

מבוא למחשבים

Pipeline Processing

השקפים מבוססים על הספרים

פרק 9 – Mano

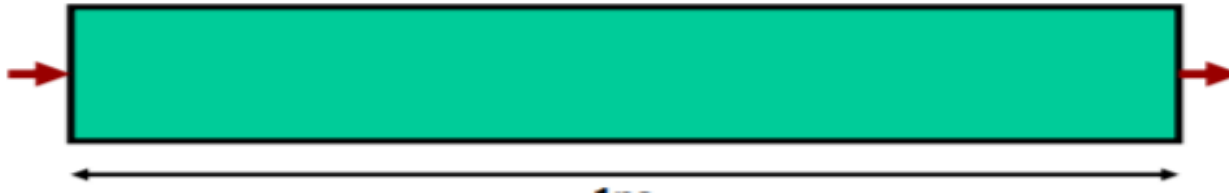
Hennessy & Patterson

ד"ר רון שמואלי

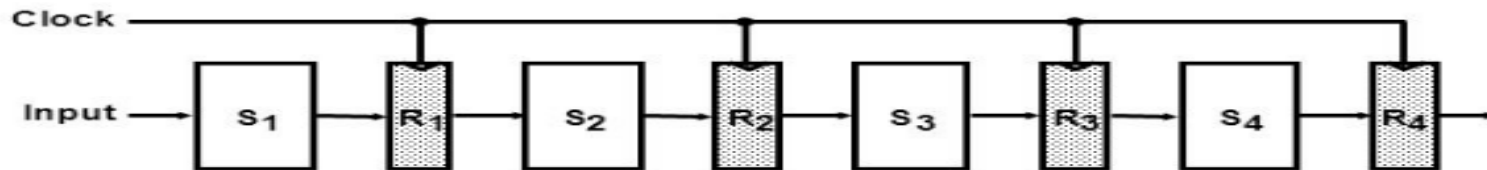
rshmueli@bgu.ac.il

Pipeline Processing

- Key idea: break big computation into pieces



General Structure of a 4-Segment Pipeline



Space-Time Diagram

	1	2	3	4	5	6	7	8	9	
Segment 1	T1	T2	T3	T4	T5	T6				→ Clock cycles
2		T1	T2	T3	T4	T5	T6			
3			T1	T2	T3	T4	T5	T6		
4				T1	T2	T3	T4	T5	T6	

חישוב גורם האצה של ה- PIPELINE

- **n - מספר המשימות לביצוע / K - מספר הסיגמנטים.**
- **במחשב ללא Pipelined.**
 - T - הזמן להשלים משימה.
 - $n * T$ - הזמן הנדרש להשלמת n משימות.

- **Pipelined במכונת**

- t_p - זמן מחזור של השעון (זמן לסיום כל סיגמנט)
- $K t_p$ - זמן להשלמת המשימה הראשונה.
- $(n-1) t_p$ - הזמן להשלמת n-1 המשימות הנותרות.
- $(k+n-1) t_p$ - הזמן להשלמת n משימות.

- **גורם ההאצה Speedup**

$$S_k = n * T / (k + n - 1) * t_p$$

$$\lim_{n \rightarrow \infty} S_k = \frac{T}{t_p} = k, \quad (\text{if } t_n = k * t_p)$$

בהנחה $T = k t_p$

תוצאות מיטביות

ה- PIPELINE תמיד מלא
גורם האצה תיאורטי הוא K
(K מספר הסגמנטים)

הגבול המעשי של Pipeline

- Pipeline - לא יכול להגיע למהירות התאורטית

– סיגמנטים שונים דורשים זמן ביצוע שונה \leftarrow תאום תדר שעון לסגמנט האיטי \leftarrow זמן המתנה בסגמנטים האחרים.

