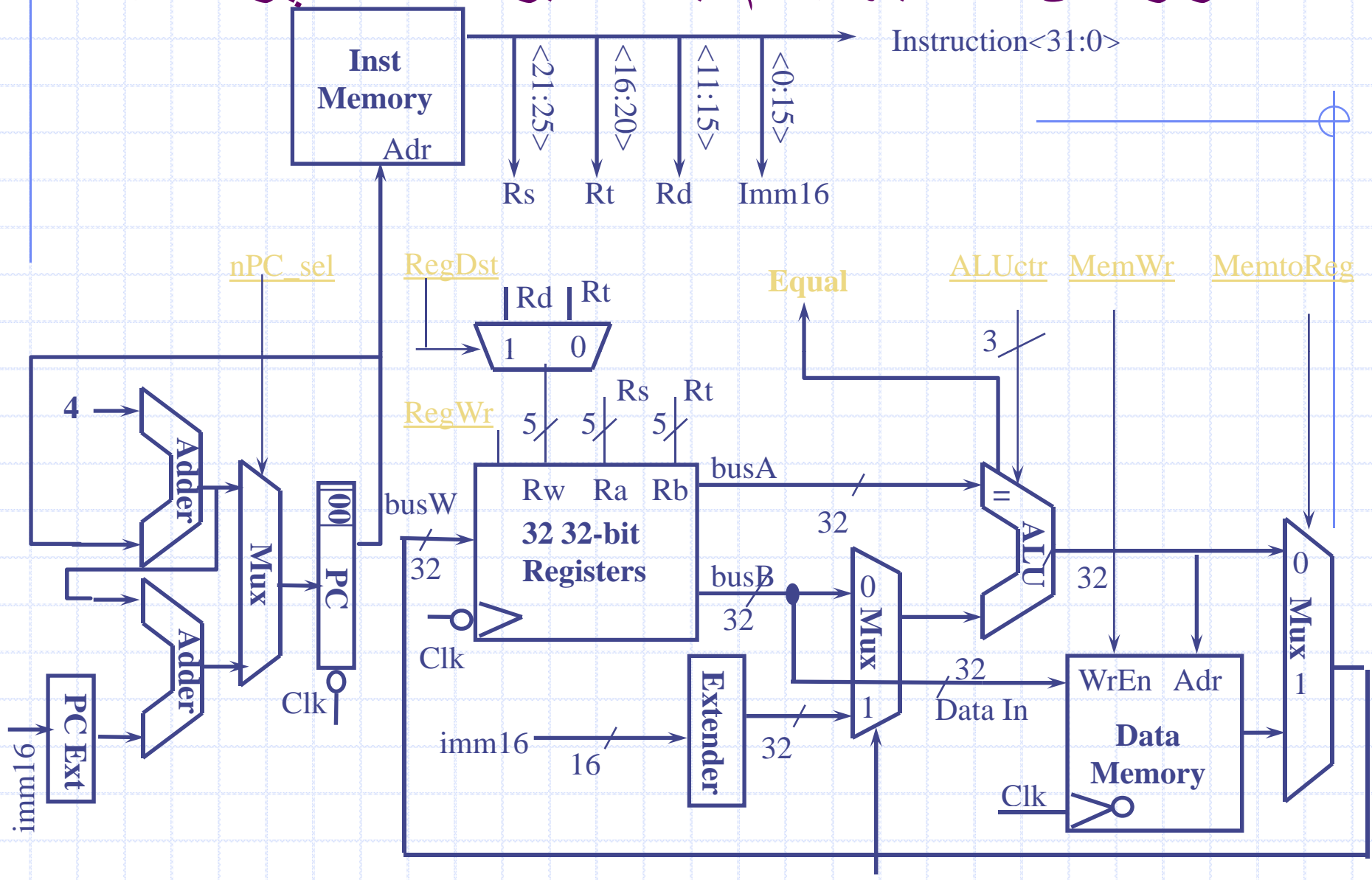
مسیر داده برای عمل انشعاب

beq rs, rt, imm16

ایجاد کردن مسیر داده شرطی (تساوی)



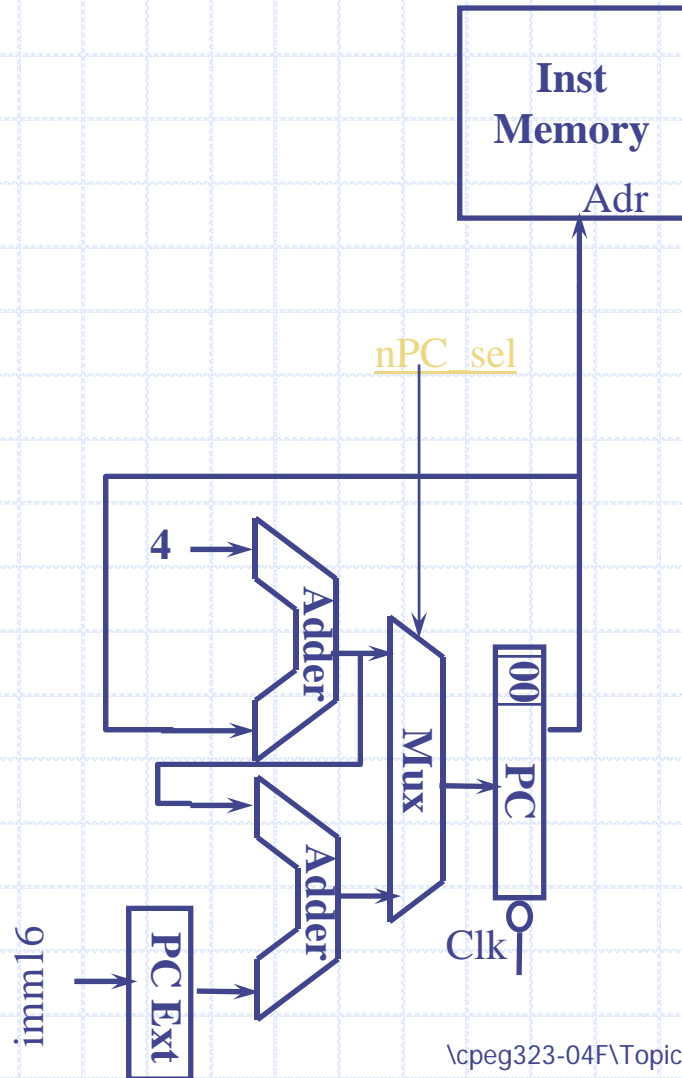
قرار دادن همه اینها باهم: یک مسیر داده تک چرخه ای



مفهوم سیگنال های کنترلی

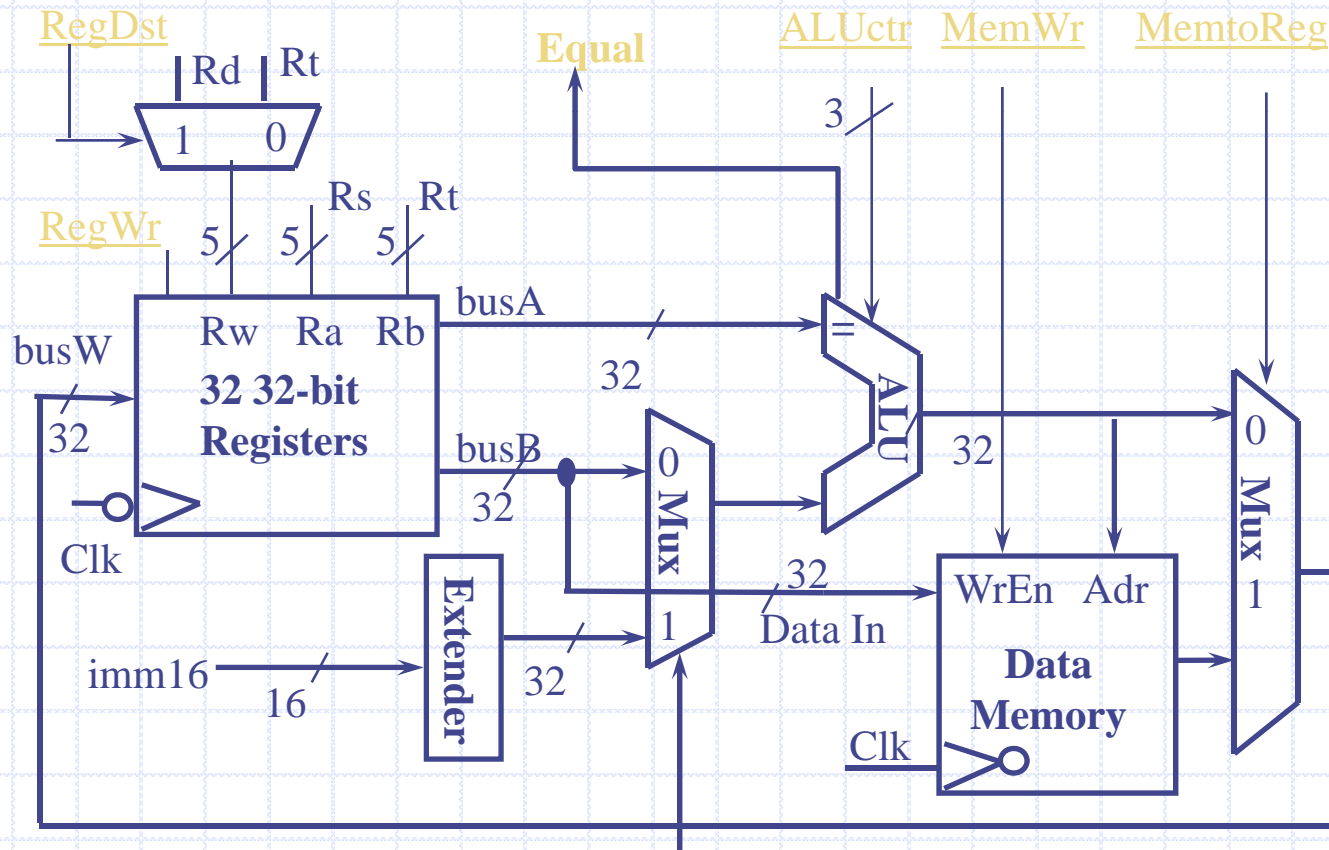
در این مسیر داده Rs, Rt, Rd و $Imed16$ سیم بندی شده اند

$0 \Rightarrow PC \leftarrow PC + 4; 1 \Rightarrow PC \leftarrow PC + 4 + \text{SignExt}(Im16) \parallel 00nPC_sel:$



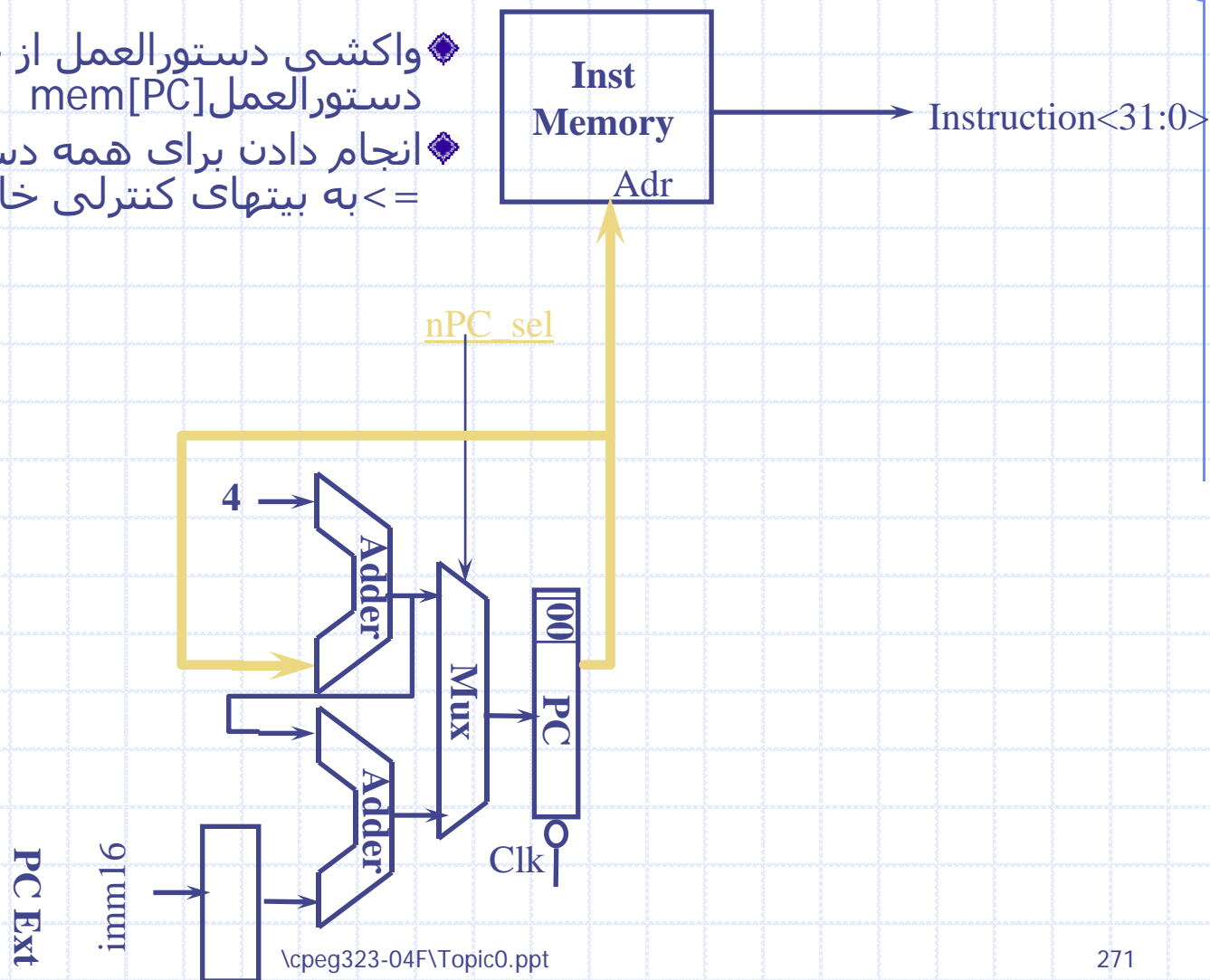
مفهوم سیگنال های کنترلی

- MemWr: نوشتن در حافظه
- MemtoReg: 0 => ALU; 1 => Mem
- ALUctr: "add", "sub", "or"
- RegDst: 0 => "rt"; 1 => "rd"
- RegWr: نوشتن در ثبات مقصد



واحد واکنشی دستورالعمل در ابتدای جمع

- ◆ واکنشی دستورالعمل از حافظه
- دستورالعمل $instruction \leftarrow mem[PC]$
- ◆ انجام دادن برای همه دستورالعمل ها
- \leq به بیت های کنترلی خاصی نیاز نیست

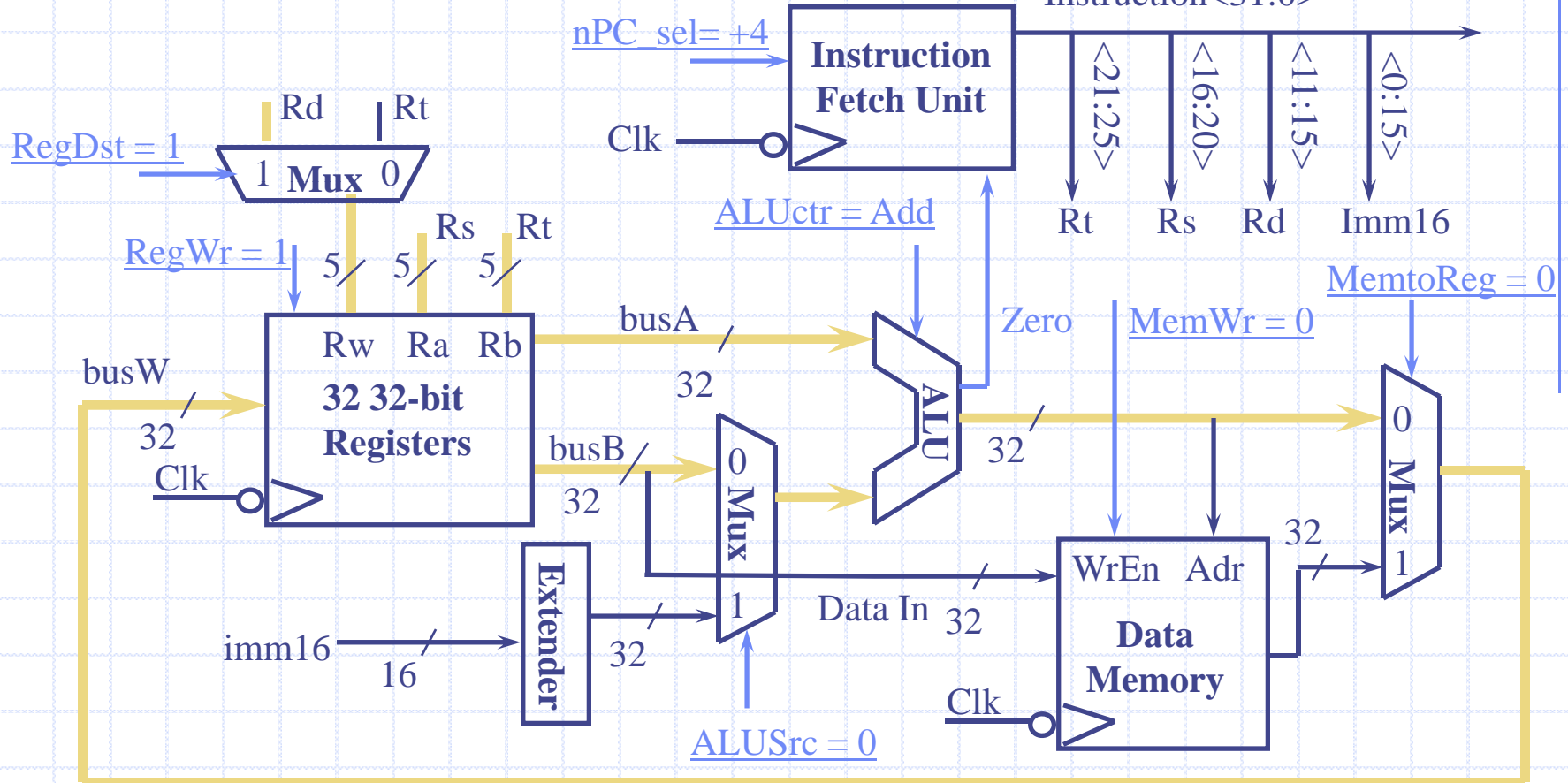


مسیر داده تک چرخه ای در طول عمل جمع و تفریق



$$R[rd] \leftarrow R[rs] \text{ op } R[rt]$$

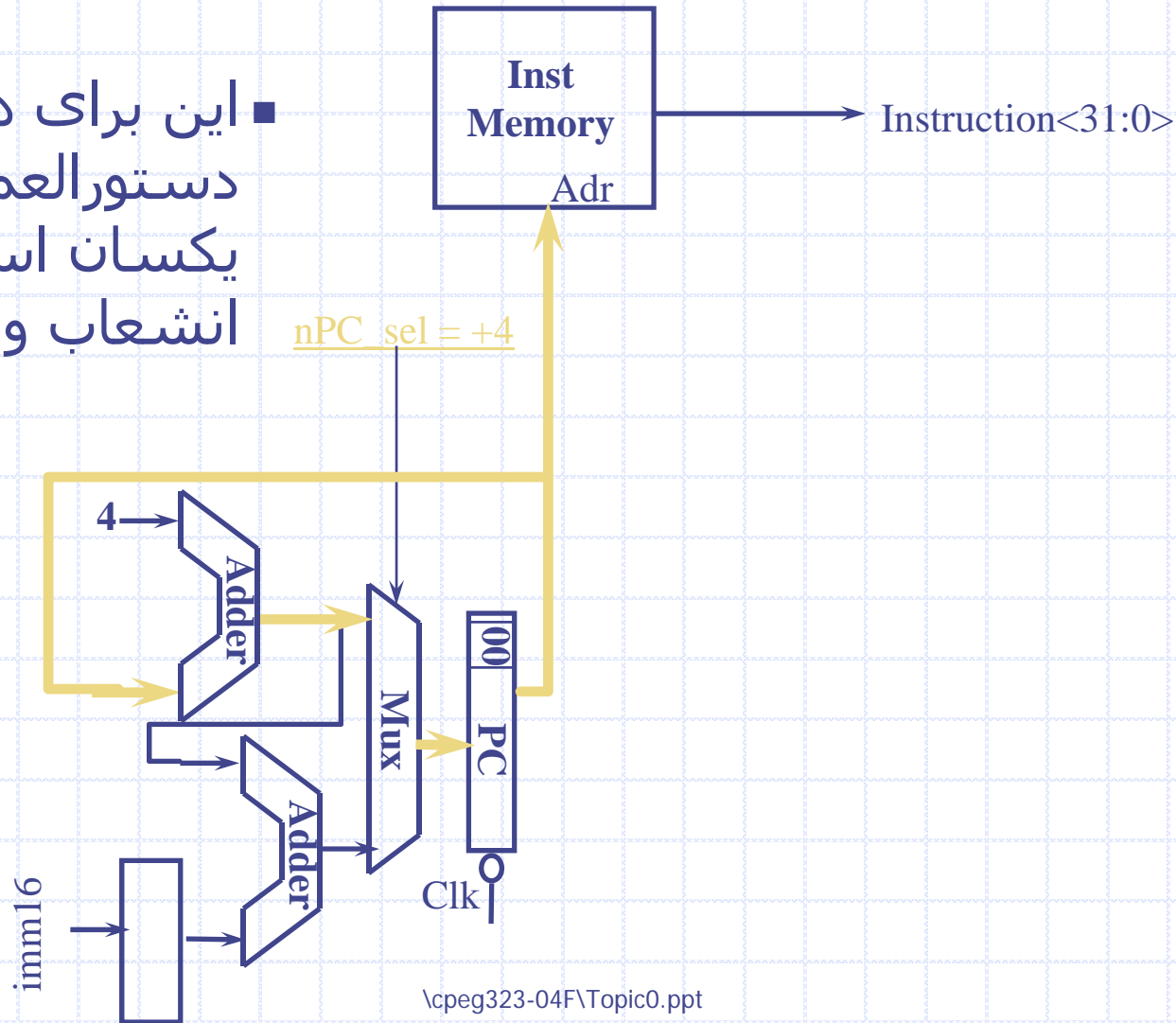
Instruction<31:0>



واحد واکنشی دستورالعمل در انتهای جمع

$$PC \leftarrow PC + 4$$

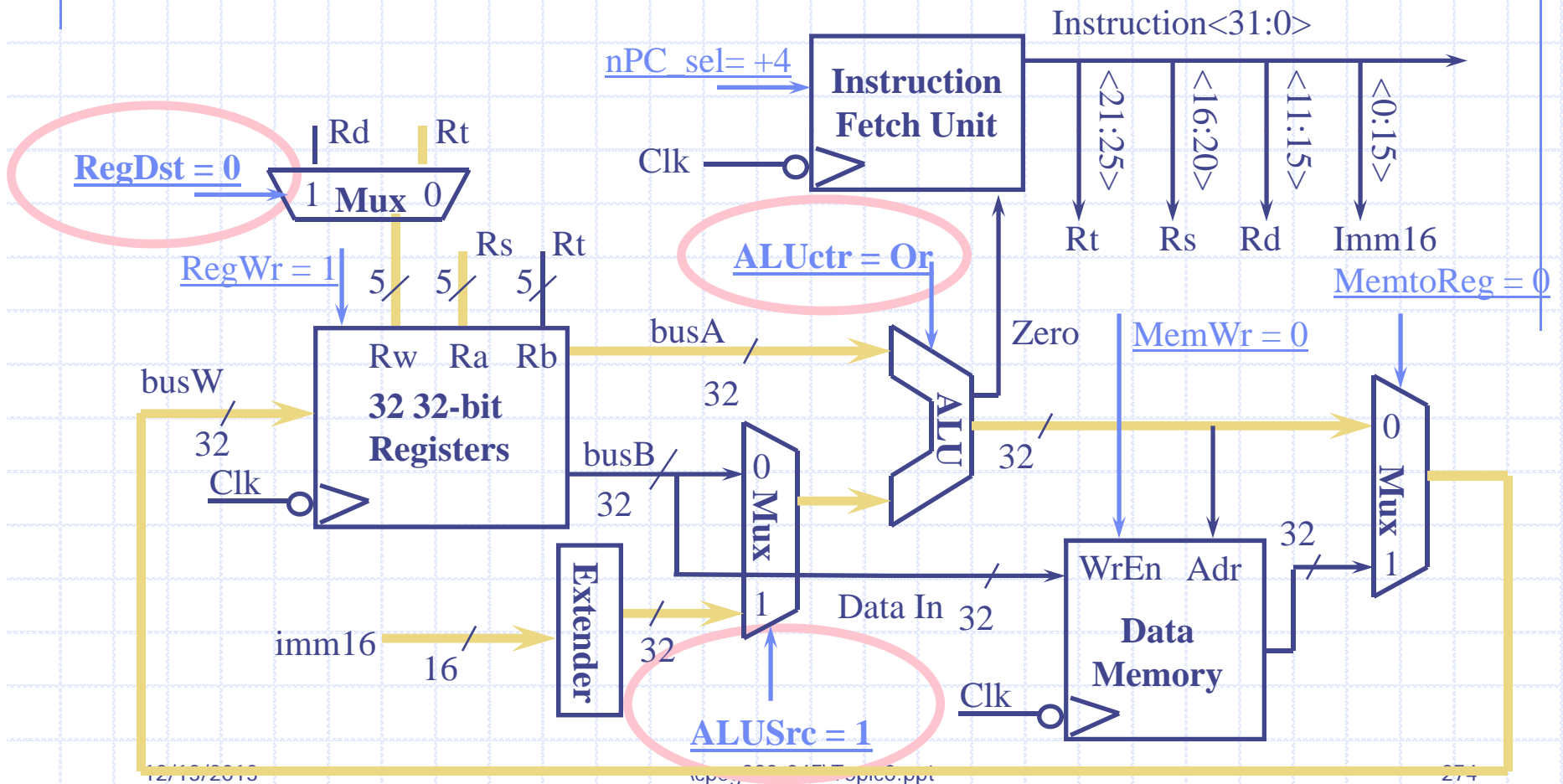
- این برای همه دستورالعمل‌ها یکسان است به جز انشعاب و پرش



OR مسیر داده تک چرخه ای در طول عمل بلافاصله



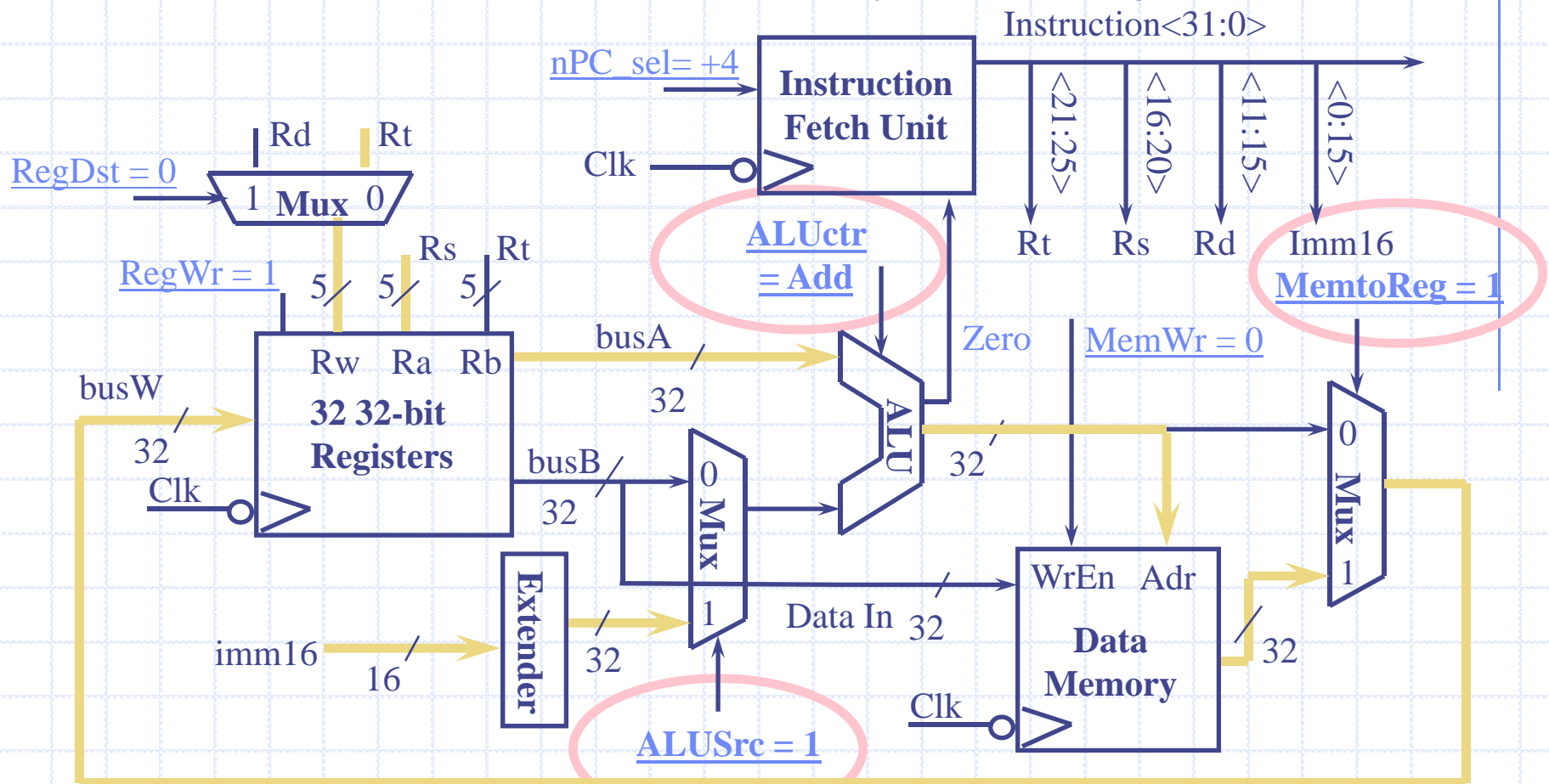
$R[rt] \leftarrow R[rs] \text{ or } [Imm16]$ علامت



مسیر داده تک چرخه ای در طول بار کردن



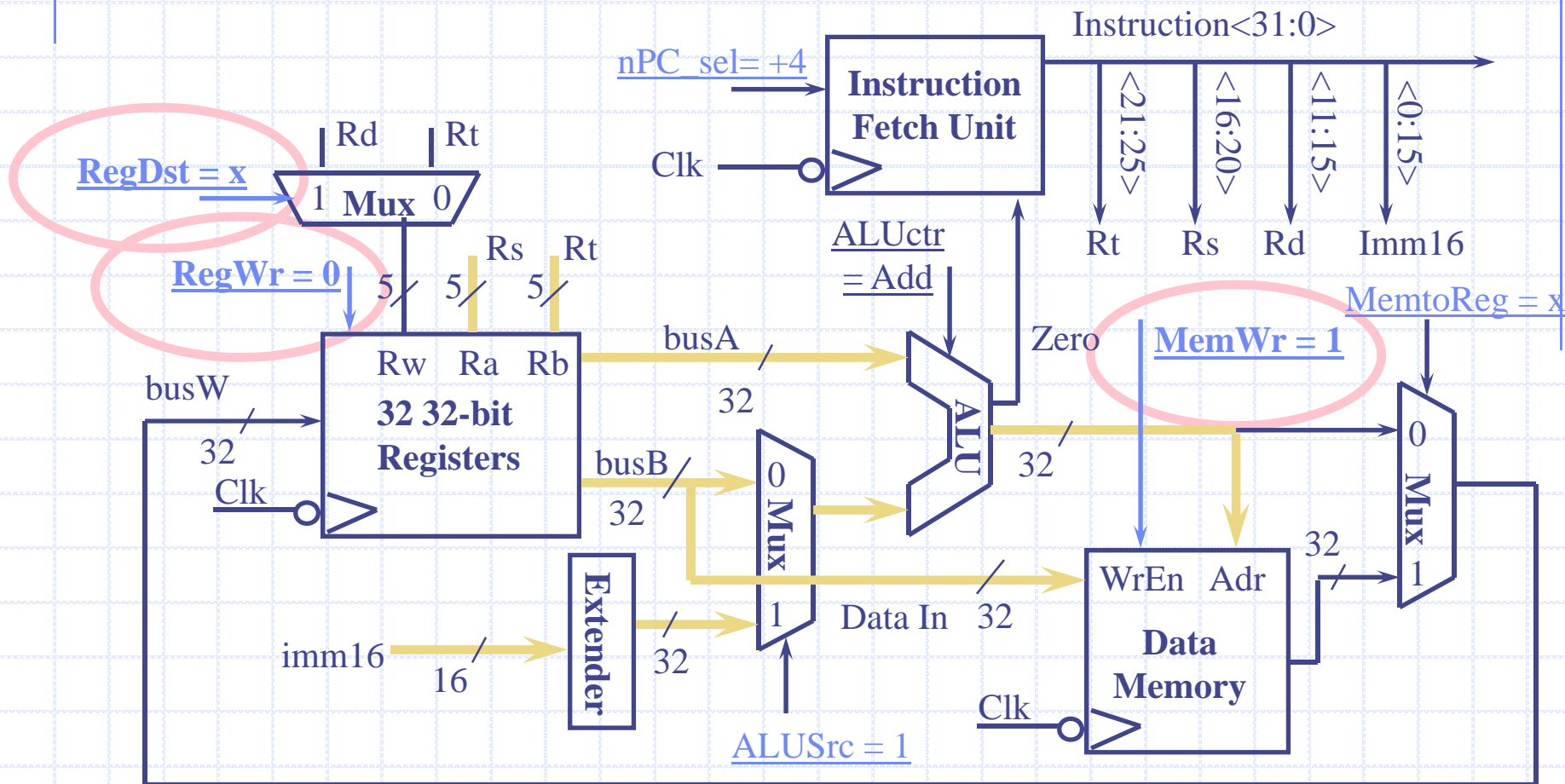
$$R[rt] \leftarrow \text{Data Memory} \{R[rs] + \text{SignExt}[\text{imm16}]\}$$



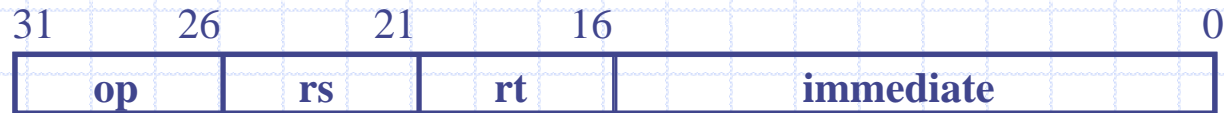
مسیر داده تک چرخه ای در طول عمل ذخیره سازی



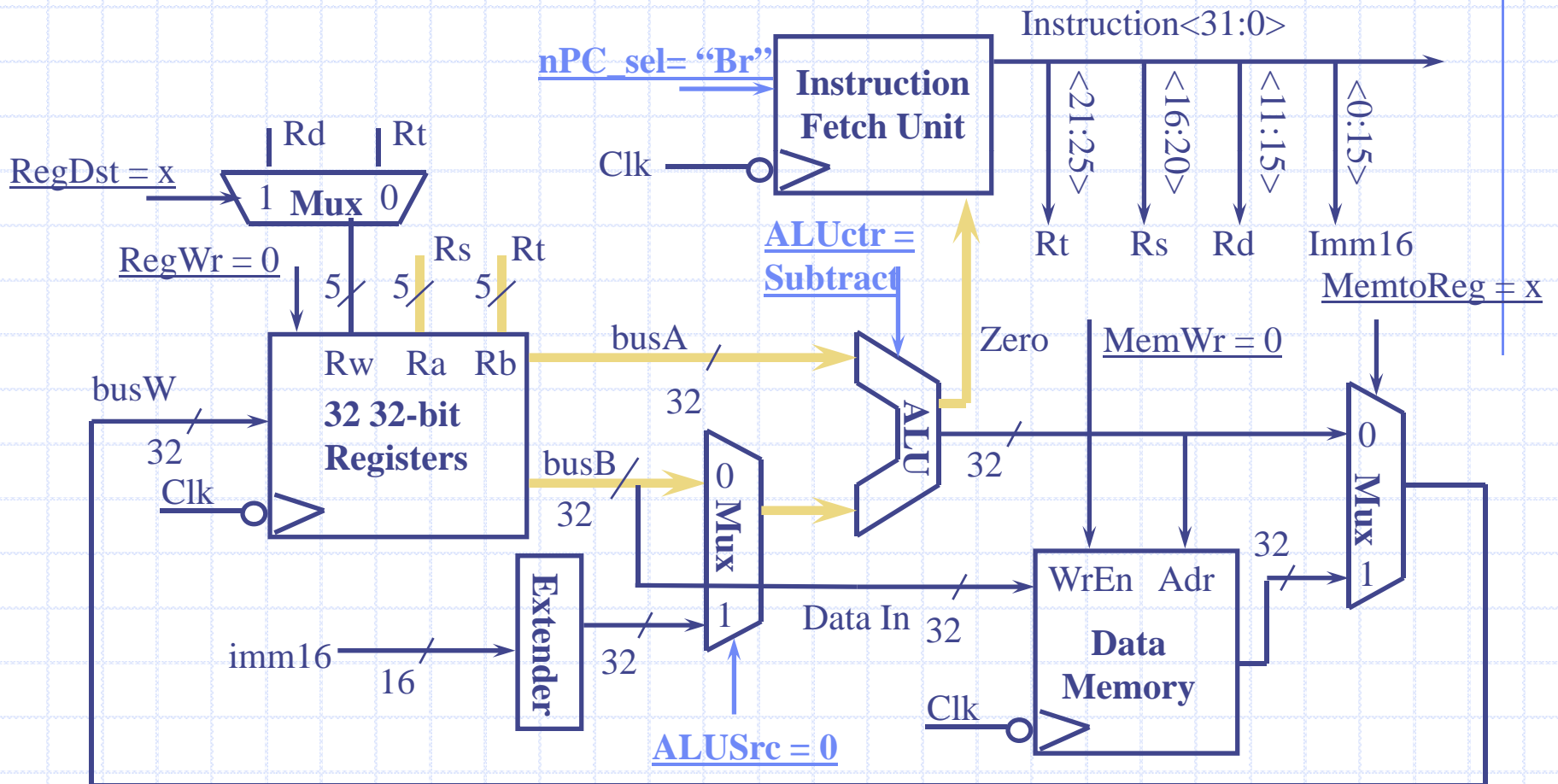
حافظه داده $\{R[rs] + \text{SignExt}[imm16]\} \leftarrow R[rt]$



مسیر داده تک چرخه ای در طول عمل انشعاب



if (R[rs] - R[rt] == 0) then Zero <- 1 ; else Zero <- 0

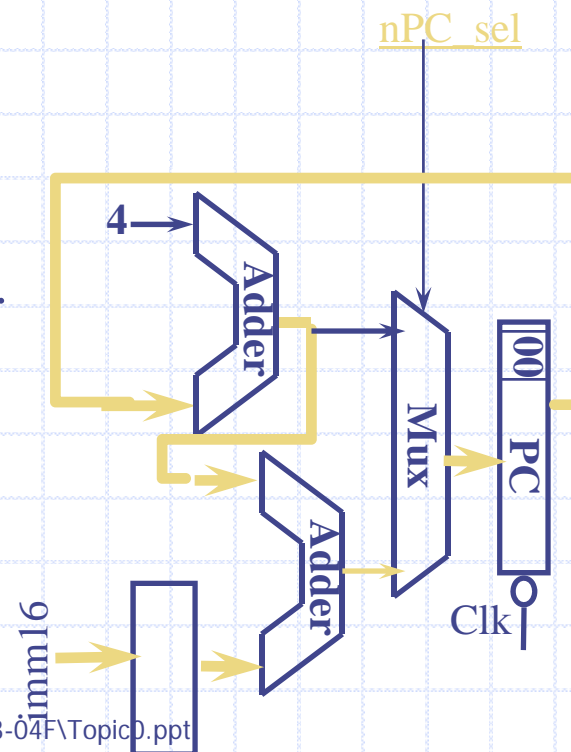


واحد واکنشی دستورالعمل در انتهای انشعاب

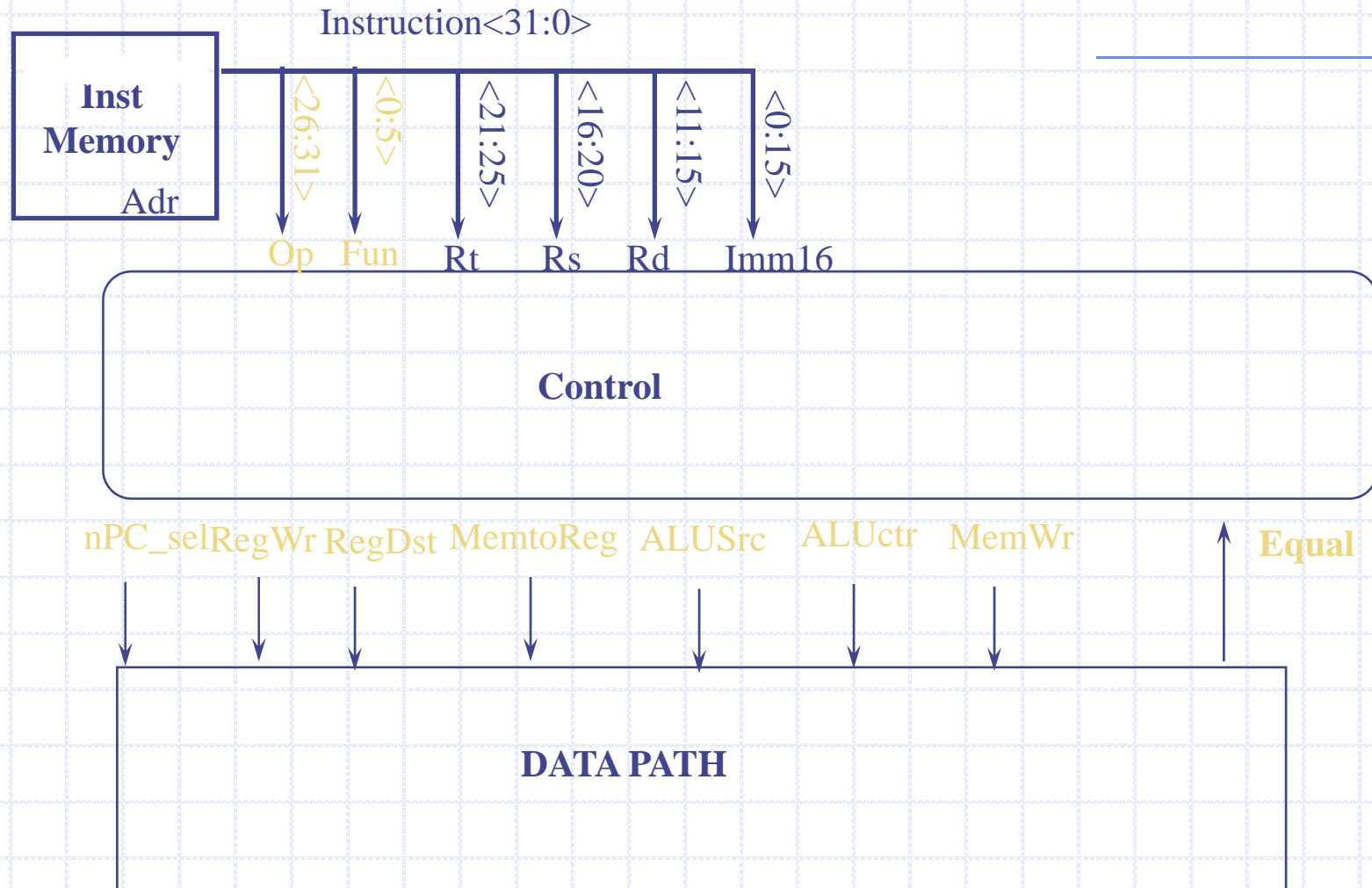


See book for what the datapath and control looks like for jump instructions.

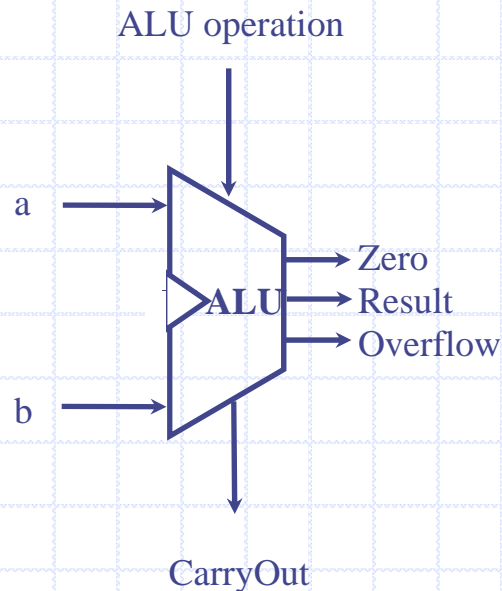
Compared to book our processor also supports the ORI instructions.



مرحله 4: مسیر داده داده شده RTL->CONTROL:



کنترل ALU



این نمادی است که به طور معمول برای نمایش یک ALU مورد استفاده قرار می گیرد. این نماد همچنین برای نمایش یک جمع کننده نیز استفاده می شود، بنابراین بوسیله یکی از دو نماد ALU یا Adder برچسب گذاری می شود. خطوط کنترلی دارای برچسب ALUOperation می باشد. مقدار آنها و عمل ALU در شکل بعد معرفی می شوند.

کنترل ALU

ALU Control lines	Function
000	And
001	Or
010	Add
110	Subtract
111	Set-on-less-than

مقادیر سه خط کنترلی ALU به همراه عمل متناظر ALU

Instruction opcode	ALUOp	Instruction operation	Function Code	Desired ALU action	ALU control input
LW	00	Load word	XXXXXX	add	010
SW	00	store word	XXXXXX	add	010
Branch equal	01	branch equal	XXXXXX	subtract	110
R-type	10	add	100000	add	010
R-type	10	subtract	100010	subtract	110
R-type	10	AND	100100	and	000
R-type	10	OR	100101	or	001
R-type	10	set-on-less-than	101010	set-on-less-than	111

این جدول چگونگی تنظیم بیت‌های کنترلی ALU و وابسته به بیت‌های کنترلی ALUOp و کدهای تابعی متفاوت برای دستورات نوع ثباتی را نشان می‌دهد.

Opcode، که در اولین ستون لیست شده است، وضعیت بیت‌های ALUOp را تعیین می‌کند. تمام کد گذارها به صورت دودویی نمایش داده می‌شوند. توجه کنید زمانی که کد ALUOp برابر 00 یا 01 باشد، فیلد‌های خروجی به فیلد کدهای تابعی وابسته نیست. در این مورد، گفته می‌شود که مقدار کدهای تابعی برای ما "بدون توجه" (don't care) می‌باشد و فیلد کد تابعی به صورت xxxxxx نمایش داده می‌شود. زمانی که مقدار ALUOp برابر 10 باشد، کد تابعی برای تنظیم ورودی کنترل ALU به کار می‌رود.

دستورالعمل‌های که از ALU استفاده می‌کنند

محاسبه آدرس - جمع

تفریق

جمع / تفریق

and/or

set-on-less-than

Load/store:

Branch eq:

R-type:

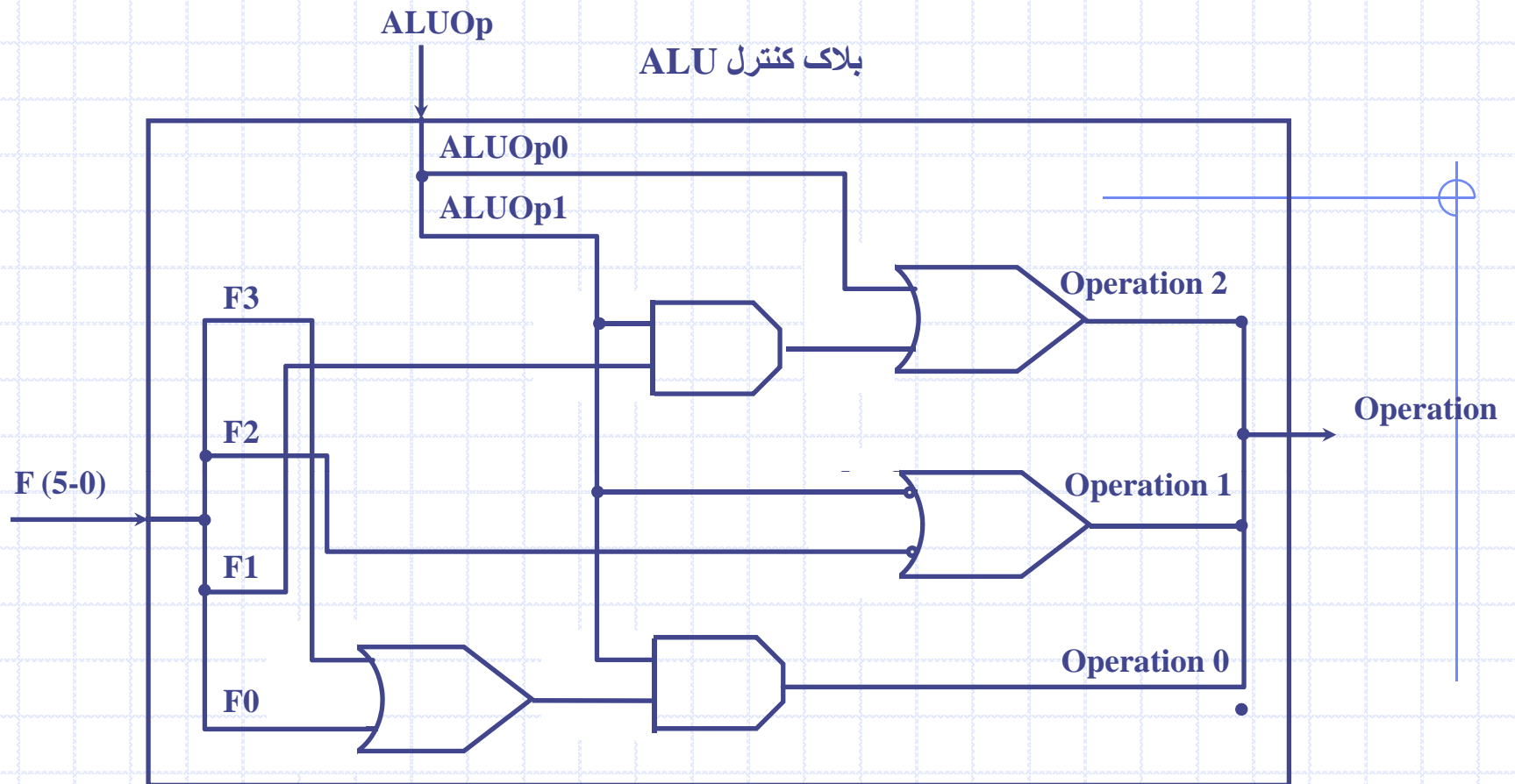
نیاز به کدهای تابعی

طراحی کنترلر ALU

ALUOp		Function code						ALU control input
ALUOp1	ALUOp0	F5	F4	F3	F2	F1	F0	
0	0	X	X	X	X	X	X	010
X	1	X	X	X	X	X	X	110
1	X	X	X	0	0	0	0	010
1	X	X	X	0	0	1	0	110
1	X	X	X	0	1	0	0	000
1	X	X	X	0	1	0	1	001
1	X	X	X	1	0	1	0	111

جدول درستی برای سه بیت کنترلی ALU به عنوان تابعی از ALUOp و فیلد کد تابعی

بسیاری از ورودیهای اضافه شده "بدون توجه" (don't care) هستند. به عنوان مثال، ALUOp از کد گذاری 11 استفاده ای نمی کند، بنابراین جدول درستی می تواند ترجیحاً شامل ورودیهای 1x و x1 به جای 10 و 01 باشد.



بلاک کنترلی ALU سه بیت کنترلی ALU را مبنی بر کد تابعی و بیت‌های ALUop تولید می‌کند.

چهره قالب دستورالعملهای MIPS

Field	0	rs	rt	rd	shamt	funct
Bit positions	31-26	25-21	20-16	15-11	10-6	5-0

a. R-type instruction

Field	35 or 43	rs	rt	address
Bit positions	31-26	25-21	20-16	10-6

b. Load or store instruction

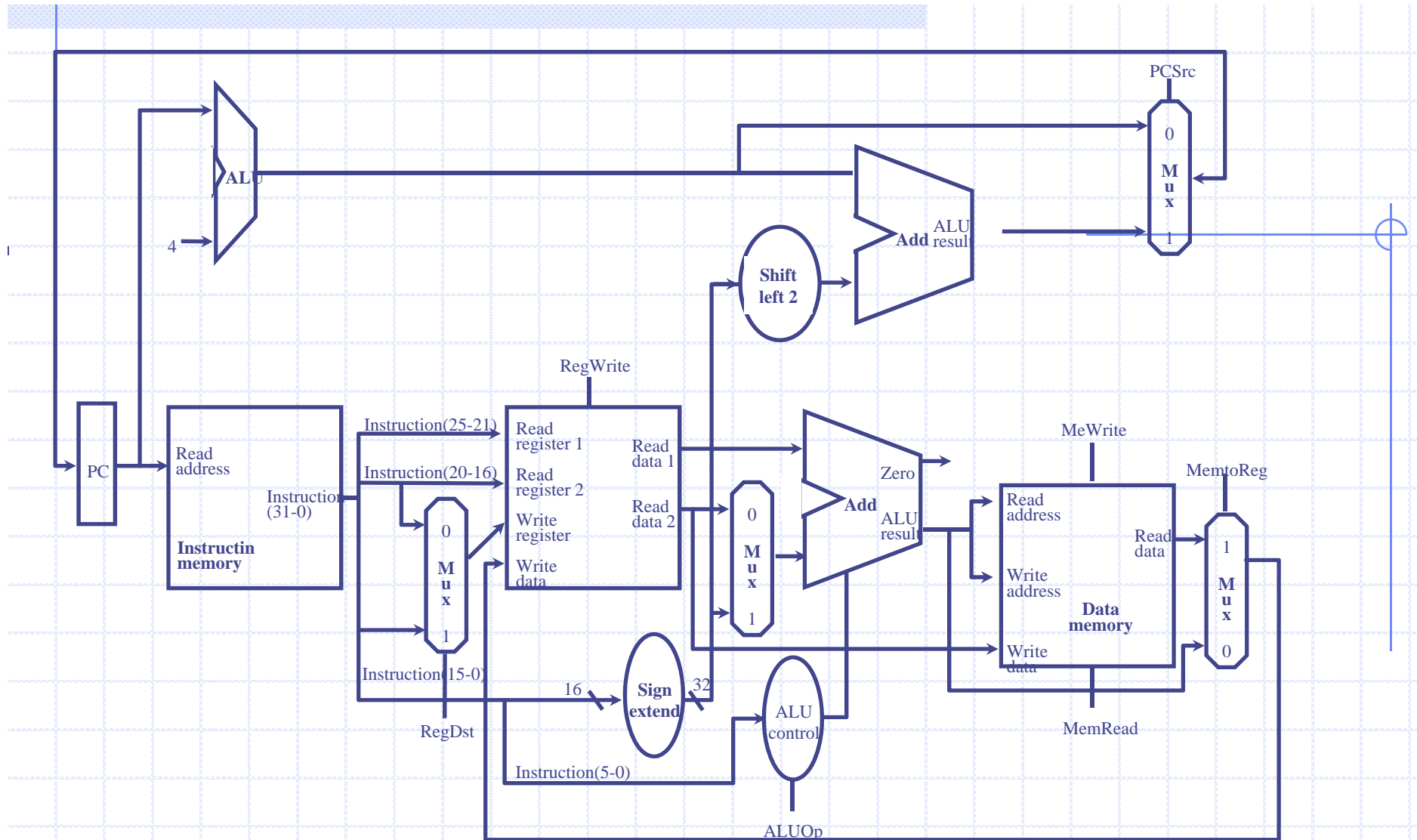
Field	4	rs	rt	address
Bit positions	31-26	25-21	20-16	15-0

c. branch instruction

چهره قالب دستورالعملهای MIPS

cont'd

- فیلد op ، همچنین معروف به opcode ، همیشه در بیت‌های 26 تا 31 قرار می‌گیرد. ما به این فیلد به صورت Op[5-0] اشاره می‌کنیم.
- دو رجیستر خواندن همیشه با فیلدهای rs و rt در موقعیت 21 تا 25 و 16 تا 20 مشخص می‌شوند. این مطلب برای دستورات ثباتی، انشعاب در صورت تساوی و ذخیره‌سازی صدق می‌کند.
- ثبات پایه برای دستورات load و store همیشه در موقعیت‌های بی‌تی 21 تا 25 (rs) قرار دارد.
- 16 بیت افسست برای دستورات انشعاب در صورت تساوی، load و store همیشه در موقعیت‌های 0 تا 15 قرار دارد.
- ثبات مقصد در یکی از این دو مکان است. برای دستورات load در بیت‌های موقعیت 16 تا 20 (rt) و در حالی که برای یک دستورالعمل ثباتی در موقعیت‌های بی‌تی 11 تا 15 (rt) است. بنابراین ما به افزودن یک مالتی پلکسر برای انتخاب اینکه کدام فیلد دستورالعمل به شماره رجیستر نوشتن اشاره می‌کند.