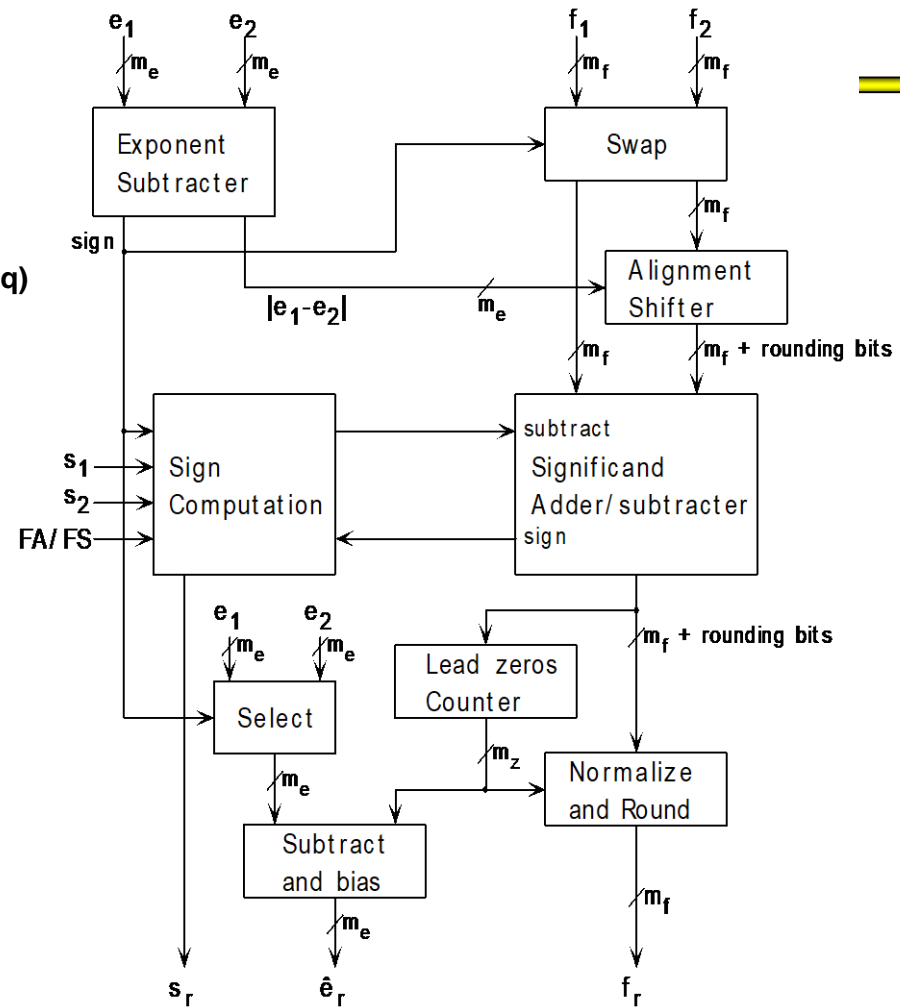
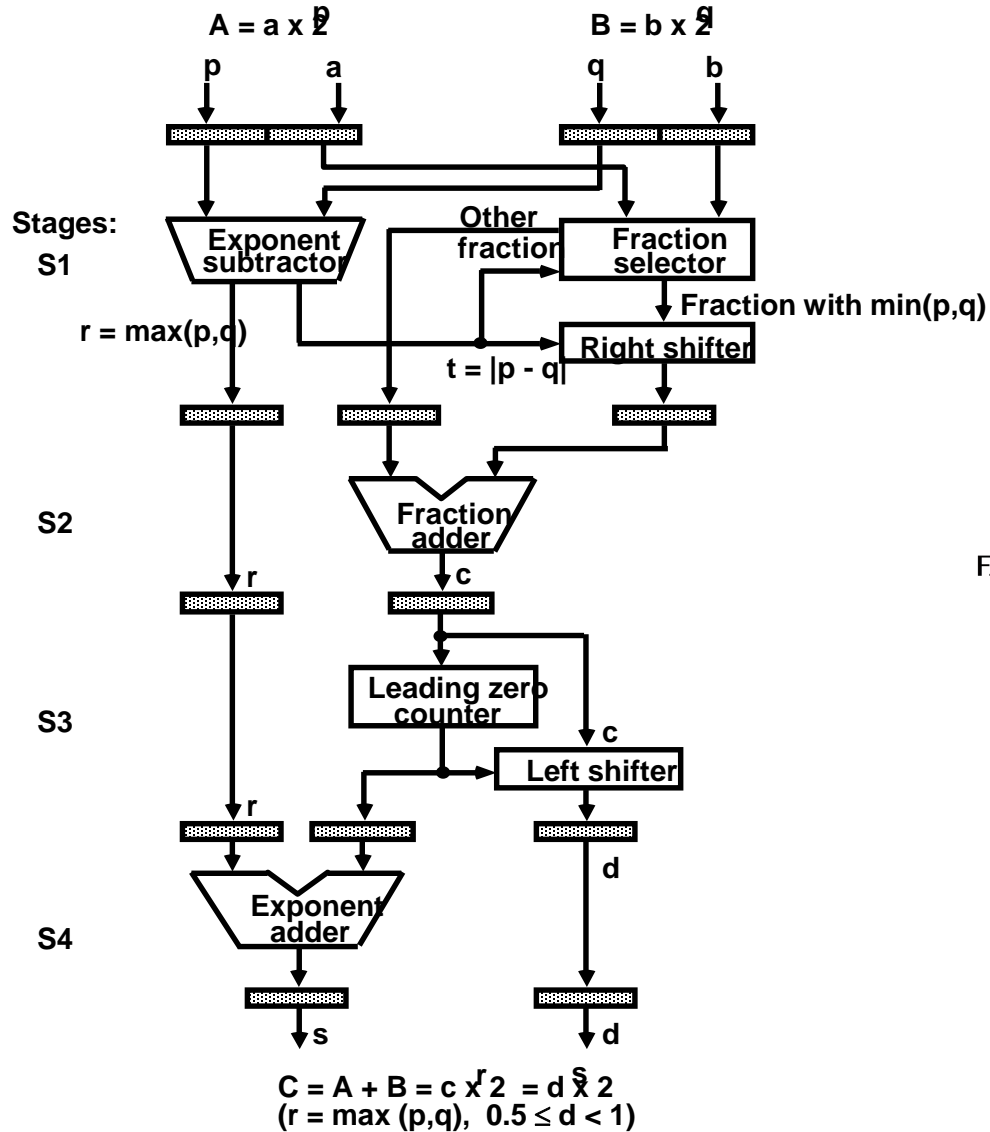
- בזבוז זמן עם אוגרי הביניים (לא נדרשים במערכת רגילה).

- Pipeline – טוב לשימוש:

- פעולות אריתמטיות

- ביצוע במקביל של מחזור פקודה (fetch-decode-execute)

4-STAGE FLOATING POINT ADDER



חישוב גורם האצה

- בהנחה שהזמן הנדרש לסיגמנטים השונים הוא:

$$t_1=60\text{ns}, t_2=70\text{ns}, t_3=100\text{ns}, t_4=80\text{ns} \quad -$$

$$- \text{ההשהיה של האוגרים } t_r=10\text{ns}$$

- זמן מחזור (כל מחזור מתקבלת תפוקה):

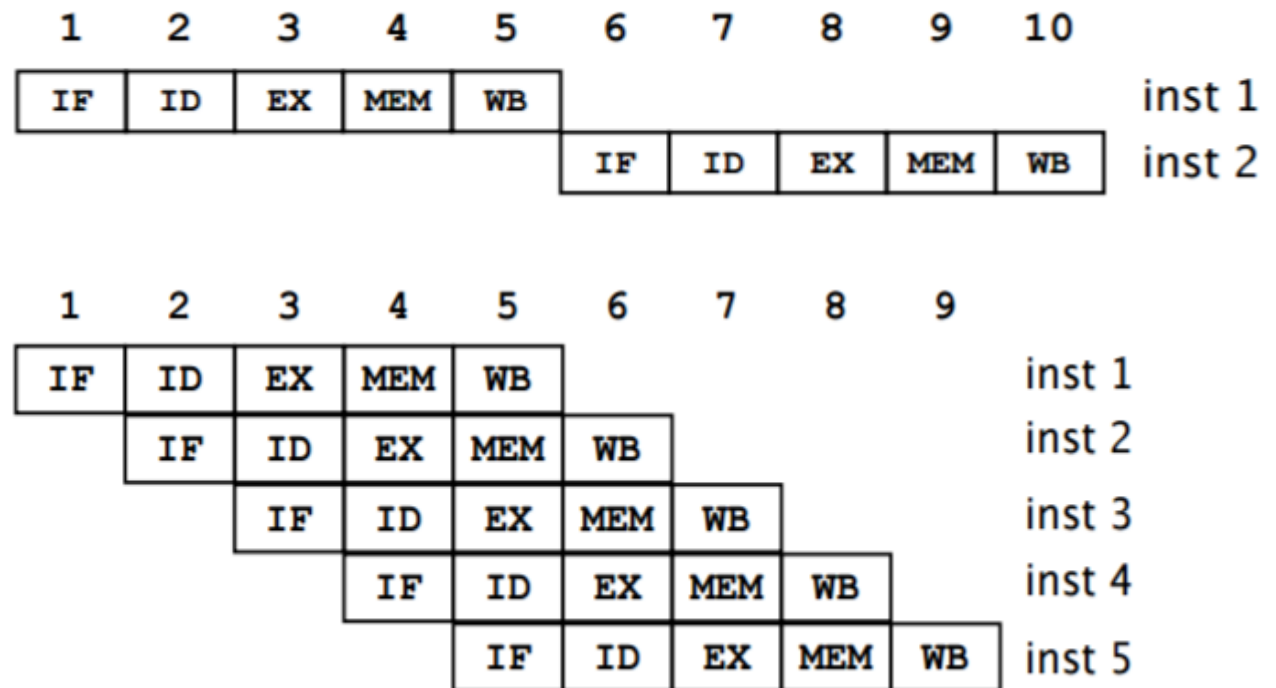
$$t_p = t_3+t_r=100+10=110$$

- זמן ביצוע ללא Pipeline (סכום הזמנים):