مسیر داده به همراه تمامی مالتی پلکسرها و تمامی خطوط کنترلی شناخته شده

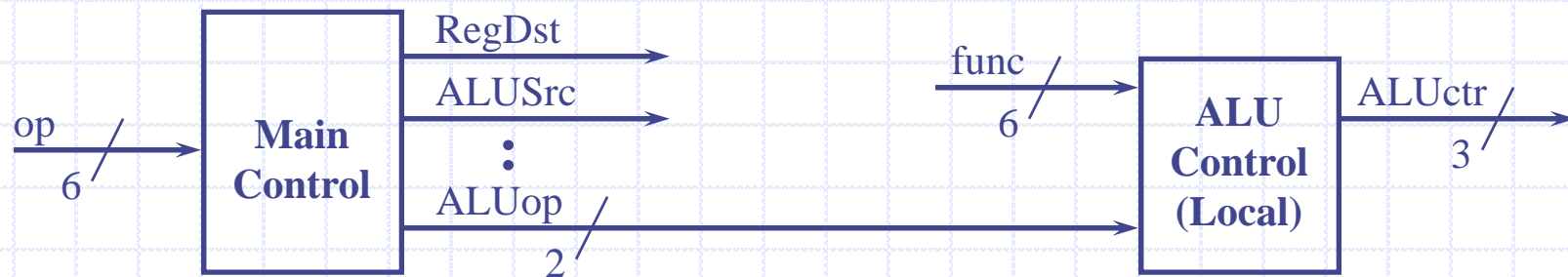
طراحی واحد کنترل برای اجرای دستورالعمل تک چرخه ای

Signal name	Effect when deasserted	Effect when asserted
MemRead	None	Data memory contents at the read address are put on read data output
MemWrite	None	Data memory contents at address given by write address is replaced by value on write data input.
ALUSrc	The second ALU operand comes from the second Register file output.	The second ALU operand is the sign-extended lower 16-bits of the instruction.
RegDst	The register destination number for the Write register comes from the rt field	The register destination number for the Write register comes from the rd field.
RegWrite	None	The register on the Write register input is written into with the value on the write data input.
PCSrc	The PC is replaced by the output of the adder That computes the value of PC + 4.	The PC is replaced by the output of the adder that computes the branch target.
MemtoReg	The value fed to the register write data input comes from the ALU	The value fed to the register write data input comes from the data memory.

عمل هر کدام از هفت سیگنال کنترلی. زمانیکه بیت کنترل (1) به سمت مالتی پلکسر هدایت شود، مالتی پلکسر ورودی متناظر با پایه 1 را انتخاب می کند. در غیر اینصورت در صورتی که بیت کنترل آزاد نگردد (0 باشد) مالتی پلکسر ورودی مربوط به پایه 0 را انتخاب می کند. به خاطر بسپارید که همه عناصر حالت دارای یک ورودی ساعت به عنوان یک ورودی ضمنی می باشد و ساعت به عنوان کنترل کننده نوشتن ها به کار می رود.

طراحی واحد کنترل





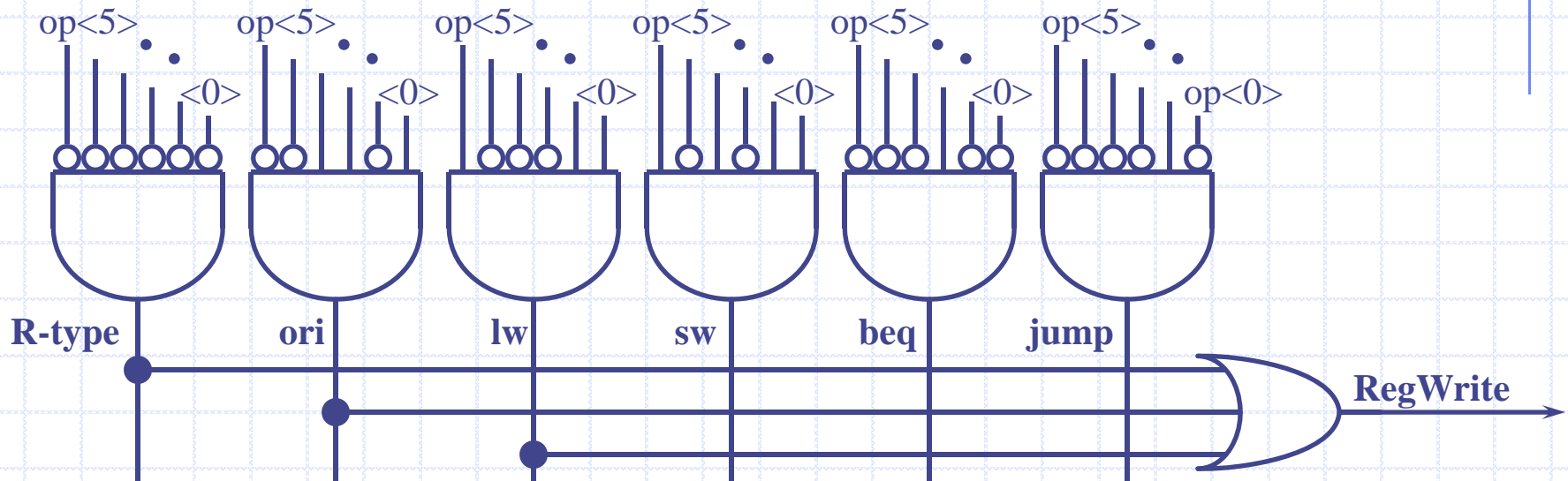
R-format		lw	sw	beq	
Inputs	Op5	0	1	1	0
	Op4	0	0	0	0
	Op3	0	0	1	0
	Op2	0	0	0	1
	Op1	0	1	1	0
	Op0	0	1	1	0
Outputs	RegDst	1	0	X	X
	ALUSrc	0	1	1	0
	MemtoReg	0	1	X	X
	RegWrite	1	1	0	0
	MemRead	0	1	0	0
	MemWrite	0	0	1	0
	Branch	0	0	0	1
	ALUOp1	1	0	0	0
ALUPp0	0	0	0	1	

عملکرد کنترل برای پیاده سازی تک چرخه ای ساده کاملاً به وسیله جدول درستی تعیین می شود. نیمه بالای جدول ترکیب سیگنالهای ورودی که مرتبط با 4 کد عملیاتی (opcode) می باشد و تعیین کننده وضع خروجی کنترل هستند را به دست می دهد. (به خاطر داشته باشید که Op(0-5) مرتبط با بیت‌های 26 تا 31 دستورالعمل است، که فیلد کد عمل (opcode) می باشند.) بخش پایین جدول خروجیها را به دست می دهد.

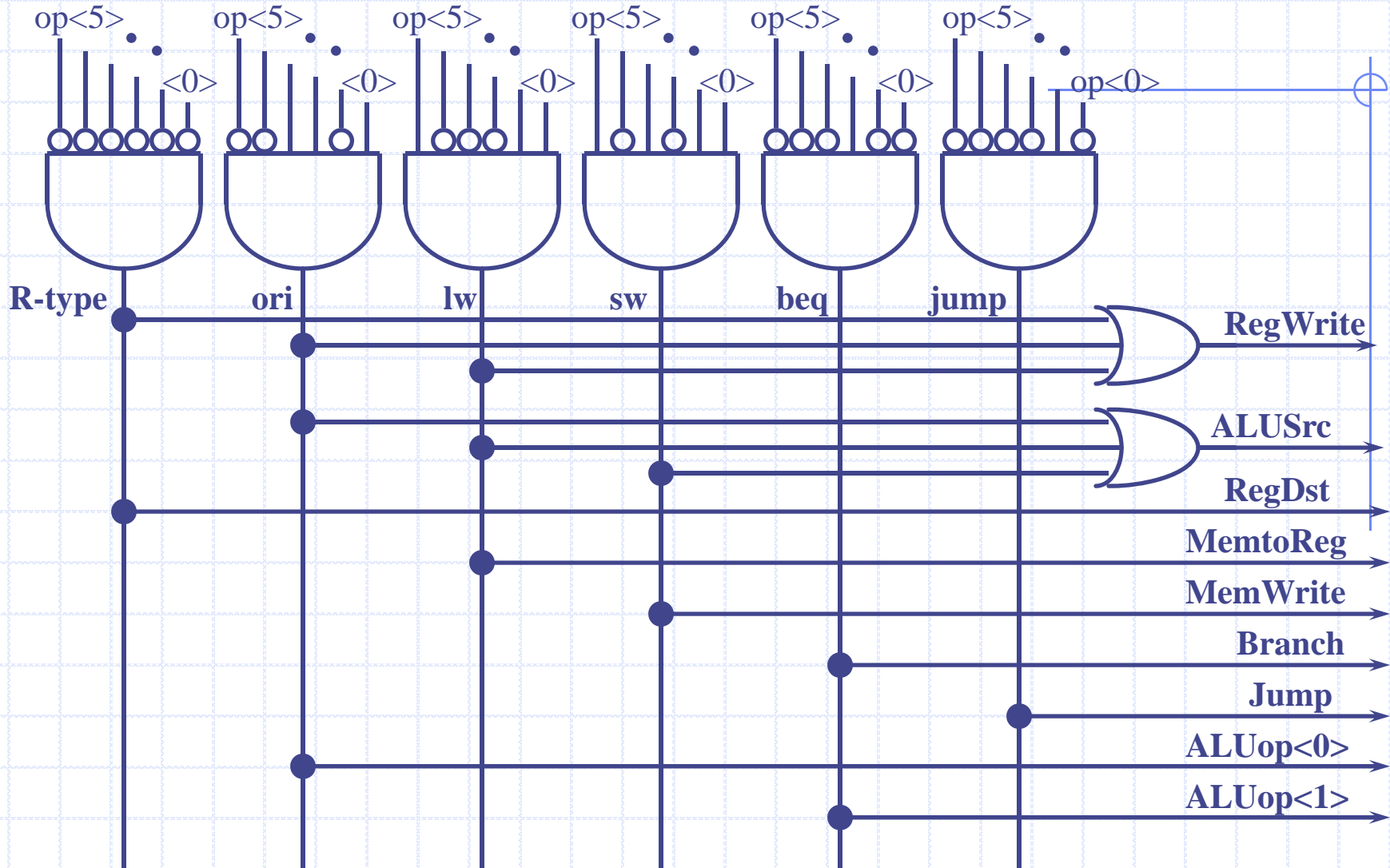
"جدول درستی" برای نوشتن ثباتی

	op	00 0000	00 1101	10 0011	10 1011	00 0100	00 0010
		R-type	ori	lw	sw	beq	jump
RegWrite		1	1	1	0	0	0

$$\begin{aligned}
 \text{RegWrite} &= \text{R-type} + \text{ori} + \text{lw} \\
 &= !\text{op}<5> \& !\text{op}<4> \& !\text{op}<3> \& !\text{op}<2> \& !\text{op}<1> \& !\text{op}<0> \text{ (R-type)} \\
 &\quad + !\text{op}<5> \& !\text{op}<4> \& \text{op}<3> \& \text{op}<2> \& !\text{op}<1> \& \text{op}<0> \text{ (ori)} \\
 &\quad + \text{op}<5> \& !\text{op}<4> \& !\text{op}<3> \& !\text{op}<2> \& \text{op}<1> \& \text{op}<0> \text{ (lw)}
 \end{aligned}$$



PLA روشهای کنترل اصلی



پیاده سازی یک دستورالعمل چند چرخه ای

توضیحاتی از پیاده سازی تک چرخه ای

- دستورالعملهای ماشین ممکن است دارای طولهای متفاوت مسیر بحرانی باشند
 - دستورالعمل Load
 - دستورات ممیز شناور
 - روشهای آدرس دهی متفاوت
- زمان چرخه ساعت بوسیله بدترین مسیر بحرانی معین می شود.
- تکثیر FU ممکن است هزینه دار باشد.

چه چیزی باعث مشکل می شود زمانی که $CPI=1$ پردازشگر است

Arithmetic & Logical



Load



Store



Branch



♦ طولانی بودن زمان چرخه

♦ همه دستورالعمل ها به اندازه کندترین دستور زمان می برند

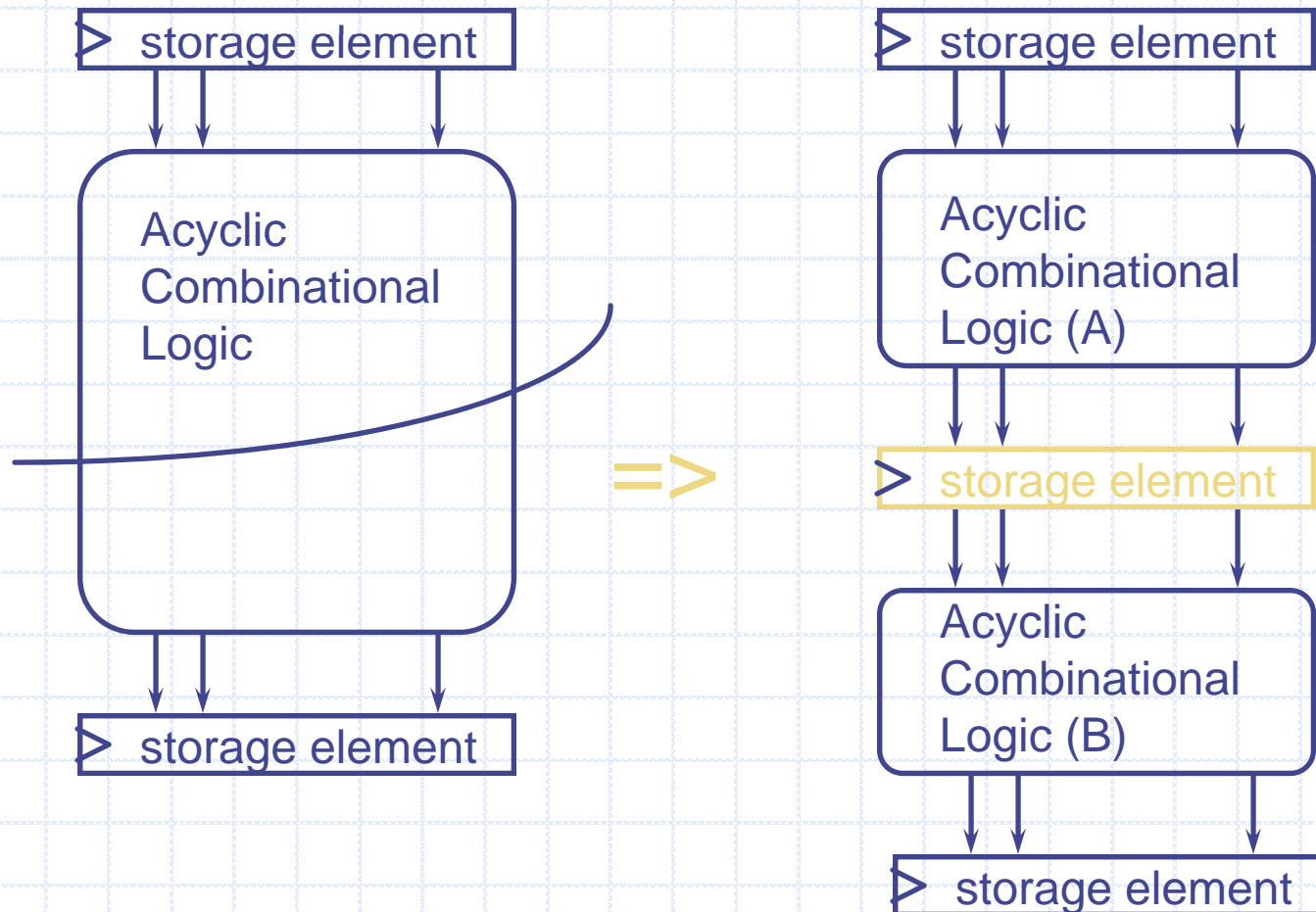
♦ Real memory is not so nice as our idealized memory

cannot always get the job done in one (short) cycle ■

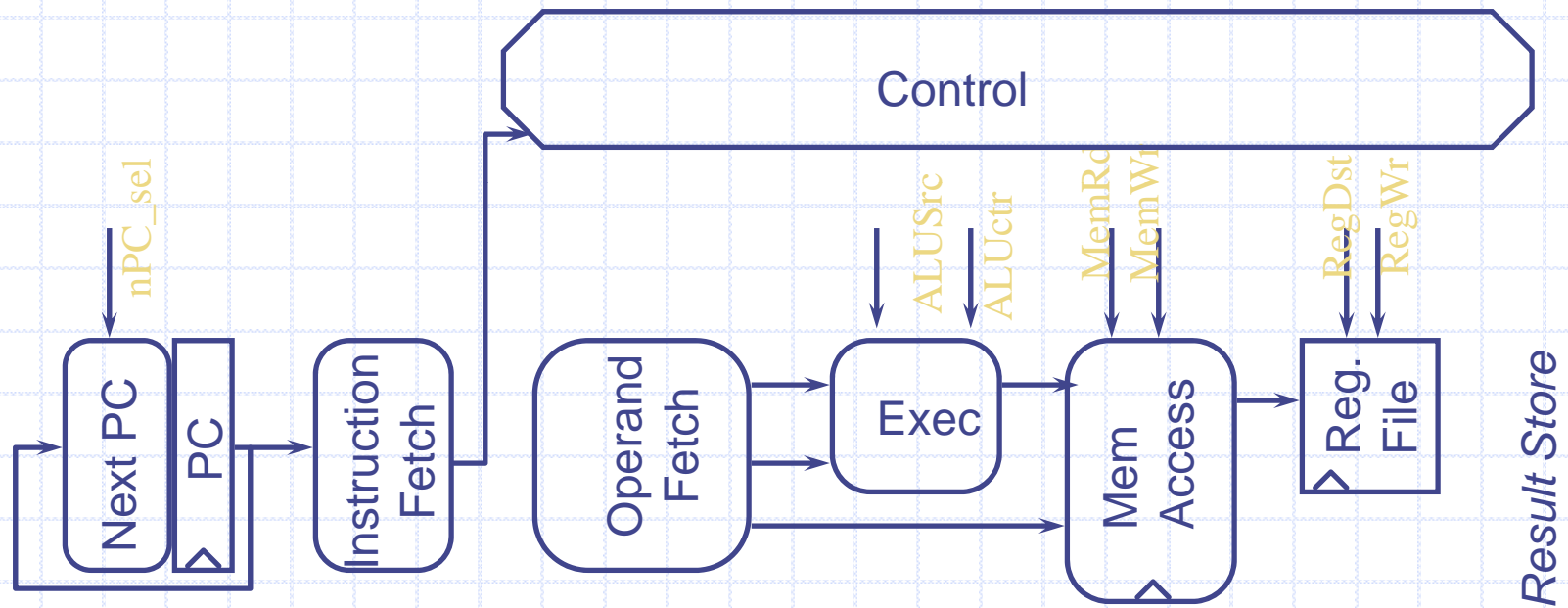
کاهش زمان چرخه

Cut combinational dependency graph and insert register / latch

درست همان کار رادر 2 چرخه سریع نسبت به آن که 1 چرخه کند بود انجام میدهد

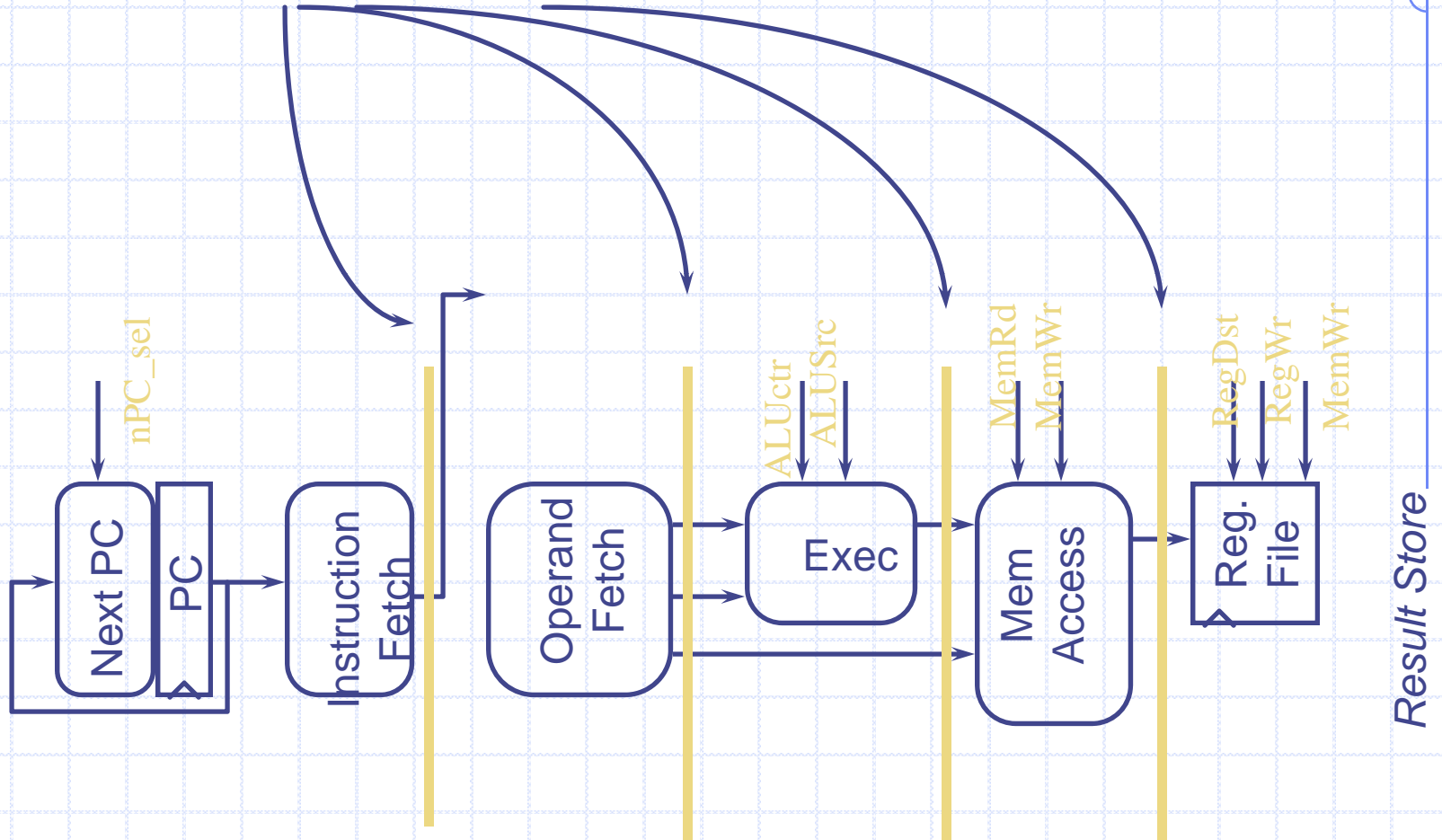


محدودیت های بنیادی زمان چرخه



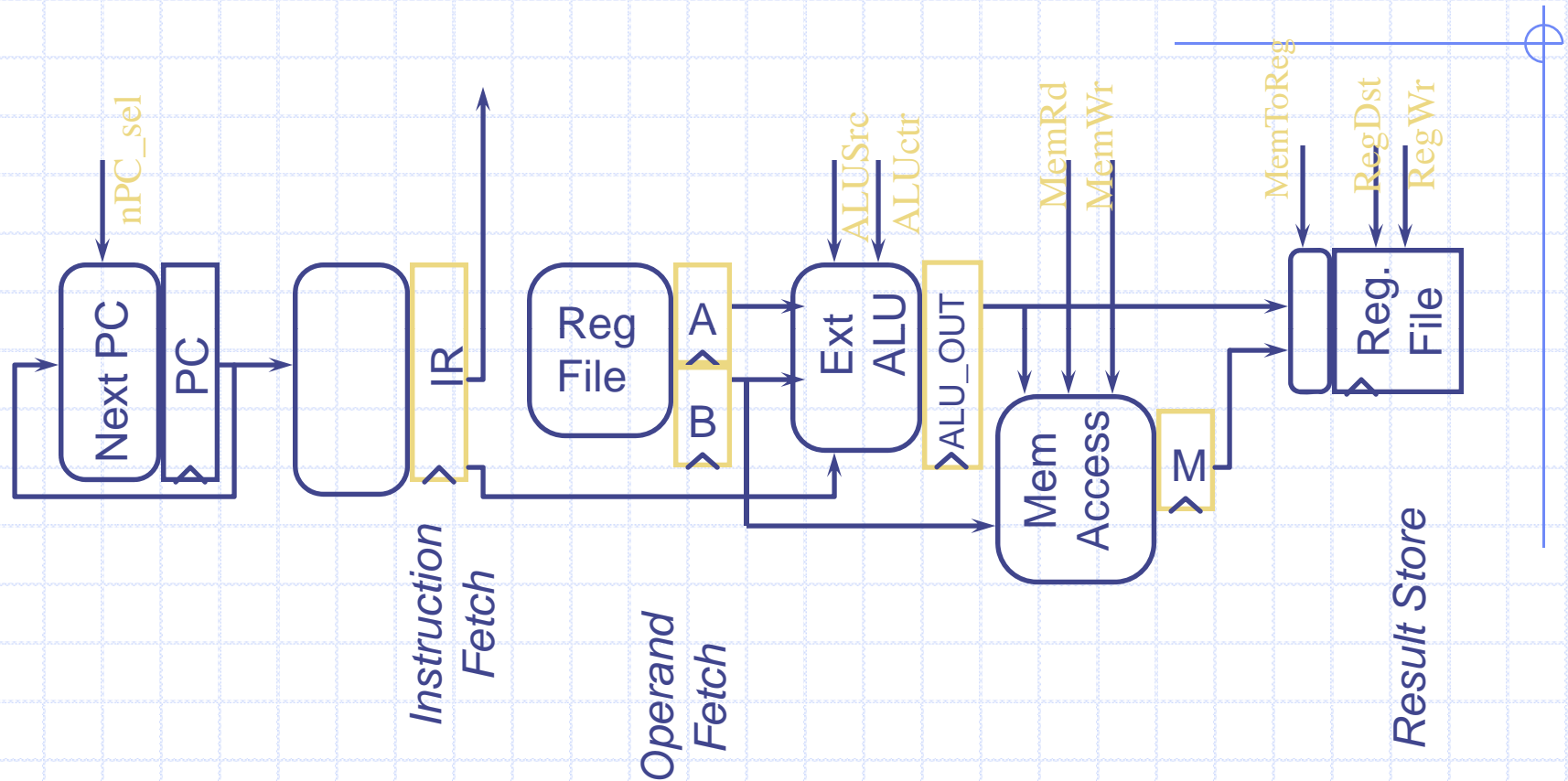
تقسیم بندی $CPI=1$

♦ اضافه کردن ثباتها بین مسیرهای کوچکترین مرحله ها



امکان انجام دستورالعمل در چند چرخه

مثال مسیر داده چند چرخه ای



◆ ثبتهای اضافی اضافه میشوند برای ذخیره کردن مقادیر بین مراحل

خصوصیات

- یک واحد حافظه یکسان برای دستورات و داده ها استفاده می شود.
- یک ثبات دستورالعمل به کار می رود. (IR)
- یک ALU واحد

R-type دستور العمل های (add, sub, . . .)

inst Logical Register Transfers

ADDU $R[rd] \leftarrow R[rs] + R[rt]; PC \leftarrow PC + 4$

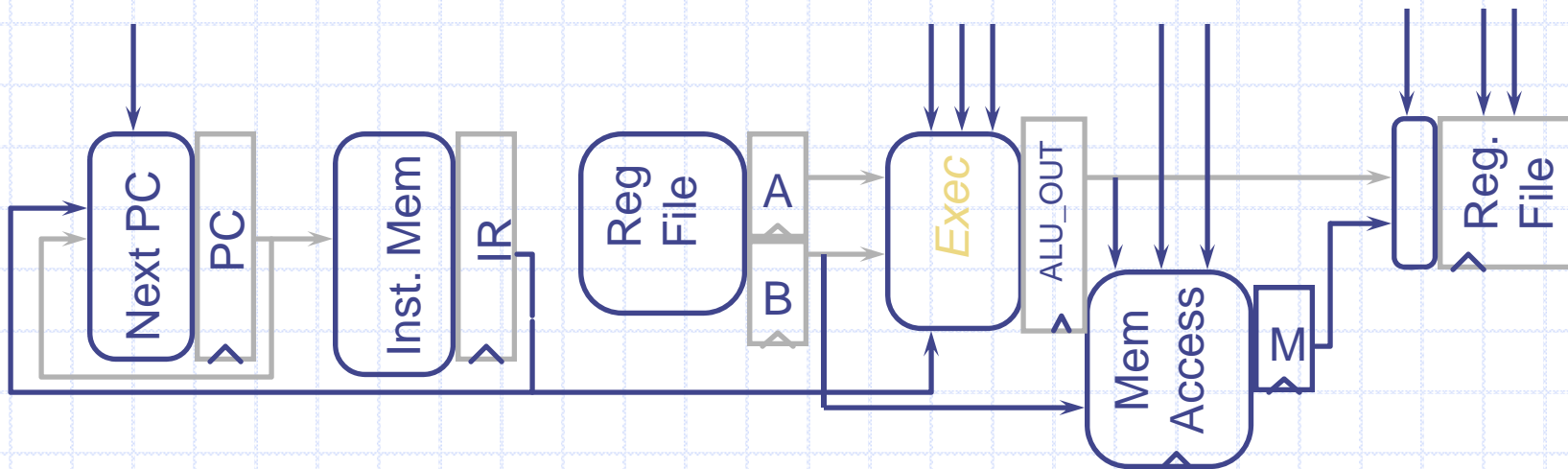
inst Physical Register Transfers

$IR \leftarrow MEM[pc]$

ADDU $A \leftarrow R[rs]; B \leftarrow R[rt]$

$S \leftarrow A + B$

$R[rd] \leftarrow S; PC \leftarrow PC + 4$



دستور العمل های فوری منطقی

inst Logical Register Transfers

ADDU $R[rt] \leftarrow R[rs] \text{ OR } sx(Im16); PC \leftarrow PC + 4$

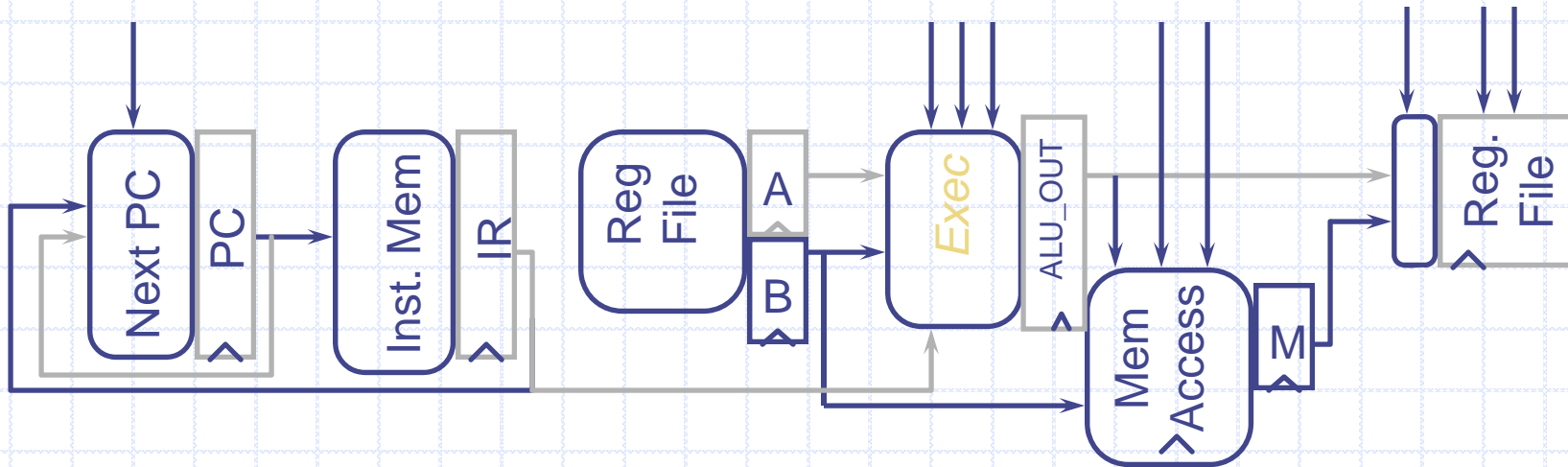
inst Physical Register Transfers

$IR \leftarrow MEM[pc]$

ADDU $A \leftarrow R[rs]; B \leftarrow R[rt]$

$S \leftarrow A \text{ or } SignExt(Im16)$

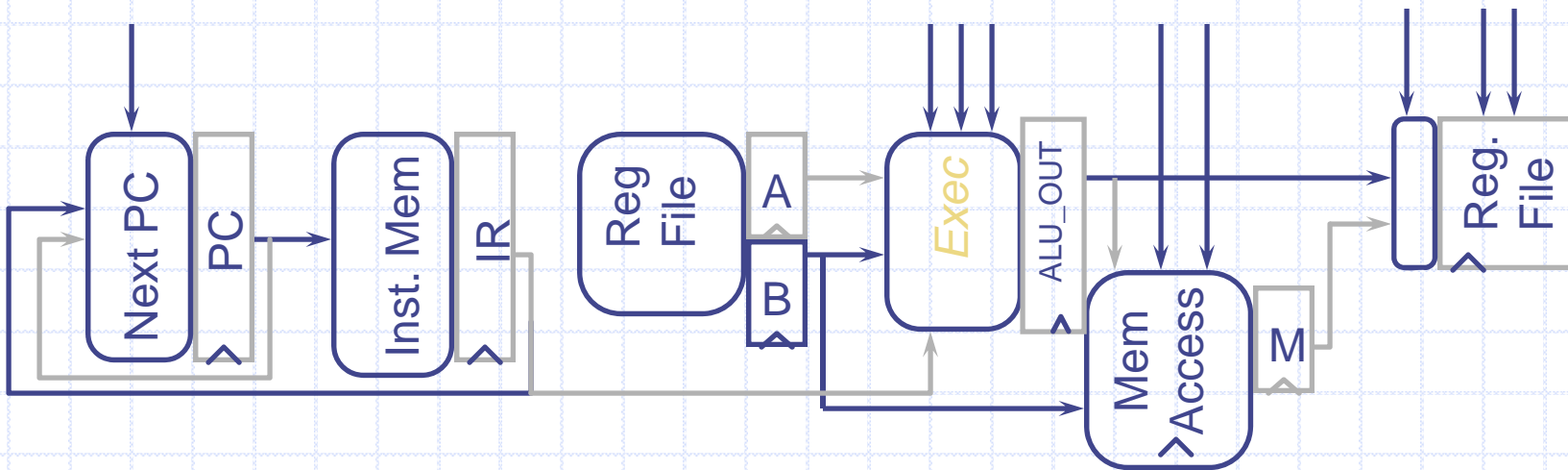
$R[rt] \leftarrow S; PC \leftarrow PC + 4$



دستور العمل بار کردن

inst انتقال دادن منطقی ثبات
LW $R[rt] \leftarrow MEM(R[rs] + sx(Im16));$
 $PC \leftarrow PC + 4$

inst انتقال دادن فیزیکی ثبات
 $IR \leftarrow MEM[pc]$
LW $A \leftarrow R[rs]; B \leftarrow R[rt]$
 $S \leftarrow A + SignEx(Im16)$
 $M \leftarrow MEM[S]$
 $R[rd] \leftarrow M; PC \leftarrow PC + 4$



دستور العمل ذخیره کردن

inst انتقال منطقی ثبات

SW $MEM(R[rs] + sx(Im16)) \leftarrow R[rt];$
 $PC \leftarrow PC + 4$

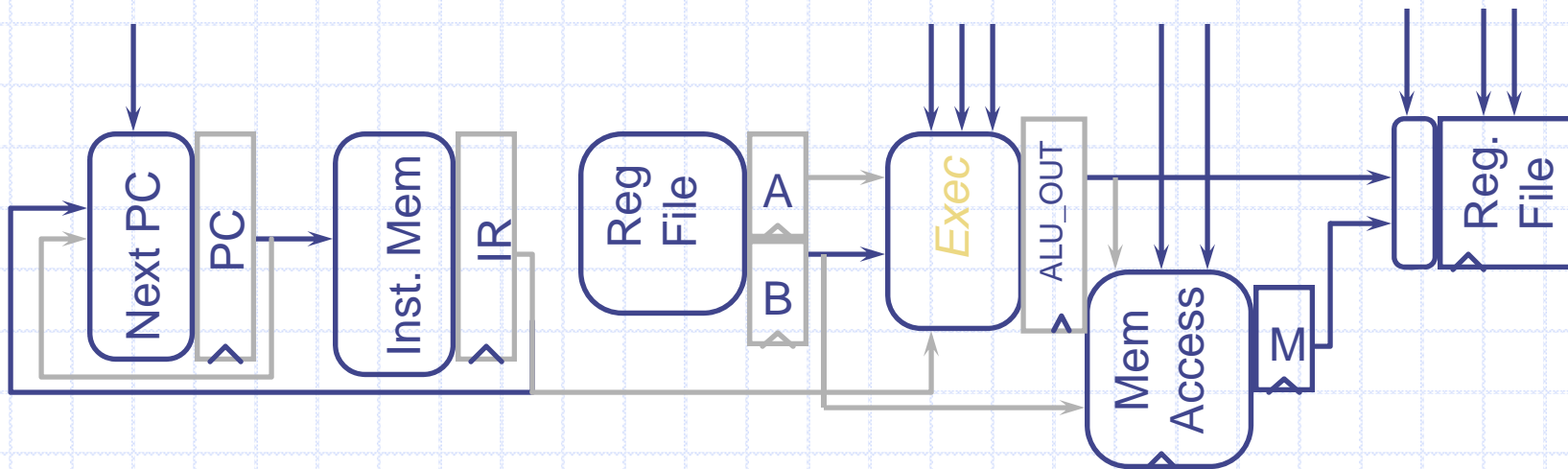
inst انتقال فیزیکی ثبات

IR $IR \leftarrow MEM[pc]$

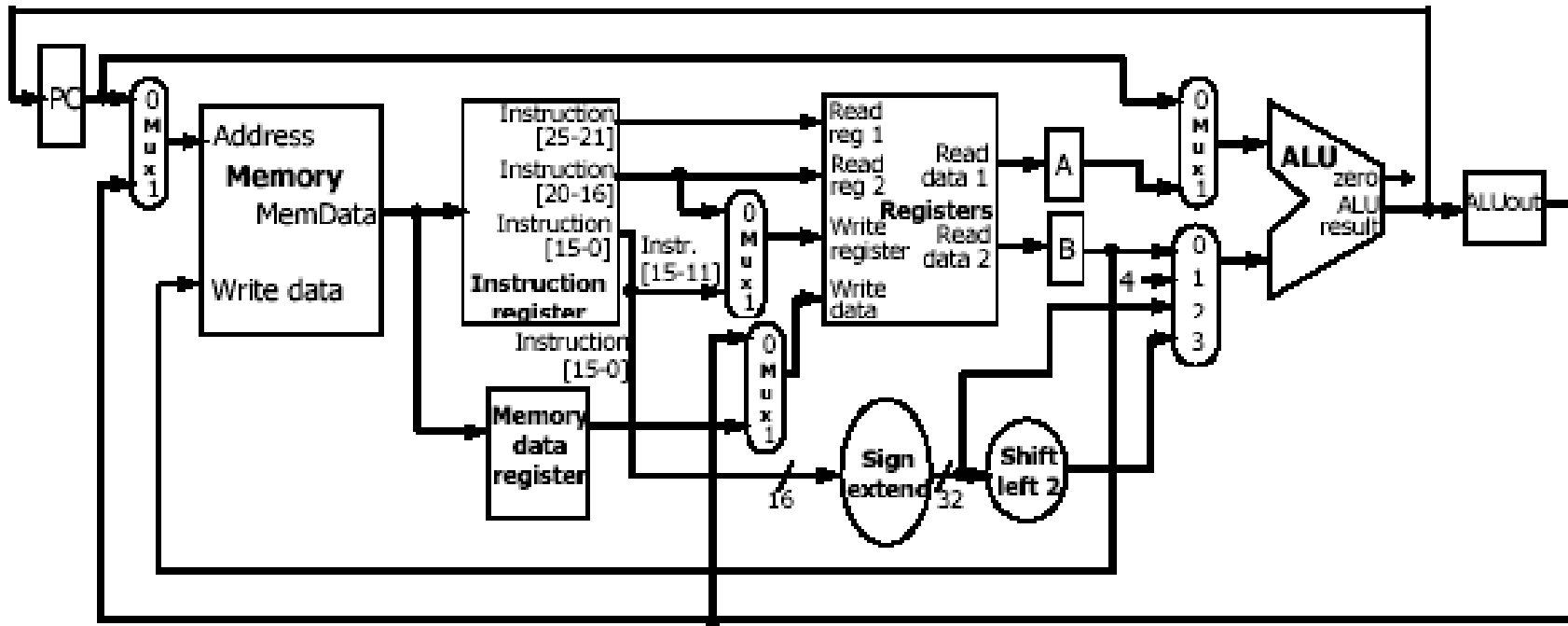
SW $A \leftarrow R[rs]; B \leftarrow R[rt]$

$S \leftarrow A + SignEx(Im16);$

$MEM[S] \leftarrow B$ $PC \leftarrow PC + 4$

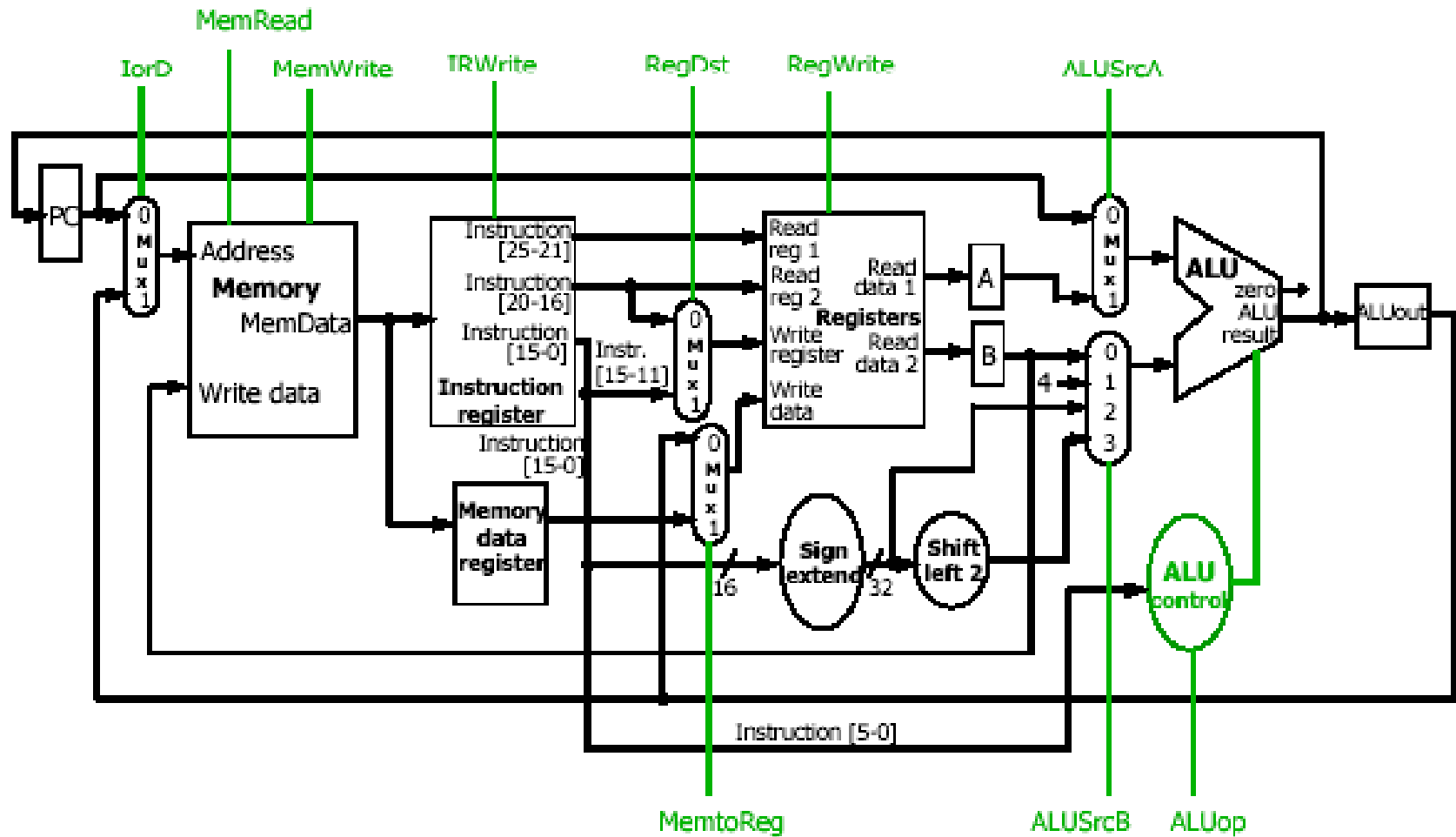


مسیر داده چند چرخه ای (Figure 5.26, p.327)



تفاوت های بین مسیر داده تکچرخه ای و چند چرخه ای
یک واحد حافظه تنها برای هر دو دستورالعمل و داده ها استفاده می شود

مسیر داده چند چرخه ای به همراه نمایش خطوط کنترلی



Signal name	Effect when deasserted	Effect when asserted
MemRead	None	Contents of memory at the read address are put on read data output.
MemWrite	None	Memory contents at the write address is replaced by value on write data input.
ALUSelA	The first ALU operand is the PC.	The first ALU operand comes from the register given by the rs field.
RegDst	The register destination number for the register write comes from the rt field.	The register destination number for the register write comes from the rd field.
RegWrite	None	The register given by Write register number is written with the value on the write data input.
MemtoReg	The value fed to the register write data input comes from the ALU.	The value fed to the register write data input comes from the Data memory.
lorD	The PC is used to supply the address to the memory unit.	The output of the ALU is used to supply the address to the memory unit.
IRWrite	None	The value from the memory unit is written into the instruction register (IR).

عمل سیگنالهای کنترلی 1 بیتی مشخص می باشد.

اعمال سیگنالهای کنترلی

Signal name	Value	Effect
ALUSelB	00	The second input to the ALU comes from the register given by the rt field.
	01	The second input to the ALU is the constant 4.
	10	The second input to the ALU is the sign-extended lower 16 bits of the IR.
	11	The second input to the ALU is the sign-extended and shifted lower 16 bits of the IR.
ALUOp	00	The ALU performs an add operation.
	01	The ALU performs an subtract operation.
	10	The function code field of the instruction determine the ALU operation.

عمل سیگنالهای کنترلی 2 بیتی مشخص می باشد.

اعمال سیگنالهای کنترلی

طرح ساعت زنی

- ما نیاز به یک ثبات موقتی داریم زمانیکه:

1. یک سیگنال در یک چرخه محاسبه و در یک چرخه دیگر استفاده می شود و
2. ورودیهای بلاکهای کاری (functional block) که این سیگنال را تولید می کنند می توانند قبل از اینکه این سیگنال در عناصر حالت ذخیره شود تغییر کنند.

- مثال:

- ما به IR نیاز داریم

- خروجی ALU

- ثباتهای A و B به ورودیهای ALU

مراحل پایه محاسبه دستور العمل

1. IF مرحله:	واکشی دستور العمل IR= Memory [PC] PC= PC + 4
2. ID step	کد برداری و واکشی عملوندها ثبات: [IR(25-21)] ثبات: [IR(20-16)] PC + (علامت گسترش یافته: [IR(15-0) << 2])
=A	
=B	
=Target	

یادداشت: ما می خواهیم عملیات را "خوشبینانه" انجام دهیم. those are common to all, or at least do no hurt any.

(cont'd)

مراحل پایه محاسبه دستورالعمل

محاسبه . محاسبه آدرس حافظه ، یا تکمیل انشعاب

EX .3

نمونه:

ارجاع به حافظه

ALU = A + sign-extend [IR(15-0)]; خروجی

دستورالعمل های ALU :

ALU خروجی; = A op B

انشعاب:

If (A ==B) PC= Target;

(cont'd)

مراحل پایه محاسبه دستور العمل

MEM .4 دسترسی به حافظه یا تکمیل دستور العمل R__TYPE

نمونه:

ارجاع به حافظه:

MDR = Memory [ALUoutput];
Memory[ALUoutput] = B; or

دستور العمل های ALU:

Reg[IR(5-11)] = ALUoutput;

WB .5 : بازنوشتن

Reg[IR(20-16)] =MDR

کنترل برای PC بعدی

PC بعدی تعیین می شود به وسیله:

- زمانی که PC برای واکنشی دستورالعمل ترتیبی یک واحد افزوده شده خروجی ALU یک منبع است.
- ثبات هدف ، منبع است زمانی که دستورالعمل داده شده یک انشعاب شرطی باشد ما همچنین به یک سیگنال برای نوشتن در ثبات نیاز خواهیم داشت که به آن Target Write گفته می شود.

PCSource: یک سیگنال کنترلی 2 بیتی برای بالا

Step name	Action for R-type instructions	Action for memory-reference instructions	Action for branches
Instruction fetch		$IR = \text{Memory}[PC]$ $PC = PC + 4$	
Instruction decode/ register fetch		$A = \text{Registers}[IR(25-21)]$ $B = \text{Register}[IR(20-16)]$ $\text{Target} = PC + (\text{sign-extend } [IR(15-0)] \ll 2)$	
Execution, address computation, or branch completion	$ALUoutput = A \text{ op } B$	$ALUoutput = A + \text{sign-extend } [IR(15-0)]$	if (A == B) then PC = Target
Memory access or R-type completion	$\text{Reg}[IR(15-11)] = ALUoutput$	memory-data = Memory[ALUoutput] or Memory [ALUoutput] = B	
Write-back		$\text{Reg}[IR(20-16)] = \text{memory-data}$	

خلاصه ای از مراحل که برای اجرای هر نوع دستورالعمل انجام می شود.

تعداد مراحل اجرای دستورالعمل ها 3 تا 5 مرحله است. دو مرحله اول مستقل از نوع دستورالعمل هستند. بعد از این مراحل، یک دستورالعمل از 1 تا حداکثر 3 چرخه برای تکمیل شدن زمان می برد که این وابسته به نوع دستورالعمل است.

Signal name	Effect when deasserted	Effect when asserted
PCWrite	None	The PC is written; the source is controlled by PCSource.
PCWriteCond	None	The PC is written; if the Zero output from the ALU is also active.
TargetWrite	None	The output of the ALU is written into the register Target.

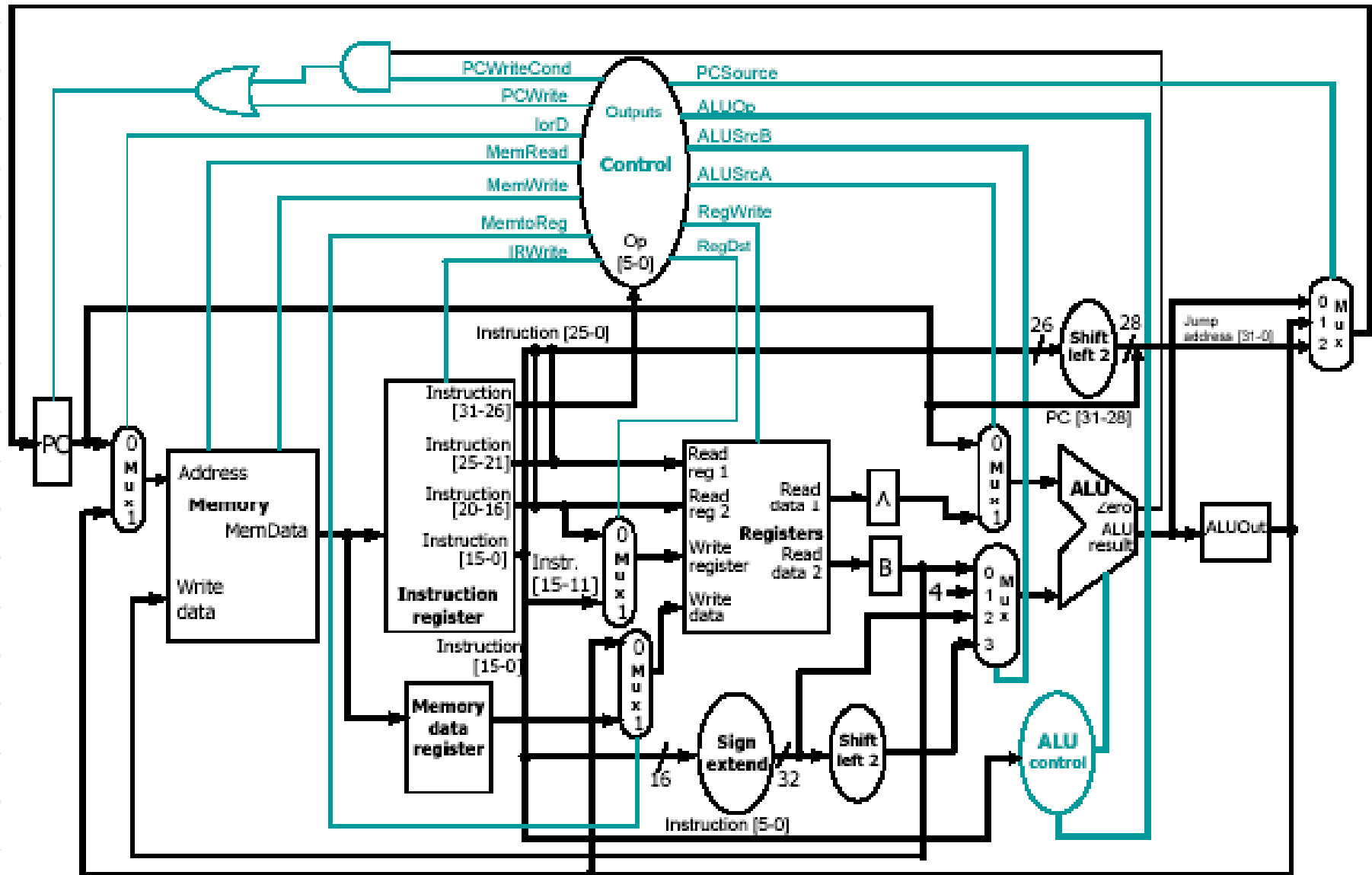
عمل ناشی از تنظیم هر یک از سیگنال های کنترلی 1 بیتی

Signal name	Value	Effect
	00	The ALU output is sent to the PC for writing.
PCSource	01	The contents of the register Target are sent to the PC for writing.
	10	The jump target address [PC + 4(29-26)] concatenated with IR (25-0) and shifted left two bits) is sent to the PC for writing.