$$t_n = t_1+ t_2+ t_3+t_4+t_r= 60+70+100+80\text{ns}= 320$$

- $S=t_n/t_p=320/110=2.9$

CPU cycle



The MIPS CPU

MIPS - Microprocessor without Interlocked Pipeline Stages

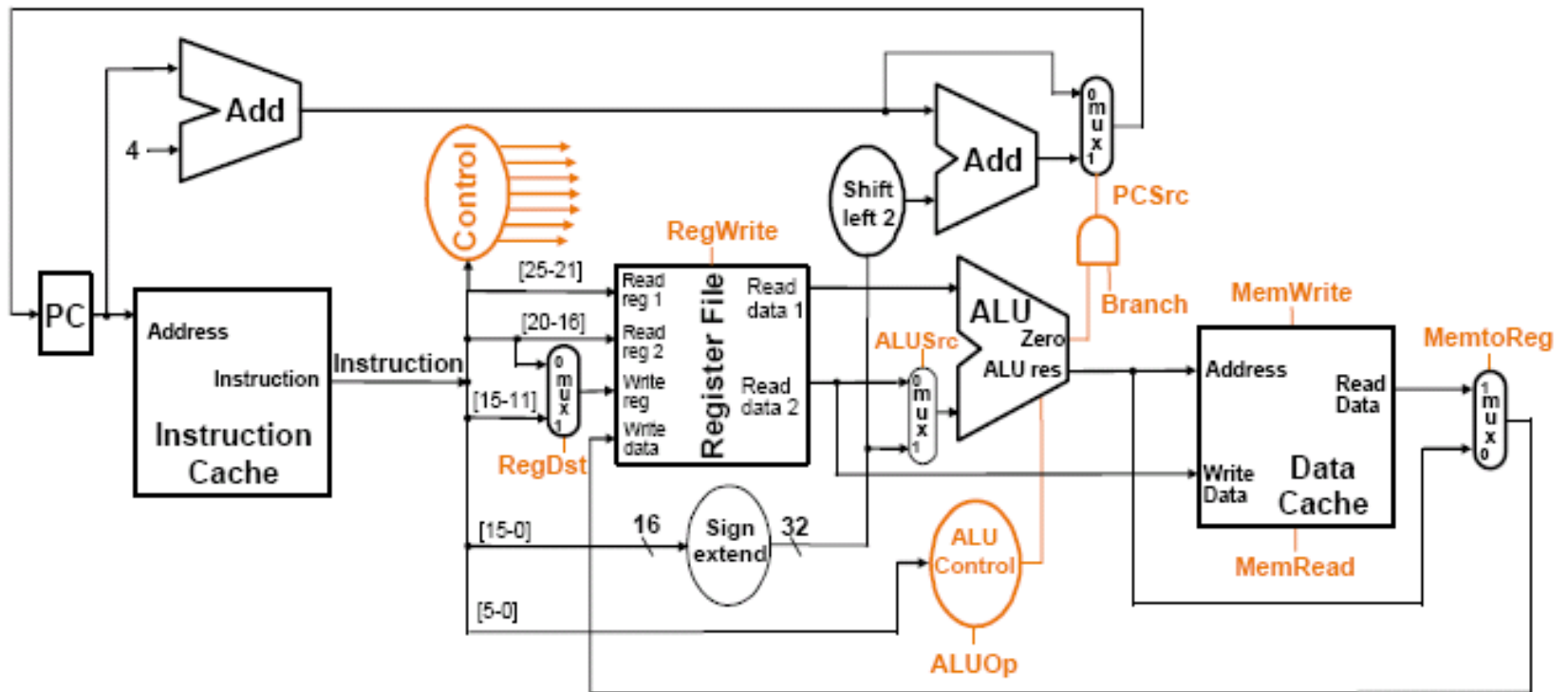
Instruction
fetch

Instruction
Decode /
register fetch

Execute /
address
calculation

Memory
access

Write
back



MIPS Instruction Formats

- R-type (register insts)

31	26	21	16	11	6	0
op	rs	rt	rd	shamt	funct	
6 bits	5 bits	5 bits	5 bits	5 bits	6 bits	
- I-type (Load, Store, Branch, inst's w/imm data)

31	26	21	16	0
op	rs	rt	immediate	
6 bits	5 bits	5 bits	16 bits	
- J-type (Jump)

31	26	0
op	target address	
6 bits	26 bits	

op: operation of the instruction

rs, rt, rd: the source and destination register specifiers

shamt: shift amount

funct: selects the variant of the operation in the “op” field

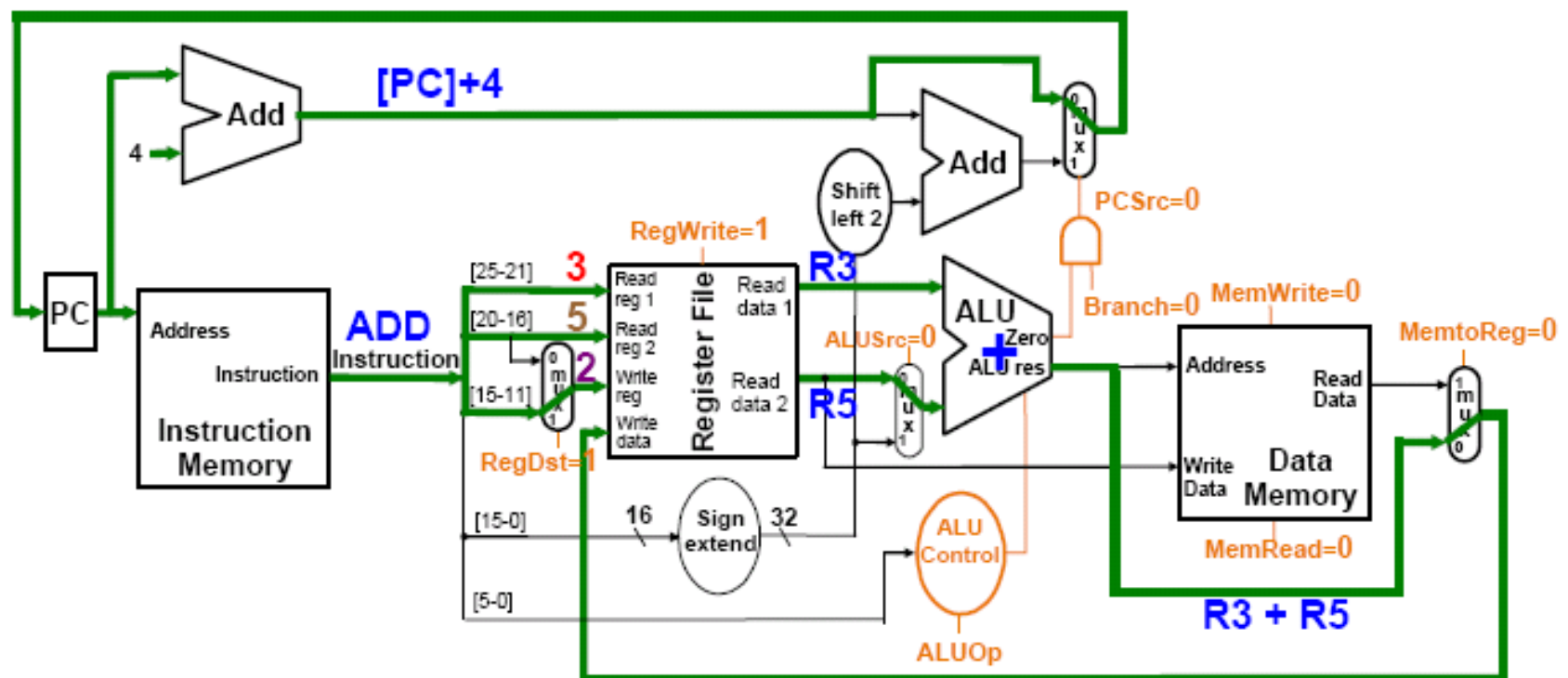
address / immediate: address offset or immediate value