عمل ناشی از تنظیم هر یک از سیگنال های کنترلی 2 بیتی

کنترل کردن pc بعدی

شکل 5.28: مسیر داده و سیگنال های کنترلی کامل برای چند چرخه ای (همراه پرش)



روشهای طراحی کنترل برای اجرای دستورات چند چرخه ای

- مبنا قرار دادن ماشین با حالات محدود

- مبنا قرار دادن Microprogramming

ساختار FSM

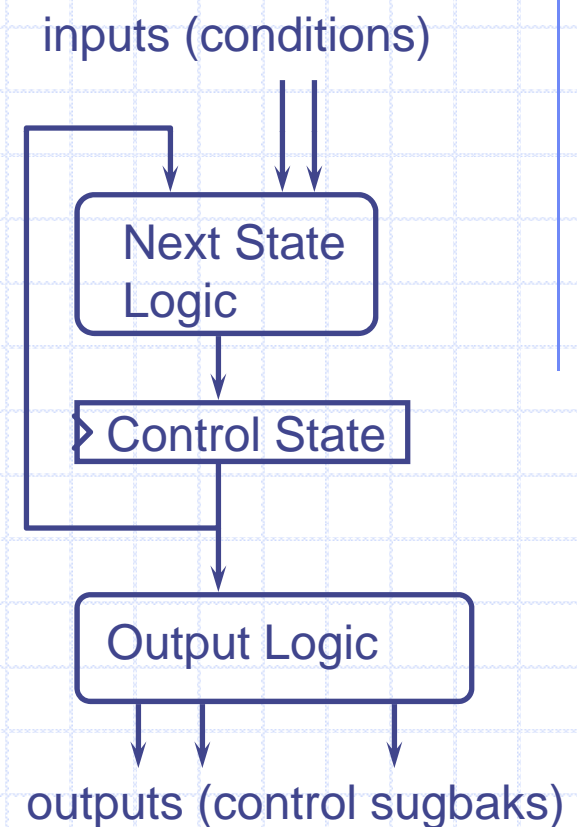
- مراحل:

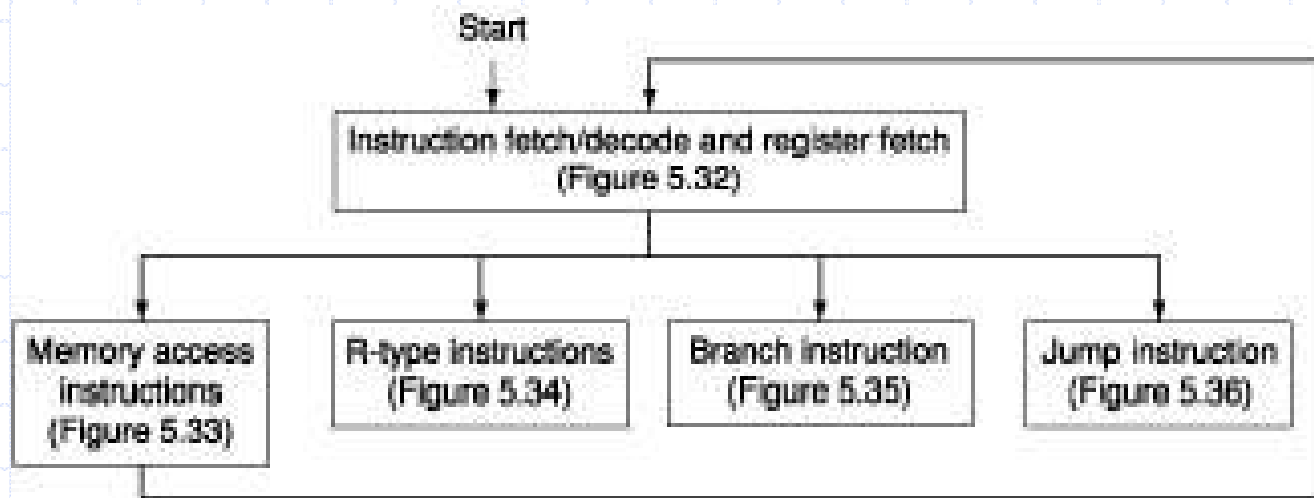
- Specify a set of output to be asserted

- Assume signals not asserted are left
disserted by default

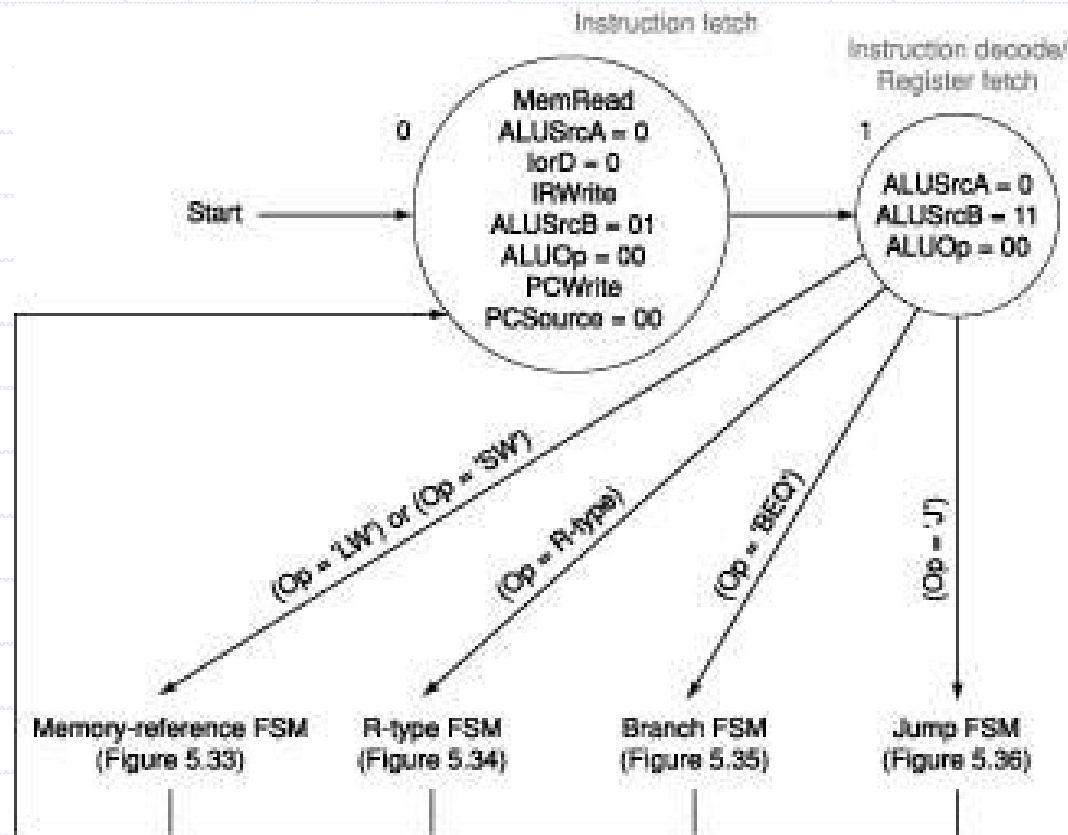
- کنترل کردن تسهیم کننده ها همیشه صراحتا مشخص
شده است

- تابع حالت بعدی

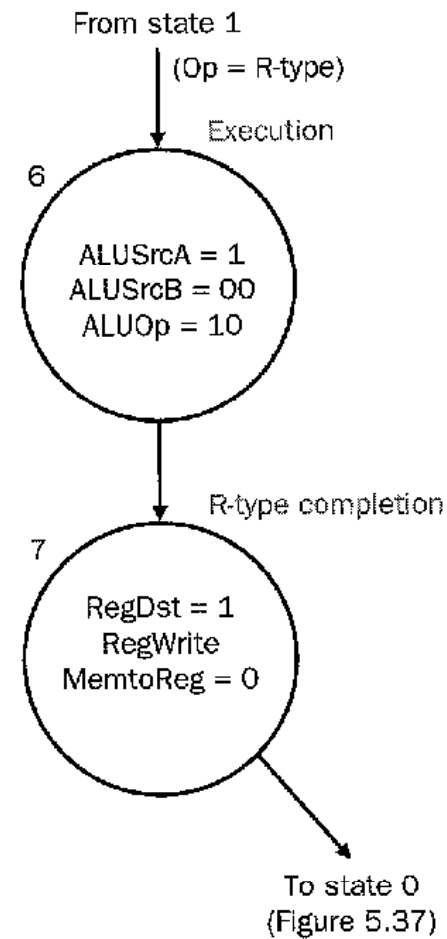




دید سطح بالایی از کنترل ماشین حالات متناهی (FSM)

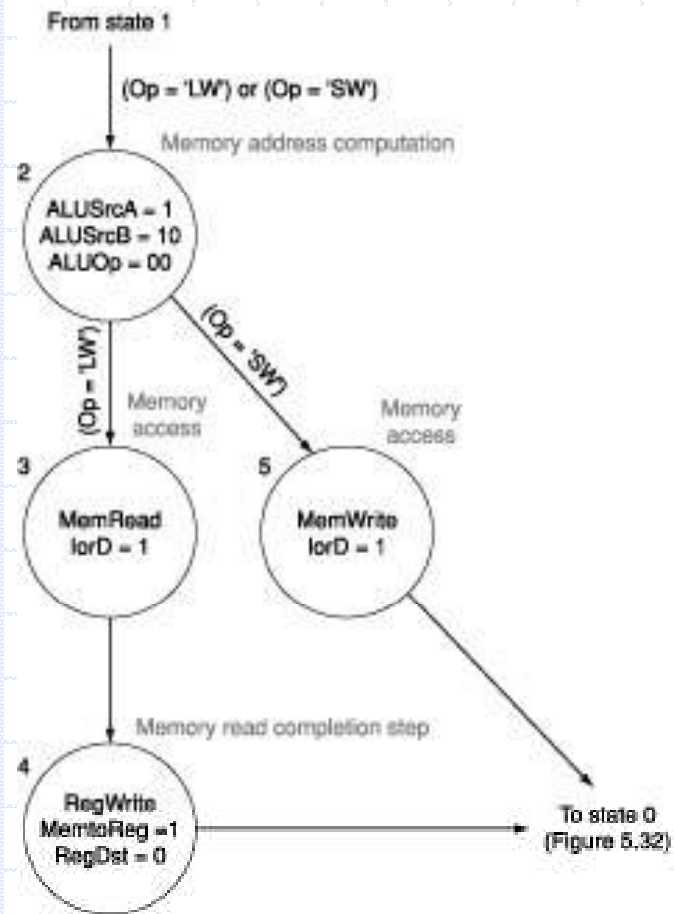


دستور العمل واكشى و كدبرداری قسمت یکسانی برای همه دستور العمل ها هستند (شکل 5.37_جدید 5.32)

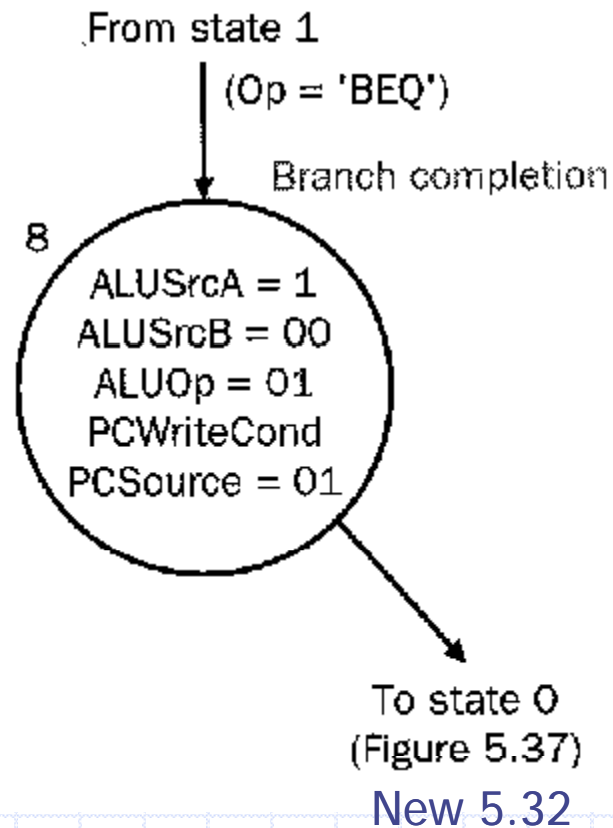


New 5.32

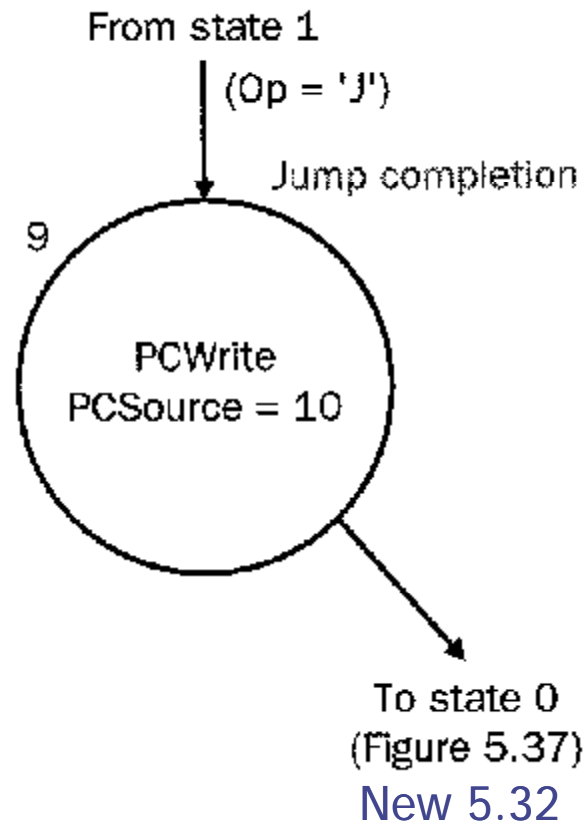
دستور العمل های R_type می توانند با یک FSM 2 مرحله ای ساده پیاده سازی شوند



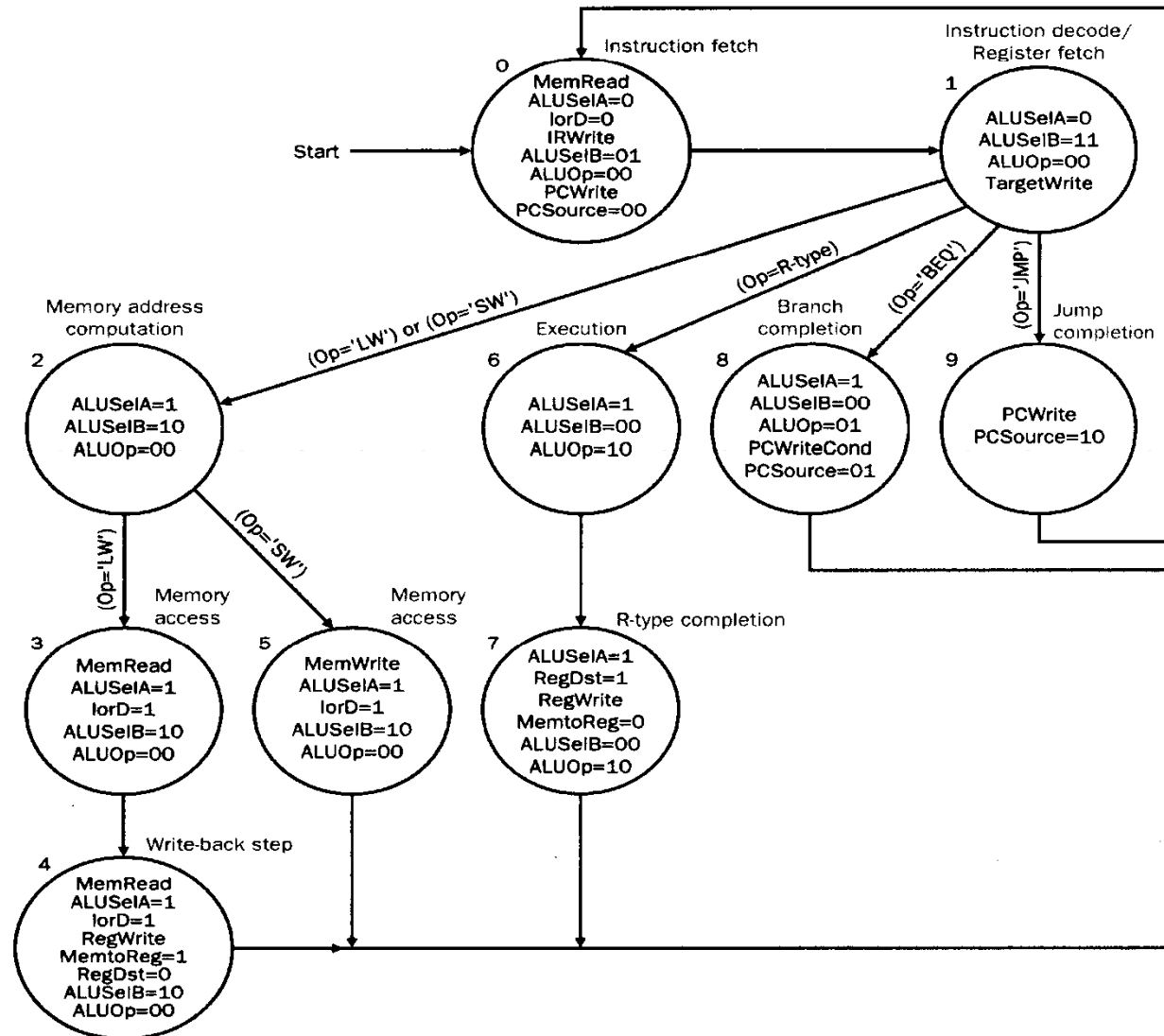
FSM برای کنترل کردن دستور العمل های ارجاع به حافظه 4 مرحله دارد (شکل 5.38_جدید 5.33)



دستور العمل انشعاب فقط به یک ماشین تک مرحله ای
نیاز دارد (شکل 5.40 جدید 5.35)

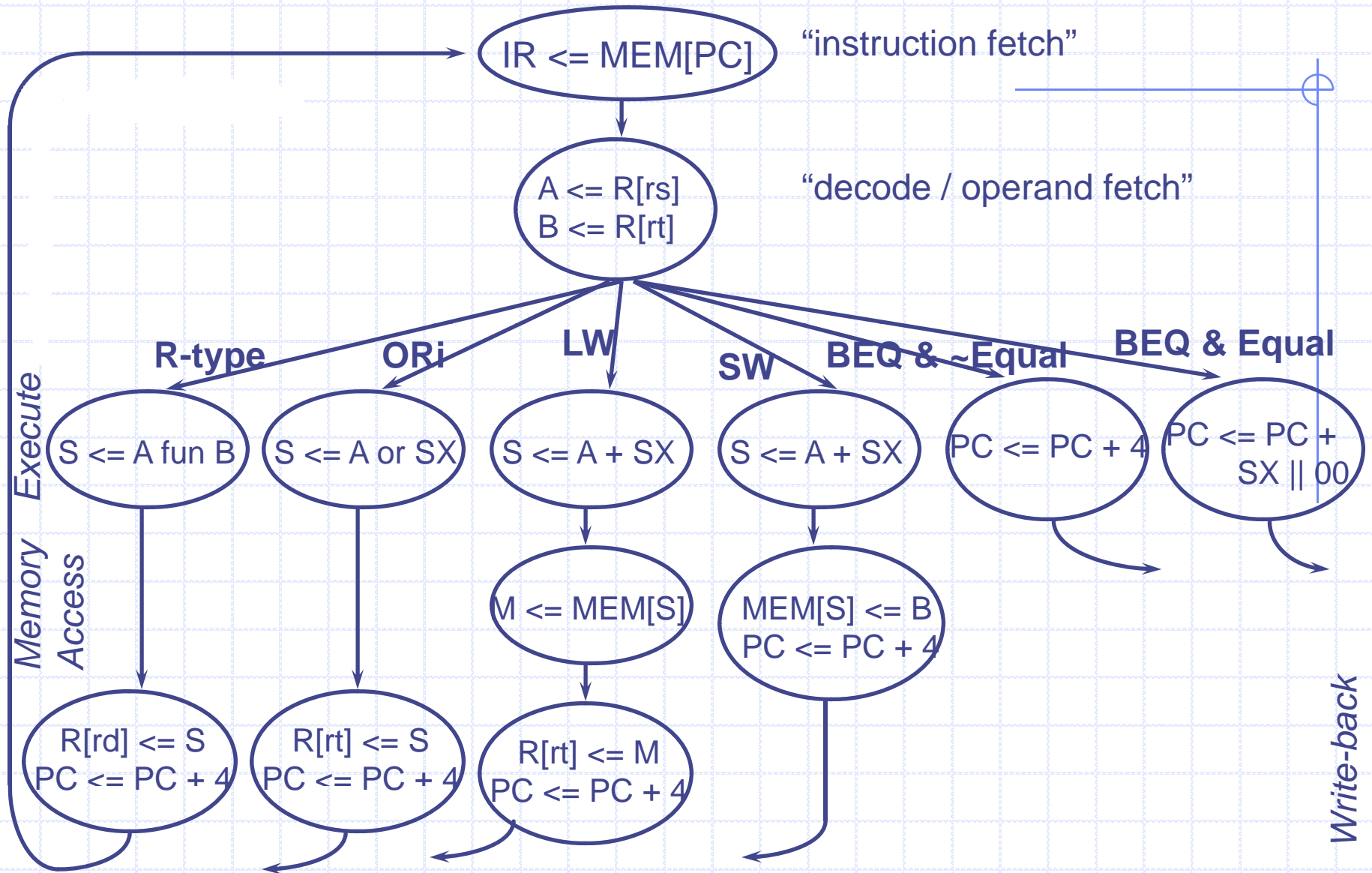


The jump instruction requires a single state that asserts two control signals to write the PC with the lower 26 bits of the instruction register shifted left two bits. (Figure 5.41 – new 5.36))

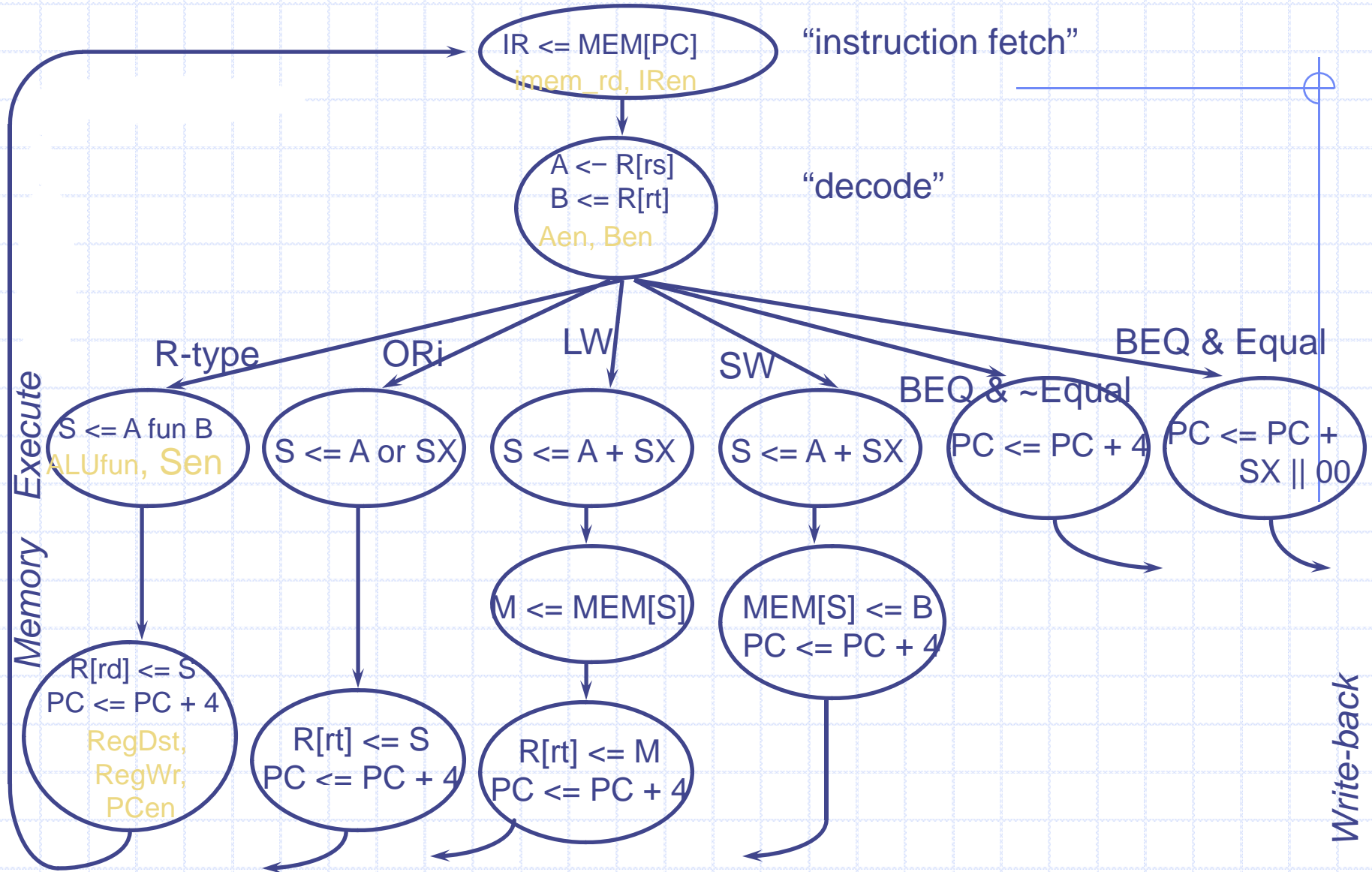


واحد کنترل ماشین حالت منتهای به صورت کامل بر روی مسیره داده ها (Figure 5.42 – New 5.38)

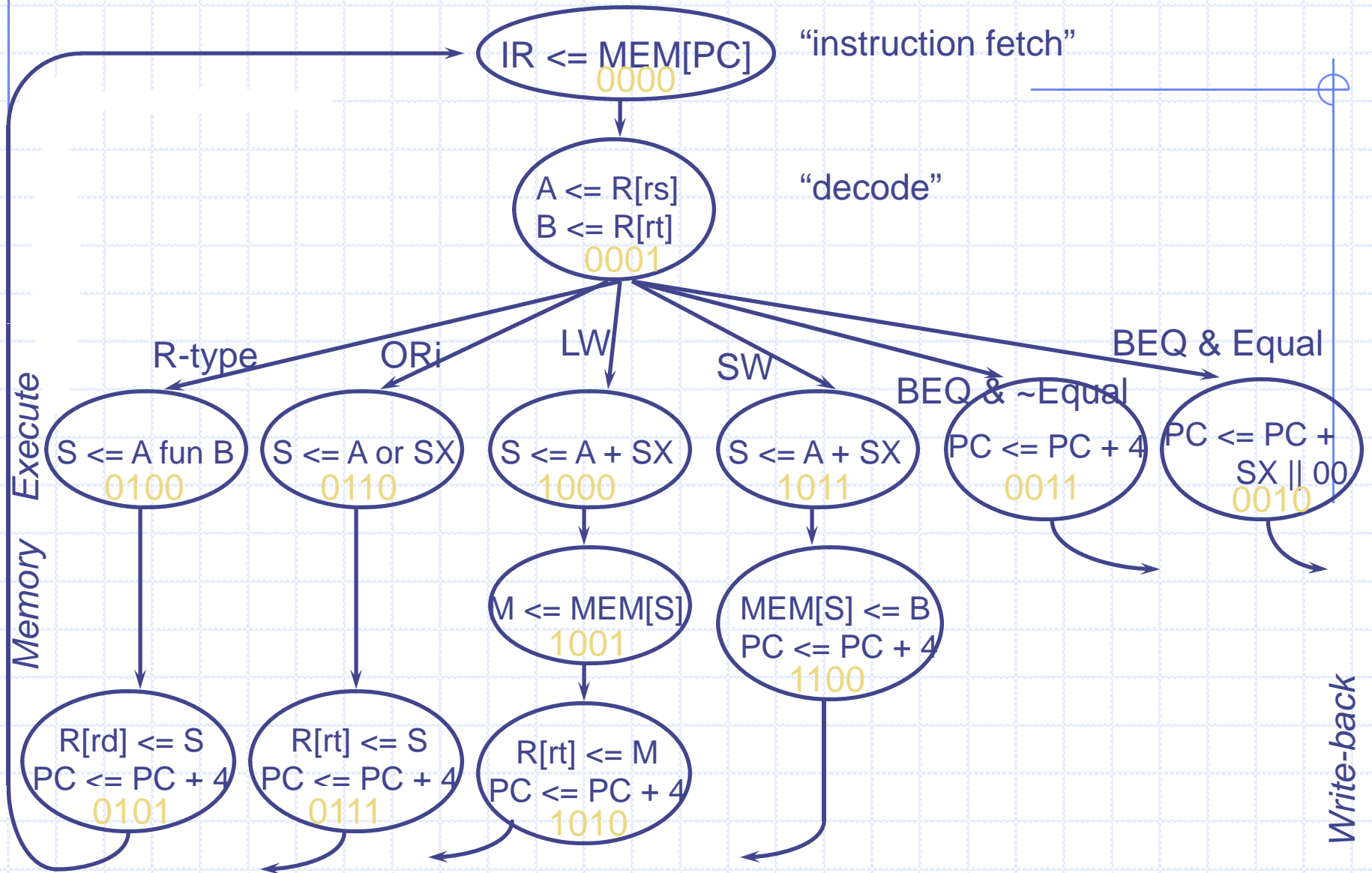
کنترل اختصاصی برای چند چرخه ای



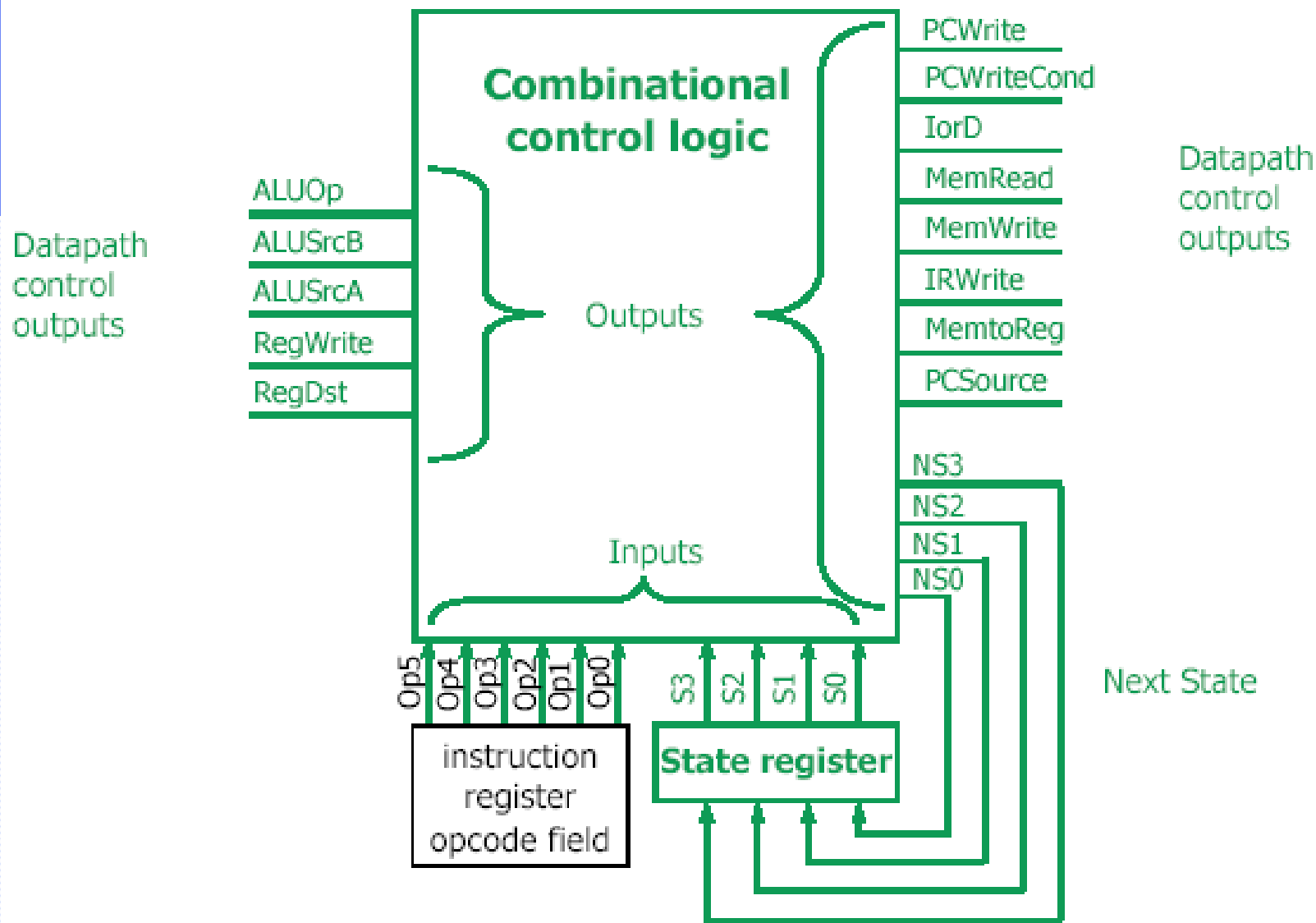
Mapping RTs to Control Points



تخصیص حالات



اجرای دستگاه کنترل ماشین حالت متناهی



تعادل منطقہ، براہ، خروہ، ہای، سگنال کنٹرول

Output	Current States
PCWrite	state0 + state9
PCWriteCond	state8
IorD	state3 + state5
MemRead	state0 + state3
MemWrite	state5
IRWrite	state0
MemtoReg	state4
PCSource1	state9
PCSource0	state8
ALUOp1	state6
ALUOp0	state8
ALUSrcB1	state1 + state2
ALUSrcB0	state0 + state1
ALUSrcA	state2 + state6 + state8
RegWrite	state4 + state7
RegDst	state7

For Example:

$$PCWrite = \overline{S3} \cdot \overline{S2} \cdot \overline{S1} \cdot \overline{S0} + S3 \cdot \overline{S2} \cdot \overline{S1} \cdot S0$$

تعادل منطقی برای خروجی های حالت بعدی

Output	Current States	Op
NextState0	state4 + state5 + state7 + state8 + state9	
NextState1	state0	
NextState2	state1	(Op = 'lw') + (Op = 'sw')
NextState3	state2	(Op = 'lw')
NextState4	state3	
NextState5	state2	(Op = 'sw')
NextState6	state1	(Op = 'R-type')
NextState7	state6	
NextState8	state1	(Op = 'beq')
NextState9	state1	(Op = 'jump')

For Example:

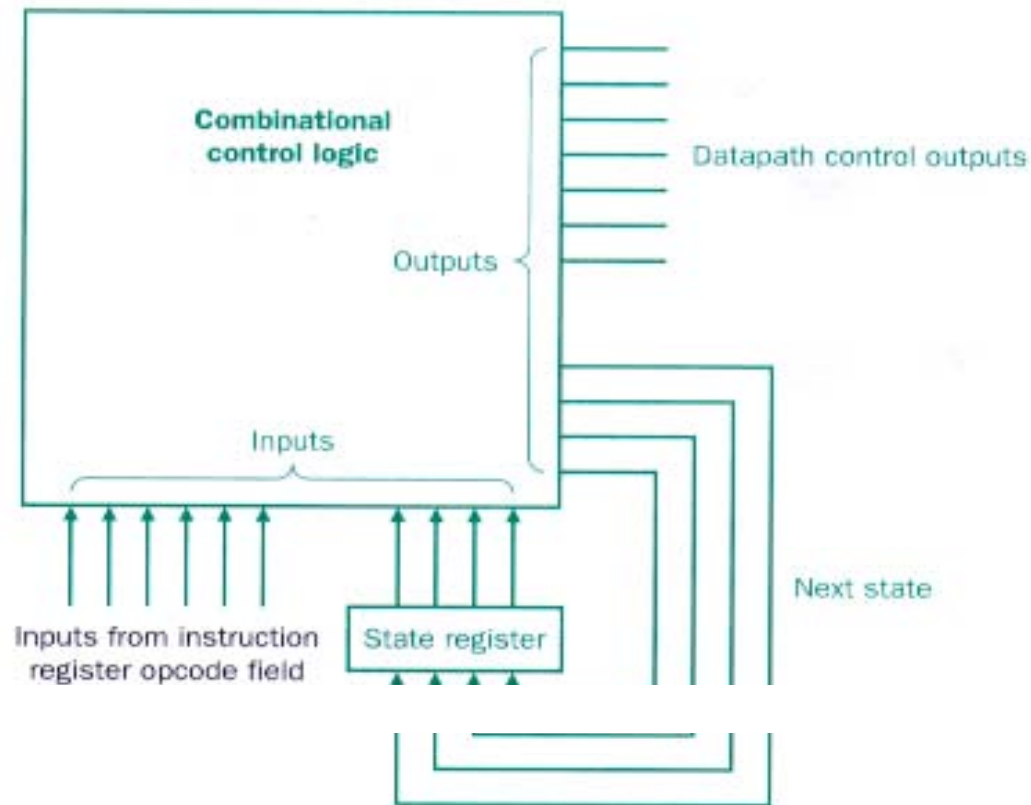
$$\text{NextState1} = \text{State0} = \overline{S3} \cdot \overline{S2} \cdot \overline{S1} \cdot \overline{S0}$$

$$\text{NextState3} = \text{State2} \cdot (\text{Op}[5-0] = \text{'lw'})$$

$$= \overline{S3} \cdot \overline{S2} \cdot S1 \cdot \overline{S0} \cdot \text{Op5} \cdot \overline{\text{Op4}} \cdot \overline{\text{Op3}} \cdot \text{Op2} \cdot \text{Op1}$$

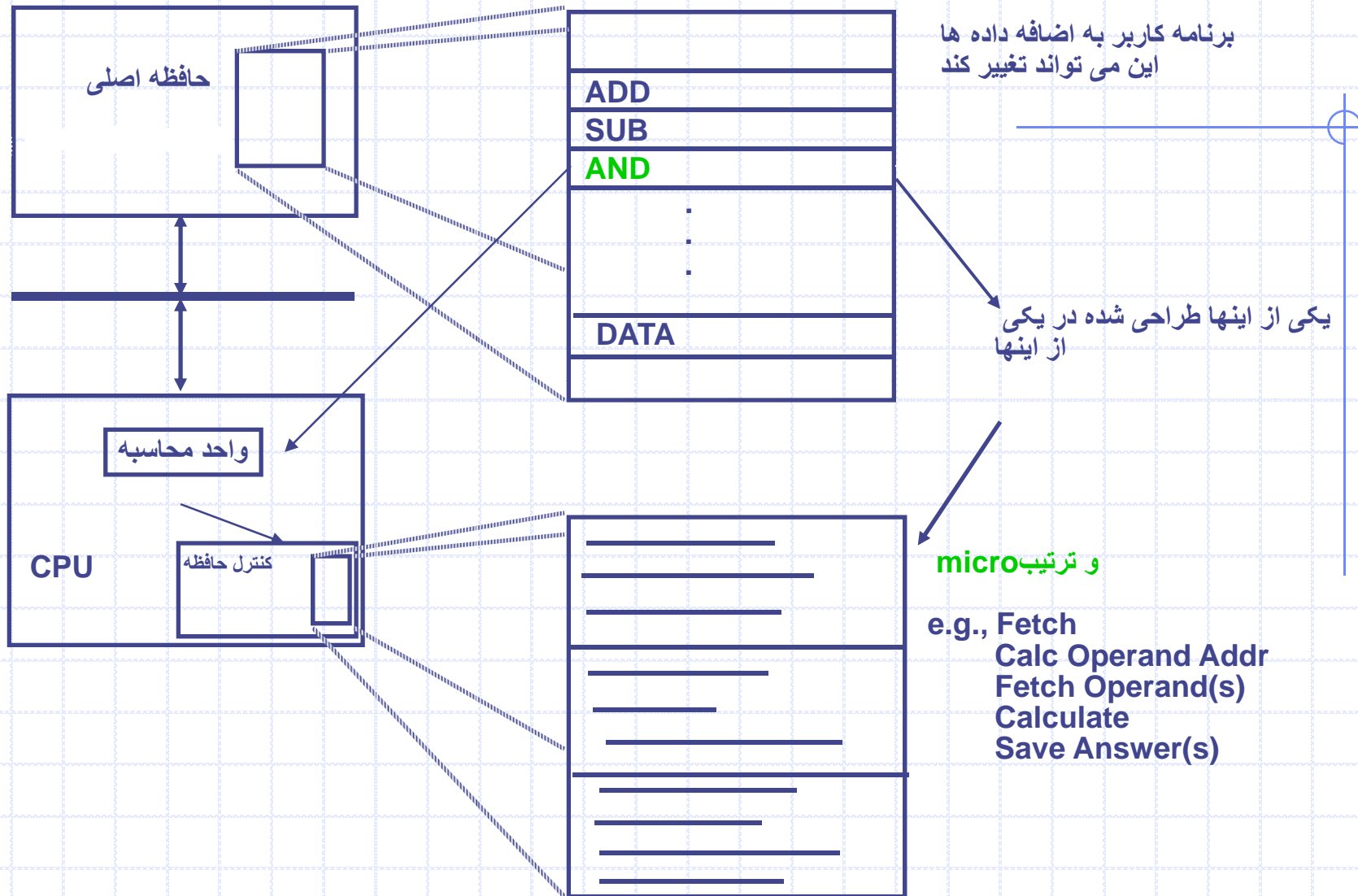
سوال:

- چه طور میشود CPI را از روی FSA تعیین کرد؟
- FSA چه طور انجام می شود؟



Finite state machine controllers are typically implemented using a block of combinational logic and a register to hold the current state.

تفسیر "دستور العمل Macro"



Field name	Values for field	Function of field with specific value
Label	Any string	Used to specify labels to control microcode sequencing. Labels that end in a 1 or 2 are used for dispatching with a jump table that is indexed based on the opcode. Other labels are used as direct targets in the microinstruction sequencing. Labels do not generate control signals directly but are used to define the contents of dispatch tables and generate control for the Sequencing field.
ALU control	Add	Cause the ALU to add.
	Subt	Cause the ALU to subtract; this implements the compare for branches.
	Func code	Use the instruction's funct field to determine ALU control.
SRC1	PC	Use the PC as the first ALU input.
	A	Register A is the first ALU input.
SRC2	B	Register B is the second ALU input.
	4	Use 4 for the second ALU input.
	Extend	Use output of the sign extension unit as the second ALU input.
	Extshft	Use the output of the shift-by-two unit as the second ALU input.
Register control	Read	Read two registers using the rs and rt fields of the IR as the register numbers, putting the data into registers A and B.
	Write ALU	Write the register file using the rd field of the IR as the register number and the contents of ALUOut as the data.
	Write MDR	Write the register file using the rt field of the IR as the register number and the contents of the MDR as the data.
Memory	Read PC	Read memory using the PC as address; write result into IR (and the MDR).
	Read ALU	Read memory using ALUOut as address; write result into MDR.
	Write ALU	Write memory using the ALUOut as address; contents of B as the data.
PCWrite control	ALU	Write the output of the ALU into the PC.
	ALUOut-cond	If the Zero output of the ALU is active, write the PC with the contents of the register ALUOut.
	Jump address	Write the PC with the jump address from the instruction.
Sequencing	Seq	Choose the next microinstruction sequentially.
	Fetch	Go to the first microinstruction to begin a new instruction.
	Dispatch i	Dispatch using the ROM specified by i (1 or 2).

واکشی بر نامه میکرو

Label	ALU control	SRC1	SRC2	Register control	Memory	PCWrite control	Sequencing
Fetch	Add	PC	4		Read PC	ALU	Seq
	Add	PC	Extshft	Read			Dispatch 1

دستور العمل ارجاع به حافظه برنامه می‌کرو

Label	ALU control	SRC1	SRC2	Register control	Memory	PCWrite control	Sequencing
Mem1	Add	A	Extend				Dispatch 2
LW2					Read ALU		Seq
				Write MDR			Fetch
SW2					Write ALU		Fetch

R-type برنامہ میکر و

Label	ALU control	SRC1	SRC2	Register control	Memory	PCWrite control	Sequencing
Rformat1	Func code	A	B				Seq
				Write ALU			Fetch

انشعاب برنامه میکرو

Label	ALU control	SRC1	SRC2	Register control	Memory	PCWrite control	Sequencing
BEQ1	Subt	A	B			ALUOut-cond	Fetch

پیش برنامه میکرو

Label	ALU control	SRC1	SRC2	Register control	Memory	PCWrite control	Sequencing
JUMP1						Jump address	Fetch

همه میکرو کد

Label	ALU control	SRC1	SRC2	Register control	Memory	PCWrite control	Sequencing
Fetch	Add	PC	4		Read PC	ALU	Seq
	Add	PC	Extshft	Read			Dispatch 1
Mem1	Add	A	Extend				Dispatch 2
LW2					Read ALU		Seq
				Write MDR			Fetch
SW2					Write ALU		Fetch
Rformat1	Func code	A	B				Seq
				Write ALU			Fetch
BEQ1	Subt	A	B			ALUOut-cond	Fetch
JUMP1						Jump address	Fetch

فصل ششم

خط لوله ای کردن و معماری خط لوله ای شده

خط لوله ای

◆ خط لوله ای چیست – مفاهیم پایه

◆ مسیر داده خط لوله ای شده

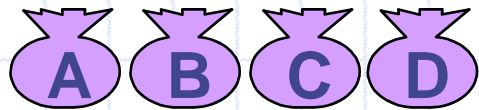
یک بررسی موردی از MIPS

◆ کنترل خط لوله

◆ رفع hazard در خط لوله

مفاهیم پایه خط لوله ای

خط لوله ای کردن، معمول است



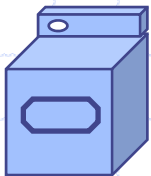
◆ خط لوله ای کردن، روشی را برای اجرای همزمان چند دستورات عمل بوجود می آورد.



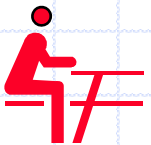
◆ مثال شستن لباسها

◆ Ann, Brian, Dave, Cathy هر کدام توده ای از

لباس را برای شستن، خشک کردن و تا کردن در اختیار دارند.



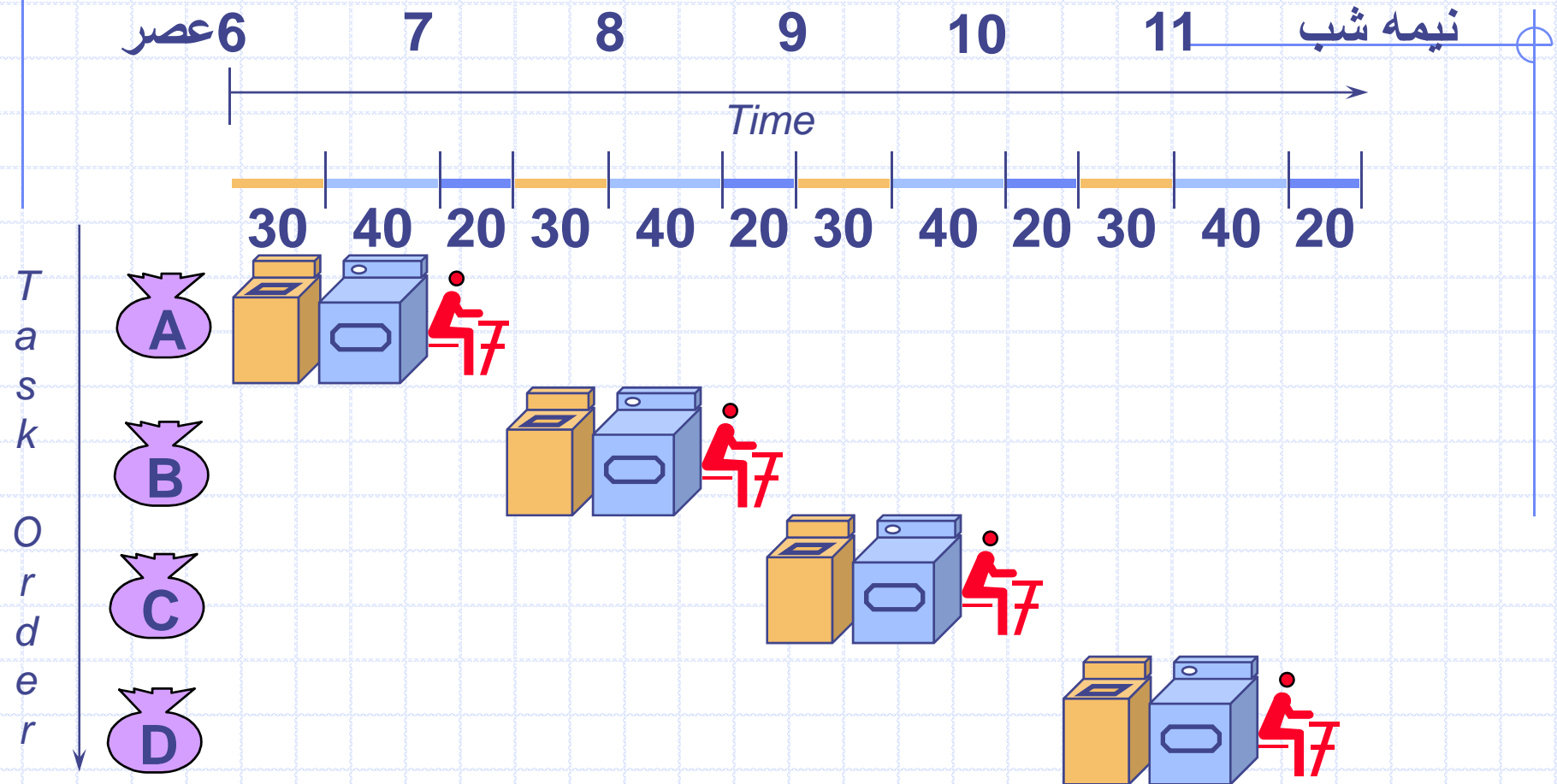
◆ شستن 30 دقیقه زمان می برد.



◆ خشک کن 40 دقیقه زمان می برد.

◆ تا کردن 20 دقیقه زمان نیاز دارد.

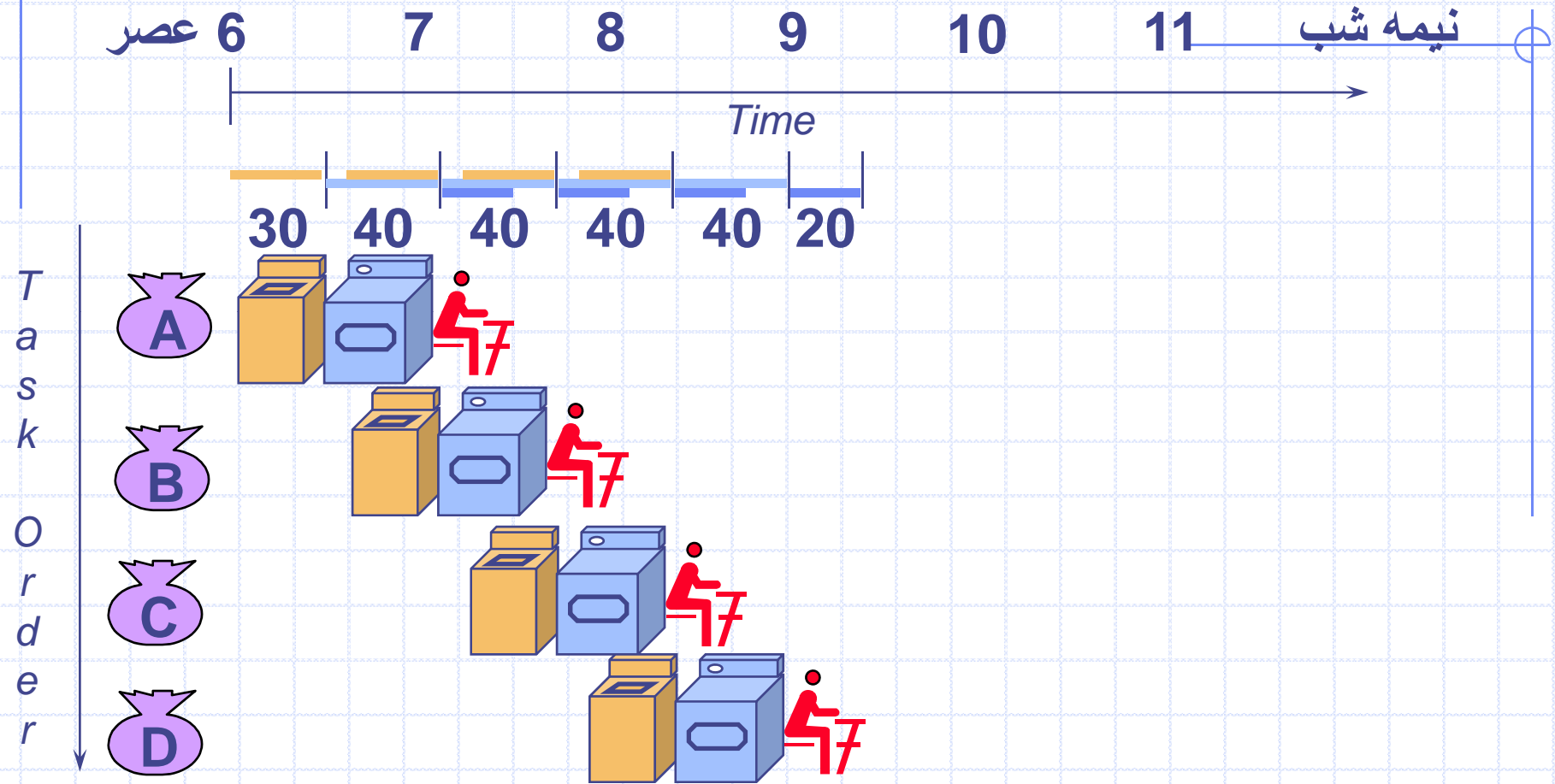
شستن ترتیبی



◆ شستن ترتیبی 4 توده لباس، 6 ساعت زمان می برد.