target address: target address of the jump instruction

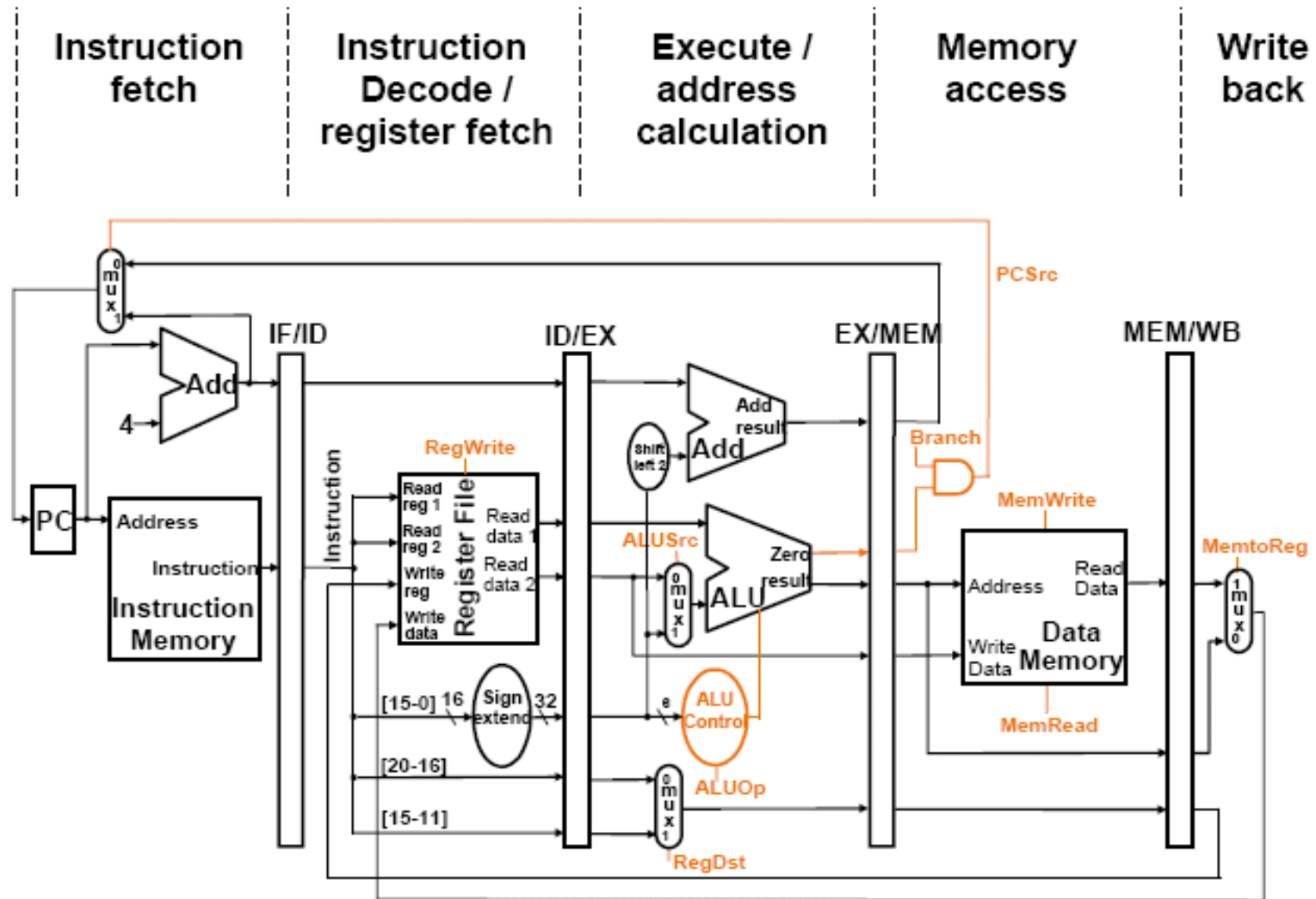
Executing an Add Instruction

Add R2, R3, R5 ; $R2 \leftarrow R3 + R5$