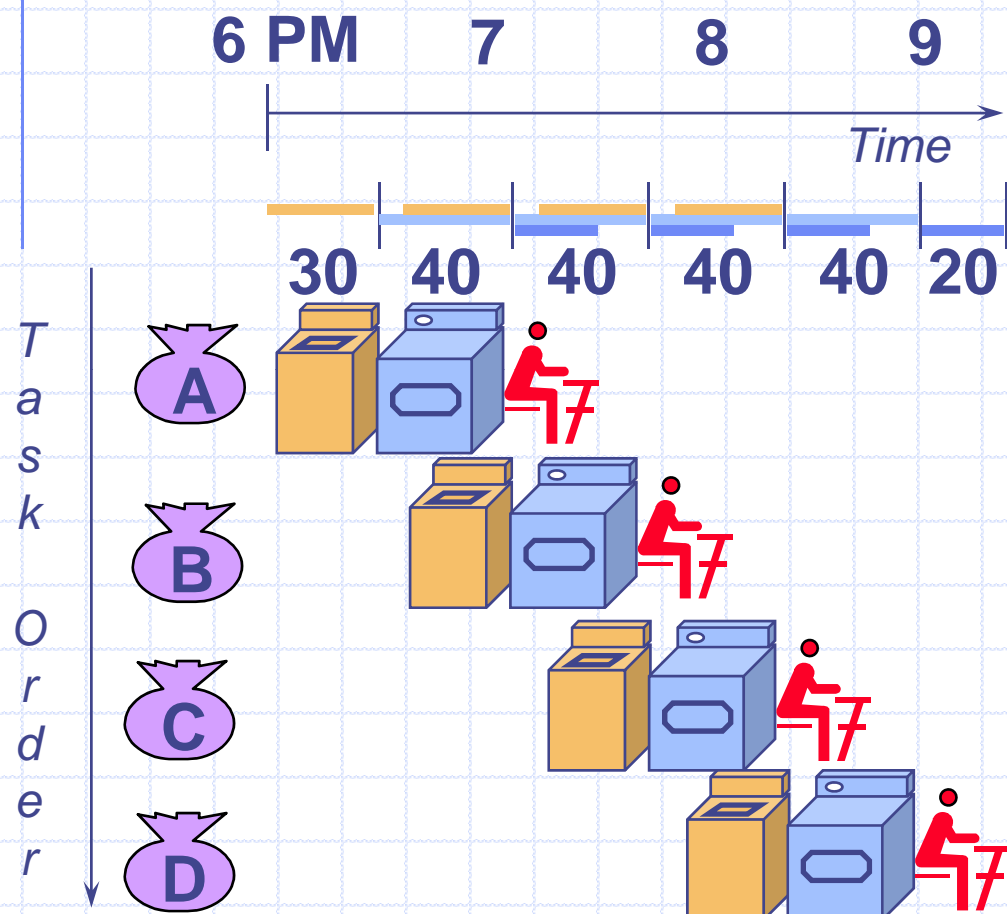
◆ اگر آنها خط لوله ای را فرا گرفته بودند، چقدر طول می کشید؟

شستن خط لوله ای شده:



◆ شستن خط لوله ای 4 توده لباس، 3 ساعت و نیم زمان می برد.

دروس خط لوله ای کردن



◆ خط لوله ای به زمان تاخیر یک کار

کمکی نمی کند، بلکه توان عملیاتی کل

بار کاری را بهبود می بخشد.

◆ نرخ خط لوله محدود به کندترین مرحله

است.

◆ اجرای همزمان چند کار نیاز به منابع

متفاوتی دارد.

◆ پتانسیل تسریع = تعداد مراحل خط لوله

طول نامتوازن مراحل خط لوله باعث

کاهش تسریع می شود.

◆ زمان مصرفی برای پر شدن و خالی

شدن خط لوله باعث کاهش تسریع می

شود.

◆ متوقف شدن به خاطر وابستگی

ها (STALL)

مفهوم پایه

◆ خط لوله: چندین دستور العمل به طور همزمان در حال اجرایند.

◆ خط لوله به بخشها یا قطعات تقسیم می شود.

◆ چرخه ماشین:

■ زمان مورد نیاز برای گذر از یک مرحله

■ چرخه ماشین بوسیله کندترین مرحله خط لوله معین می گردد.

■ معمولا $\text{چرخه ماشین} = \text{چرخه ساعت}$

◆ در یک ماشین خط لوله ای شده کاملاً متوازن:

$$(1) \quad \frac{\text{زمان غیر خط لوله ای}}{\text{تعداد بخشهای خط لوله}} = \text{زمان} \quad \blacksquare$$

◆ در یک ماشین معمولی

(1) درست نیست،

زمان مراحل برابر نیست، یک سربار وجود دارد.

اما می تواند تا اختلاف 10% به (1) نزدیک شود.

خط لوله ای کردن به عنوان یک تکنیک معماری

◆ بطور کلی، تقریباً برای کاربر نامحسوس است.

Scalar pipelined machine vs. vector
machine

◆ توان عملیاتی خط لوله ای :

تعداد دستورالعملهای کامل شده

چرخه

◆ بنابراین خط لوله ای کردن توان عملیاتی (ظرفیت پذیرش) را افزایش می دهد، اما زمان اجرای هر دستورالعمل بدون تغییر می ماند.

◆ نرخ ساعت یک ماشین خط لوله ای شده محدود شده است به:

■ latch time

■ clock skew

■ زمان تاخیر مورد نیاز برای انتشار سیگنال ساعت در یک تراشه

محدودیت‌های کارایی در یک خط لوله

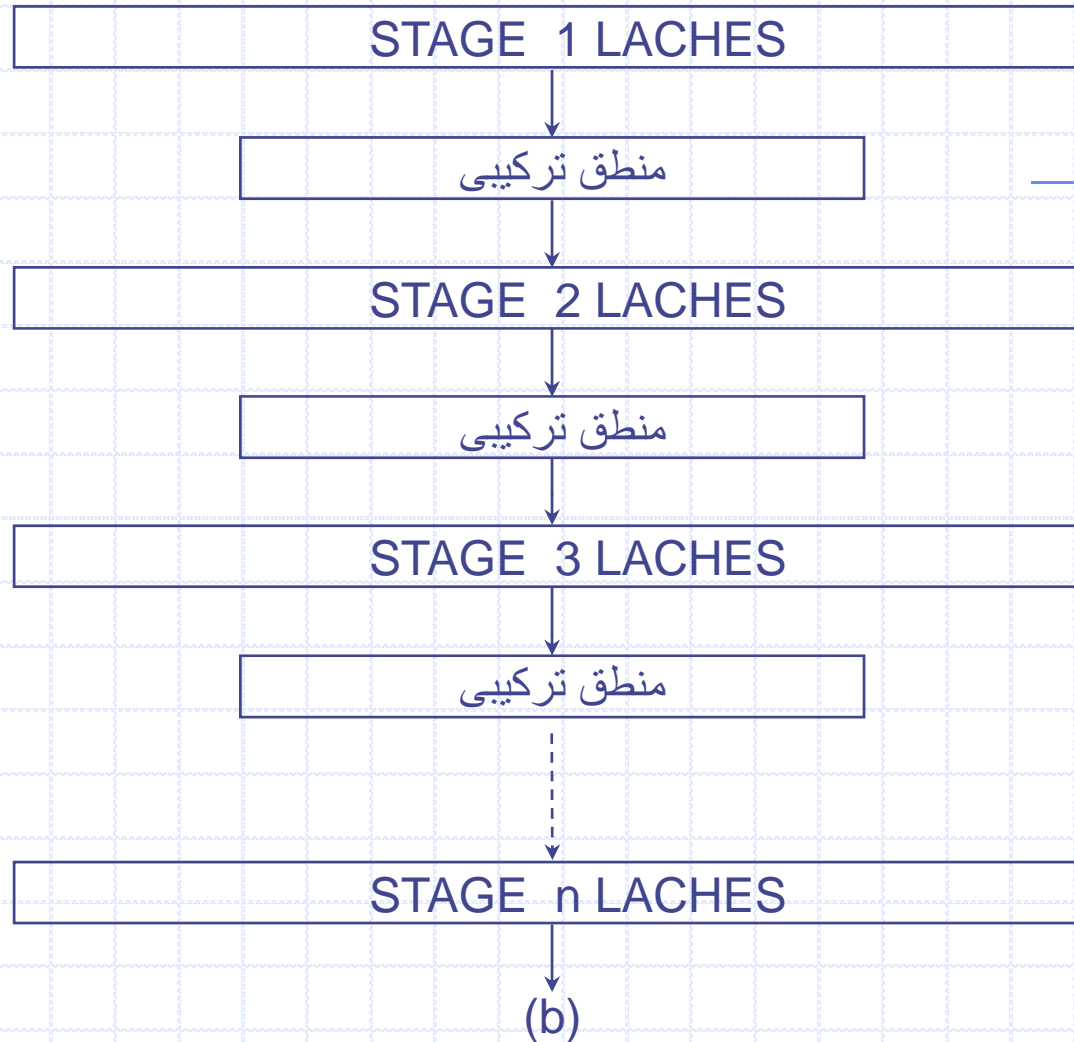
◆ نمی‌توان از کندترین مرحله انتظار سرعت داشت

◆ پیچیدگیها در واقعیت، زمانهای پردازش متفاوت برای

مراحل مختلف، فعل و انفعالات/وابستگیها بین

مراحل – ممکن است باعث وابستگی داده ای شود.

(پویا)



فاکتورهای مهم در مسیر لوله (مدت چرخه) ، تاخیر latch و clock skew

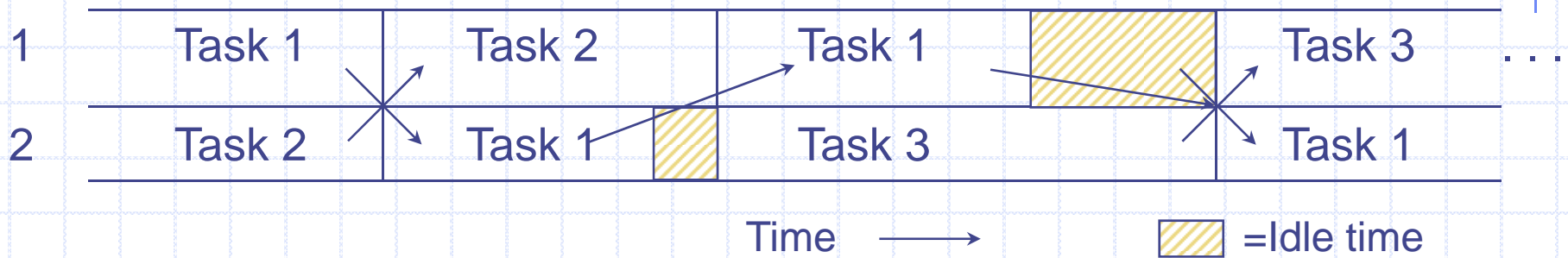
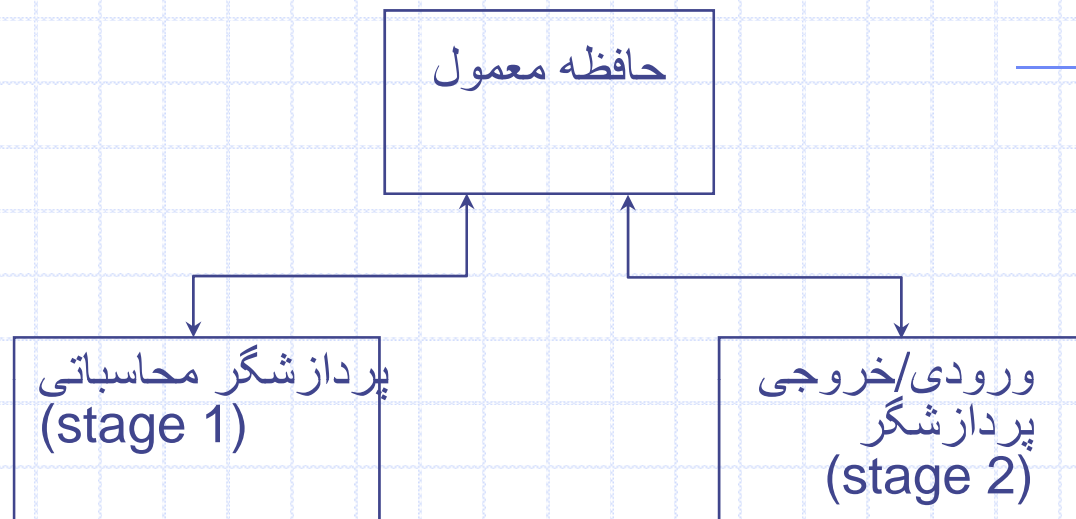
رویهم اندازی (Overlap) در برابر خط لوله ای

◆ خط لوله ای

- ارتباط قوی زیر بخشها (sub functions)
- زمان پایه ثابت مرحله
- ارزیابی تابع پایه مستقل

◆ مشترکات

- پیوند ضعیف زیر تابع ها
- زمان متغیر مرحله
- ارزیابی وابستگی بین تابعی



رویهم اندازی CPU/I/O

خط لوله پویا و ایستا

◆ ایستا:

- تنها ارزیابی های مکرر همان کار با داده های متفاوت انجام می شود

- - no dynamic data dependencies between initiations

- طرح ثابت آغازها

◆ پویا: پویا متضاد است با ایستا

- (اسنکرون)

- (روی هم افتادگی)

◆ مثالی از خط لوله ایستا:

- جمع ممیز شناور در خط لوله

یک تابع واحدلوله _ چند تابع لوله

◆ یک تابع واحد لوله

- مثال: جمع لوله

◆ چند تابع لوله

- مثال

- لوله حسابی

◆ بردار لوله

- کنترل لوله قابل برنامه ریزی

- (دستورالعمل برداری وظیفه و ورودی ها را تعیین می کند)

خط لوله MIPS

مراحل خط لوله:

- واکنشی
- ID (کد برداری + واکنشی ثبات ها)
- محاسبات
- دستیابی به حافظه
- باز نویسی

Instruction number	Clock number								
	1	2	3	4	5	6	7	8	9
Instruction i	IF	ID	EX	MEM	WB				
Instruction $i+1$		IF	ID	EX	MEM	WB			
Instruction $i+2$			IF	ID	EX	MEM	WB		
Instruction $i+3$				IF	ID	EX	MEM	WB	
Instruction $i+4$					IF	ID	EX	MEM	WB

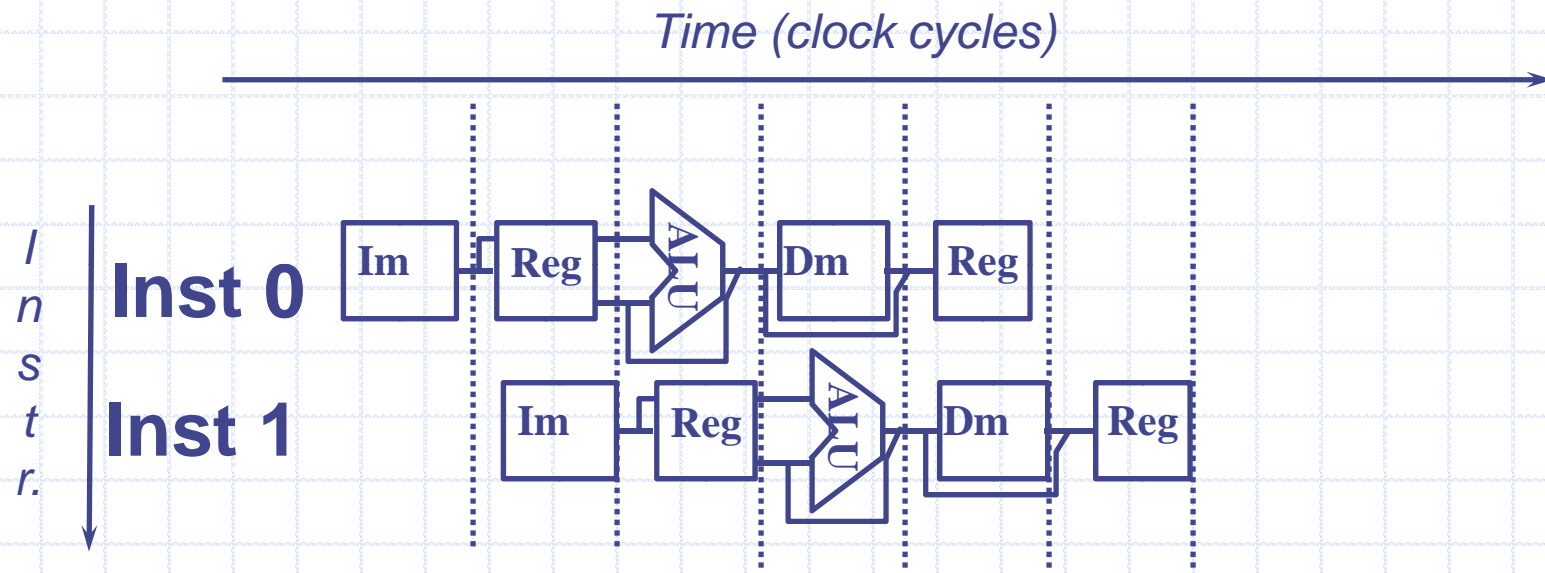
در هر چرخه ساعت دستورالعمل دیگری واکنشی شده و اجرا در 5 مرحله آغاز می شود. اگر یک دستورالعمل شروع شود در هر چرخه ساعت started every clock cycle, the performance will be five times that of a machine that is not pipelined.

خط لوله MIPS

نمایش دیگر

	1	2	3	4	5	6	7	8	9
IF	1	2	3	4	5				
ID		1	2	3	4	5			
EX			1	2	3	4	5		
MEM				1	2	3	4	5	
WB					1	2	3	4	5

نمایش منقوش خط لوله



◆ با پاسخ دادن به سوالات مشابه می تواند کمک کند

- چند چرخه زمان می برد که این کد اجرا شود؟
- ALU در طول 4 چرخه چه کاری انجام می دهد؟
- آیا 2 دستورالعمل سعی در استفاده همزمان از یک منبع مشترک دارند؟

چرا خط لوله ای؟

◆ فرض کنید

- 100 دستورالعمل را می خواهیم اجرا کنیم
- ماشین تک چرخه ای که یک چرخه ساعت دارد NS45 زمان می خواهد
- ماشین چند چرخه ای و خط لوله ای که چند چرخه زمان می برند NS10 زمان می خواهد

■ CPI ماشین چند چرخه ای 4.6 است

◆ ماشین تک چرخه ای

$$45 \text{ ns/cycle} \times 1 \text{ CPI} \times 100 \text{ inst} = 4500 \text{ ns}$$

◆ ماشین چند چرخه ای

$$10 \text{ ns/cycle} \times 4.6 \text{ CPI} \times 100 \text{ inst} = 4600 \text{ ns}$$

◆ ماشین خط لوله ای ایده ال

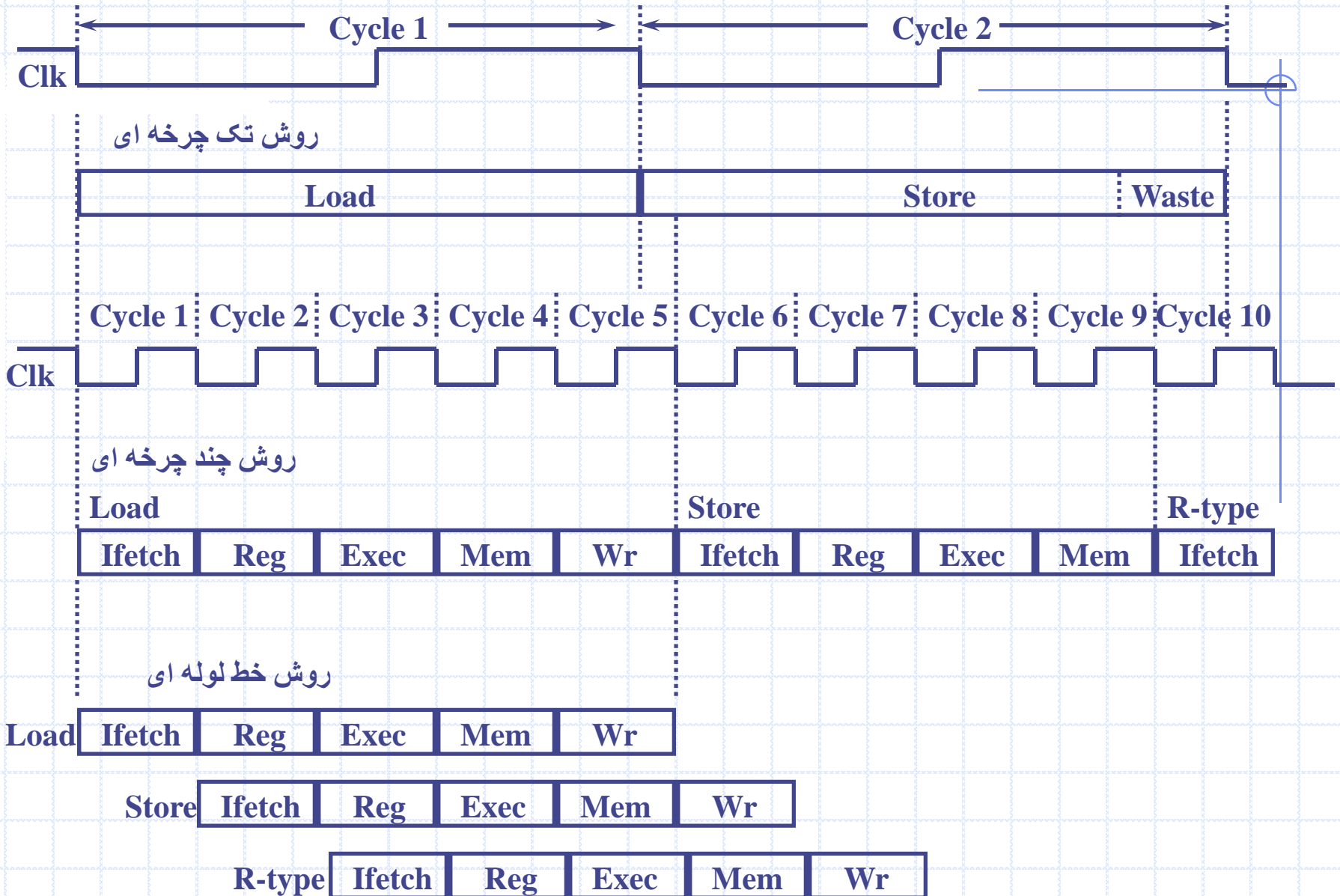
$$10 \text{ ns/cycle} \times (1 \text{ CPI} \times 100 \text{ inst} + 4 \text{ cycle drain}) = 1040 \text{ ns}$$

◆ میزان تسریع خط لوله ایده ال و ماشین تک چرخه ای

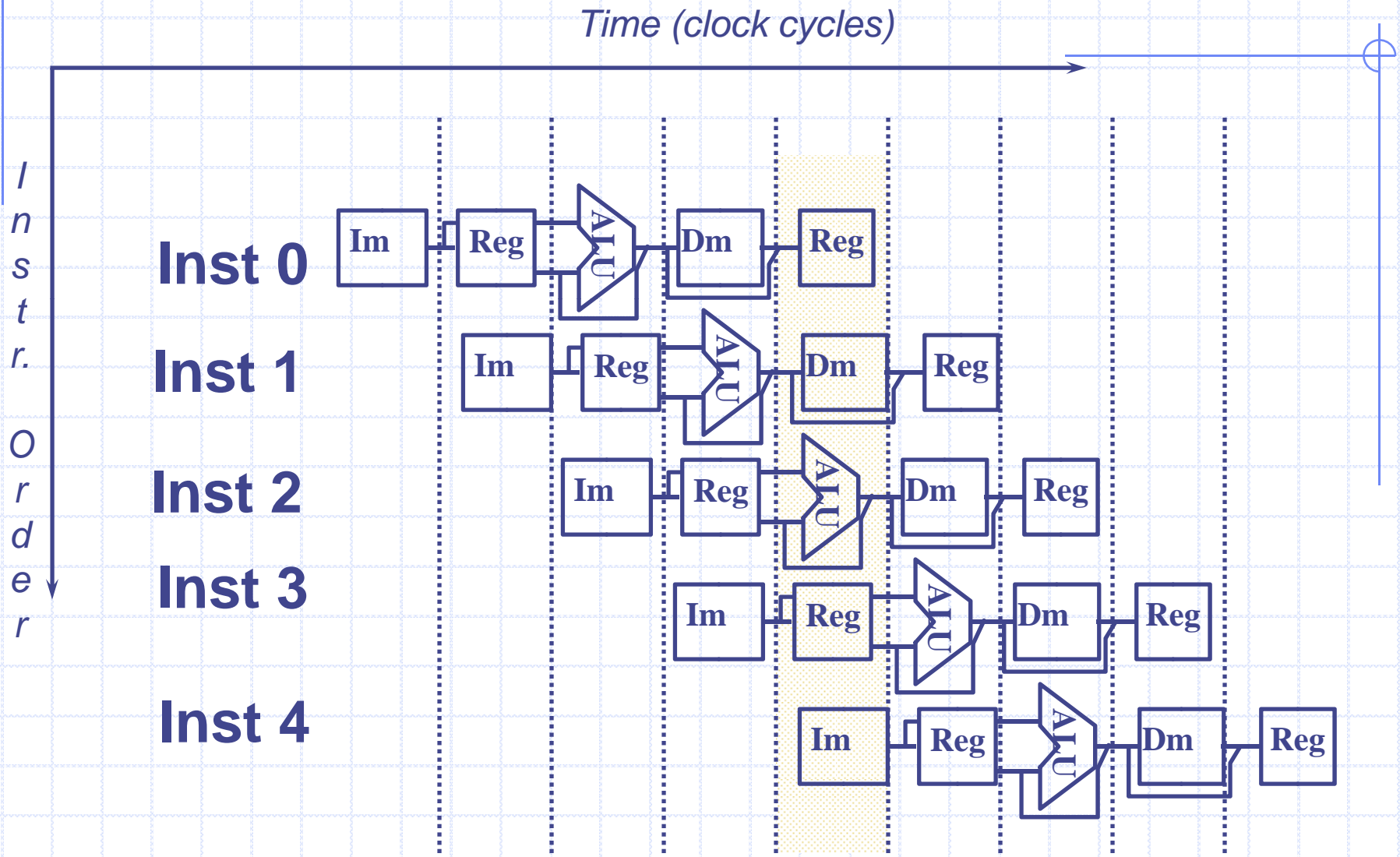
$$4500 \text{ ns} / 1040 \text{ ns} = 4.33$$

◆ What has not yet been considered?

تک چرخه ای_ چند چرخه ای_ خط لوله ای



چرا خط لوله؟ برای اینکه منابع را در اینجا داریم!



آیا خط لوله می تواند برای ما مشکل ایجاد کند؟

◆ بله: هزاردهای خط لوله ای

■ هزاردهای ساختاری: تلاش برای استفاده از یک منبع یکسان برای 2 کار متفاوت در یک زمان

◆ مثال: سعی در خواندن همزمان دو دستورالعمل از یک حافظه

■ هزاردهای داده ای: تلاش برای استفاده از یک آیتیم قبل از اینکه آماده شود

◆ دستورالعمل وابسته است به نتیجه دستورالعمل قبلی که هنوز در خط لوله است

```
add r1, r2, r3
```

```
sub r4, r2, r1
```

■ هزاردهای کنترلی: تلاش برای ظاهر کردن یک تصمیم قبل از اینکه شرط اجرا شود

◆ دستورالعمل های انشعاب

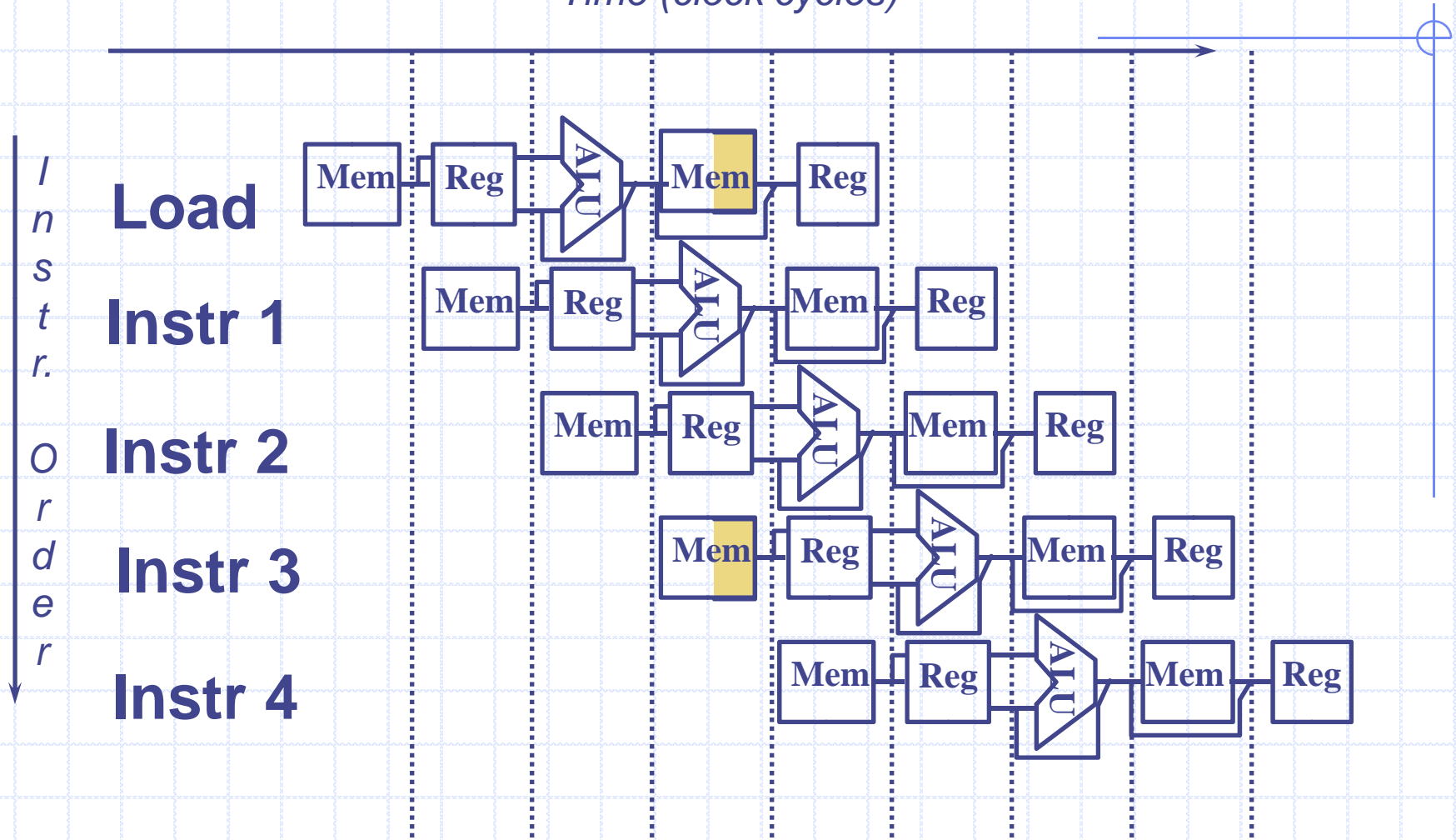
```
beq r1, loop
```

```
add r1, r2, r3
```

- ◆ همیشه میتوان هزاردها را با انتظار رفع کرد
- کنترل خط لوله ای باید هزارد را کشف کند
 - برداشتن عمل (تاخیر عمل) می تواند هزارد را حل کند

حافظه تک یک هزارد ساختاری است

Time (clock cycles)



تشخیص در این نمونه آسان است (نیمه روشن سمت راست یعنی خواندن و نیمه چپ یعنی نوشتن)

هزاردهای ساختاری کارایی را محدود می کند

◆ مثال: اگر 1.3 حافظه در هر دستورالعمل دستیابی شود و فقط یک حافظه مورد دستیابی قرار می گیرد در هر چرخه
 $CPI = 1.3$ میانگین

■ منبع بیشتر از 100% مورد استفاده واقع شده

◆ راه حل اول: جدا کردن حافظه داده و دستورالعمل از یکدیگر

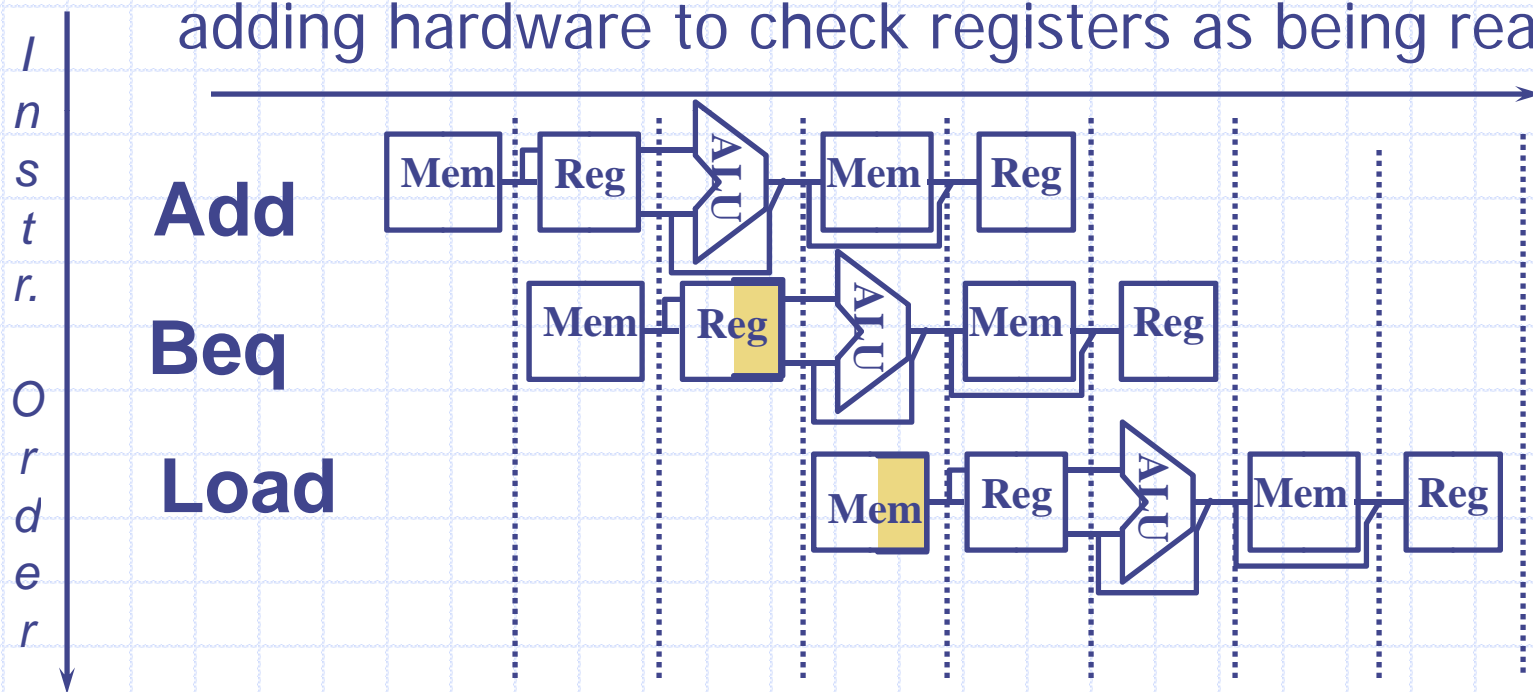
◆ راه حل دوم: اجزه دهیم حافظه بیش از یک خواندن و نوشتن در یک چرخه انجام دهد

◆ راه حل سوم: توقف

راه حل های هزارد کنترلی

◆ توقف: تا زمانی که تصمیم مشخص شود

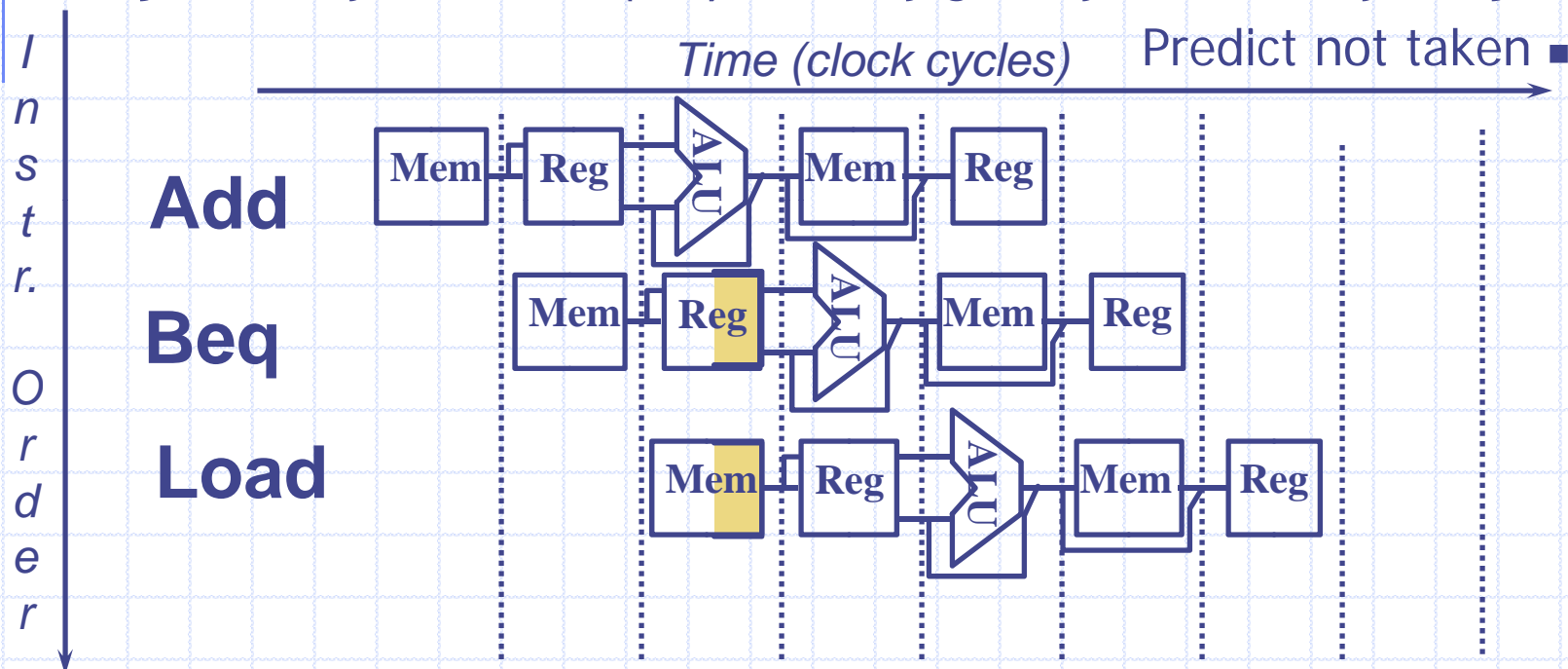
Its possible to move up decision to 2nd stage by ■ adding hardware to check registers as being read



برخورد: هر دستورالعمل انشعاب دو چرخه ساعت زمان می برد <= کند شدن

راه حل های هزار د کنترلی

پیشگویی کردن: یک مسیر حدس زدن بعد تهیه پشتیبان ذ اگر اشتباه بود

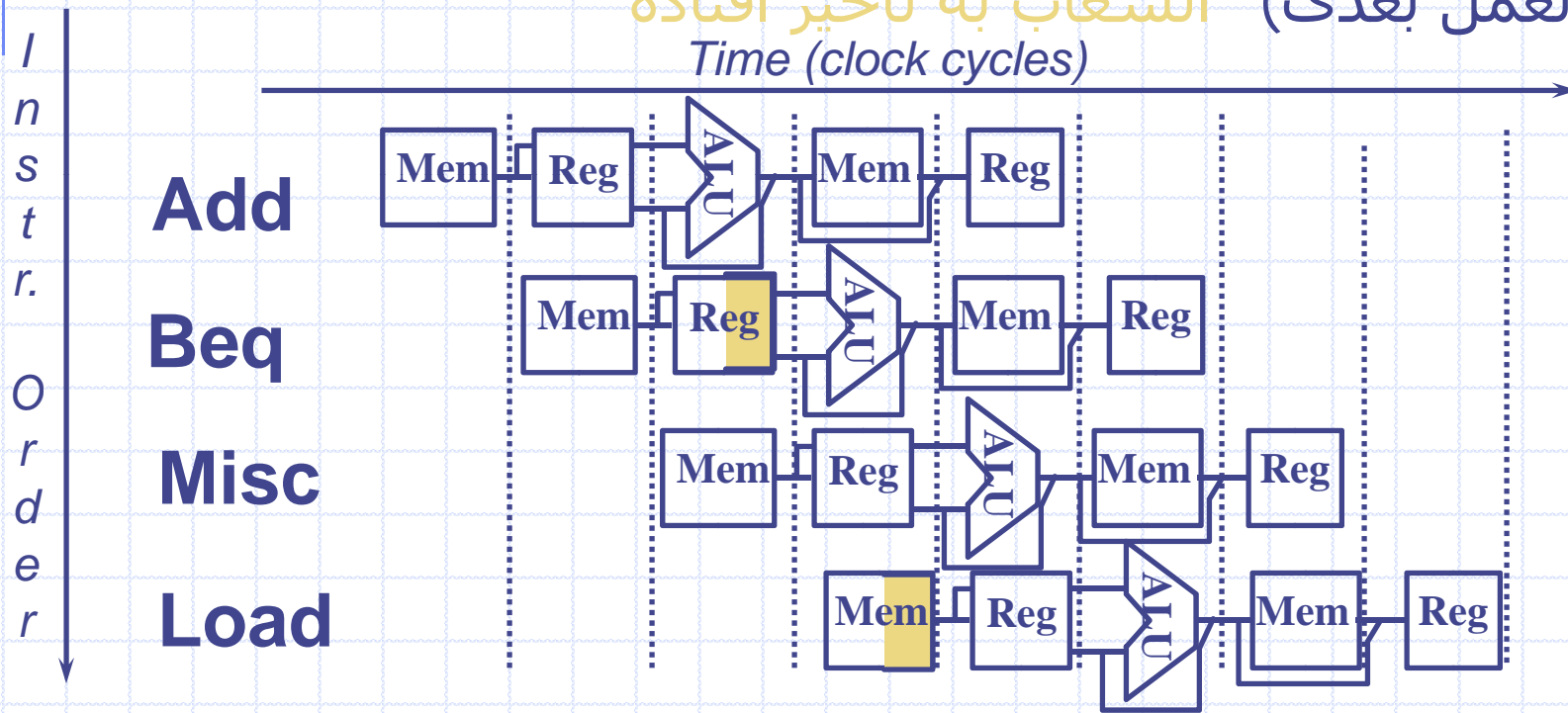


برخورد: هر دستورالعمل انشعاب اگر درست باشد یک چرخه ساعت و اگر اشتباه باشد 2 چرخه زمان می برد (50% موارد درست است)

More dynamic scheme: history of 1 branch (- 90%)

راه حل های هزارد کنترلی

دوباره تعریف کردن رفتار انشعاب (اتفاق می افتد بعد از دستورالعمل بعدی) "انشعاب به تاخیر افتاده"



Impact: 1 clock cycles per branch instruction if can find instruction to put in "slot" (- 50% of time)

Launch more instructions per clock cycle => less useful

هزاردهای داده ای روی r1

مشکل: r1 نمی تواند به وسیله دستورالعمل های دیگر خوانده شود قبل از اینکه به وسیله add نوشته شود

add r1, r2, r3

sub r4, r1, r3

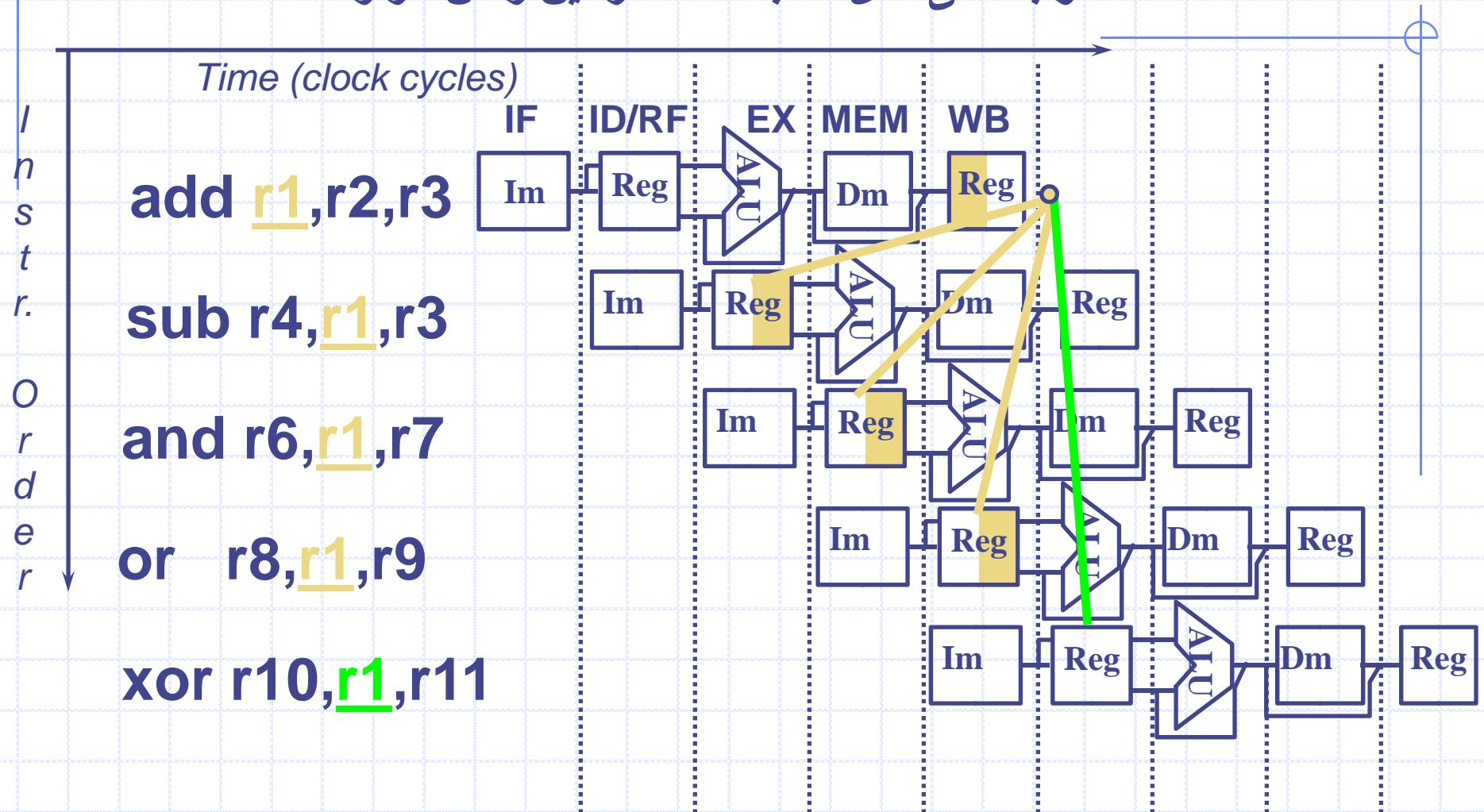
and r6, r1, r7

or r8, r1, r9

xor r10, r1, r11

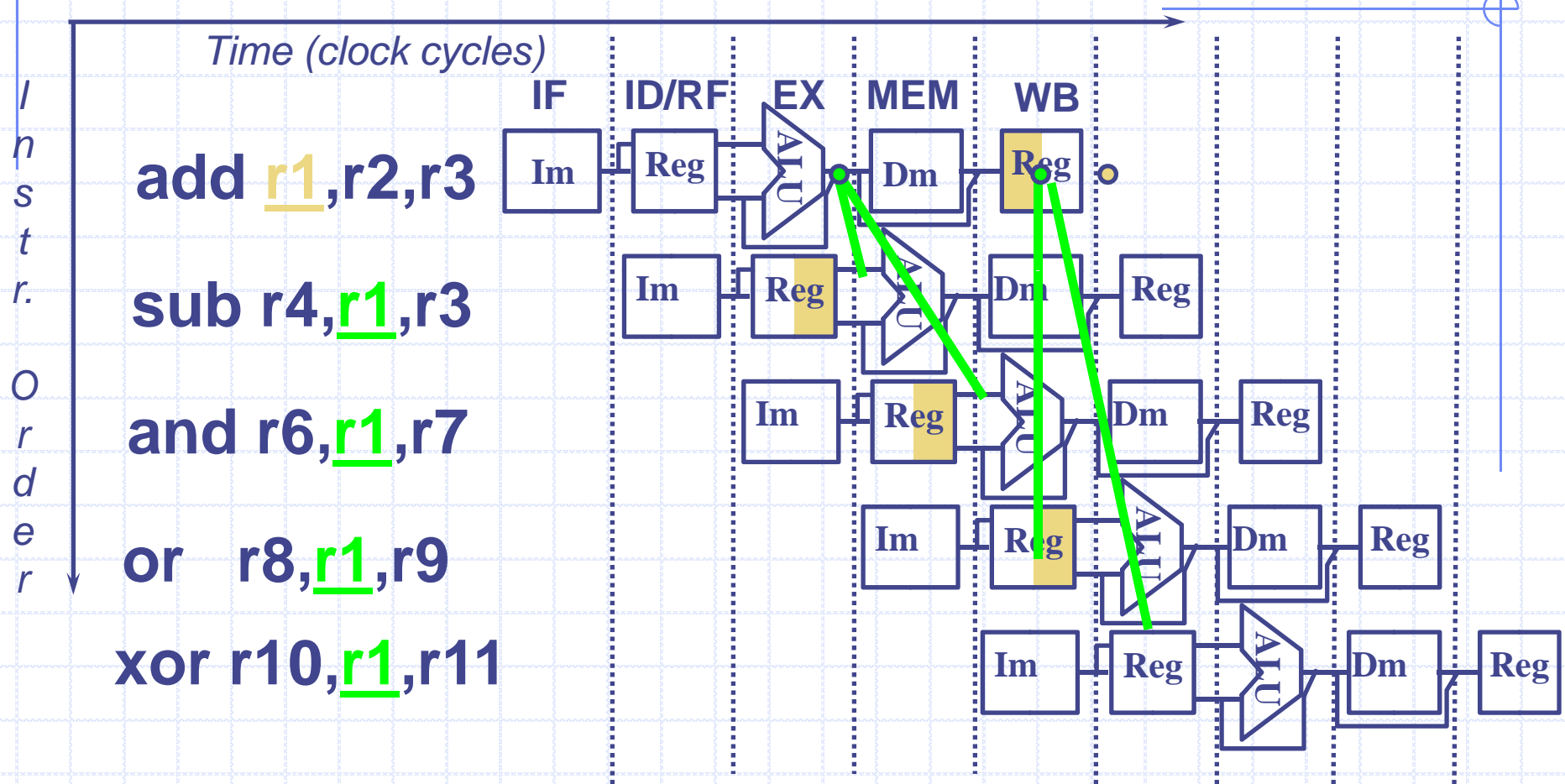
هزارداده ای روی r1:

وابستگی های عقب افتاده در این زمان هزاردها هستند



راه حل هزارد داده ای:

• پیش رو "نتیجه از یک مرحله برای دیگری

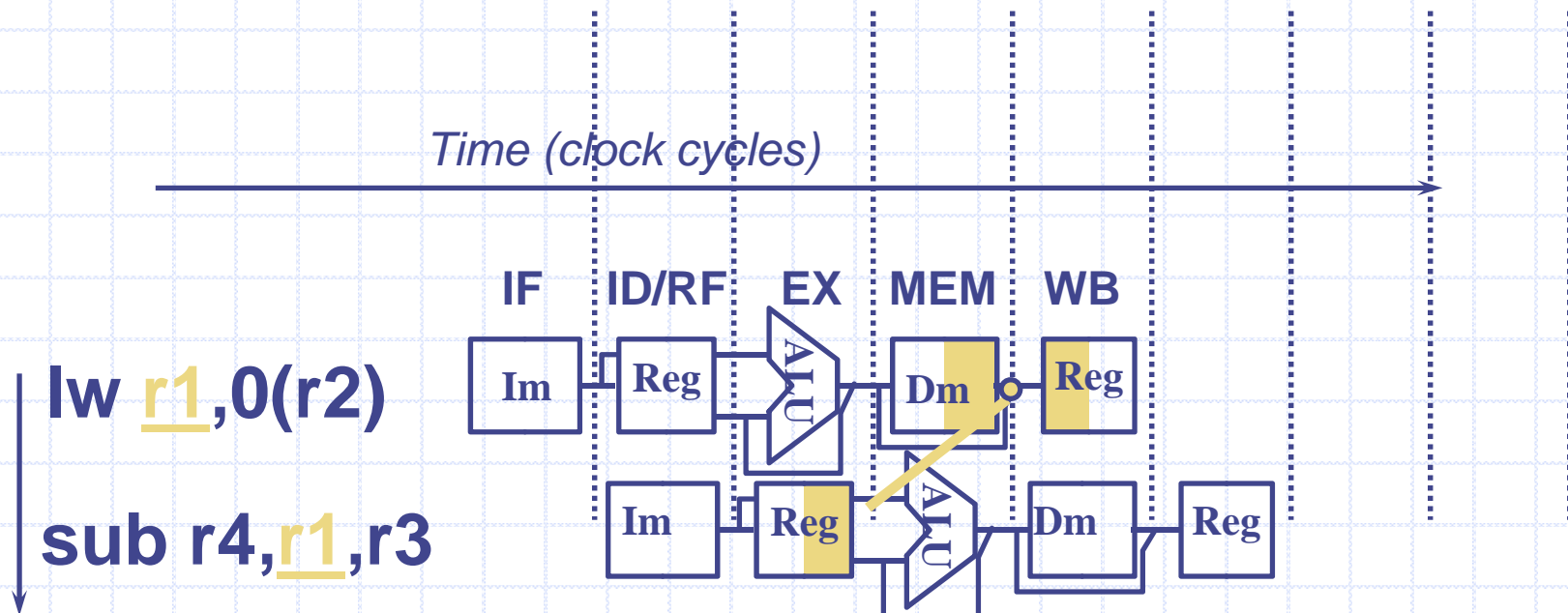


• "or" بسیار خوب اگر که تعریف نوشتن/خواندن درست است

Forwarding (یاگذشتن): در مورد بار کردن چی؟

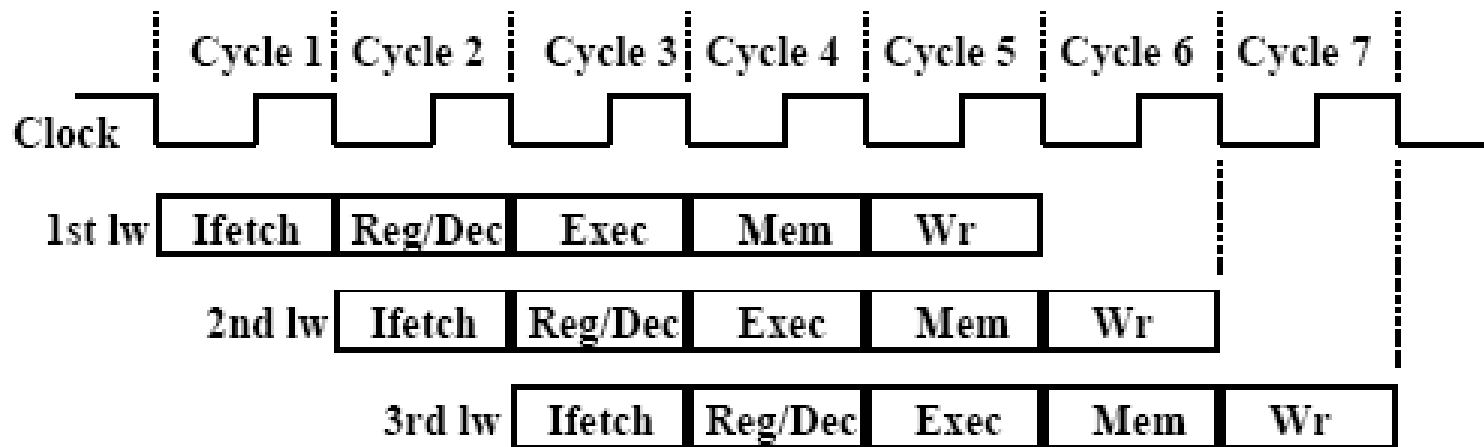
• وابستگی های به جا مانده از قبل هزارد هستند

• نمی توان با ارسال (Forwarding) این مسئله را حل کرد؟
• باید دستورالعمل وابسته به بارگذاری را توقف / تاخیر دهیم



طراحی یک مسیر داده خط لوله ای شده

خط لوله ای کردن دستور العمل بار کردن



• 5 تا واحد تابعی مستقل در مسیر داده خط لوله ای هستند.

• حافظه دستور العمل برای مرحله واکنشی

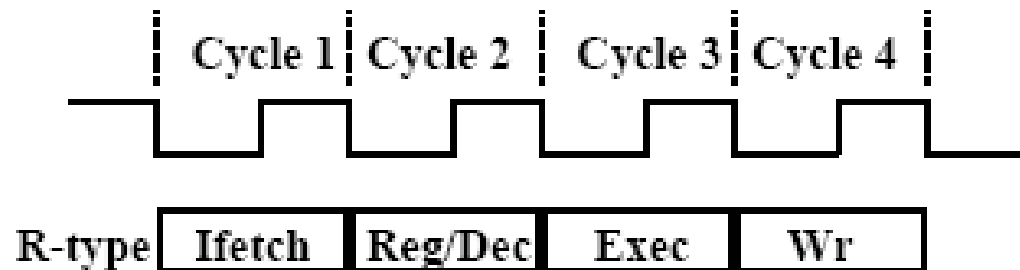
- خواندن از بانک ثباتی برای مرحله Reg/Dec (bus A and busB)

- ALU برای مرحله اجرا

- حافظه داده برای مرحله دسترسی به حافظه

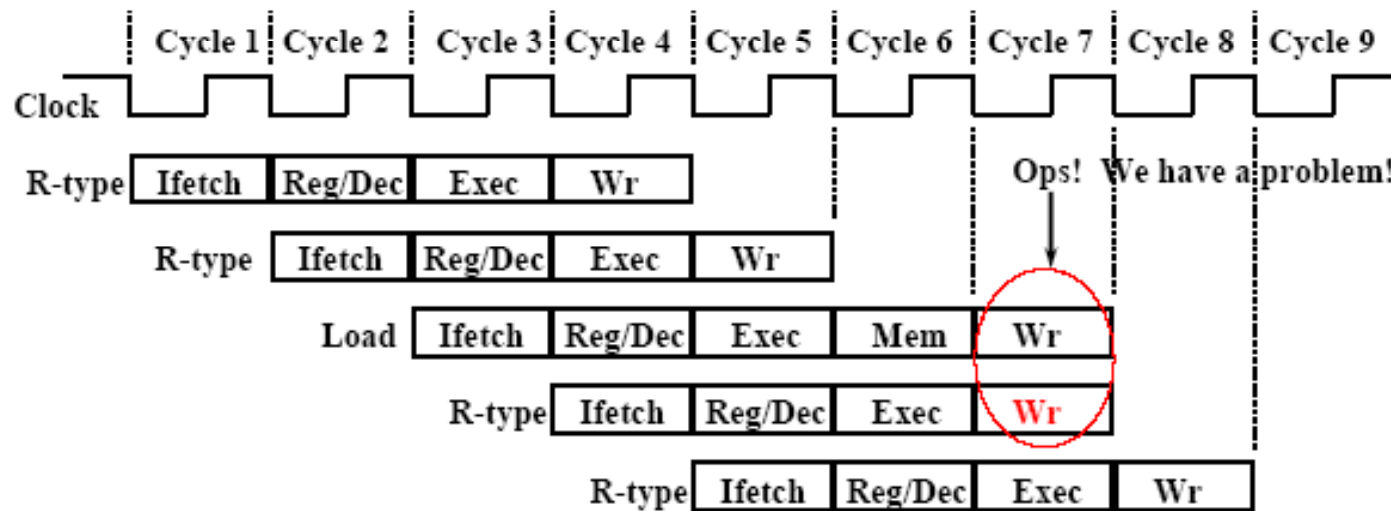
- نوشتن بر بانک ثباتی برای مرحله Wr (bus W)

4 مرحله برای R-type



- واکنشی: واکنشی دستور العمل
- واکنشی دستور العمل از حافظه دستور العمل
- جدید کردن PC
- Reg/Dec: واکنشی ثباتها و کد برداری دستور العمل ها
- Exec: عمل کردن ALU بر روی 2 تا ثبات عملوند
- Wr: باز نویسی خروجی ALU بر روی بانک ثباتی

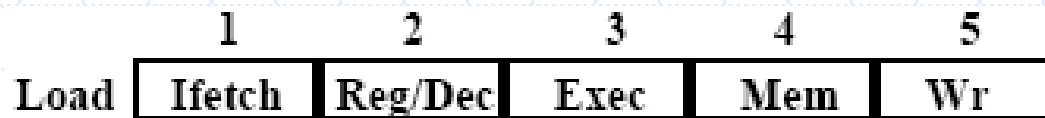
خط لوله ای کردن دستورالعمل های R-type و Load



ما در خط لوله ای تضاد داریم یا هزارد های ساختاری
- دو دستورالعمل سعی در نوشتن در بانک ثبتی به صورت همزمان دارند!
- تنها نوشتن یک پورت

مشاهدات مهم

هر واحد کاری تنها یک بار برای هر دستورالعمل می تواند استفاده شود
- هر واحد کاری باید در همان مرحله برای همه دستورالعمل ها استفاده شود:
- در دستورالعمل Load نوشتن در بانک ثباتی در مرحله پنجم انجام می شود

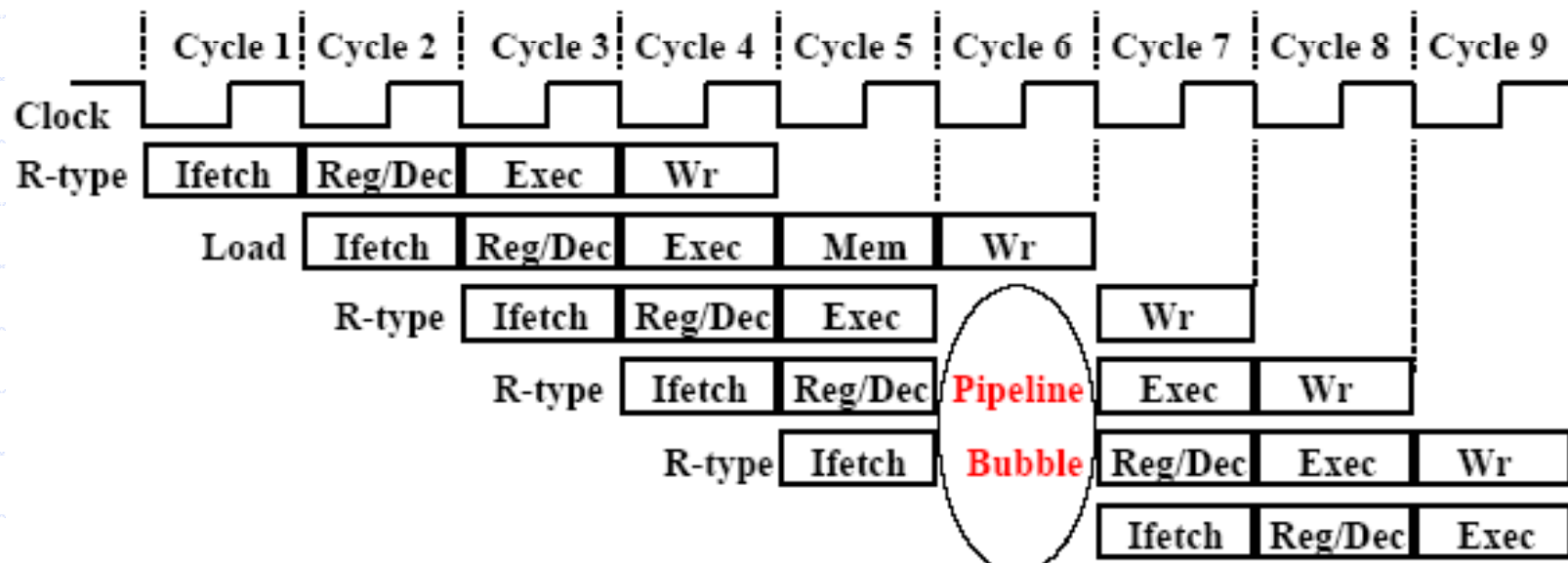


- در دستورالعمل R-type نوشتن در بانک ثباتی در مرحله چهارم انجام می شود



2 راه حل برای هزار خط لوله ای وجود دارد

راه حل اول: وارد کردن "حباب" در خط لوله



• وارد کردن حباب در خط لوله برای مانع شدن از دو نوشتن همزمان در یک چرخه

- منطق کنترل می تواند پیچیده شود.

- از دست دادن واکنشی دستورالعمل و فرصت انتشار.

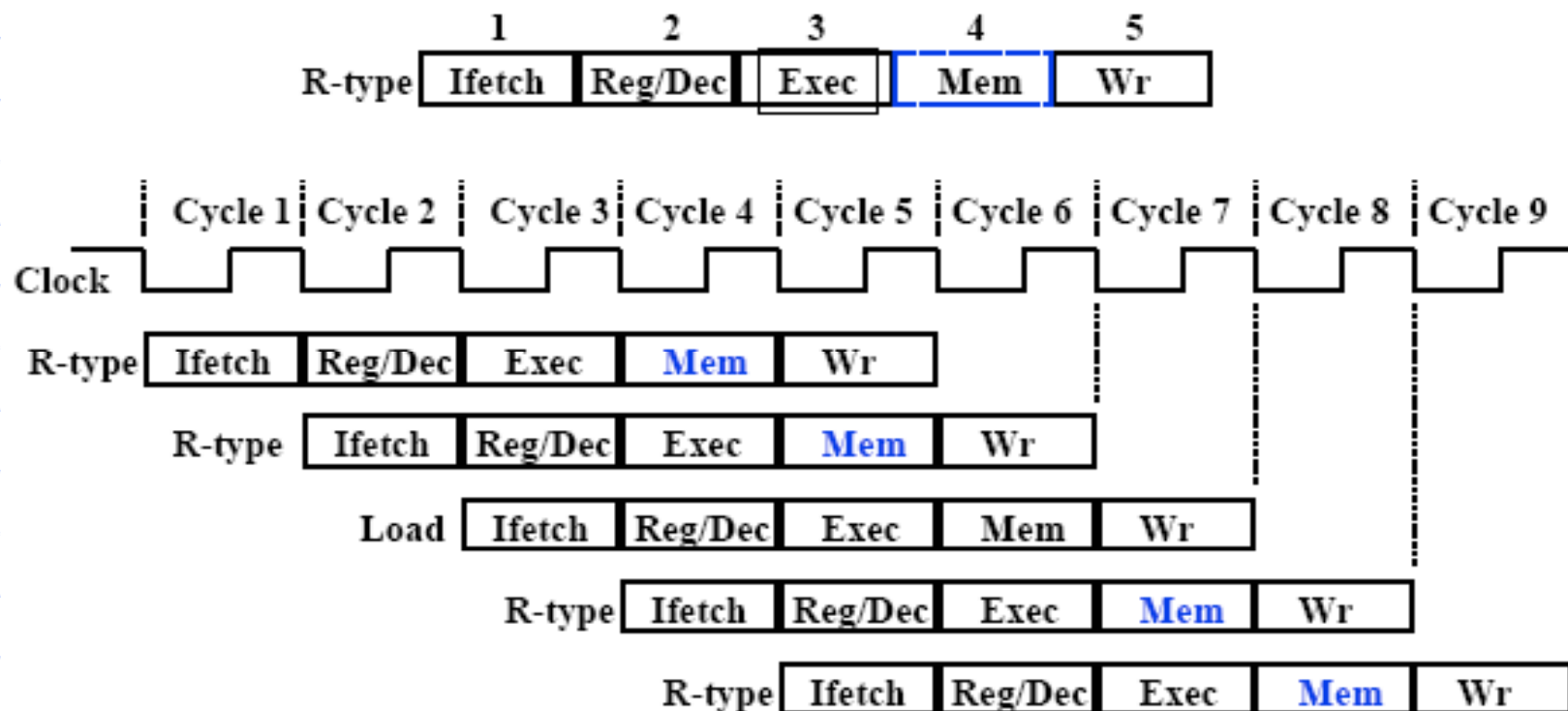
• هیچ دستورالعملی در چرخه 6 شروع نمی شود

راه حل دوم : تأخیر نوشتن R-type به وسیله یک چرخه

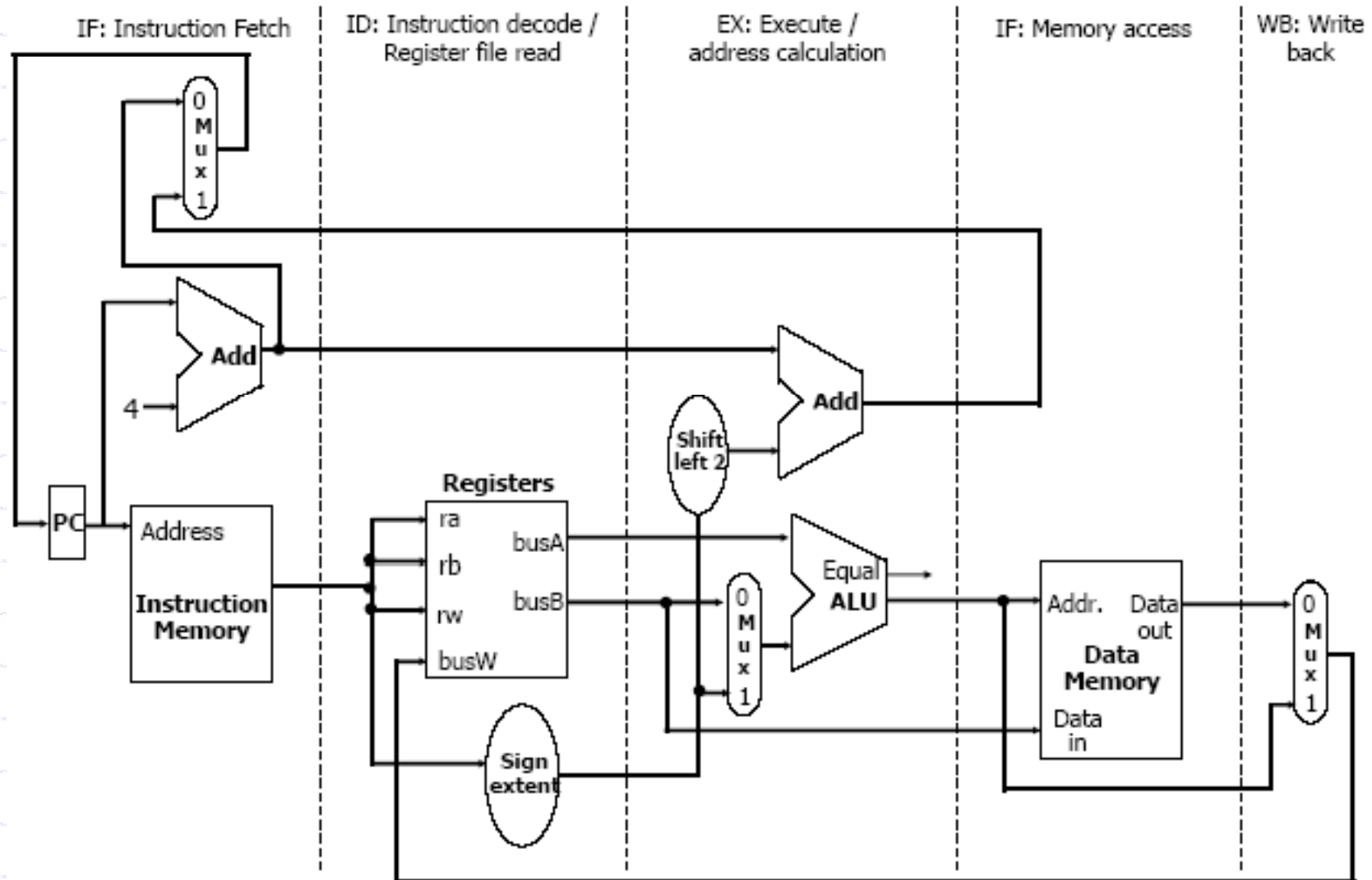
◆ تأخیر نوشتن R-type به وسیله یک چرخه

◆ - حال می توان نوشتن دستورالعمل های R-type را در مرحله پنجم انجام داد

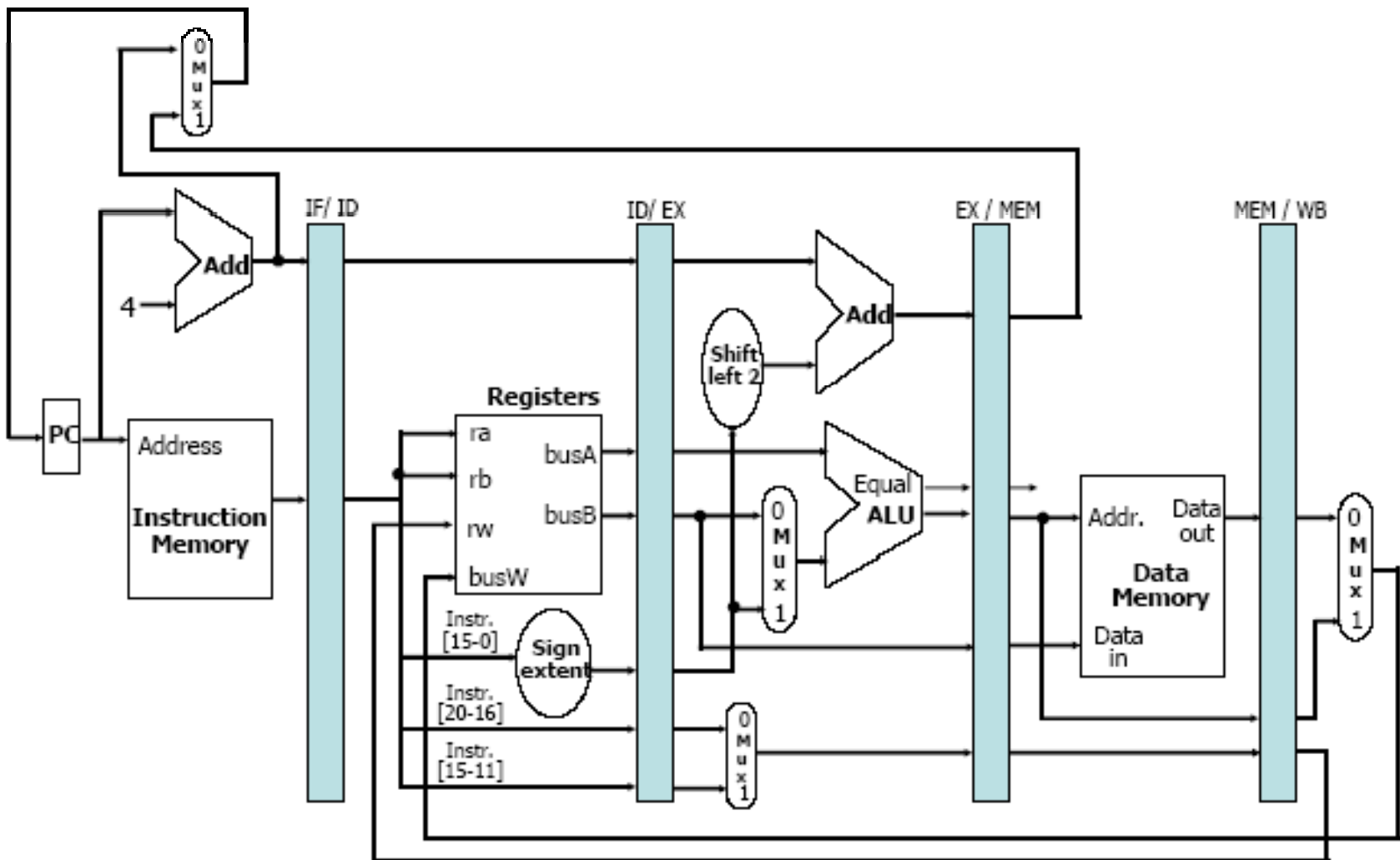
- مرحله دستیابی به حافظه یک مرحله بلااستفاده است: هیچ چیزی در آنجا انجام نمی شود.



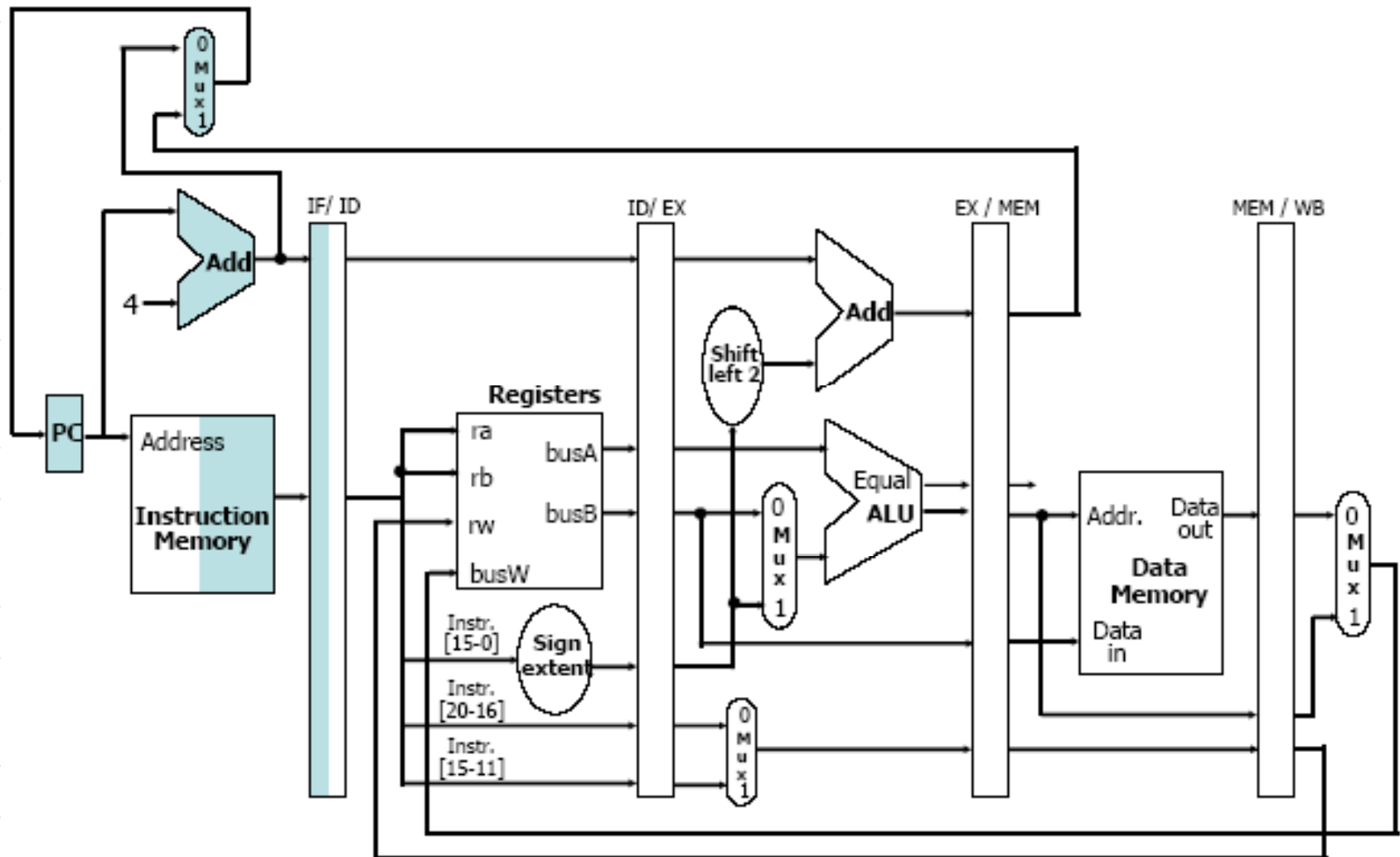
مسیر داده تک چرخه ای



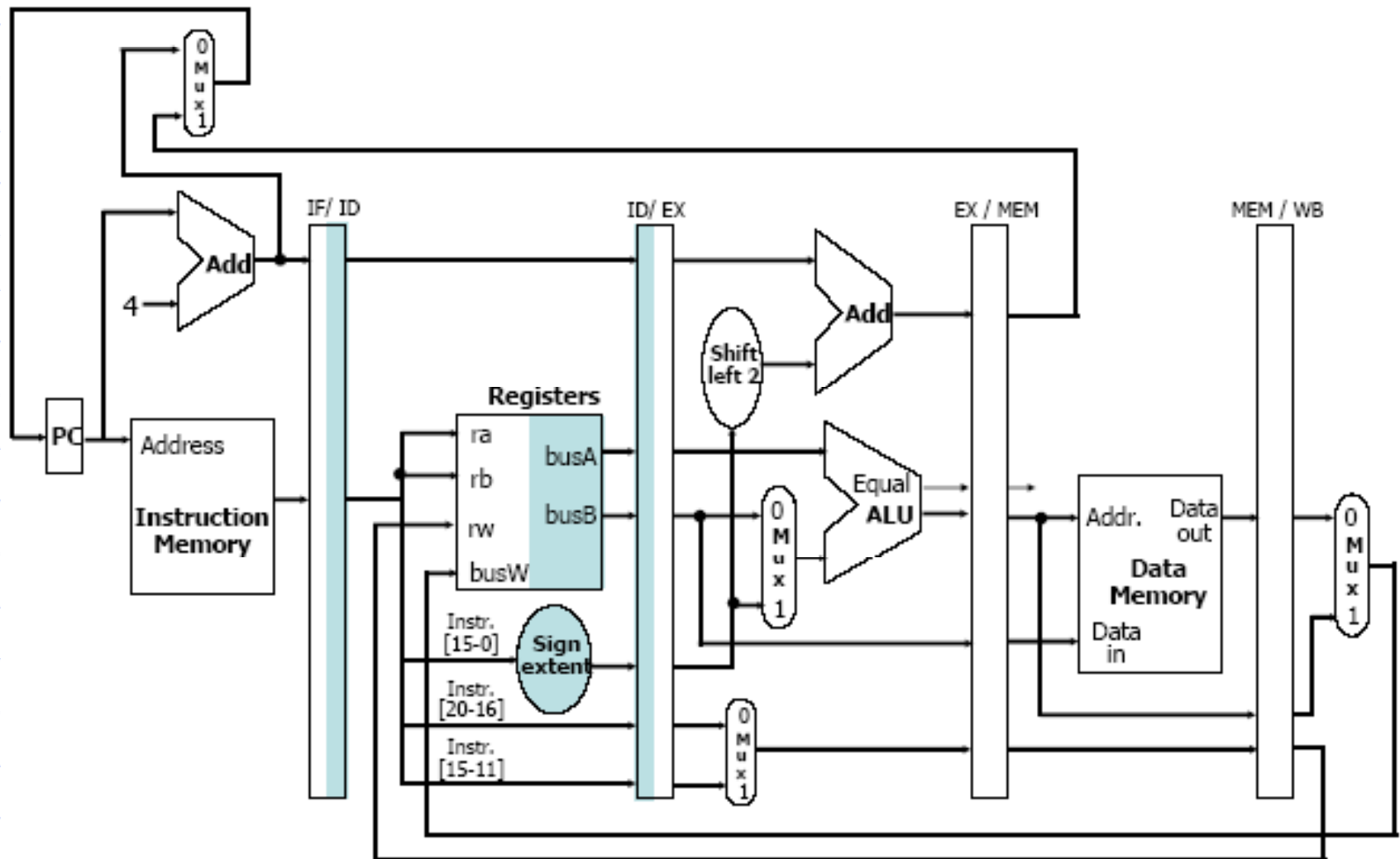
نسخه خط لوله ای شده مسیر داده



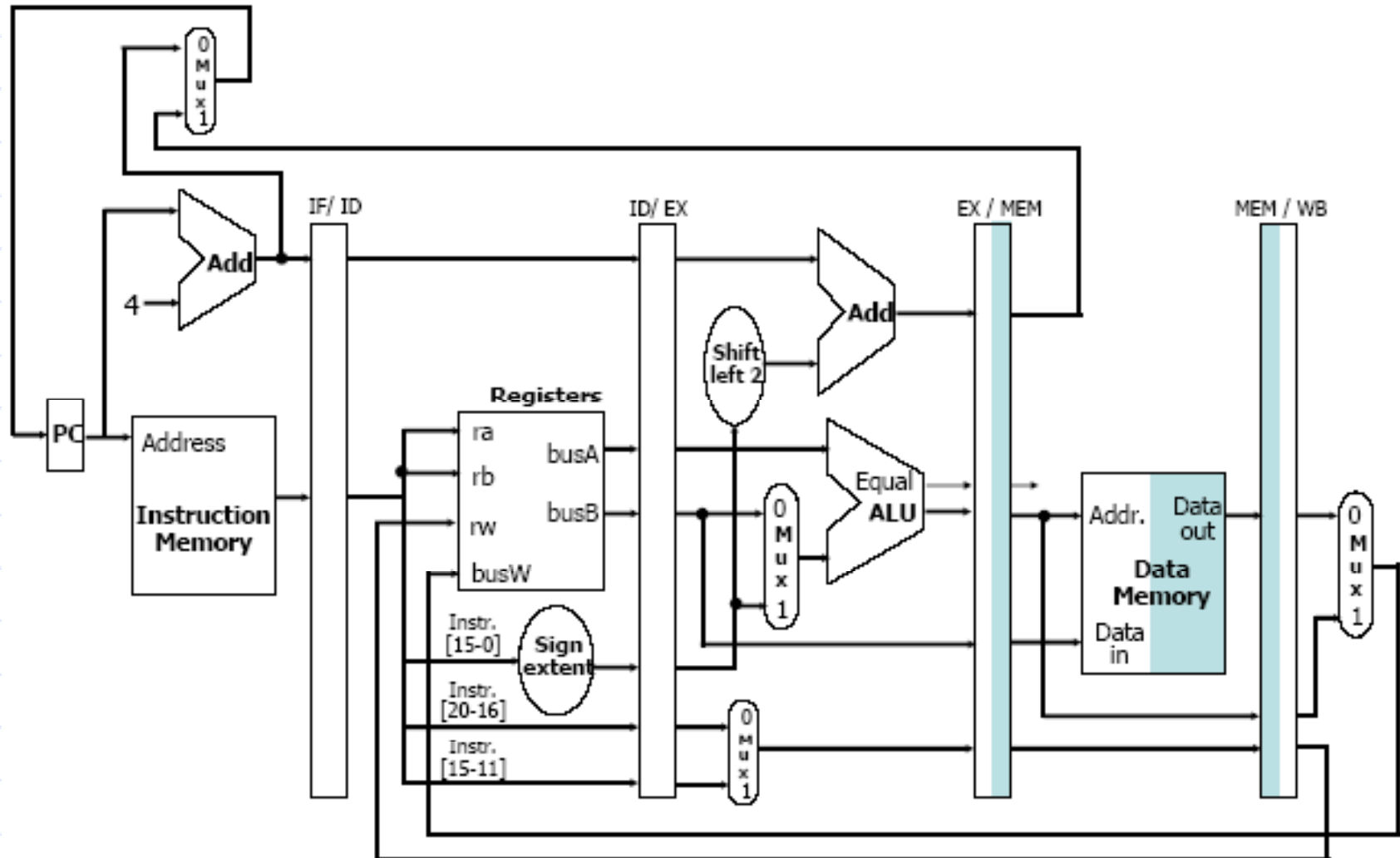
IF: اولین مرحله لوله دستور العمل بارکردن



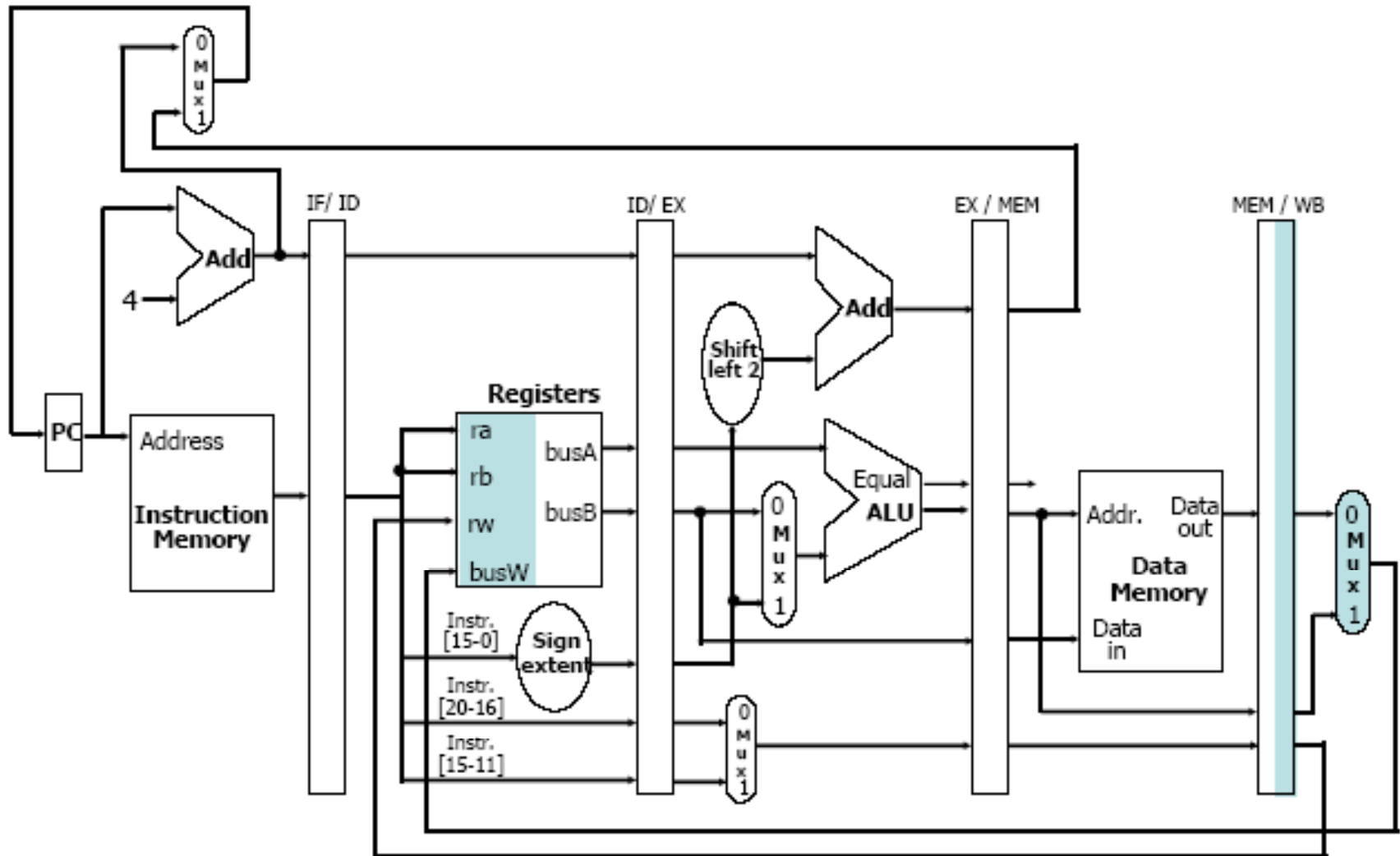
ID: دومین مرحله لوله دستورالعمل بارکردن



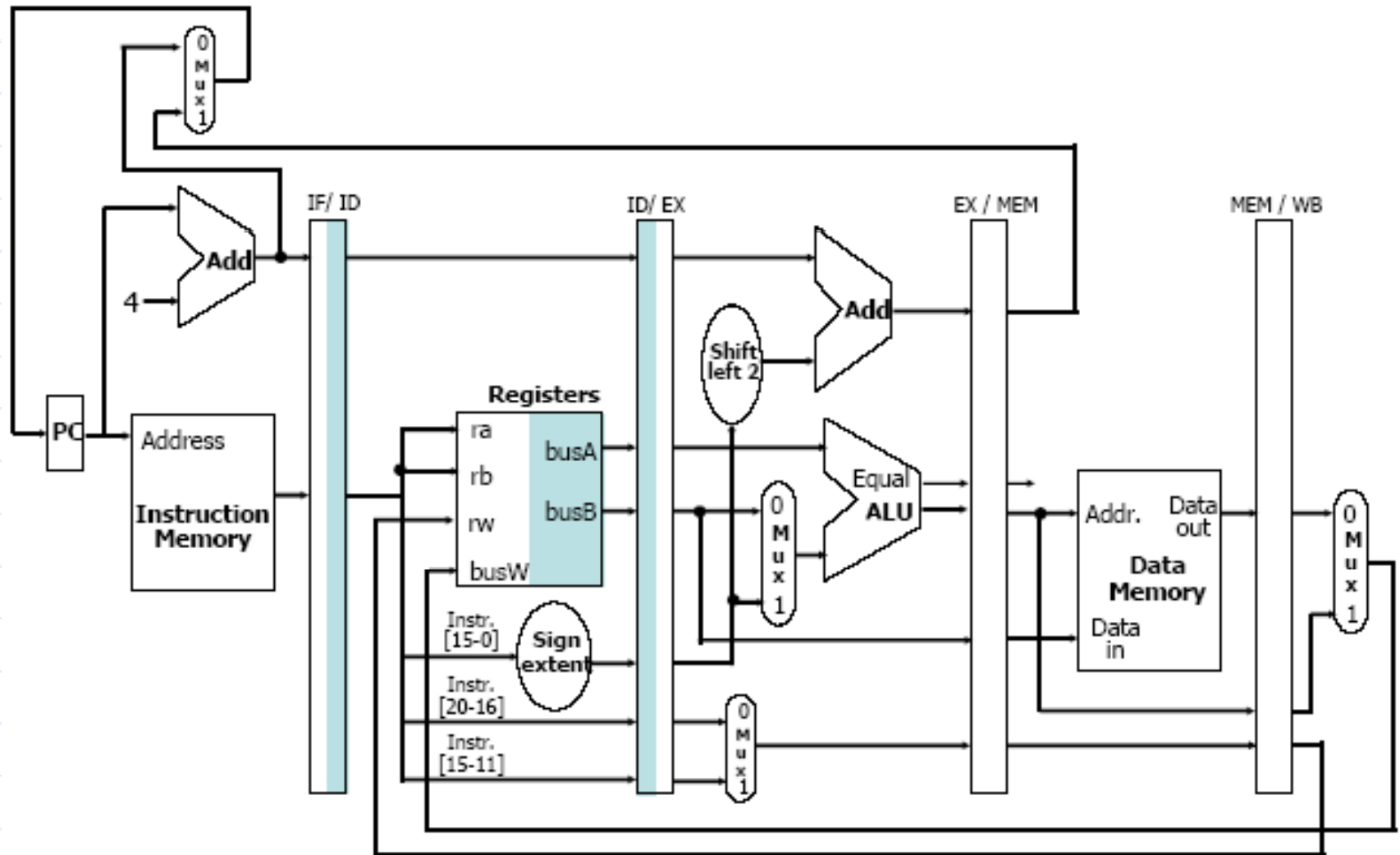
MEM: چهارمین مرحله لوله دستورالعمل بارکردن



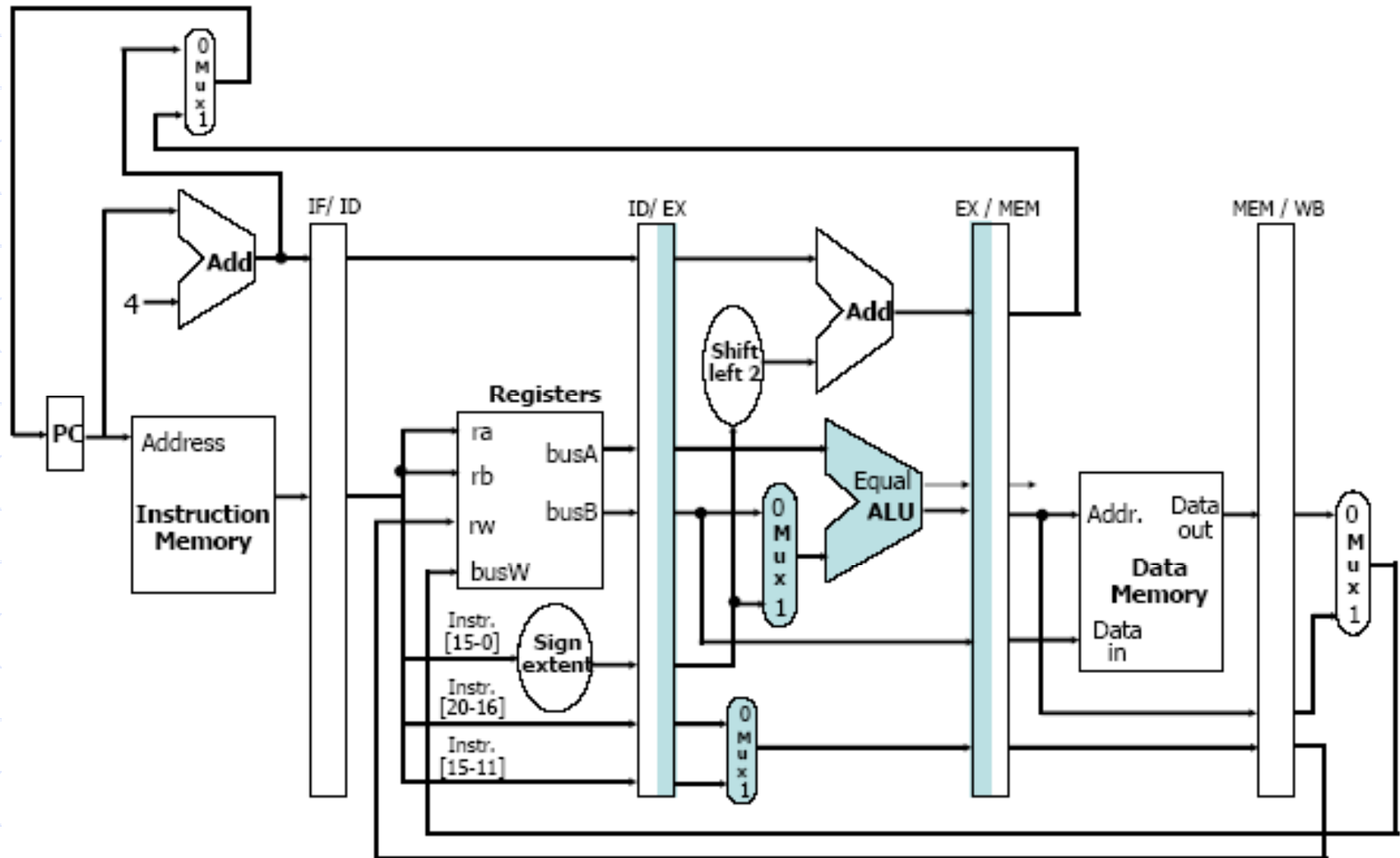
WB: پنجمین مرحله لوله دستورالعمل بارکردن



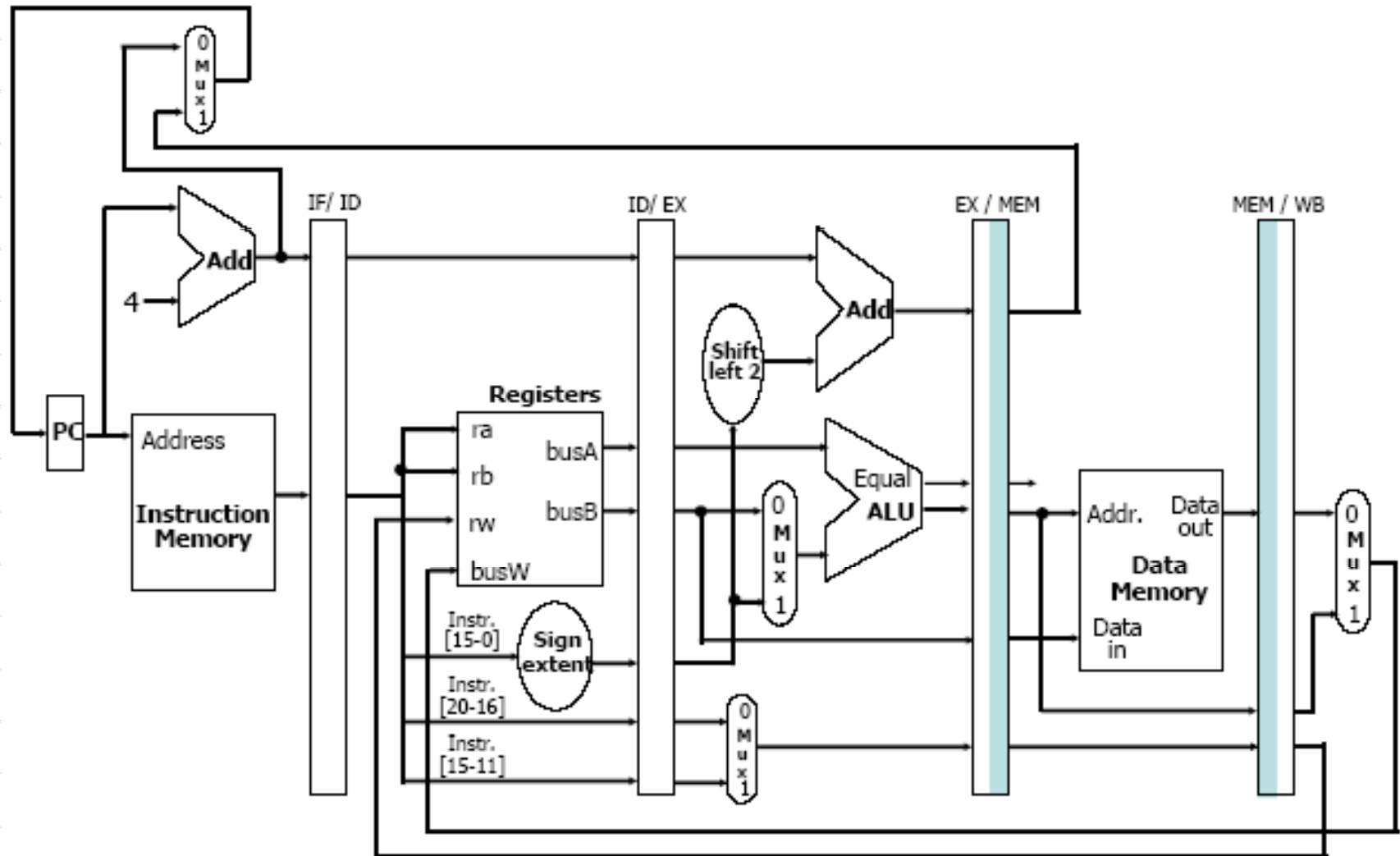
ID: دومین مرحله لوله دستورالعمل R-type



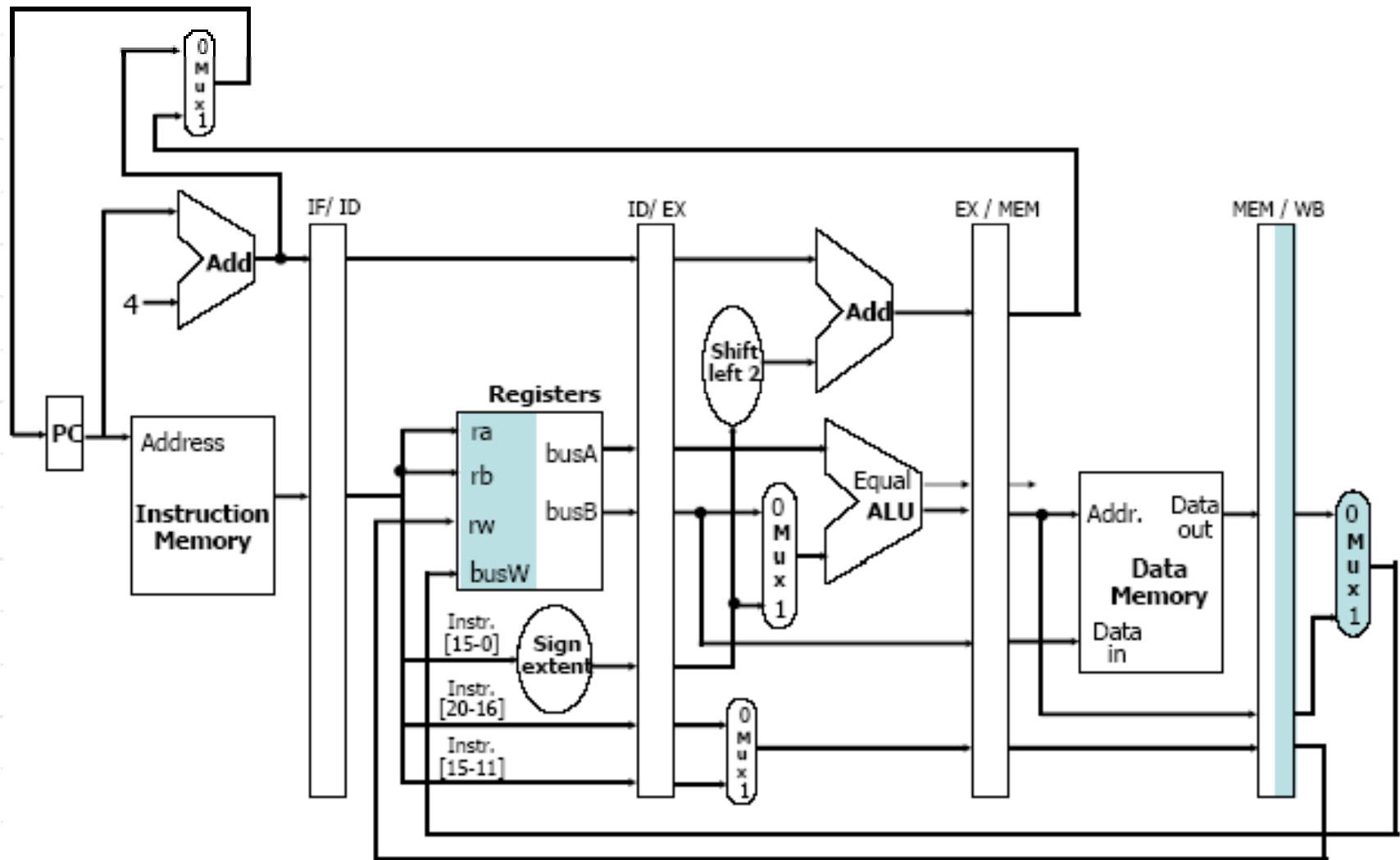
EX: سومین مرحله لوله دستورالعمل R-type



MEM: چهارمین مرحله لوله دستورالعمل R-type



WB: پنجمین مرحله لوله دستورالعمل R-type

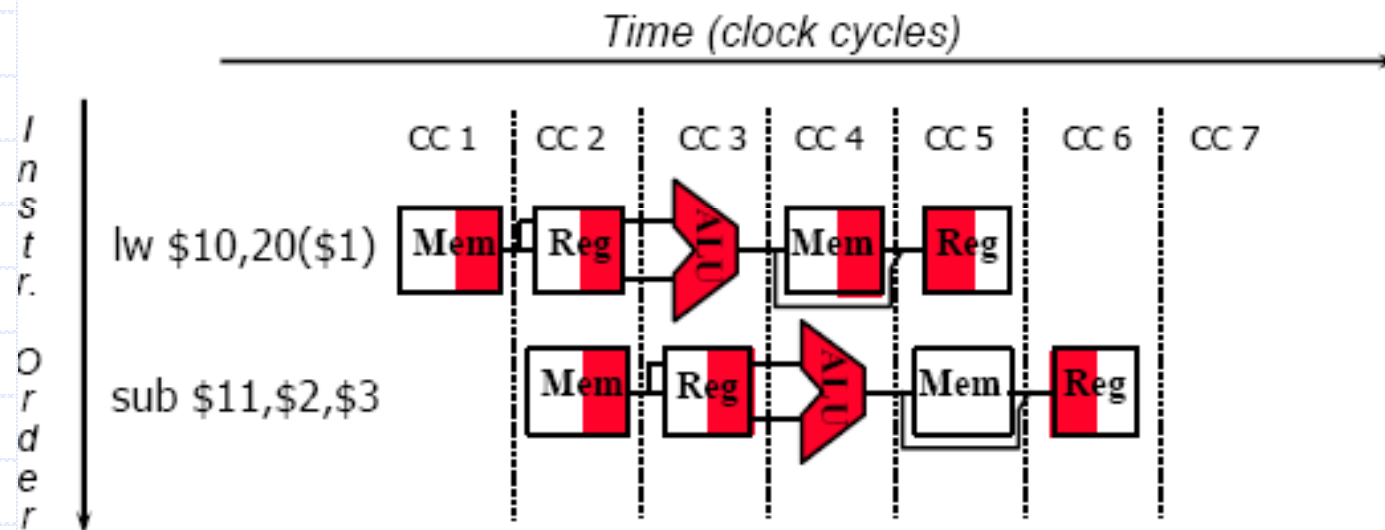


یک مثال برای روشن شدن خط لوله

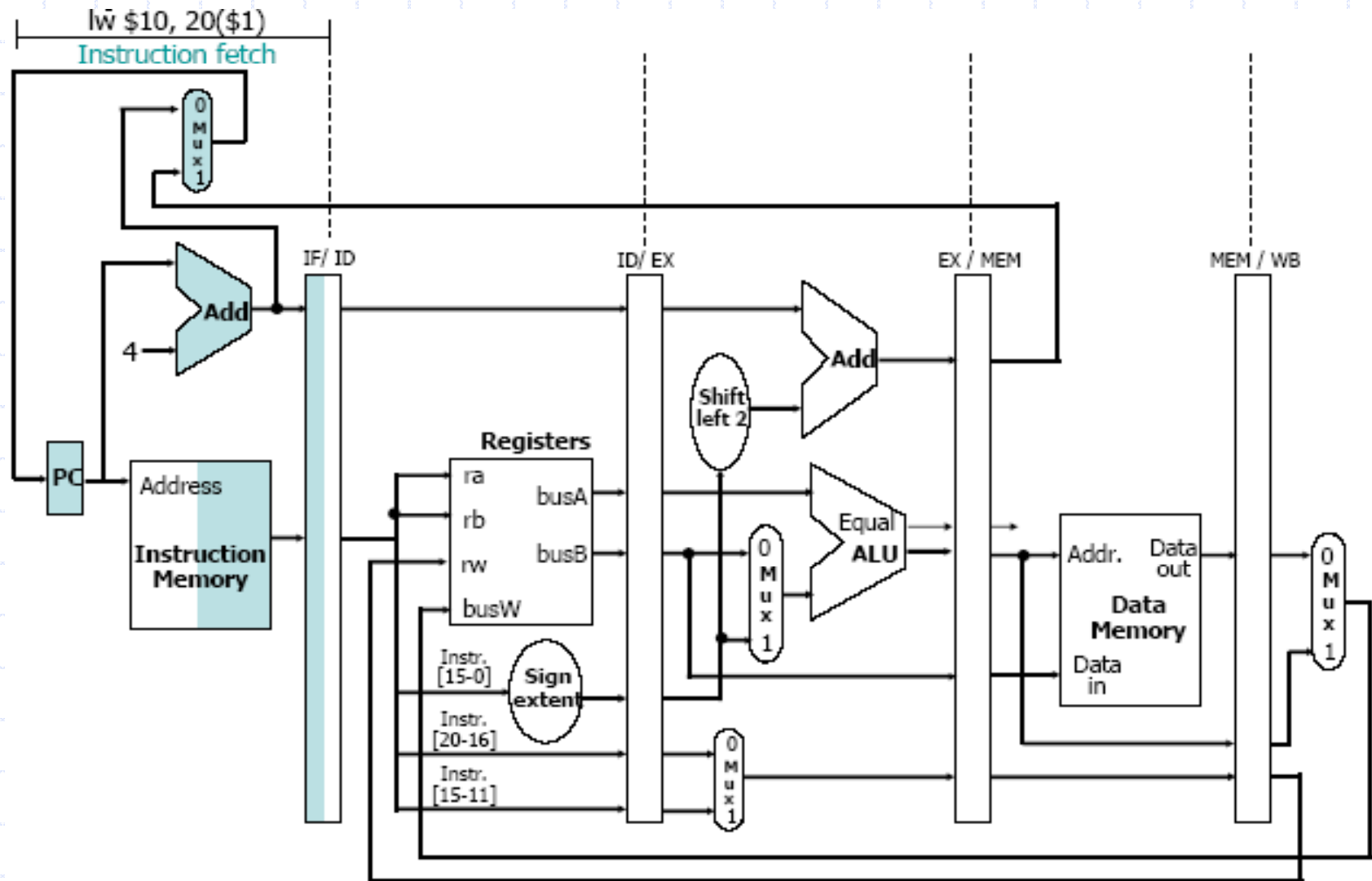
- از زمانی که بیشتر دستورالعمل ها باهم در یک مسیر داده تک چرخه ای اجرا می شوند، این می تواند برای درک کردن مشکل باشد
- دنبال کردن کد زیر یک آزمون خواهد بود.

`lw $10, 20($1)`

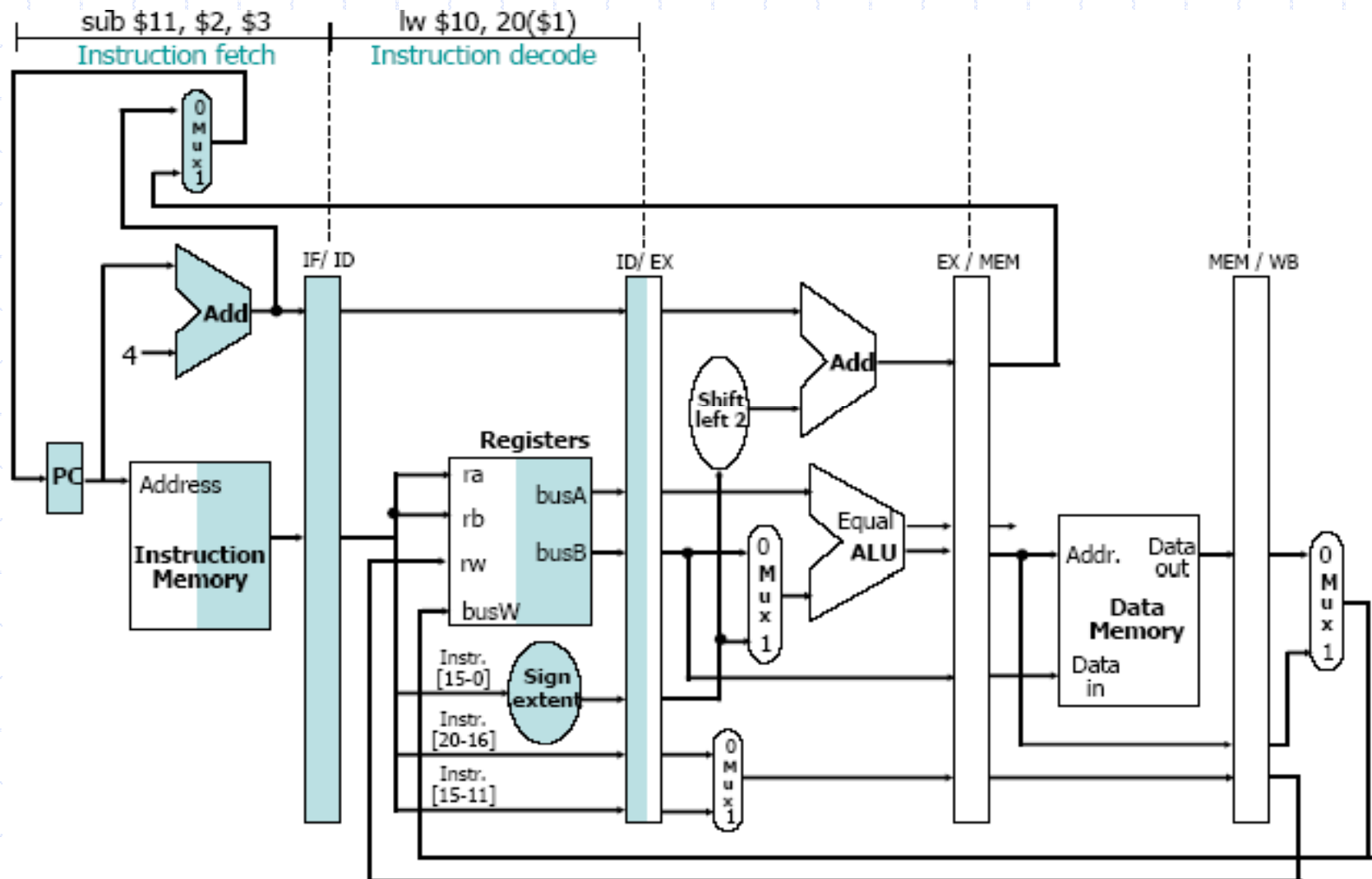
`sub $11, $2, $3`



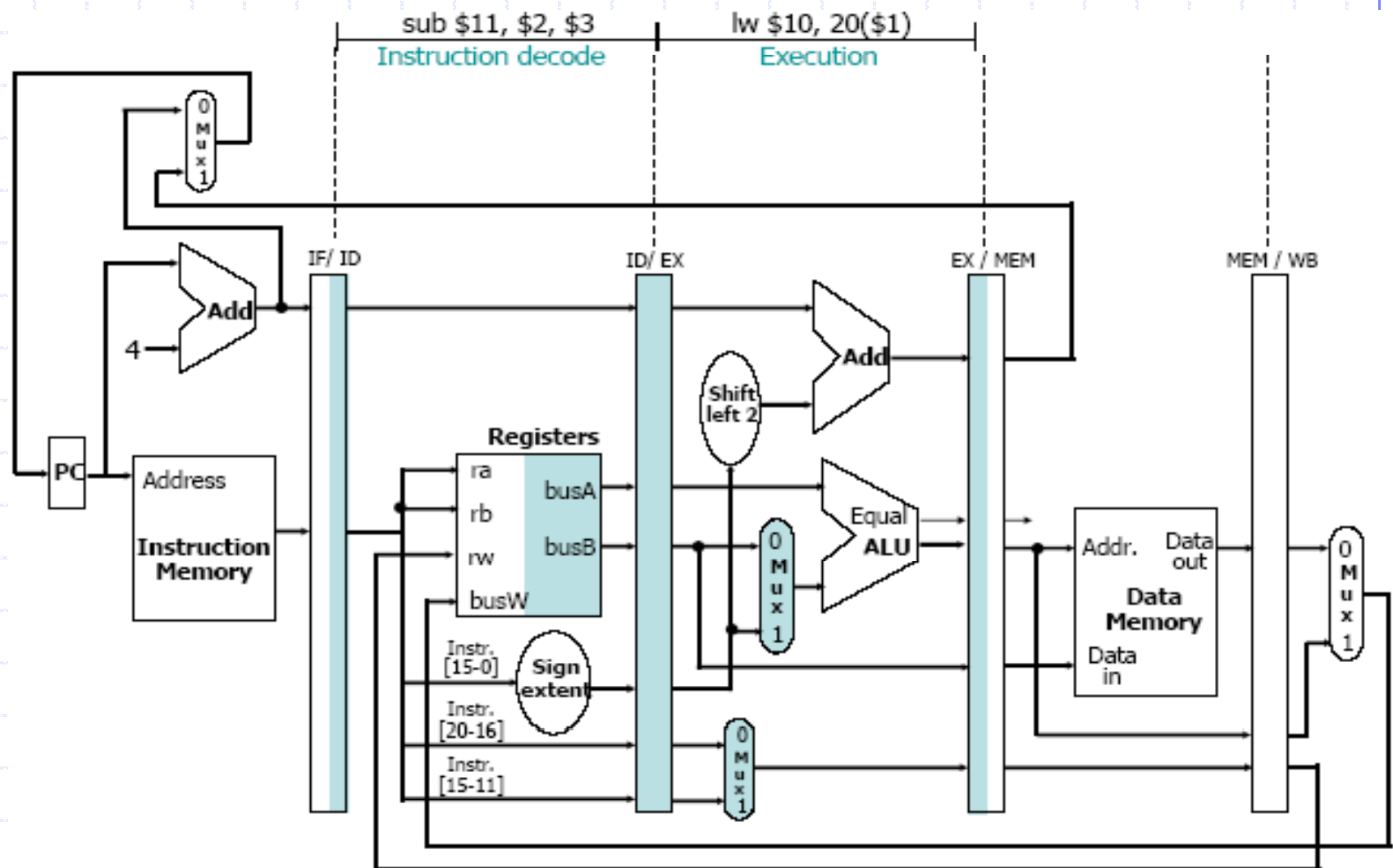
Clock 1



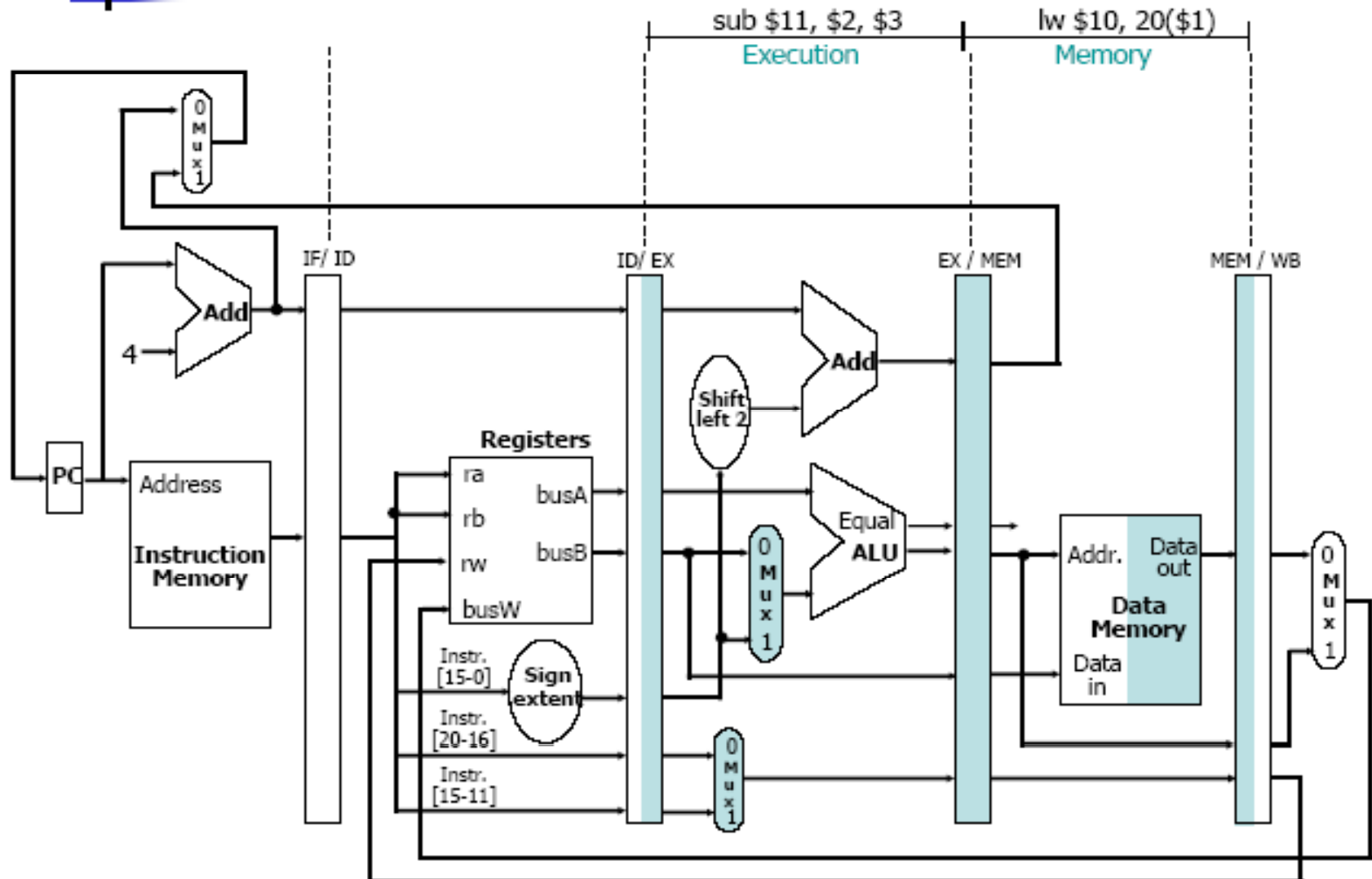
Clock 2



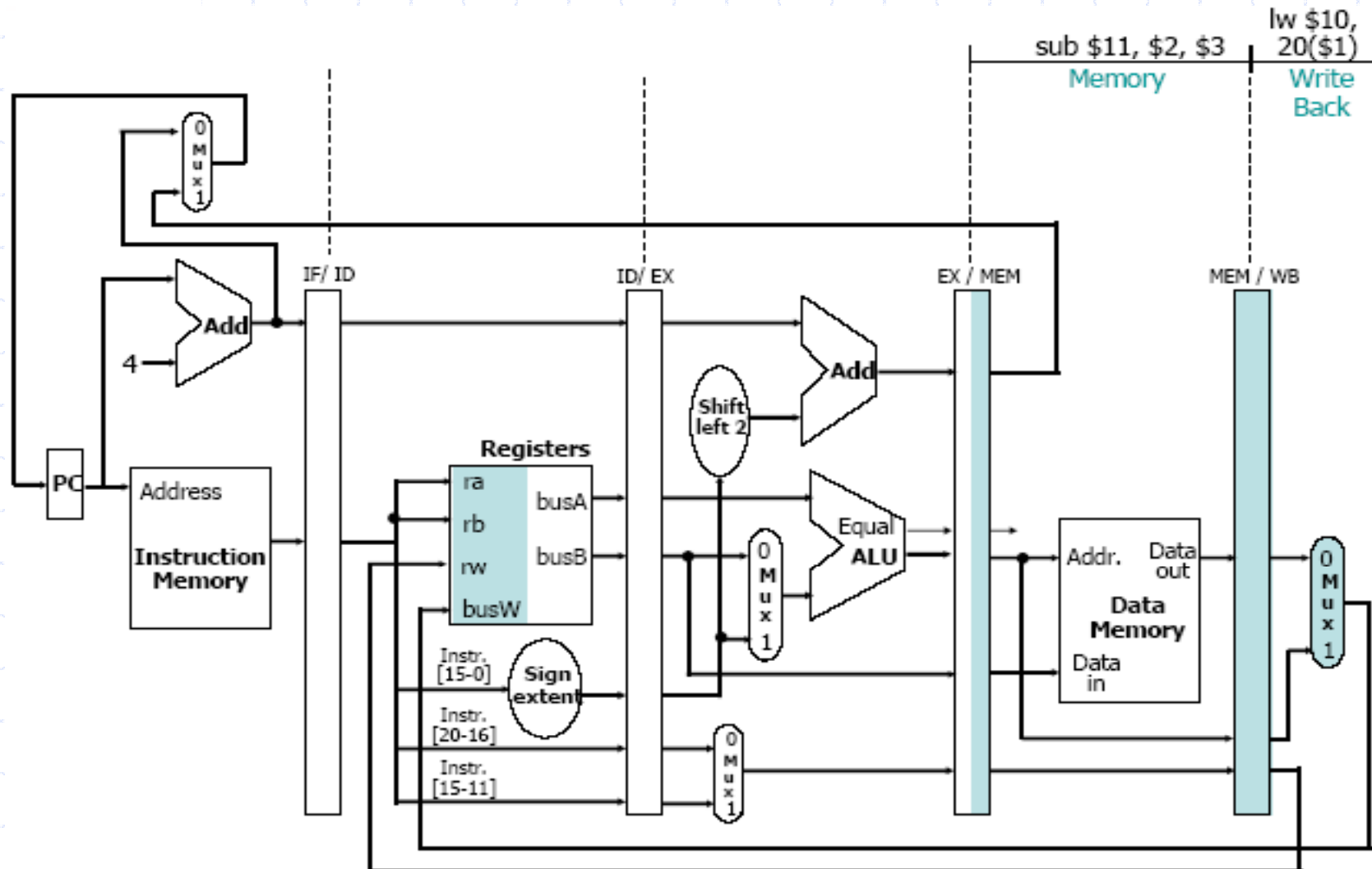
Clock 3



Clock 4



Clock 5

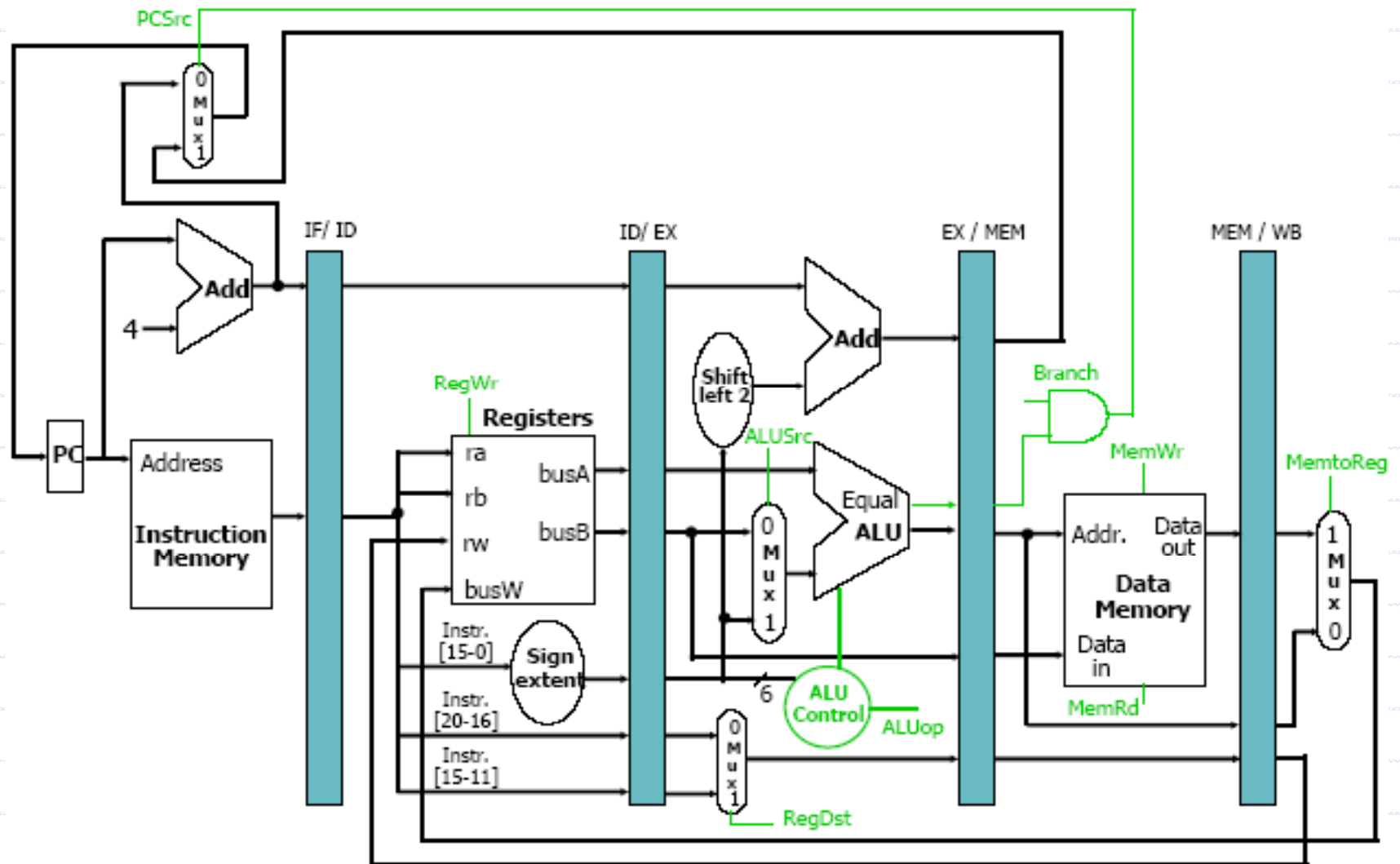


خلاصه: خط لوله ای کردن

- چه چیزی آن را آسان می سازد؟
همه دستورالعمل هایی که طول یکسان دارند فقط یک تعداد کمی از قالب های دستورالعمل عملوندهای حافظه تنها در بارکردن و ذخیره سازی ظاهر می شوند
- چه چیزی آن را سخت می سازد؟
هزاردهای سختاری: فرض کنید ما فقط یک حافظه داریم
هزاردهای کنترلی: ما باید نگران دستورالعمل های انشعاب باشیم
هزاردهای داده ای: یک دستورالعمل وابسته به دستورالعمل های قبلی است
- خط لوله ای کردن یک مفهوم بنیادی است
چندین مرحله استفاده می شود برای جدا کردن مراحل
- پردازشگرهای جدید واقعاً آن را سخت می کنند
استثنا گردانی
سعی کردن برای افزایش کارایی با یک اجرای نادرست و...

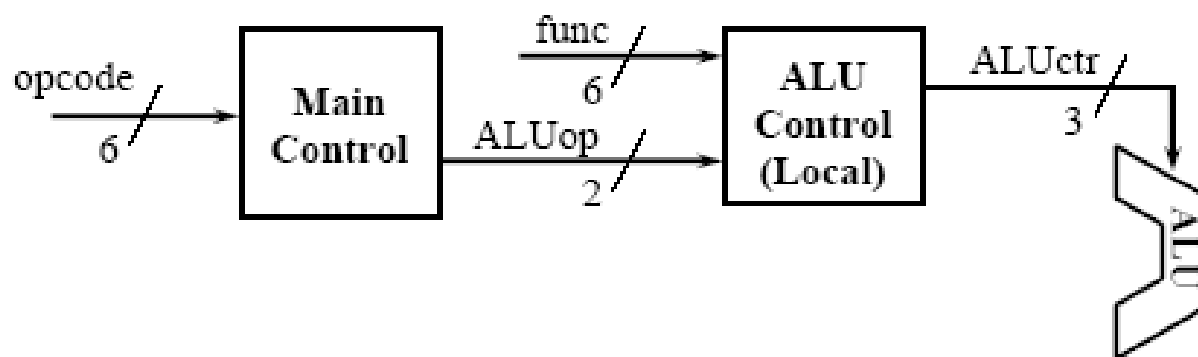
کنترل مسیر داده خط لوله ای شده

مسیر داده خط لوله ای شده با سیگنال های کنترلی



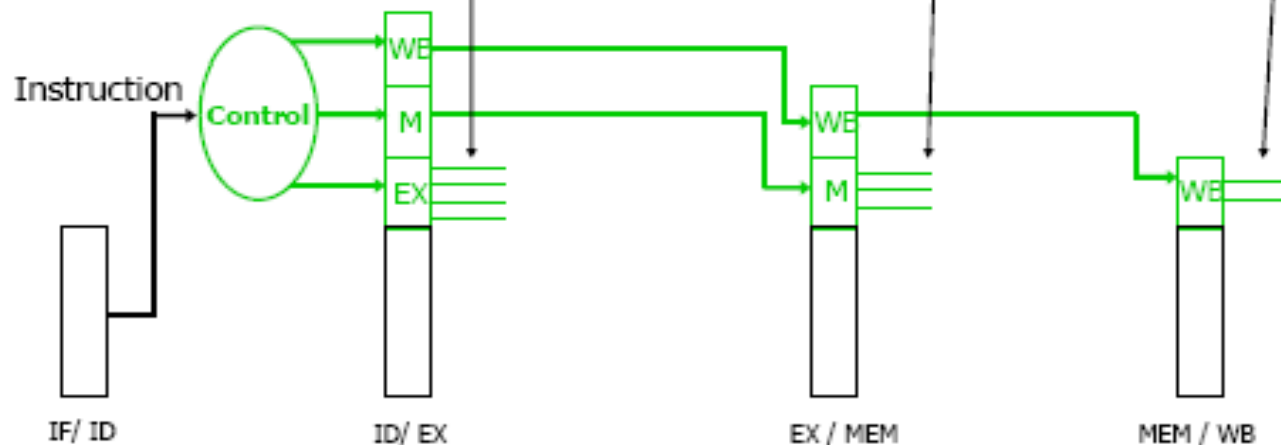
فراخوانی: بیت‌های کنترلی ALU

Instruction opcode	ALUop	Instruction operation	Function field	Desired ALU action	ALU control
LW	00	load word	XXXXXX	add	010
SW	00	store word	XXXXXX	add	010
beq	01	branch equal	XXXXXX	subtract	110
R-type	10	add	100000	add	010
R-type	10	subtract	100010	subtract	110
R-type	10	AND	100100	and	000
R-type	10	OR	100101	or	001
R-type	10	set on less than	101010	set on less than	111

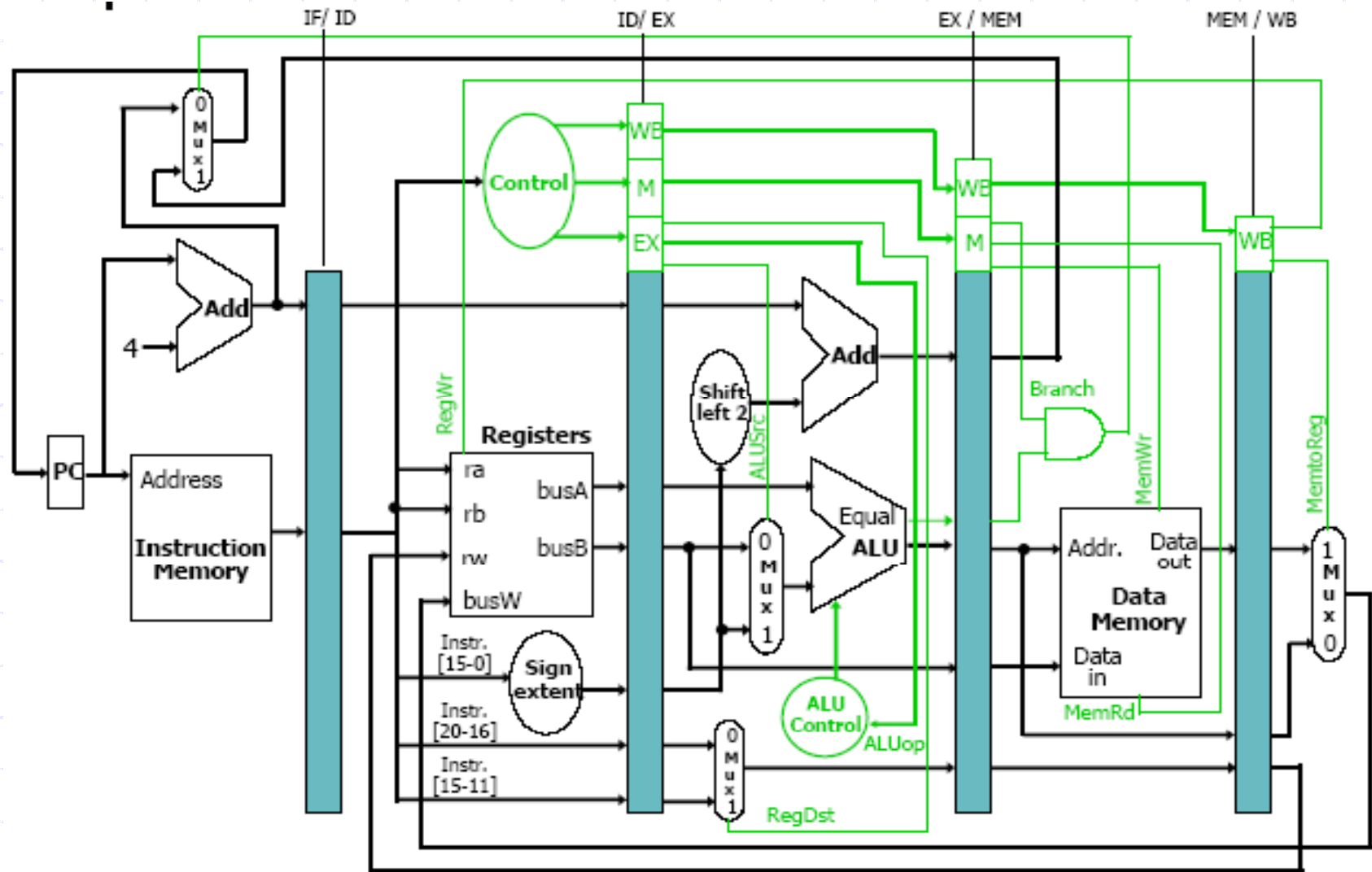


مقادیر خطوط کنترلی برای 3 مرحله آخر خط لوله

Instructions	Execution / Address Calculation stage control lines				Memory Access stage control lines			Write Back stage control lines	
	Reg Dst	ALU Op0	ALU Op1	ALU Src	Branch	Mem Rd	Mem Wr	Reg Wr	Mem to Reg
R-type	1	1	0	0	0	0	0	1	0
lw	0	0	0	1	0	1	0	1	1
sw	X	0	0	1	0	0	1	0	X
Beq	X	0	1	0	1	0	0	0	X



مسیر داده خط لوله ای شده با سیگنال های کنترلی



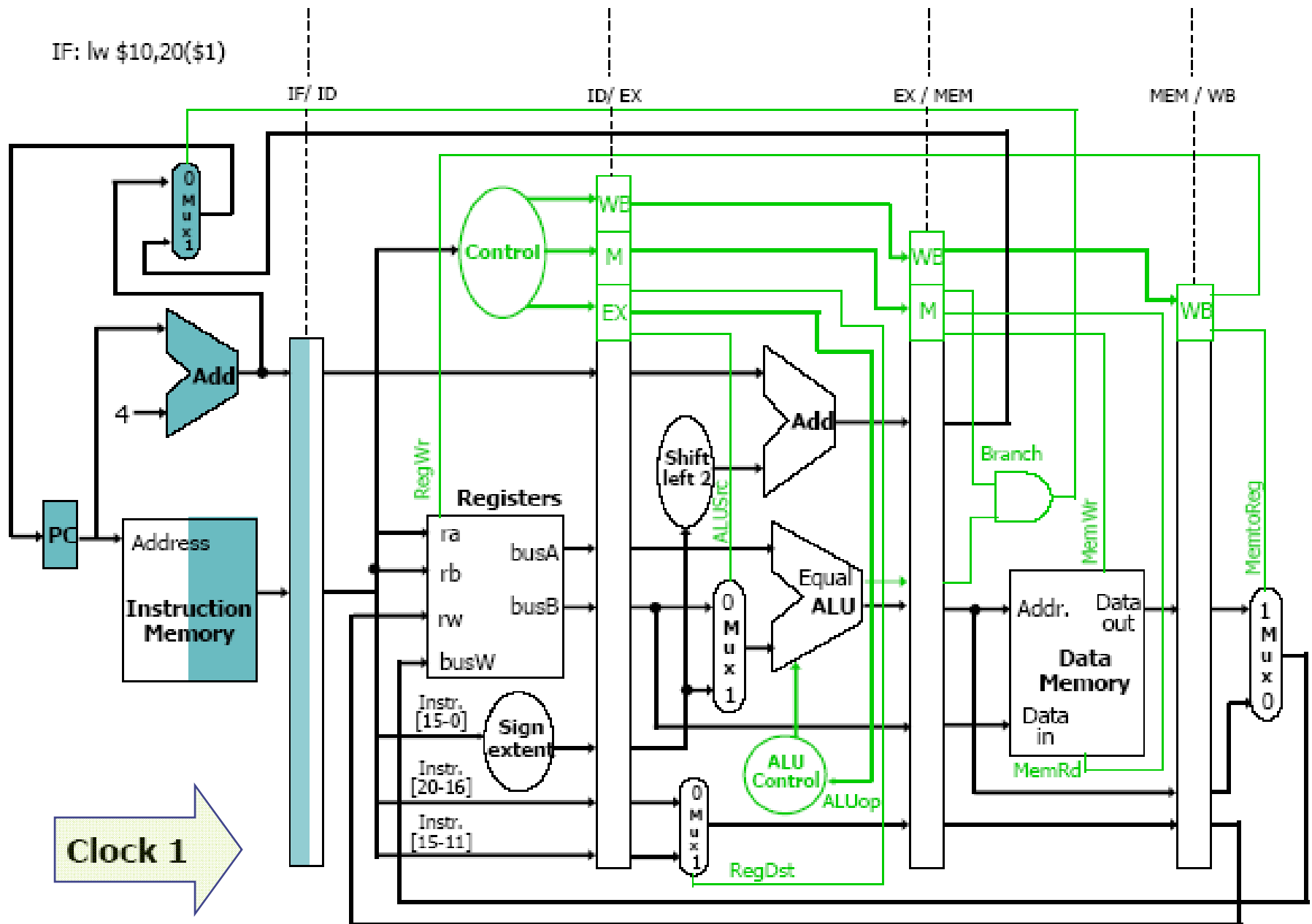
یک مثال برای روشن شدن مسئله کنترل خط لوله

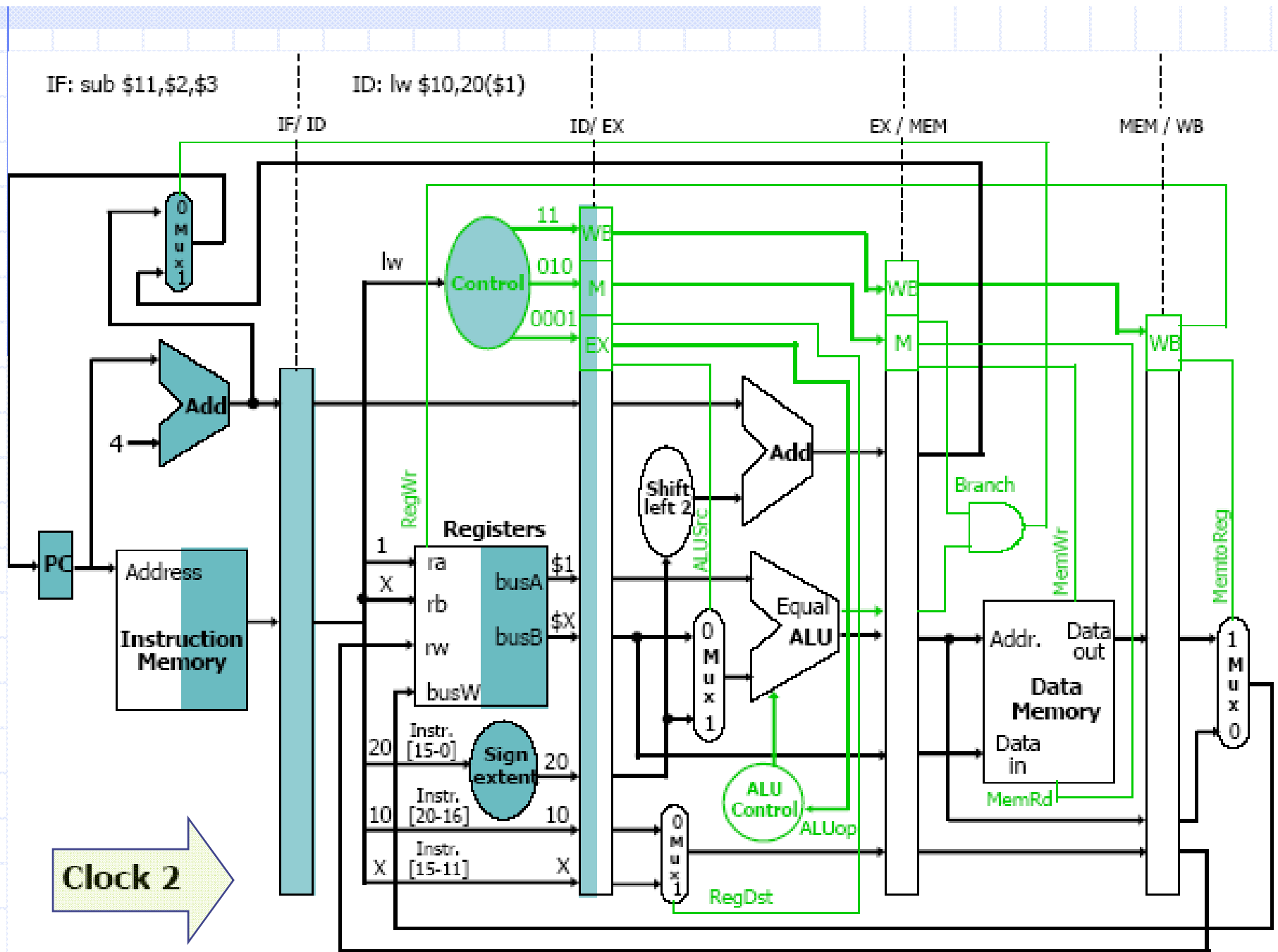
• اجازه بدهید ببینیم چه اتفاقی برای این خط لوله در برنامه زیر می افتد

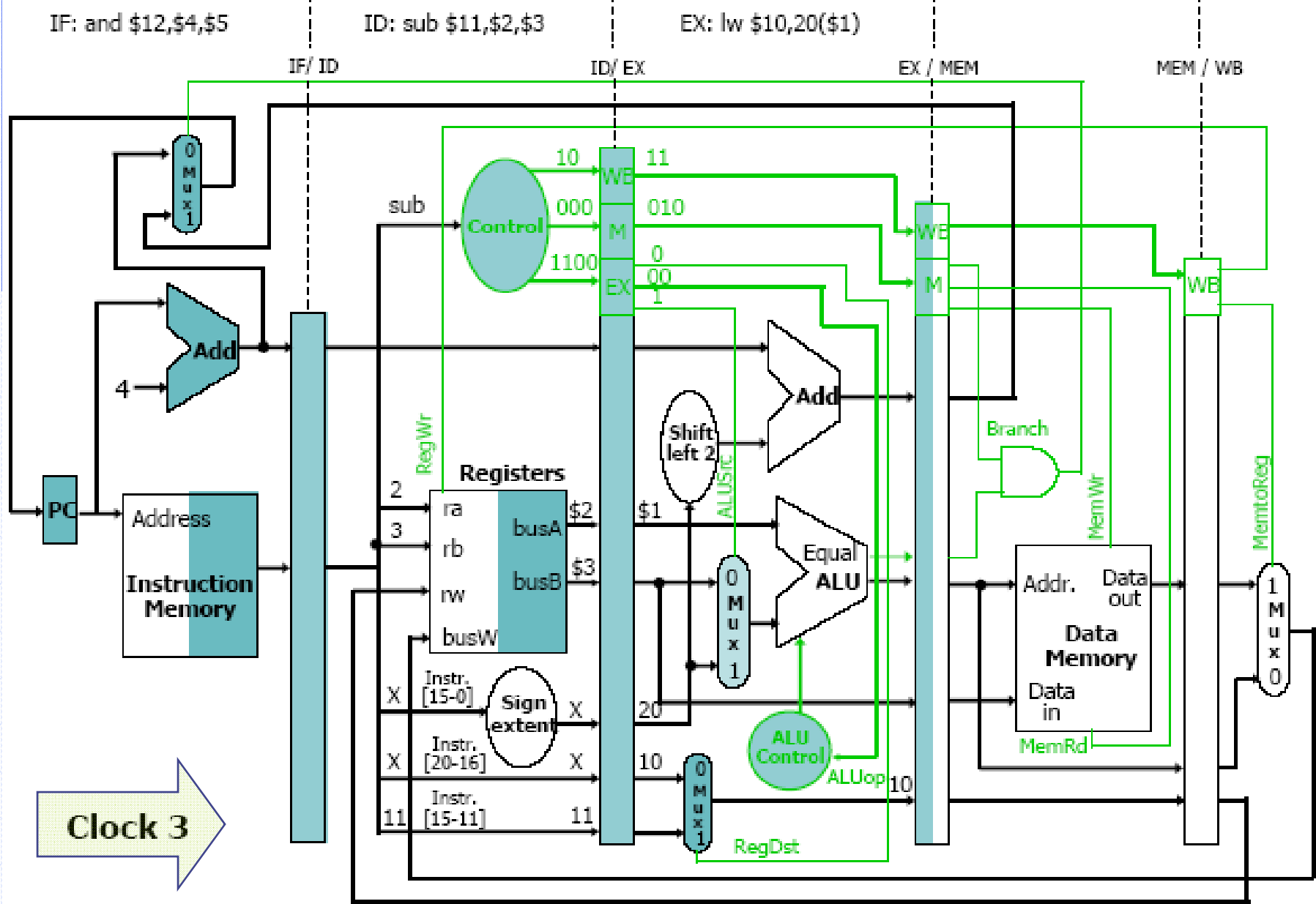
```
lw $10, 20($1)
sub $11, $2, $3
and $12, $4, $5
or $13, $6, $7
add $14, $8, $9
```

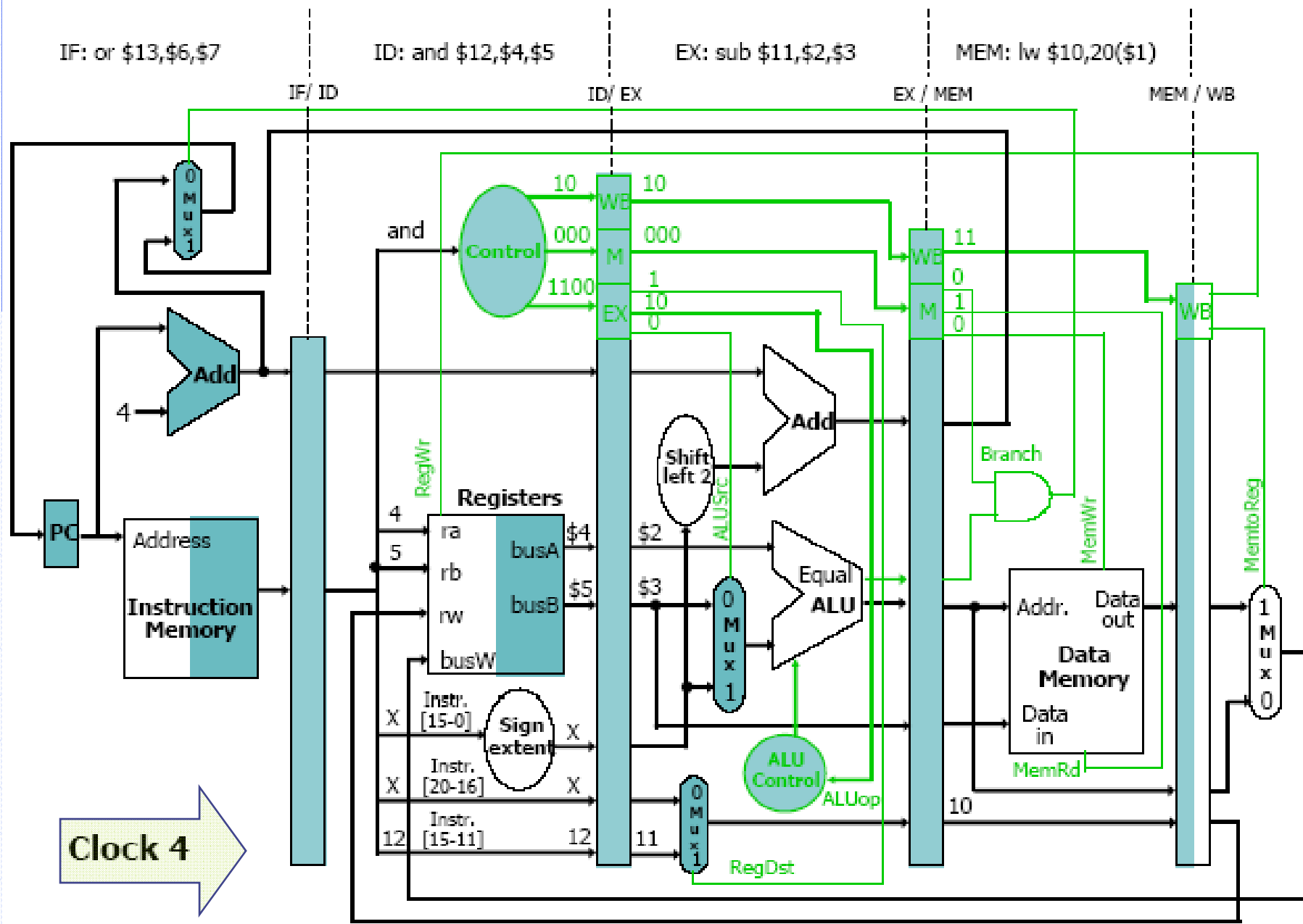
این کد هیچ هزارد داده ای، ساختاری و کنترلی ندارد

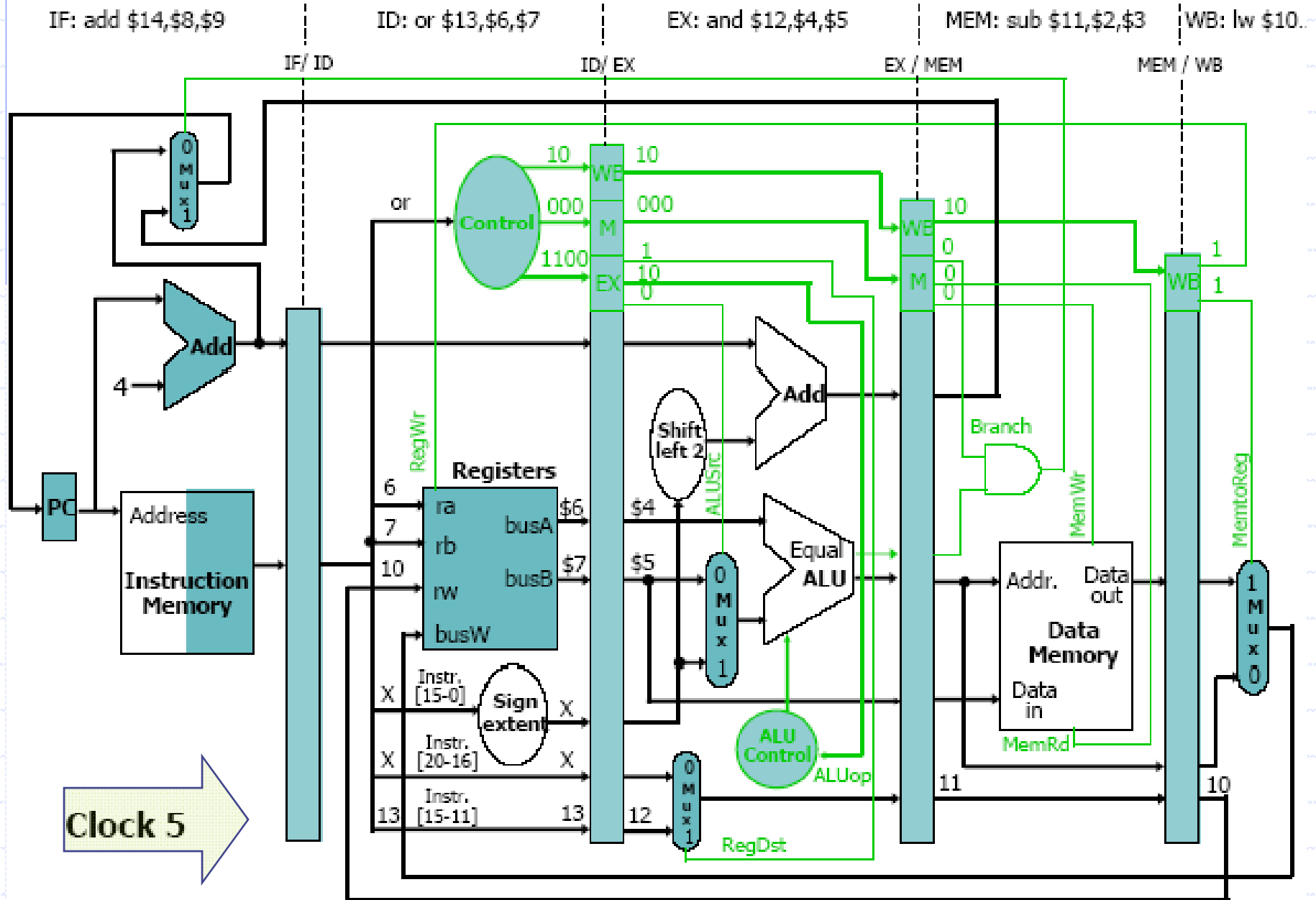
IF: lw \$t0,20(\$t1)











هزاردهای داده ای

- مثال قبل به ما نشان می دهد که دستورات عمل های مستقل که از نتایج آنها در دستورات قبل استفاده نمی شود چطور اجرا شده اند.
- این یک نمونه با برنامه های واقعی نیست.
- اجازه دهید کدهای مرتب زیر را دنبال کنیم.

sub \$2, \$1, \$3

and \$12, \$2, \$5

or \$13, \$6, \$2

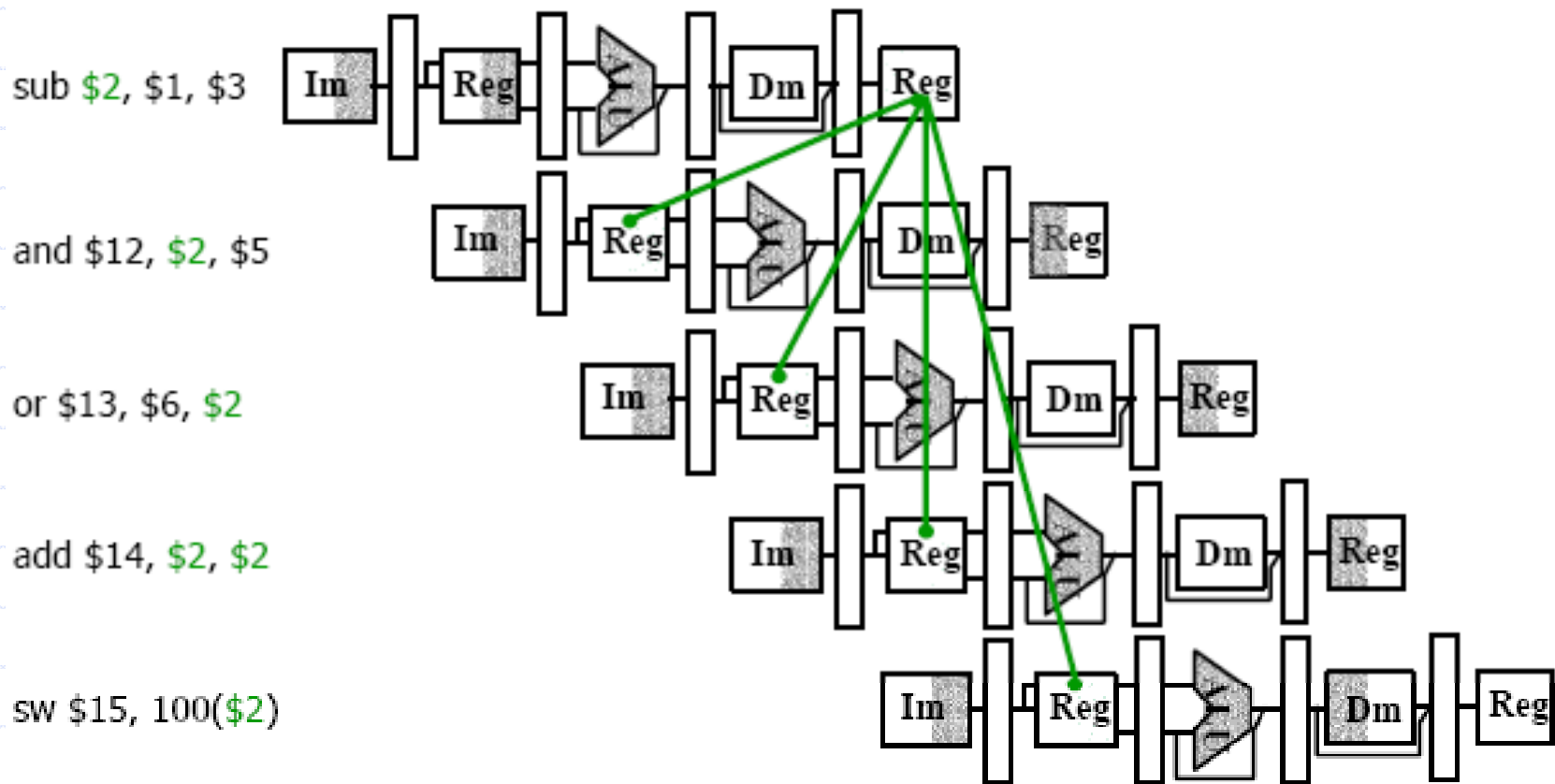
add \$14, \$2, \$2

sw \$15, 100(\$2)

- 4 دستور آخر همگی به ثبات \$2 وابسته اند که در دستور اول تولید می شود.
- فرض کنید ثبات \$2 مقدار 10 را قبل از عمل تفریق دارد و مقدار 20 را بعد از آن دارد.

هزاردهای داده ای

	CC1	CC2	CC3	CC4	CC5	CC6	CC7	CC8	CC9
The value of \$2: 10		10	10	10	10/-20	-20	-20	-20	-20



هزاردهای داده ای

راه حل ساده : کامپایلر هیچ دستورالعملی را بین دستورالعمل های sub-and وارد نمی کند

• هیچ یک از این دو دستورالعمل نه داده ای را تعریف می کنند و نه نتیجه ای را می نویسند.

```
sub $2, $1, $3
```

```
nop
```

```
nop
```

```
and $12, $2, $5
```

```
or $13, $6, $2
```

```
add $14, $2, $2
```

```
sw $15, 100($2)
```

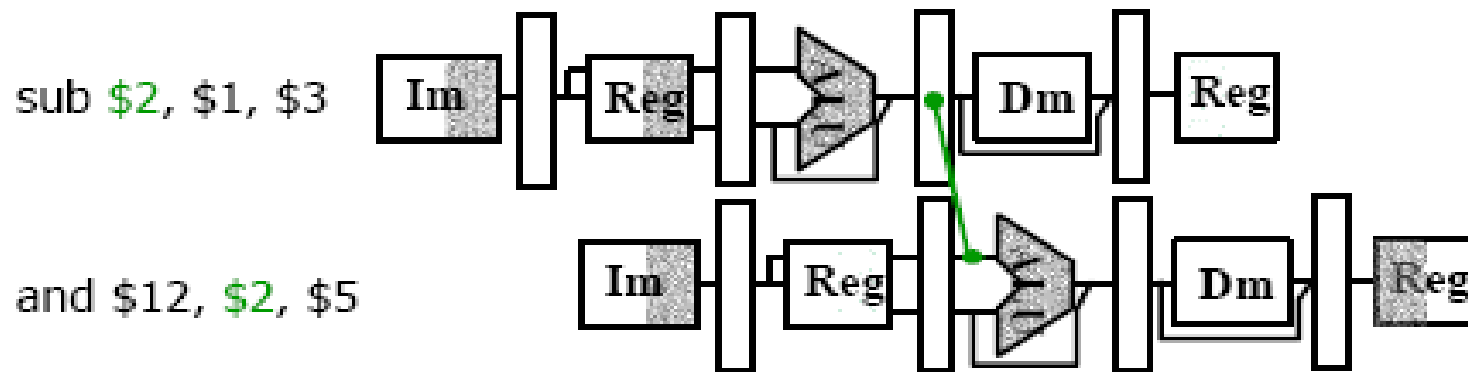
• نتیجه : این کار انجام می شود اما 2 چرخه ساعت به هدر می رود (کارایی کاهش می یابد)

کشف هزاردهای داده ای و ارسال (forwarding)

این ممکن است که هزارد داده ای کشف شود و سپس مقدار مناسب را برای حل هزارد ارسال (Forwarding) کنیم

زمانی که یک دستورالعمل سعی می کند یک ثبات را در مرحله EX بخواند دستورالعمل زودتر دیگری قصد دارد در مرحله WB بنویسد.

این یک نمونه هست بین دستورالعمل های sub-and زیر:

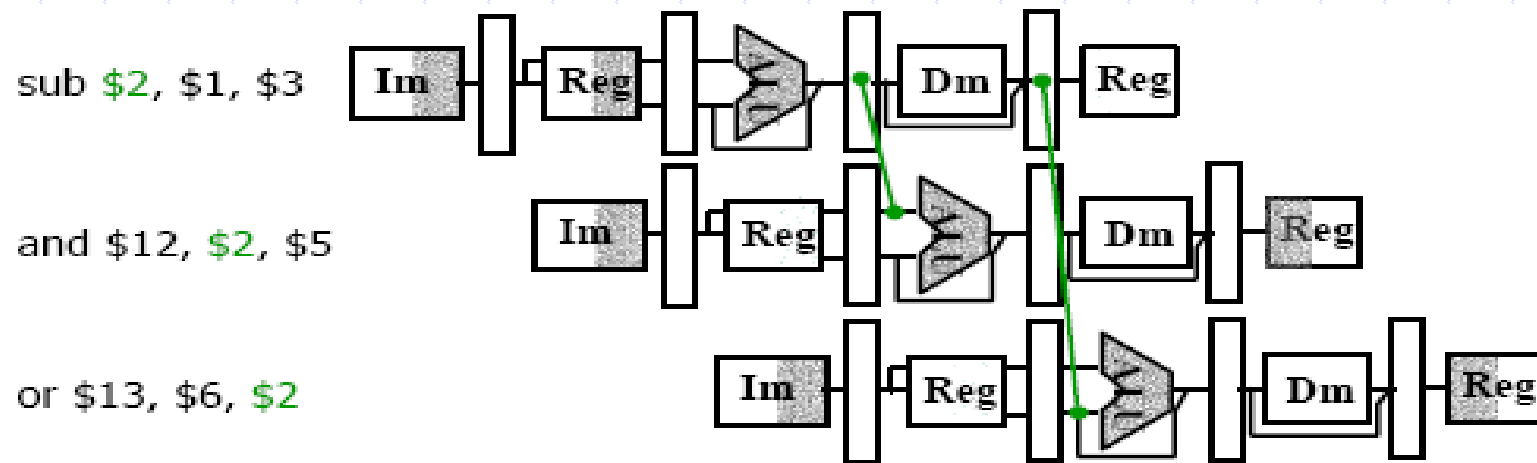


این هزارد می تواند کشف شود به وسیله یک چک کردن ساده

$$EX/MEM.RegisterRd = ID/EX.RegisterRs = \$2$$

کشف هزاردهای داده ای و ارسال (forwarding)

◆ هزارد دیگر است بین دستورالعمل های SUB_OR



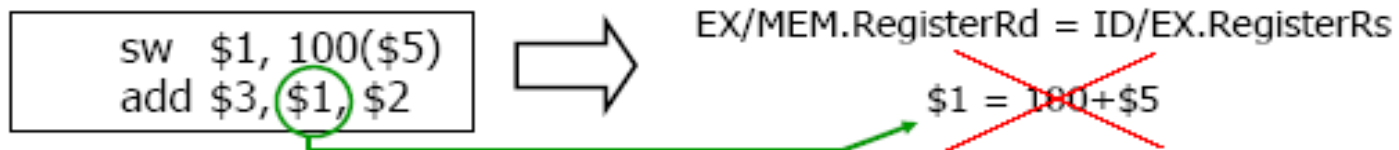
- این هزارد می تواند کشف شود به وسیله یک چک کردن ساده
 $MEM/WB.RegisterRd = ID/EX.RegisterRt = \2
- هزارد داده ای بین دستورالعمل های sub-add و sub-sw وجود ندارد

خلاصه ای از وضعیت های هزارد داده ای

- 1a. EX/MEM.RegisterRd = ID/EX.RegisterRs
- 1b. EX/MEM.RegisterRd = ID/EX.RegisterRt
- 2a. MEM/WB.RegisterRd = ID/EX.RegisterRs
- 2b. MEM/WB.RegisterRd = ID/EX.RegisterRt

This actually refers to destination field of an instruction. It is **rd** filed in **R-type** instructions and **rt** field in **I-type** instructions. Mux in the **EX** stage chooses the correct one, therefore, **EX/MEM** and **MEM/WB** pipeline registers store this information as a **rd** filed (EX/MEM.Register**Rd** and MEM/WB.Register**Rd**).

- Since some of the instructions (i.e. *sw*, *beq*) do not write to register file, the above policy is inaccurate. Consider the following code sequence:



- This problem can be solved simply by checking RegWr signal.

خلاصه ای از وضعیت های هزارد داده ای

چه اتفاقی می افتد اگر \$0 به عنوان ثبات مقصد استفاده شود؟
◆ هیچ مقدار صفری ارسال نشود
◆ بنابراین کشف هزارد باید به صورت زیر دنبال شود.

EX hazard:

if (EX/MEM.RegWr
and (EX/MEM.RegisterRd = 0)
and (EX/MEM.RegisterRd = ID/EX.RegisterRs)) ForwardA=10
if (EX/MEM.RegWr
and (EX/MEM.RegisterRd = 0)
and (EX/MEM.RegisterRd = ID/EX.RegisterRt)) ForwardB=10

MEM hazard:

if (MEM/WB.RegWr
and (MEM/WB.RegisterRd = 0)
and (MEM/WB.RegisterRd = ID/EX.RegisterRs)) ForwardA=10
if (MEM/WB.RegWr
and (MEM/WB.RegisterRd = 0)
and (MEM/WB.RegisterRd = ID/EX.RegisterRt)) ForwardB=10

خلاصه ای از وضعیت های هزارد داده ای

◆ ملاحظه کنید ترتیب زیر را :

add \$1, \$1, \$2

add \$1, \$1, \$3

add \$1, \$1, \$4

...

◆ در این نمونه ، نتیجه از مرحله MEM ارسال شده برای اینکه نتیجه در مرحله MEM نتیجه اخیر خیلی است نسبت به نتیجه در مرحله WB.

MEM hazard:

if (MEM/WB.RegWr

and (MEM/WB.RegisterRd = 0)

and (MEM/WB.RegisterRd = ID/EX.RegisterRs)) ForwardA=10

if (MEM/WB.RegWr

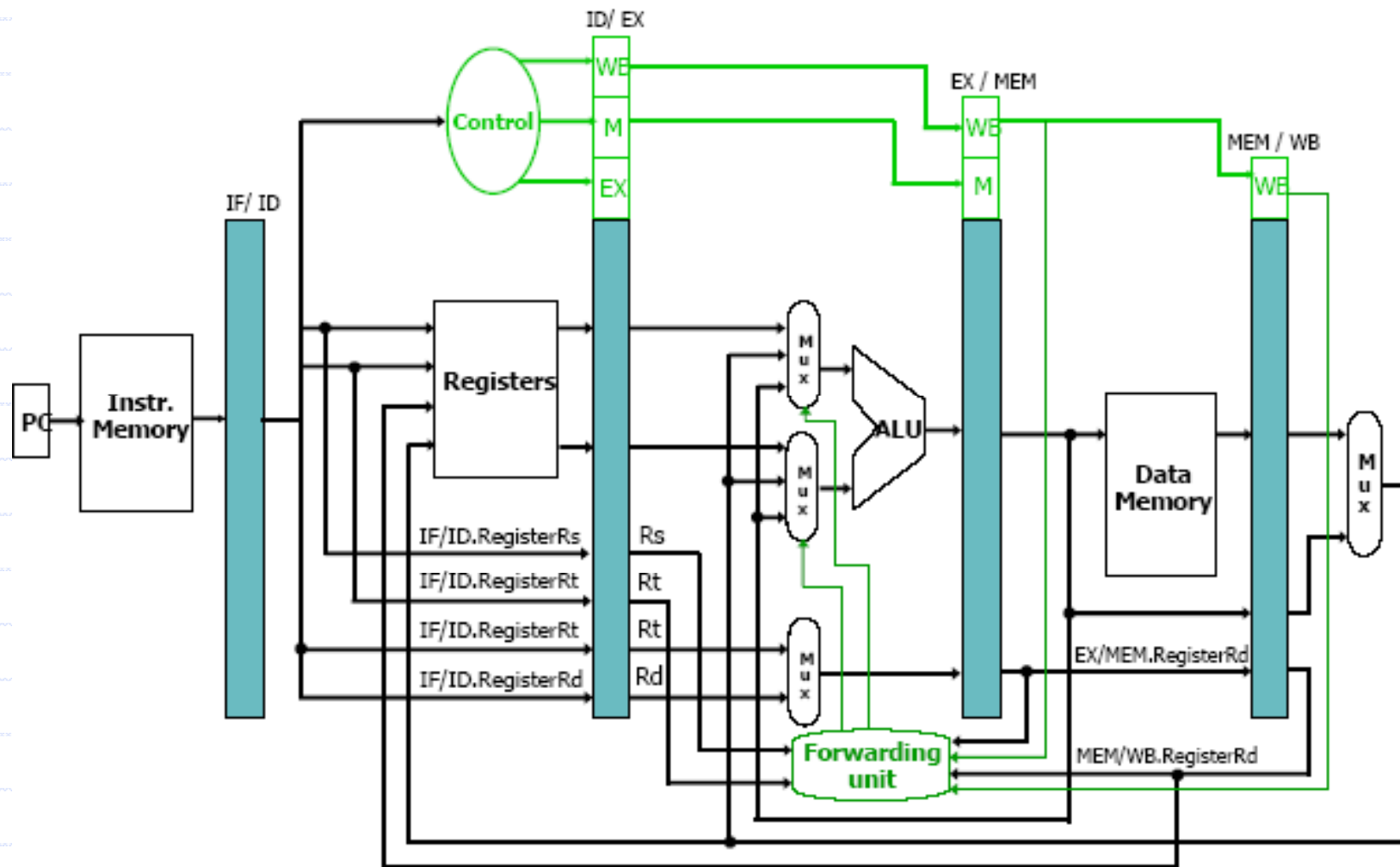
and (MEM/WB.RegisterRd = 0)

and (MEM/WB.RegisterRd = ID/EX.RegisterRt)) ForwardB=10

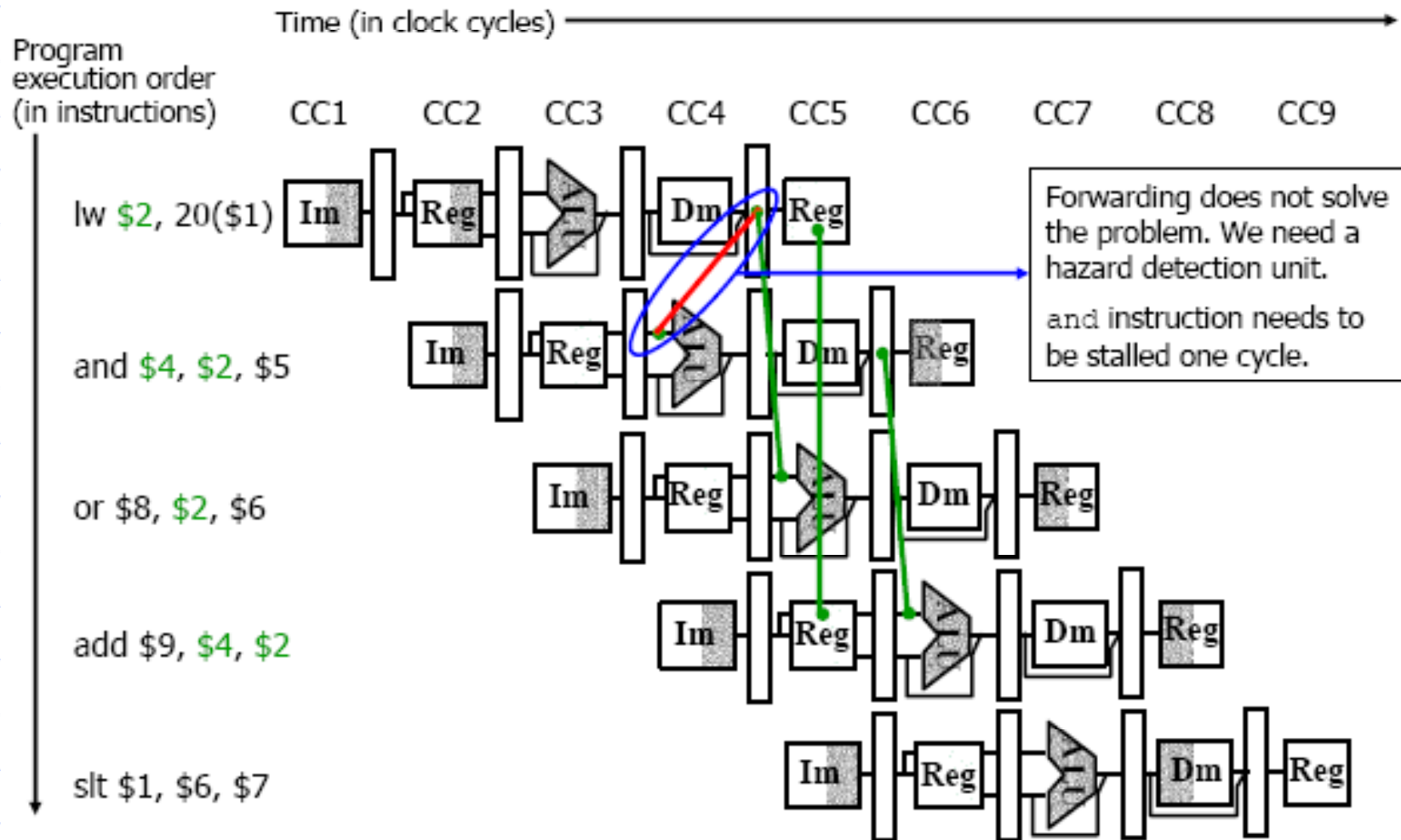
and (EX/MEM.RegisterRd = ID/EX.RegisterRs)

and (EX/MEM.RegisterRd = ID/EX.RegisterRt)

مسیر داده خط لوله ای شده با ارسال (Forwarding)



هزاردهای داده ای و توقف



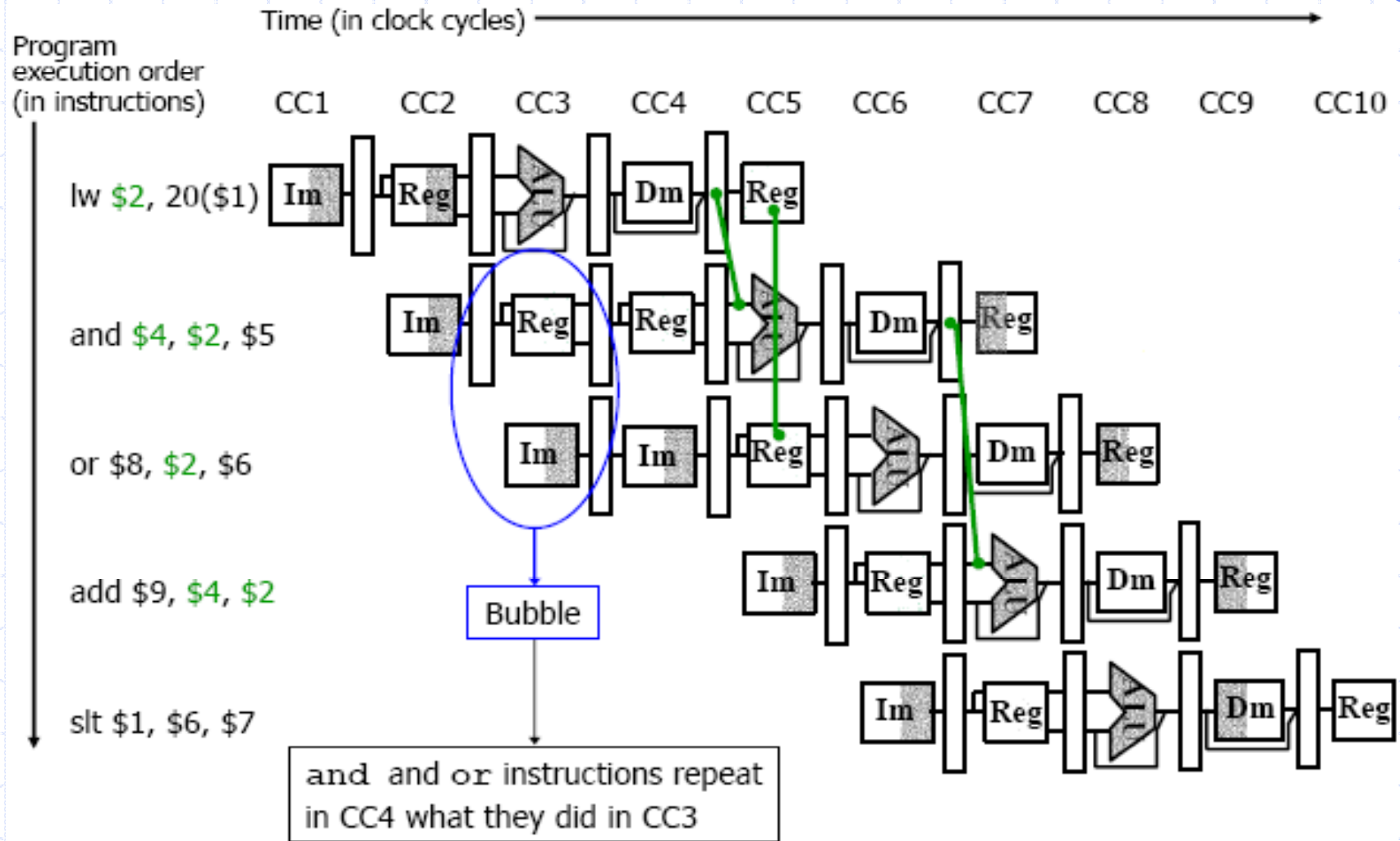
کشف هزارد

• کنترل برای کشف هزارد

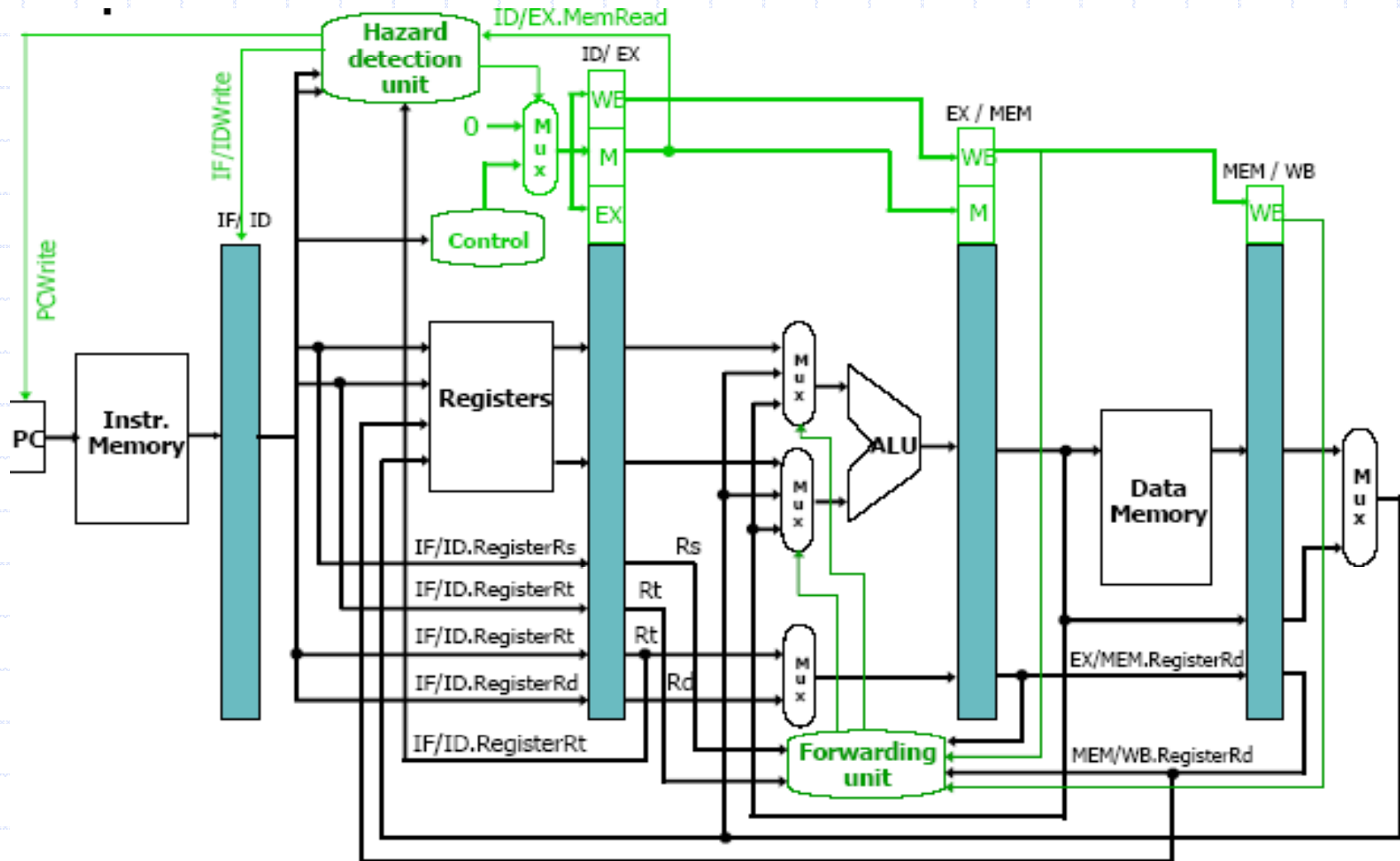
if (ID/EX.MemRd چک کن اگر دستورالعمل بار کردن است
and →
((ID/EX.RegisterRt = IF/ID.RegisterRs) or
(ID/EX.RegisterRt = IF/ID.RegisterRt)))

سپس:
توقف خط لوله

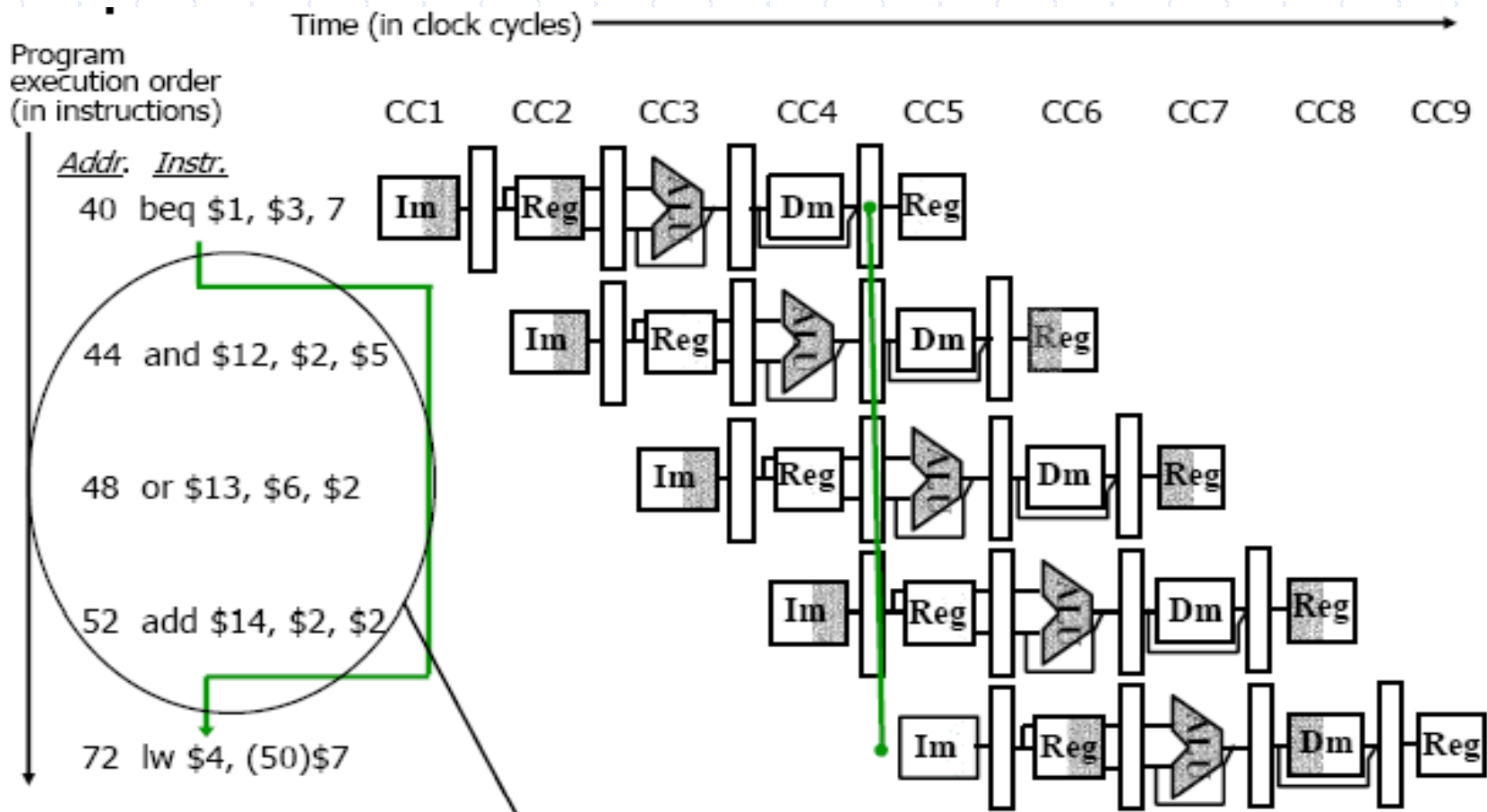
هزاردهای داده ای و توقف



مسیر داده خط لوله ای شده با (Forwarding) ارسال و واحد کشف هزارد



هزاردهای انشعاب

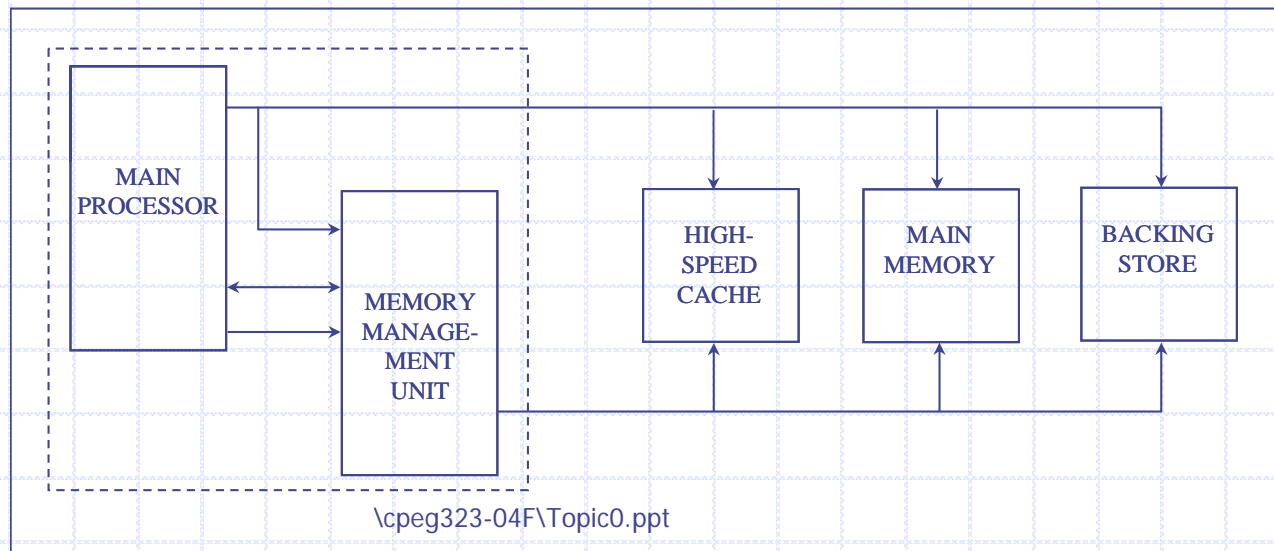
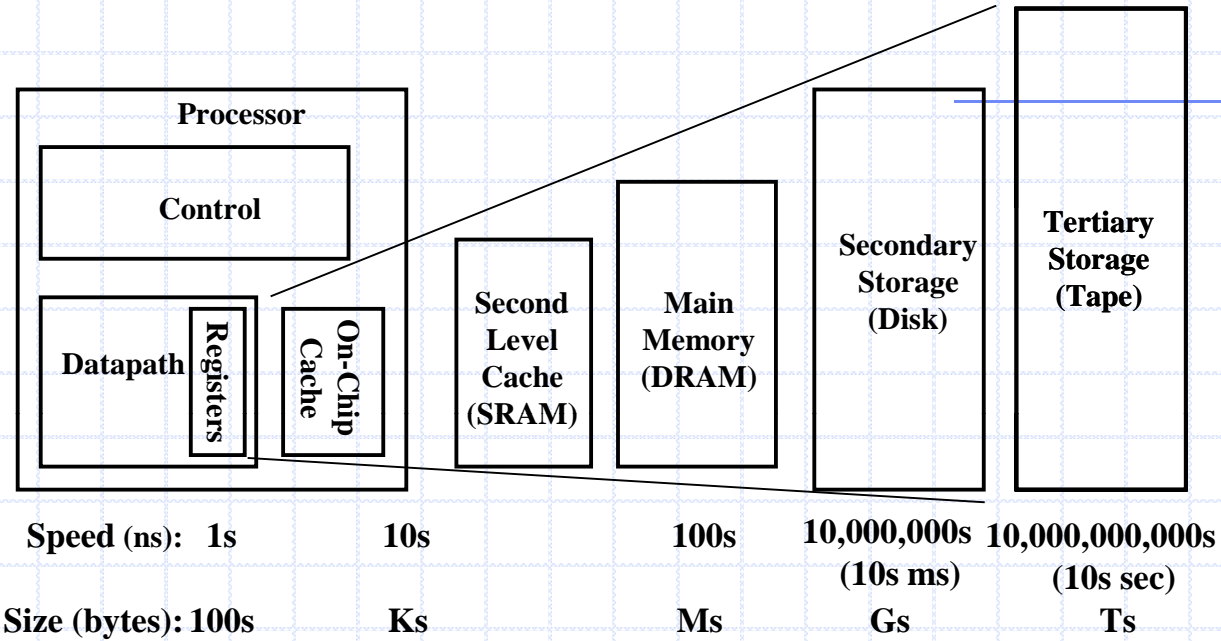


Actually, the number of instructions needs to be flushed can be reduced from 3 to 1 instruction (shown in the following slide) when the direction of branch is mispredicted.

فصل هفتم

طراحی سیستم حافظه

سلسله مراتب حافظه



مشخصات برنامه ها و سازماندهی حافظه

◆ تجزیه و تحلیل Ram در برابر دستیابی ترتیبی:

■ در مقایسه کارایی بر هزینه و تکنولوژی

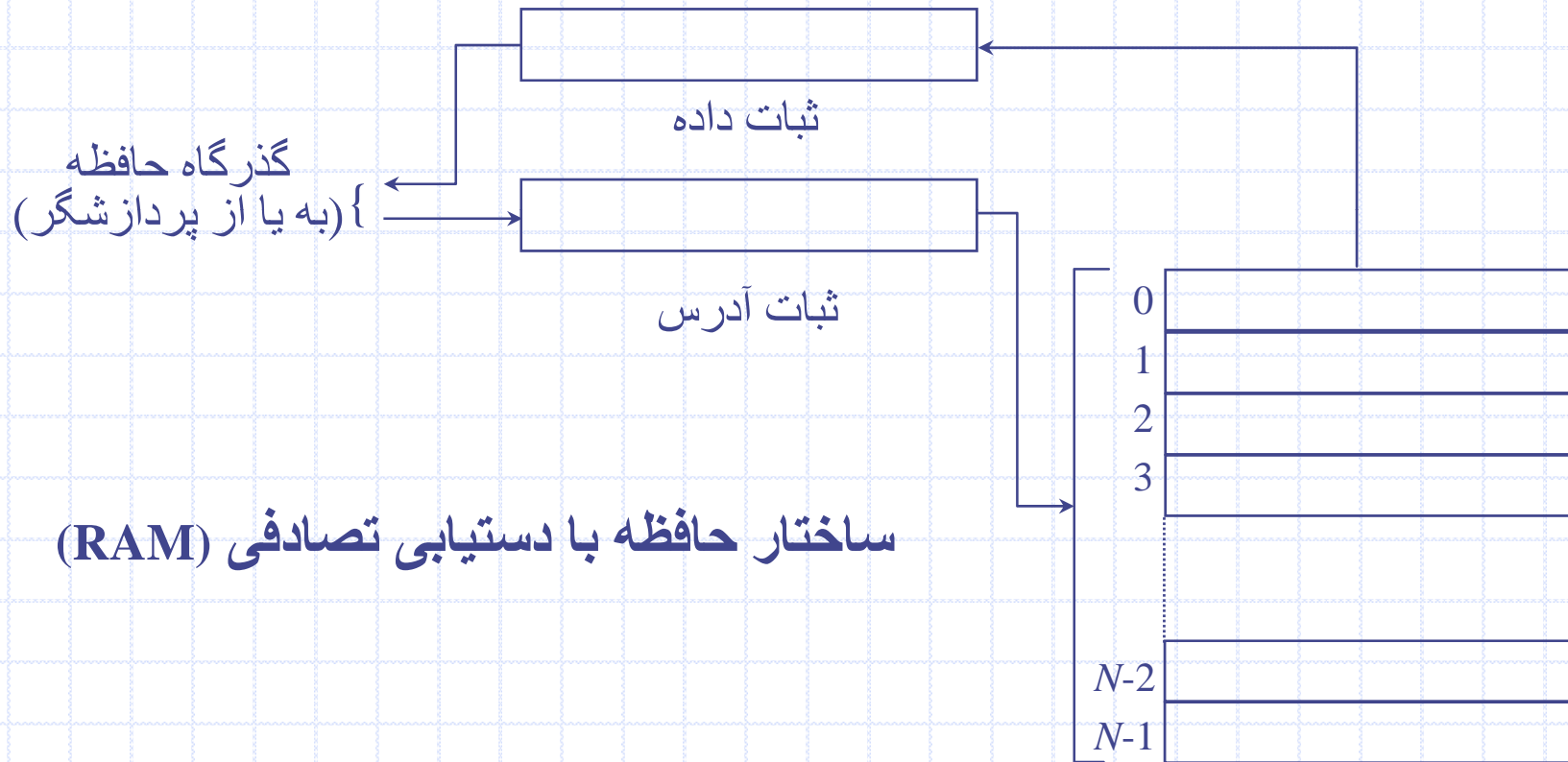
◆ همجواری در الگوهای دسترسی به حافظه

◆ سلسله مراتب در طراحی حافظه

■ حافظه نهان

■ حافظه مجازی

حافظه با دستیابی تصادفی



ساختار حافظه با دستیابی تصادفی (RAM)

نکته: دارای زمان دستیابی ثابت

آدرسها سلولهای حافظه

کارایی حافظه

پهنای باند = تعداد بیتها بر ثانیه ”که می توانند در دسترس قرار گیرند“
که برابر است با:

$$(\text{bit/word}) \times (\text{word/cycle}) \times (\text{cycle/sec})$$

بنابراین، بهبود پهنای باند؟

”گلوگاه ون نیومن“

چگونه می توان کارایی سیستم حافظه را بهبود بخشید

◆ کاهش مدت چرخه

◆ افزایش اندازه کلمه

◆ همزمانی

◆ طراحی کارآمد حافظه

حافظه نهان

- ◆ تقریباً تمام ریز پردازنده های با کارایی بالا در بازار از حافظه نهان استفاده می کنند.
- ◆ چرا معماری برداری (نوعی کامپیوتر بسیار بزرگ) Cray استفاده نشود؟

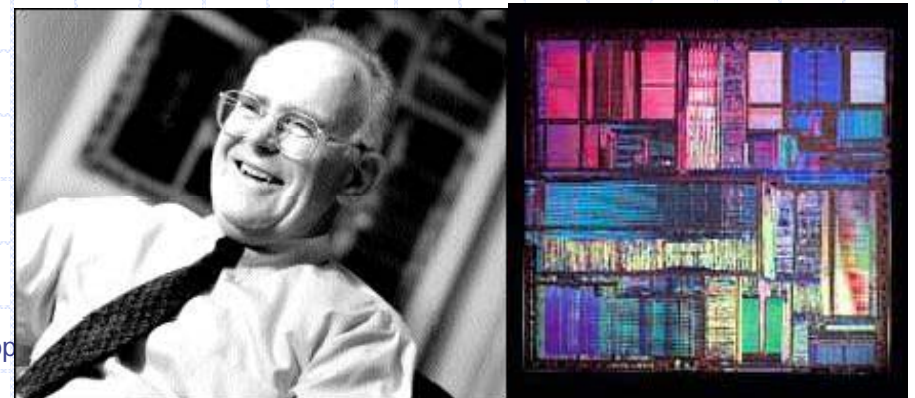
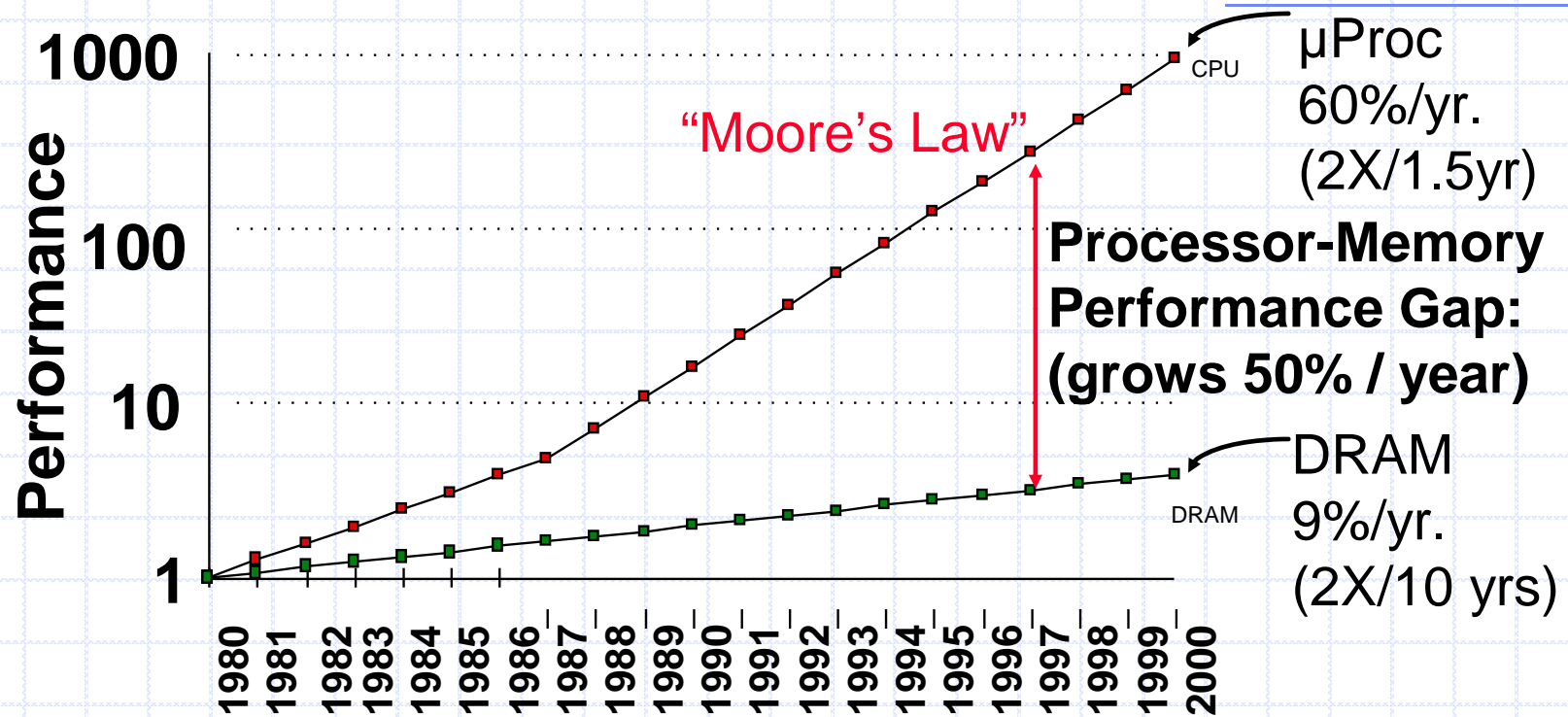
حافظه نهان

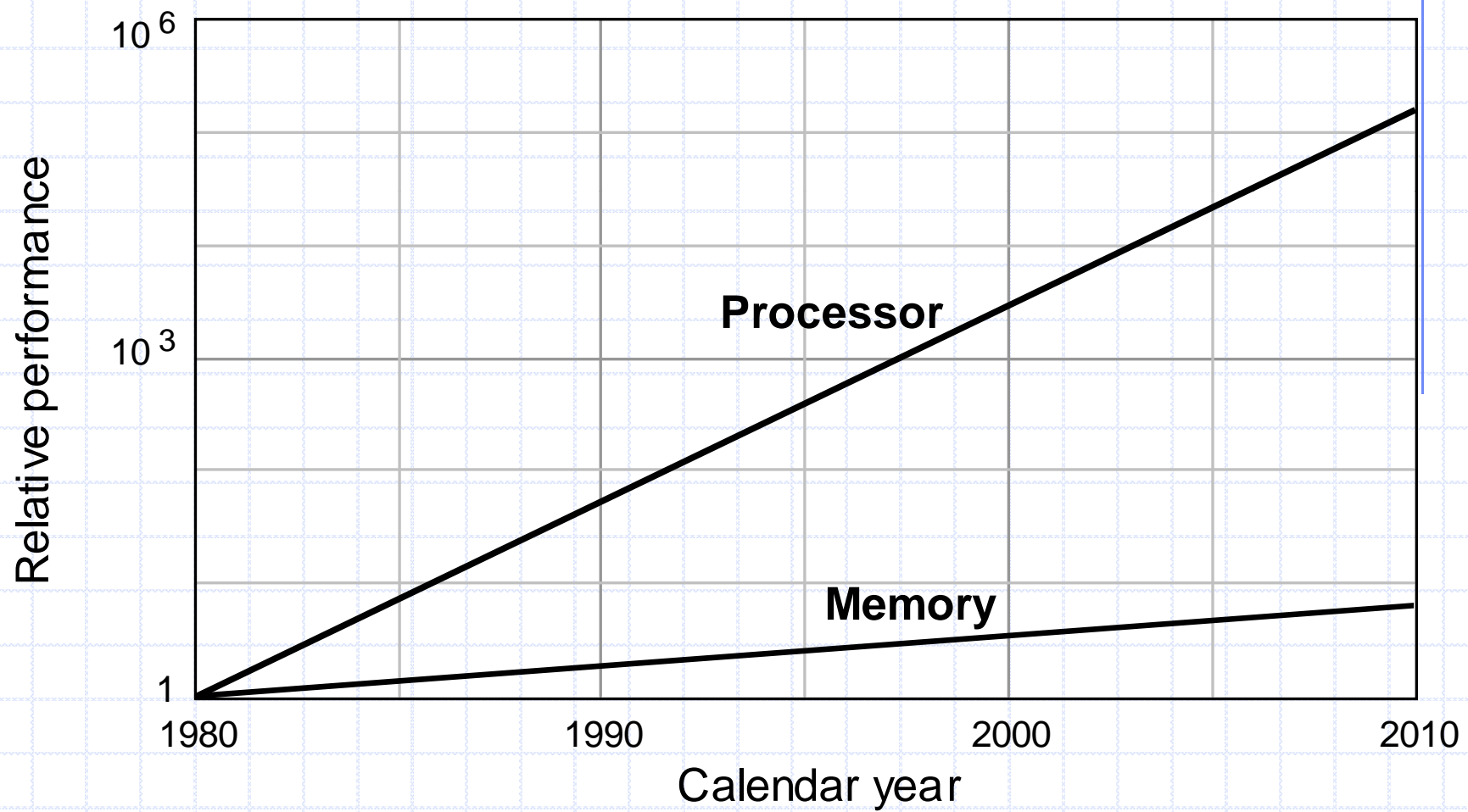
پیشرفت‌ها در تکنولوژی مدارات مجتمع نه فقط متاثر از DRAM ها، بلکه SRAM ها، هزینه حافظه نهان را بسیار کاهش داده اند. حافظه های نهان یکی از ایده های بسیار مهم در معماری کامپیوتر هستند زیرا آنها می توانند کارایی استفاده از حافظه را به طور اساسی بهبود بخشند.

افزایش شکاف بین زمان چرخه DRAM و زمان چرخه پردازنده، که در شکل بعد نمایش داده شده است، یک انگیزه برای ایجاد حافظه نهان است. اگر ما بر آنیم که پردازنده ها با سرعتی که توانایی آن را دارند اجرا شوند، ما باید حافظه هایی با سرعت بیشتر برای فراهم نمودن داده داشته باشیم.

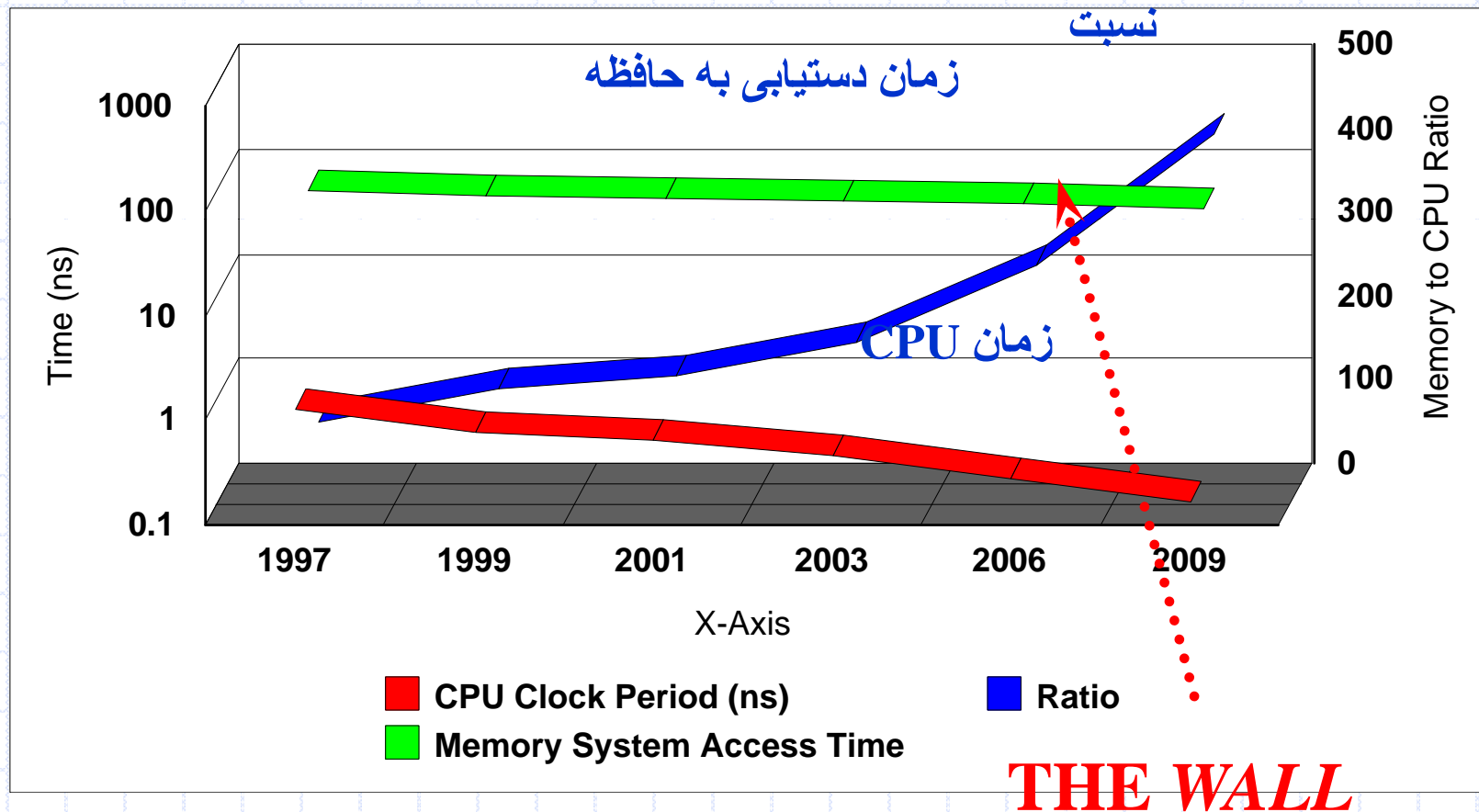
[Joupi & Hennessy 91]

قانون مور - یک فرصت از دست رفته





نهفتگی در یک سیستم واحد



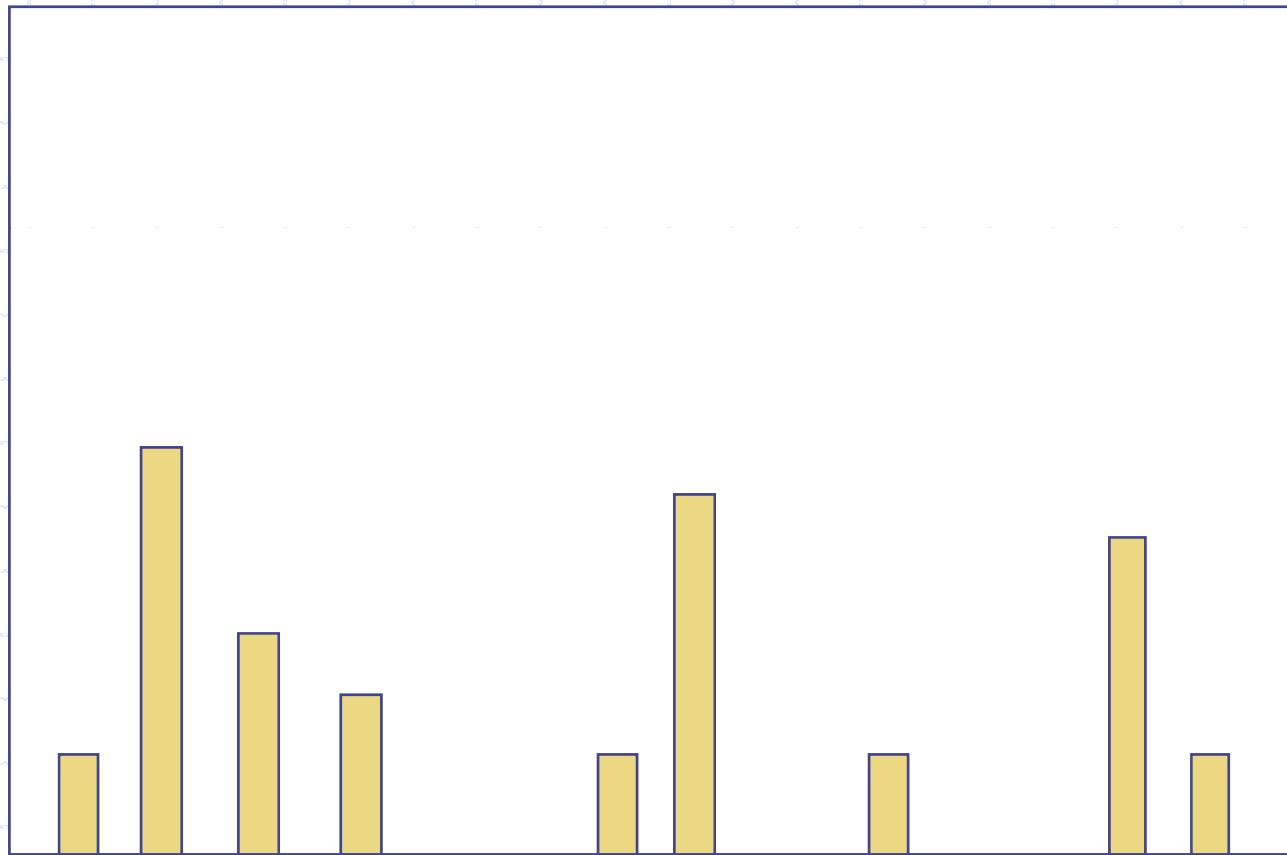
همجواری مراجع

“ مراجع برنامه تمایل دارند که به
موقع در کنار هم قرار بگیرند.”

مکانهایی با احتمال دستیابی بالا (زمان کوتاه)

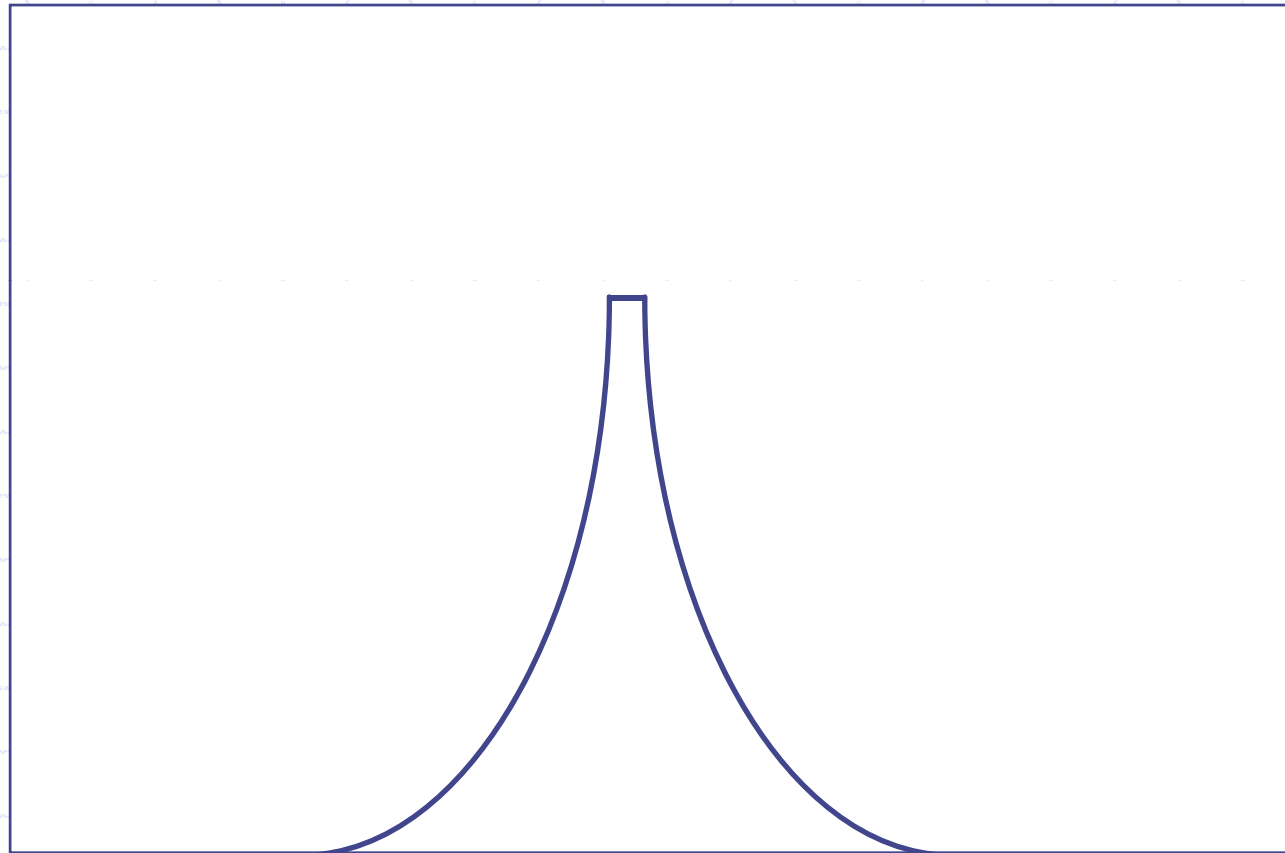
- ◆ همجواری PC
- ◆ قاب پشته (محل)
- ◆ نزدیک زیر روال ها
- ◆ داده های فعال

احتمال ارجاع



آدرس ارجاع →

احتمال ارجاع



آدرس ارجاع →

مشکل: اغلب پیش بینی یک page size می شود که بزرگتر از اندازه مورد نیاز باشد.

همجواری

موفقیت حافظه های نهان که قبلا تشریح شد برگرفته از ”**خاصیت همجواری**“ است.
خاصیت همجواری دارای دو جنبه است، *spatial* و *temporal*.

در دوره های کوتاه زمانی، یک برنامه ارجاعها به حافظه خود را، به صورت غیر
یکنواخت در فضای آدرس خود گسترش می دهد.

که این بخشهای فضای آدرس در دوره های زمانی طولانی به صورت گسترده
یکنواخت باقی می ماند.

همجواری موقتی

اولین خاصیت، معروف به همجواری موقتی، یا همجواری زمانی، به این معنی است که: اطلاعاتی که در آینده نزدیک مورد نیازند احتمالاً هم اکنون در حال استفاده اند.

این نوع رفتار می تواند از حلقه های برنامه که هر دوی داده و دستورالعمل مجدداً استفاده می شوند، انتظار رود.

همجواری مکانی

دومین خاصیت، همجواری مکانی، به این معنی است که بخشهایی از فضای آدرس که در حال استفاده اند معمولاً شامل یک تعدادی از بخشهای همجواری یکتا از فضای آدرس می باشد. همجواری مکانی، سپس، به این معناست که محل ارجاع برنامه در آینده نزدیک، احتمالاً نزدیک محل ارجاع فعلی است

[Smith82, p475]

این نوع رفتار با آگاهی معمولی از برنامه ها مورد انتظار است؛ بخشهای وابسته داده (متغیرها، آرایه ها) معمولاً با یکدیگر ذخیره می شوند، و دستورات بیشتر به صورت ترتیبی اجرا می شوند. تا زمانی که حافظه نهان سگمنتهای اطلاعاتی را که اخیراً استفاده شده اند، میانگیر (در خود نگهداری) می کند، خاصیت همجواری اشاره دارد به اینکه اطلاعات مورد نیاز احتمالاً در حافظه نهان یافت می شود.

اطلاعاتی که در آینده نزدیک مورد نیازند احتمالاً شامل اطلاعات رایجی هستند که هم اکنون در حال استفاده اند. (*locality by time*)

و آن اطلاعاتی که همجواری منطقی دارند و به طور جاری مورد استفاده اند. (*locality by space*)

موقتی

فضایی

مقدمه ای بر طراحی cache

(نوعی حافظه میانی که سرعت دستیابی به اطلاعات آن بالا است)

Cache (حافظه نهان)

مکانی امن برای پنهان سازی و ذخیره سازی.

Webster Dictionary

Cache (حافظه نهان)

◆ حتی با وجود ظرفیت cache هنوز تقریباً همه cpuها شدیداً به وسیله زمان دستیابی cache محدود شده اند. در بسیاری از گونه ها، هرچند زمان دستیابی cache کاهش داده شده باشد ماشین نیز باید تسریع شود.

-Alan Smith-

Even more so for MPs!

◆ الگوهایی که می توانند طراحی حافظه cache
حاضر اشکست دهند، باعث افزایش کارایی حافظه
cache شده، و برای هر برنامه یا بار کاری واقعاً
محاسبات مفیدی انجام دهد.

◆ بهینه کردن طراحی حافظه cache به طور گسترده که دارای 4 جنبه است:

1. به حد اعلی رساندن احتمال یافتن یک مرجع حافظه در cache (نسبت برخورد)
2. به حداقل رساندن زمان دستیابی به اطلاعاتی که در cache وجود دارند (زمان دستیابی)
3. به حداقل رساندن تأخیر به علت یک خطا و
4. به حداقل رساندن هزینه های ناشی از به روز کردن حافظه اصلی، نگهداری *cache coherence* و غیره.

فاکتورهای کلیدی در تصمیم طراحی برای VM و Cache (حافظه مجازی)

$$\frac{\text{زمان دستیابی حافظه اصلی}}{\text{cache زمان دستیابی}} = 4 \sim 20$$

$$\frac{\text{زمان دستیابی حافظه ثانویه}}{\text{زمان دستیابی حافظه اصلی}} = 10^4 \sim 10^6$$

کنترل cache معمولاً در سخت افزار اجرا می شوند

فن آوری در 1990s:

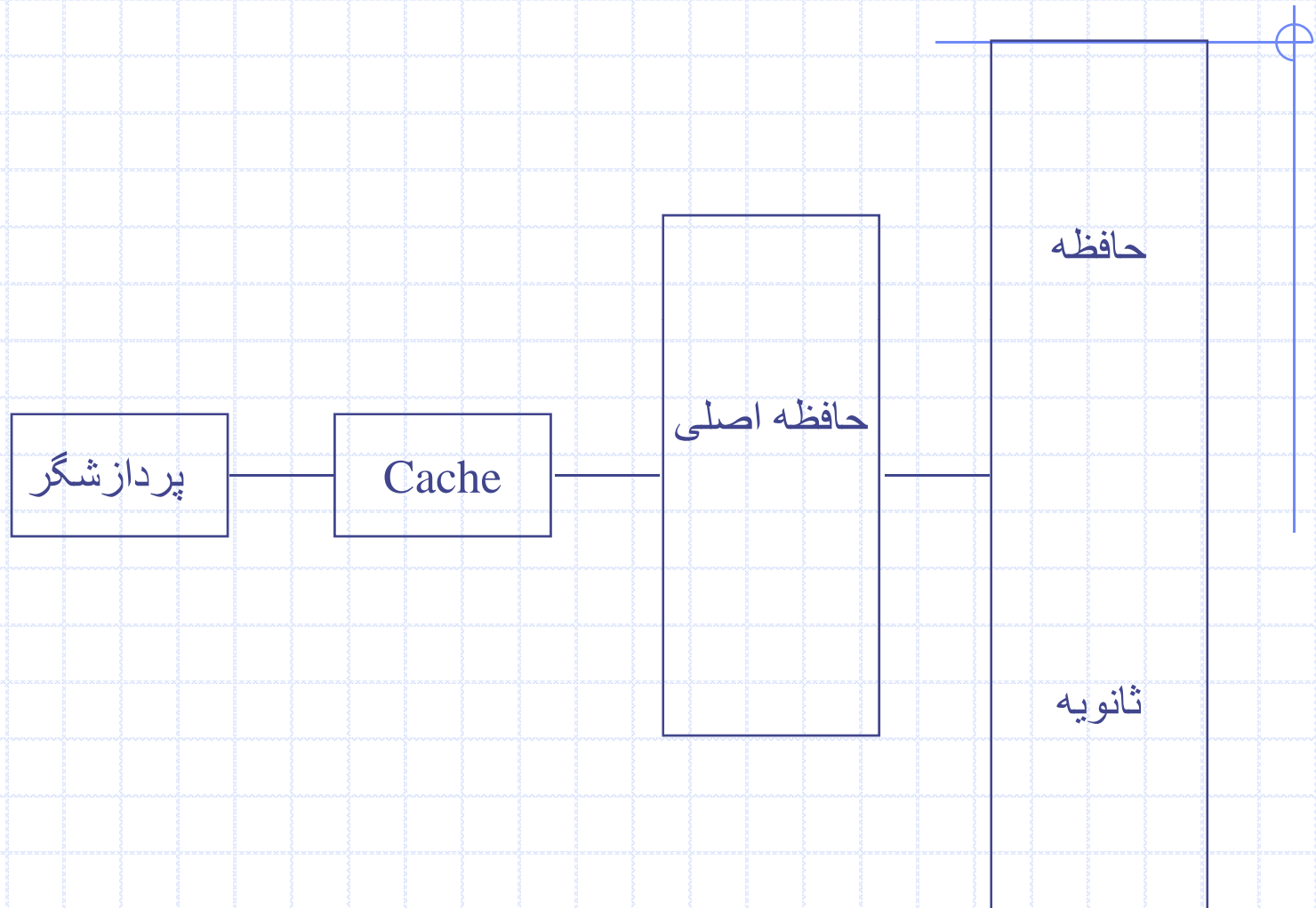
Memory Technology	Typical Access Time	\$ per Mbyte
SRAM	10-20 ns	200-400
DRAM	90-120 ns	50-100
Magnetic disk	10,000,000 - 20,000,000 ns	2-5

فن آوری در 2000s:؟

فناوری در 2000s :

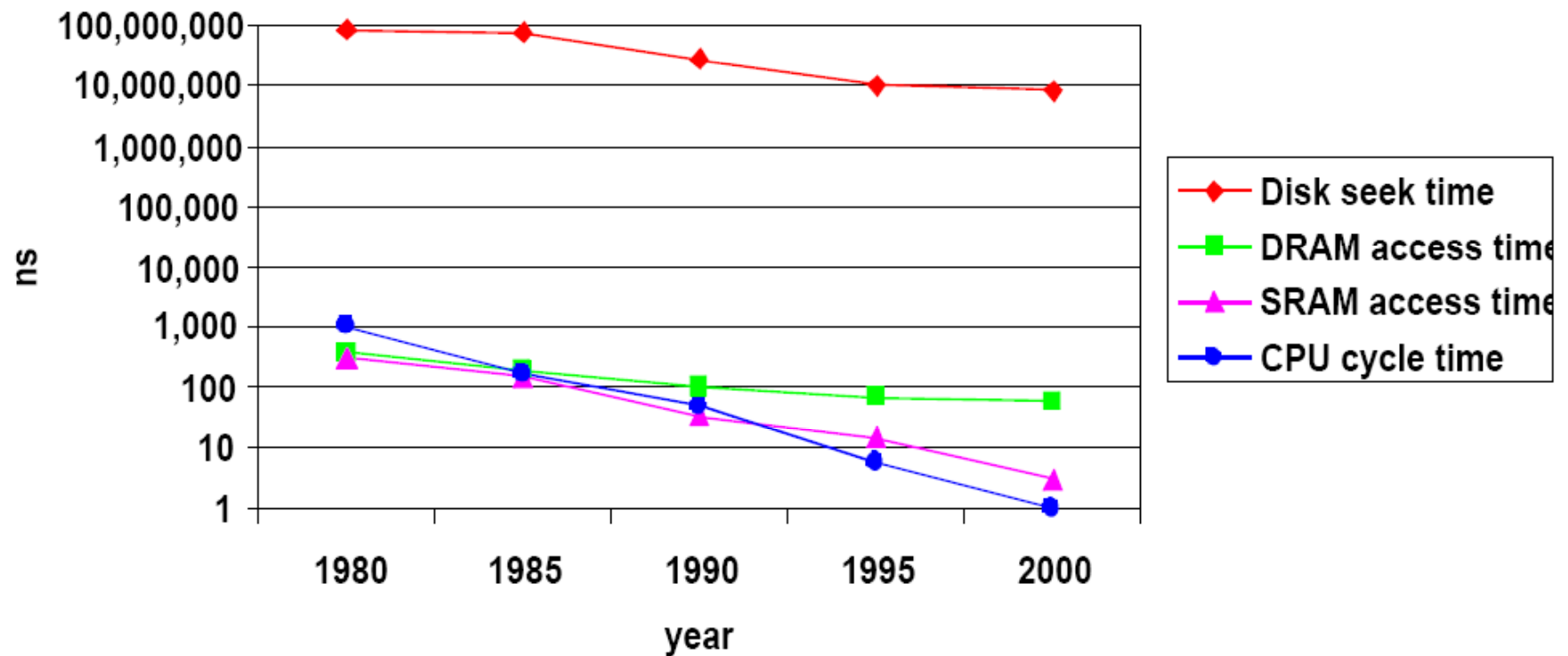
Memory Technology	Typical Access Time	\$ per Mbyte
SRAM	5-10 ns	1-20
DRAM	10-60 ns	0.1-1
Magnetic disk	2,000,000 - 5,000,000 ns	0.001-0.006

Cache در سلسله مراتب حافظه



The CPU-Memory Gap

The increasing gap between DRAM, disk, and CPU speeds.



چهار سوال برای دسته بندی سلسله مراتب حافظه

قاعده کلی و اساسی که گرداننده همه سلسله مراتب های حافظه است به ما اجازه می دهد از شرایطی که باعث برتری یافتن نسبت به سطوحی که در موردشان سخن گفتیم استفاده کنیم.

این قواعد یکسان به ما اجازه طرح چهار سوال درباره هر سطح از سلسله مراتب را می دهد:

چهار سؤال برای دسته بندی سلسله مراتب حافظه

س1: در کجا یک بلاک می تواند در سطح بالاتری قرار بگیرد؟ (مکان بلاک)

س2: چگونه وقتی یک بلاک در سطح بالاتری قرار دارد شناخته می شود؟ (هویت بلاک)

س3: کدام بلاک می تواند در یک فقدان جایگزین شود؟ (جایگذاری بلاک)

س4: در نوشتن چه روی می دهد؟ (استراتژی نوشتن)

این سوالات ما را در جهت رسیدن به یک درک کامل از مبادلات
متفاوت در خواست شده به وسیله ارتباطات حافظه ها در سطوح
متفاوت یک سلسله مراتب کمک می کند.

	0	1	2	3	4	5	6	7
TAGS	DATA							
0117X	35	72	55	30	64	23	16	14
7620X	11	31	26	22	55	...		
3656X	71	72	44	50	...			
...								
1741X			33	35	07	65	...	

خط ←

01173

آدرس

30

داده ها

مفهوم cache hit و cache miss

زمان دستیابی موثر cache	:	t_{eff}
زمان دستیابی cache	:	t_{cache}
زمان دستیابی حافظه اصلی	:	t_{main}
نسبت برخورد (موءفقیت)	:	h

$$t_{\text{eff}} = ht_{\text{cache}} + (1-h)t_{\text{main}}$$

مثال

$$\begin{aligned} &= 10 \text{ ns} - 1 - 4 \text{ clock cycles} & t_{\text{cache}} & \text{Let} \\ &= 50 \text{ ns} - 8 - 32 \text{ clock cycles} & t_{\text{main}} & \\ & & = 0.95 & h \\ & & t_{\text{effect}} & = ? \end{aligned}$$

$$\begin{aligned} &10 \times 0.95 + 50 \times 0.05 \\ &9.5 + 2.5 = 12 \end{aligned}$$

نسبت بر خورد (یا موء فقیت)

- ◆ نیاز به نیروی کافی برای رسیدن به سطح کارایی دلخواه دارد (say > 90%)
- ◆ وسعت بخشی به اثرات و نتایج تغییرات
- ◆ هرگز ثبات دائمی حتی برای ماشین های مشابه ندارد

حساسیت کارایی (hit ratio) *w.r.th*

$$t_{\text{eff}} = h t_{\text{cache}} + (1-h) t_{\text{main}}$$
$$\frac{t_{\text{main}}}{t_{\text{cache}}} = \frac{t_{\text{eff}}}{t_{\text{cache}}} [h + (1-h)]$$
$$\frac{t_{\text{main}}}{t_{\text{cache}}} = \frac{t_{\text{eff}}}{t_{\text{cache}}} [1 + (1-h)]$$

since $\frac{t_{\text{main}}}{t_{\text{cache}}} \approx 10$, the magnifactor of h changes is 10 times.

نتیجه: خیلی حساس

یادآوری

" h T "
 \approx

مثال:

Let $h = 0.90$

if $h \neq \underline{0.05}$ (0.90 \rightarrow 0.95)

then $(1 - h) = \underline{0.05}$

then $t_{\text{eff}} = t_{\text{cache}} (1 + 0.5)$

اصطلاحات بنیادی

◆ Cache line (block) - size of a room

1 ~ 16 words

“ یک مجموعه از داده های همجوار که به عنوان یک نهد ذخیره سازی cache مورد بحث واقع شده اند.”

◆ Cache directory - key of rooms



بخشی از cache که کلیدهای دستیابی را که دستیابی اشتراکی را پشتیبانی می کنند نگهداری می کنند

cache باید مشارکت دز پیدا کردن “راهنمای درست” به وسیله تطبیق داشته باشد

سازمان دهی cache

- ◆ کاملاً متحد: می تواند یک عنصر در هر بلاک باشد
- ◆ مسیر دهی بی واسطه و مستقیم: می تواند یک عنصر فقط در یک بلاک باشد
- ◆ چیدمان و قرارگیری پیوسته: می تواند یک عنصر در گروهی از بلاک ها باشد

پیوستگی کامل

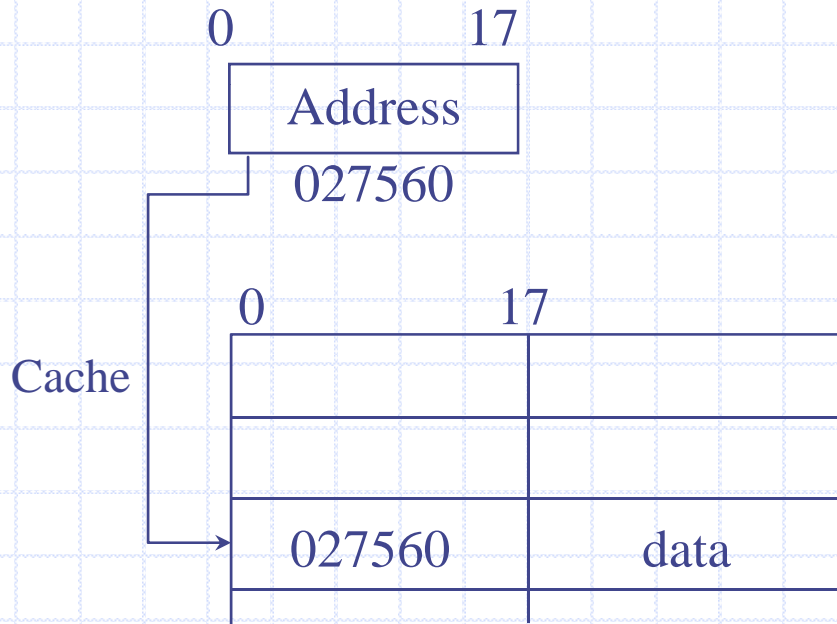
ساختار

- هر بلاک در حافظه می تواند در هر block-frame در cache باشد
- تمام ورودی های الگوی بلاک (block frame) با هم مقایسه می شوند (به وسیله جستجوی پیوسته)

یک بلاک = a word : ساده ترین مثال

تمام حافظه آدرس دهی کلمه ای شدن

“tag”

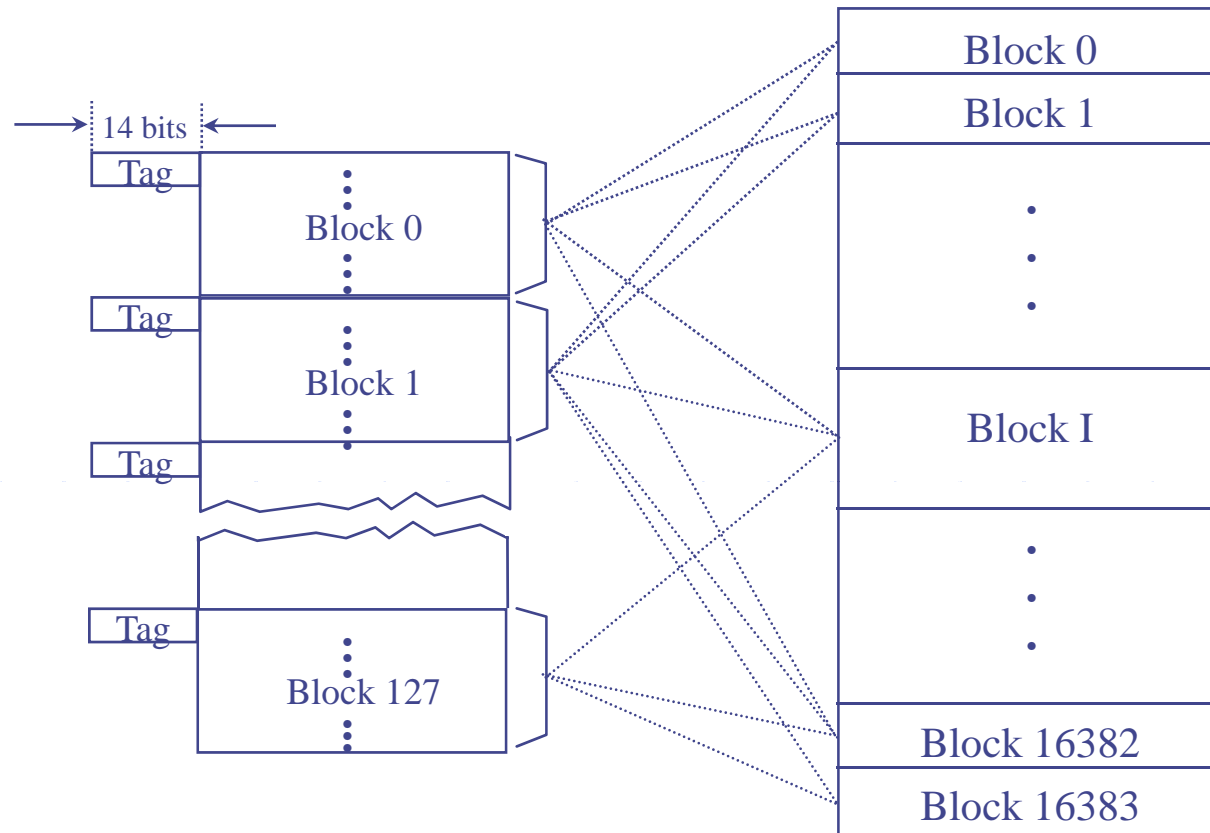


بسیار انعطاف "و با احتمال بالا به منظور ساکن شدن در cache"

مزایا: بدون اتلاف (با قابلیت سازمان دهی سریع)

معایب: overhead of associative search:

هزینه + زمان



آدرس حافظه اصلی



یدآوری: هر بلاک 16 کلمه دارد

سازمان کاملاً پیوسته cache

نگاشت مستقیم

◆ بدون پیوند پیوسته

◆ از آدرس M حافظه "به طور مستقیم" برای قاب بلاک در cache جایی

که بلاک باید قرار بگیرد شاخص دار شده. A comparison

then is to used to determine if it is a miss or hit.

نگاشت مستقیم

ساده ترین امتیاز:

سریع (با منطق کمتر)
هزینه کم: از این رو فقط یک دستگاه مقایسه کننده
لازم است که می تواند در فرم استاندارد حافظه قرار گیرد

معایب: اتلاف و بیهودگی

نگاشت مستقیم

◆ مثال:

زمانی که cache فقط 128 الگوی قاب دارد در نتیجه
درجه تسهیم چنین می باشد:

for addressing the corresponding frame or set of size 1.

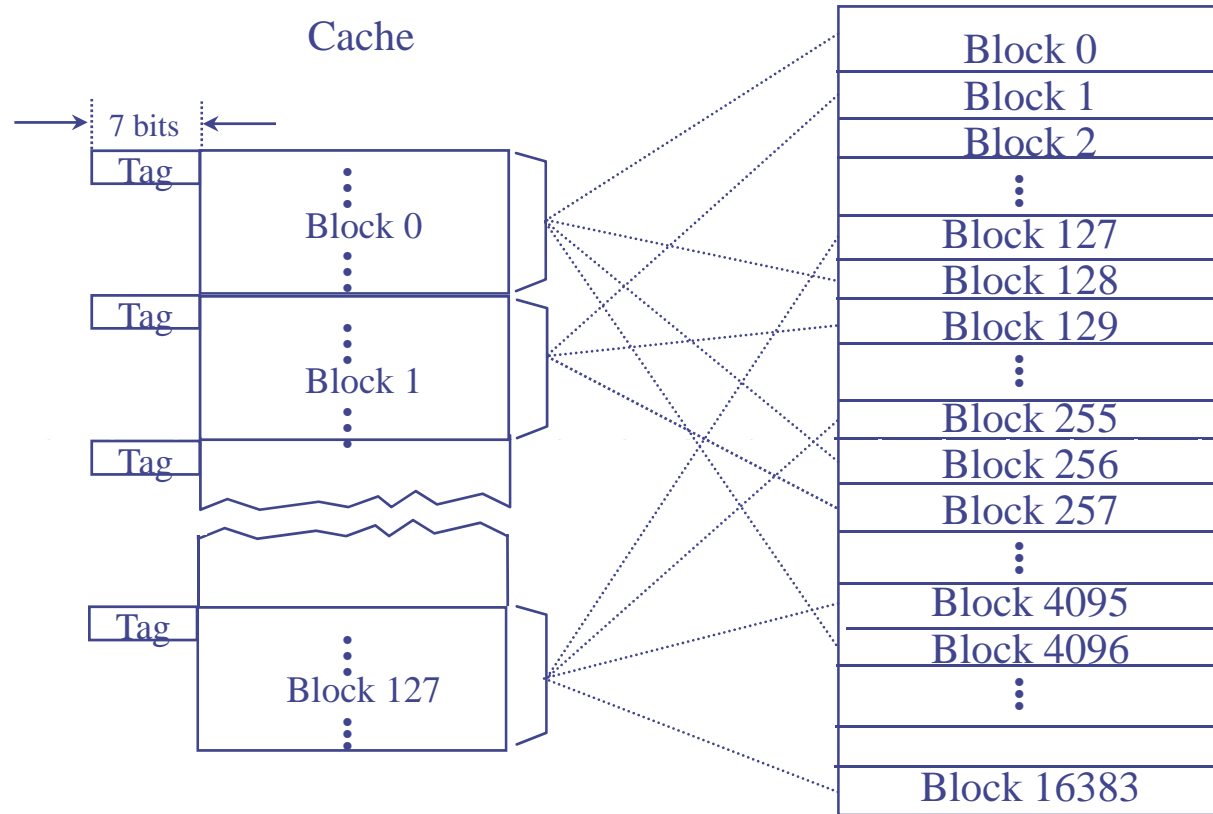
$$\frac{\text{اندازه حافظه اصلی}}{128 (2^7)} = \frac{16384 \text{ (block)}}{128} = 2^7 \text{ block/frame}$$

the high-order 7 bit is used as tag.

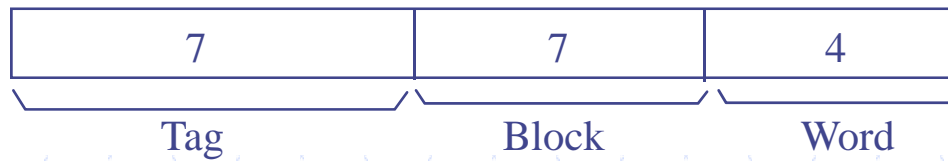
i.e. 2^7 blocks “fall” in one block frame.

Disadr: “frashing”

حافظه اصلی



آدرس حافظه اصلی



نگاشت مستقیم

نگاشت مستقیم

Cont'd

Mapping block addr mod (# of blocks in cache – in this case: mod (2^7))

مزیت: مرتبه پایین \log_2 (اندازه cache) بیت می تواند برای شاخص دار شدن استفاده شود

قرار گیری پیوسته

♦ توافق بین direct/full-associative

♦ Cache به چیدمانهای S تقسیم می شود

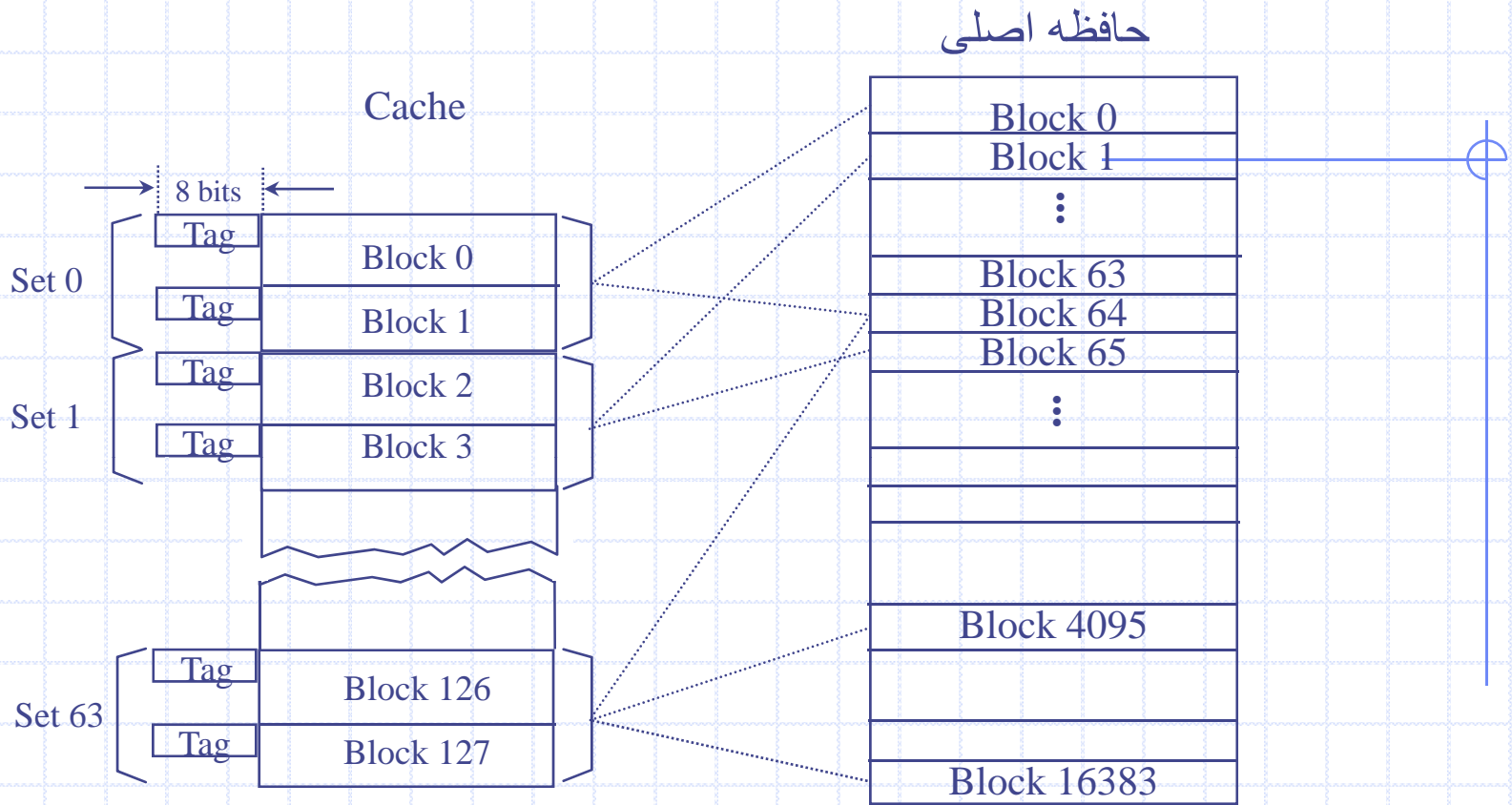
$$S = 2, 4, 8, \dots$$

♦ اگر M cache قاب داشته باشد که همه باهم باشند آنگاه داریم:

$$E = \frac{M}{S} \text{ blocks/set}$$

تعداد ساختارهای قابل دسترسی برای شاخص دار کردن

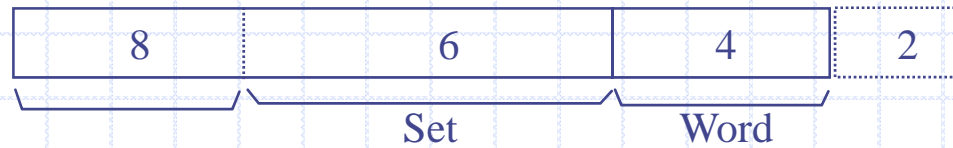
6 بیت می توانند در چیدمان سمت راست فهرست شوند و سپس 8 بیت برای پیوند پیوسته به کار خواهد رفت



آدرس حافظه اصلی



2 راه برای قرار گرفتن سازمان پیوسته



آماده دسترسی برای
شاخص دار کردن

*6 bit used to index
into the right "set"
higher order*

$$\frac{2^{14} (16k)}{2^6} = 2^8 \text{ block/set}$$

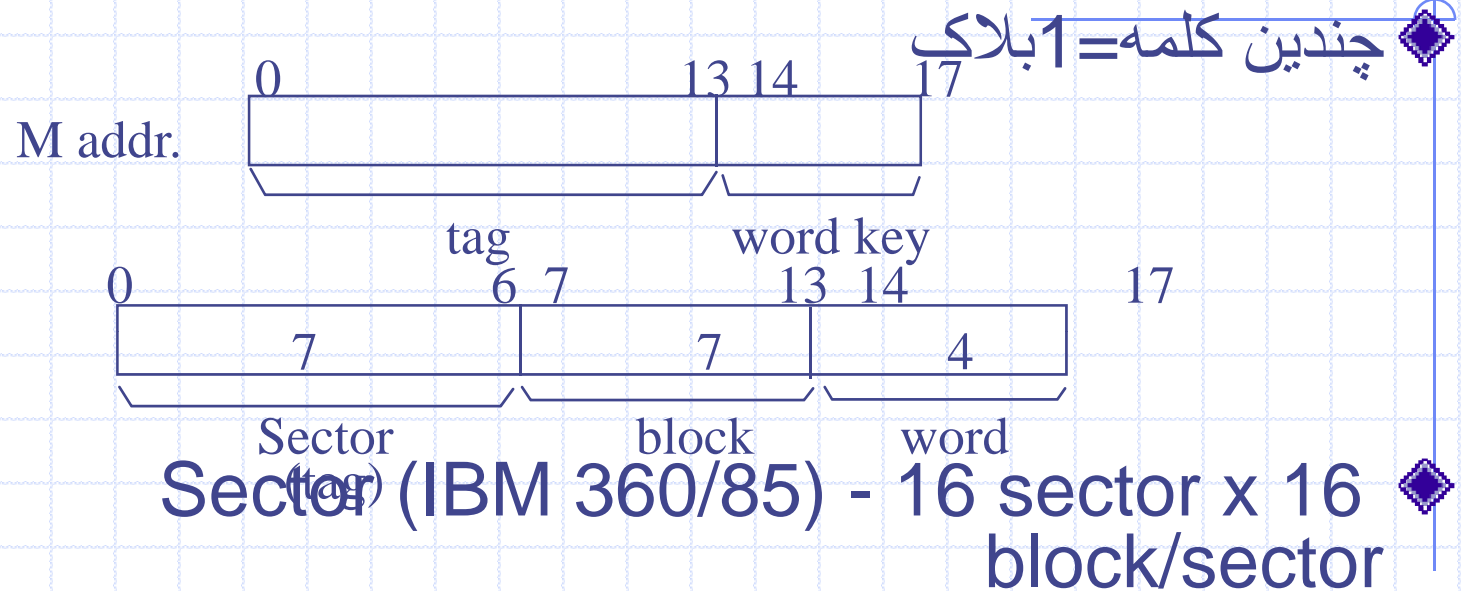
2 way

*8 bit used as tag
hence an associative
match of 8 bit with
the tags of the 2
blocks is required*

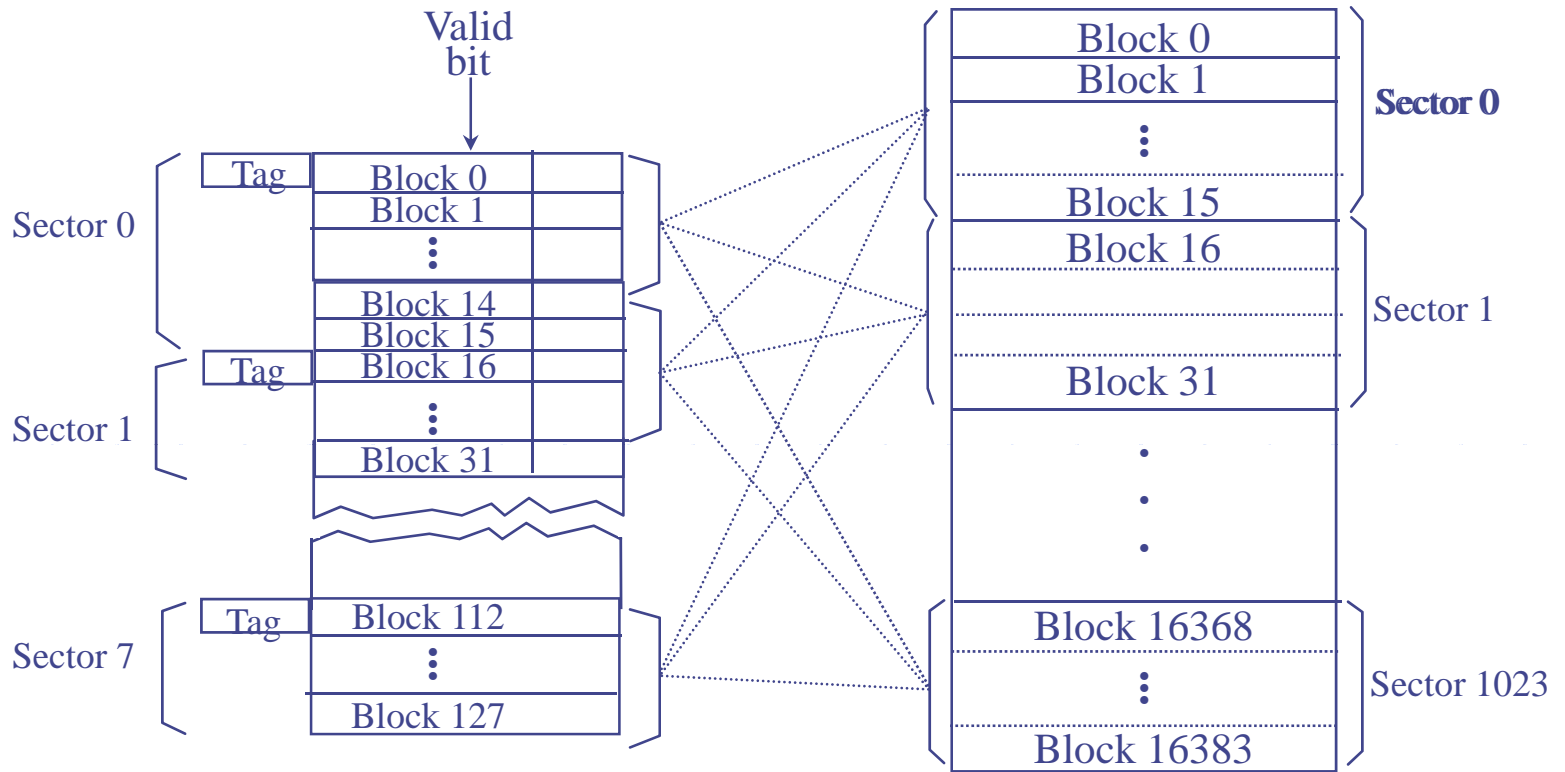
2^8 block/per set of 2 blocks

از این رو به هم پیوستن های 8 بیتی با
علامتهای 2 بیتی نیاز است

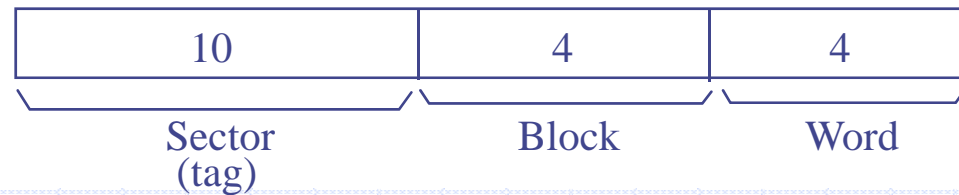
گونه های دیگر



- 1 سکتور: چندین بلاک متوالی
- Cache miss: جایگزینی سکتور
- بیت صحیح-یک بلاک بنا به درخواست انتقال داده شده



آدرس حافظه اصلی



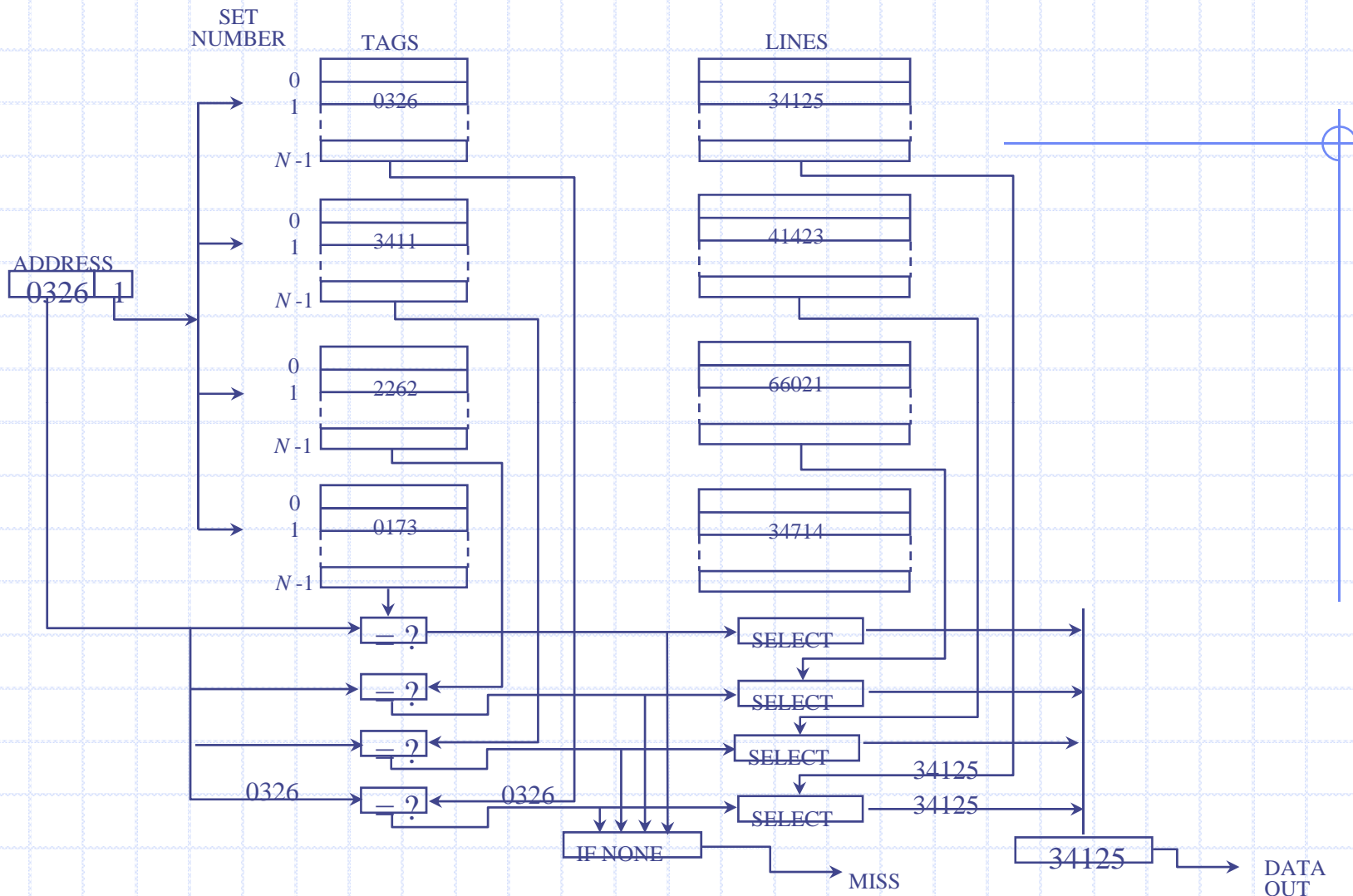
نگاشت سکتور cache

cont'd

$$\text{Cache دارد} \frac{128 \text{ blocks}}{16 \text{ blocks/sector}} = 8 \text{ sector}$$

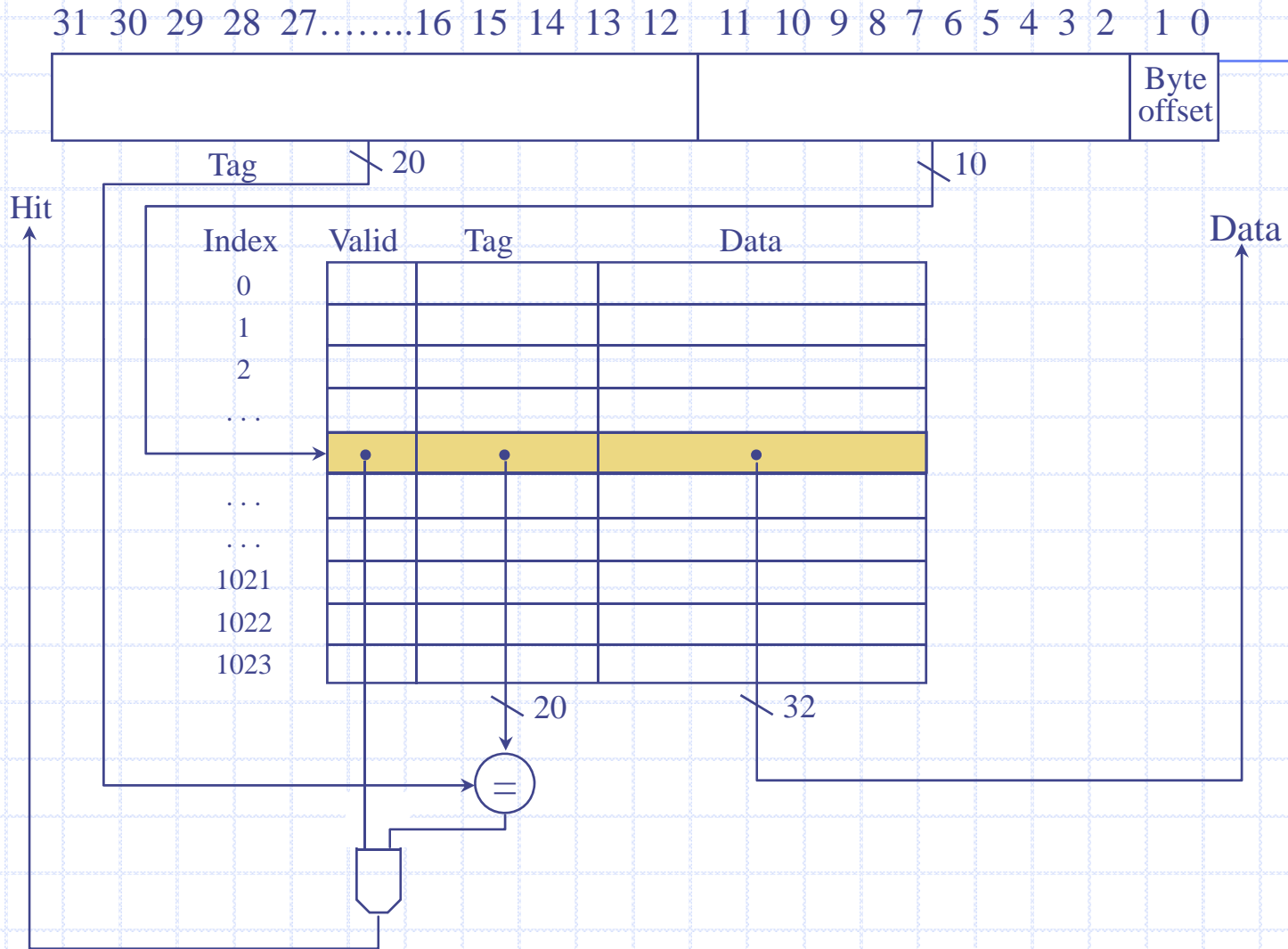
$$\text{حافظه اصلی دارد} \rightarrow \frac{16^k}{16} = 1K \text{ sectors}$$

cache نگاشت سکتور



The structure of a four-way set-associative cache with N sets

Address (showing bit positions)



جمع تمام بیت ها در cache

جمع بیتها =

(# of bits of a tag +
of bits of a block +
of bits in valid field) x Cache size

mips مثالی برای :

$$= ((32-14-2) + 32 + 1) \times 2^{14}$$

$$= 2^{14} \times 49$$

$$= 784 \text{ k bits}$$

$$\sim 100 \text{ kbytes}$$

$$= 64 \text{ K (bytes)}$$

$$= 2^{14} \text{ blocks with}$$

فرض کنید یک نگاشت مستقیم cache.

به روز کردن حافظه اصلی/خط و مشی واکنشی

• به روز کردن 8~16% are writes :

– باز نویسی

◆ یک نوشتن اشتباه شاید ربطی به حافظه اصلی نداشته

باشد

– write-through

◆ عبور و مرور بیشتر

• واکنشی

چه طور خواندن/نوشتن رابه کار ببریم

One-word-line: (DEC 3100) ◆

خواندن:آسان است ◆

- آدرس را به cache مناسب بفرستید آدرس را از یکی از این دو pc(برای خواندن یک دستورالعمل) یا از ALU(برای دستیابی داده ها)
- اگر سیگنال های Cache موفق شوند ، کلمه در خواست شده در خطوط داده ها موجود است.اگر که سیگنال های Cache دچار مشکل شوند، ما آدرس کامل را به حافظه اصلی می فرستیم. زمانی که حافظه داده ها را بر می گرداند، ما آن را در cache می نویسیم.

Read Miss is easy to be handled quickly: ◆

■ خواندن تگ و خواندن بلاک می تواند باهم انجام شود

◆ نوشتن معمولاً کندتر است

■ خواندن تگ و نوشتن بلاک نمی تواند باهم انجام شود

(one-word-line caches به جز برای)

برای خط های چند کلمه ای:

موقع نوشتن

on a write miss, it is a

خواندن



بلاک اصلی

- تغییر دادن -



یک بخش

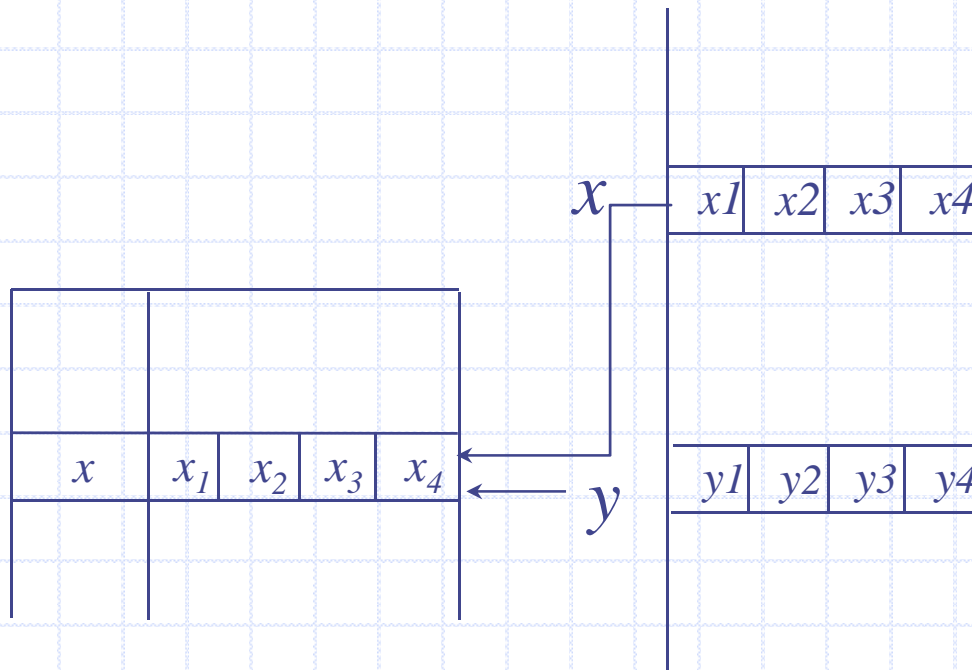
چرخه نوشتن



نوشتن بلاک

مقایسه و سنجش تگ نمی تواند به موازات آن انجام گیرد در نتیجه انجام آن آهسته تر است

چرا؟ فرض کنید x, y در یک ستون یکسان نگاشت شوند و در ابتدای x cache داریم



فرض کنید قبل از نوشتن، cache در برگیرنده سطور x است، زمانی که y نوشته می شود یک فقدان، و یک نوشته رخ می دهد

z y .3 نوشتن

Note after write miss, if not careful, we get

y	$x1$	$x2$	Z	$x4$
-----	------	------	-----	------

اما در پایان یک بازنویسی $y1, y2, y4$ ناپود می شوند

◆ – *Write through (or store through)* اطلاعاتی که نوشته می شود روی هر دو بلاک cache و سطح پایین تر حافظه

◆ – *Write back (also called copy back)* اطلاعاتی که تنها روی بلاک در cache نوشته می شود. قالب cache تغییر یافته روی حافظه اصلی تنها زمانی که جایگزین شود نوشته می شود

باز نویسی

پاک کردن و پر کردن بلاک

یک بیت پر نشان می دهد که آیا یک سطر تا زمانی که در cache است تغییر کرده است. **موقعی** که یک "dirty line" جایگزین شود، باید در حافظه اصلی بازنویسی شود.

دو تا گزینه در write miss وجود دارند

◆ تخصیص نوشتن (که به آن واکنشی به منظور نوشتن هم گفته می شود) بلاک بار می شود، به دنبال آن. This is similar to a read miss. by the write-hit actions above.

◆ *No write allocate* (also called *write around*) - The block is modified in the lower level and not loaded into the cache.

Think what you do when have a write miss!

چه موقع write through بهتر است؟

◆ زمانی که از cache دو سطحی استفاده می شود

Cpu روی یک تراشه cache کوچک است.

+

یک تراشه بزرگ cache

■ سازگاری آسانتر است

■ به وسیله دومین cache از ترافیک حافظه جلوگیری می شود

Write-back vs. Write-through

سرعت (باز نویسی سریع است) ◆

Traffic (در حالت عمومی، باز کپی بهتر است) ◆

If more than one read hit to the line ■

So attractive to multiprocessor in this sense ■

سازگاری cache (write-through بهتر است) ◆

منطقی (copy-back بسیار کامل است) ◆

Write-back vs. Write-through

◆ Buffering (4 برای write-through بهترین است)

به هر دو نیاز دارد، در حالی که copy-back فقط به یکی نیاز دارد .
مدیریت پیچیده است زیرا زمانی که یک ref ساخته می شود ، باید با
buffer همراه باشد

◆ قابلیت اطمینان (write-through بهتر است) برای اینکه حافظه اصلی

قابلیت شناسایی خطا را دارد

“انتخاب واضحی وجود ندارد” بر حسب کارایی

www.salampnu.com

سایت مرجع دانشجوی پیام نور

- ✓ نمونه سوالات پیام نور : بیش از ۱۱۰ هزار نمونه سوال همراه با پاسخنامه
- تستی و تشریحی
- ✓ کتاب ، جزوه و خلاصه دروس
- ✓ برنامه امتحانات
- ✓ منابع و لیست دروس هر ترم
- ✓ دانلود کاملاً رایگان بیش از ۱۴۰ هزار فایل مختص دانشجویان پیام نور

www.salampnu.com