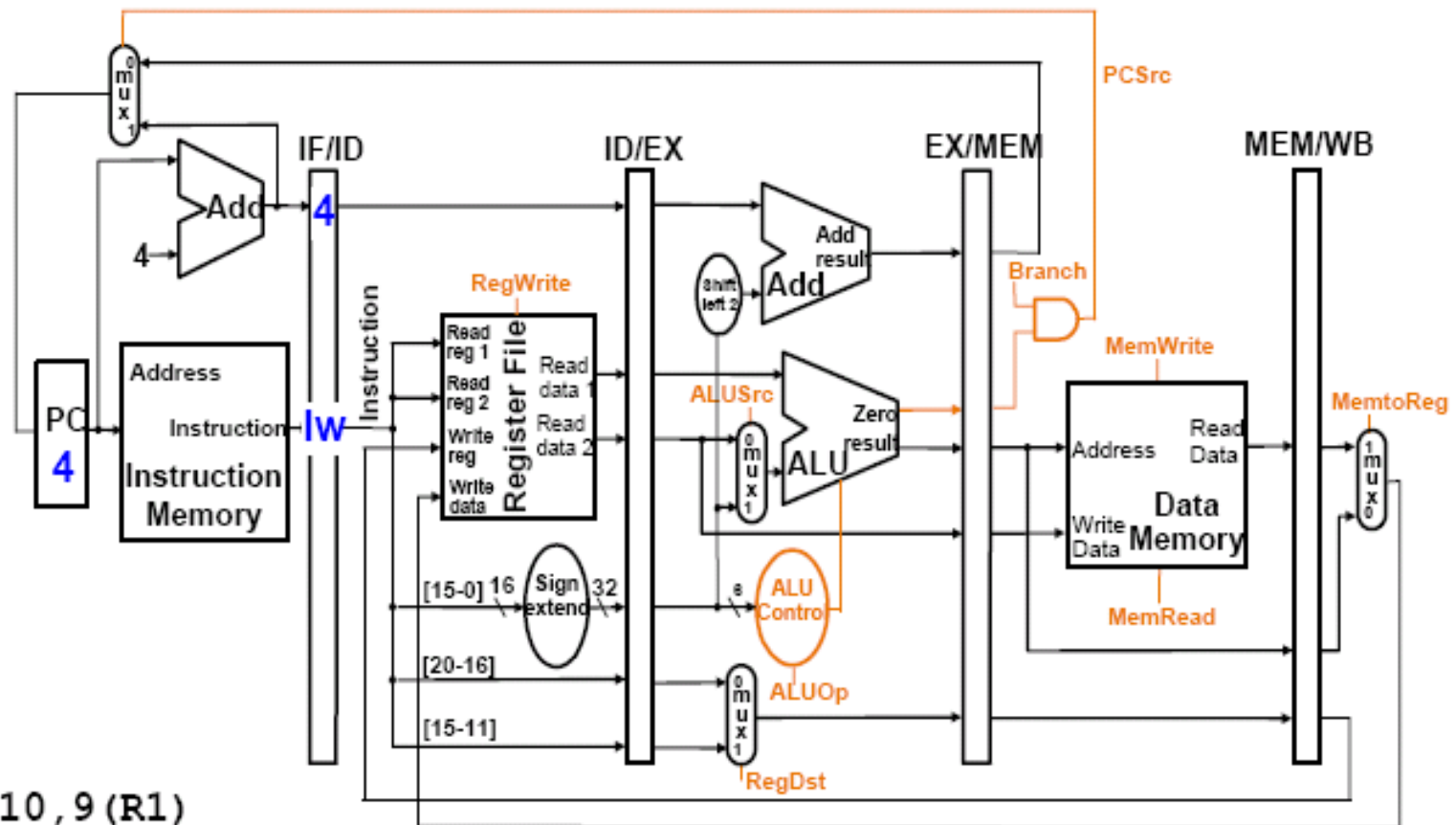
31	op	26	rs	21	rt	16	rd	11	shamt	6	funct	0
	ALU		3		5		2		0		0 = Add	



Pipelined CPU

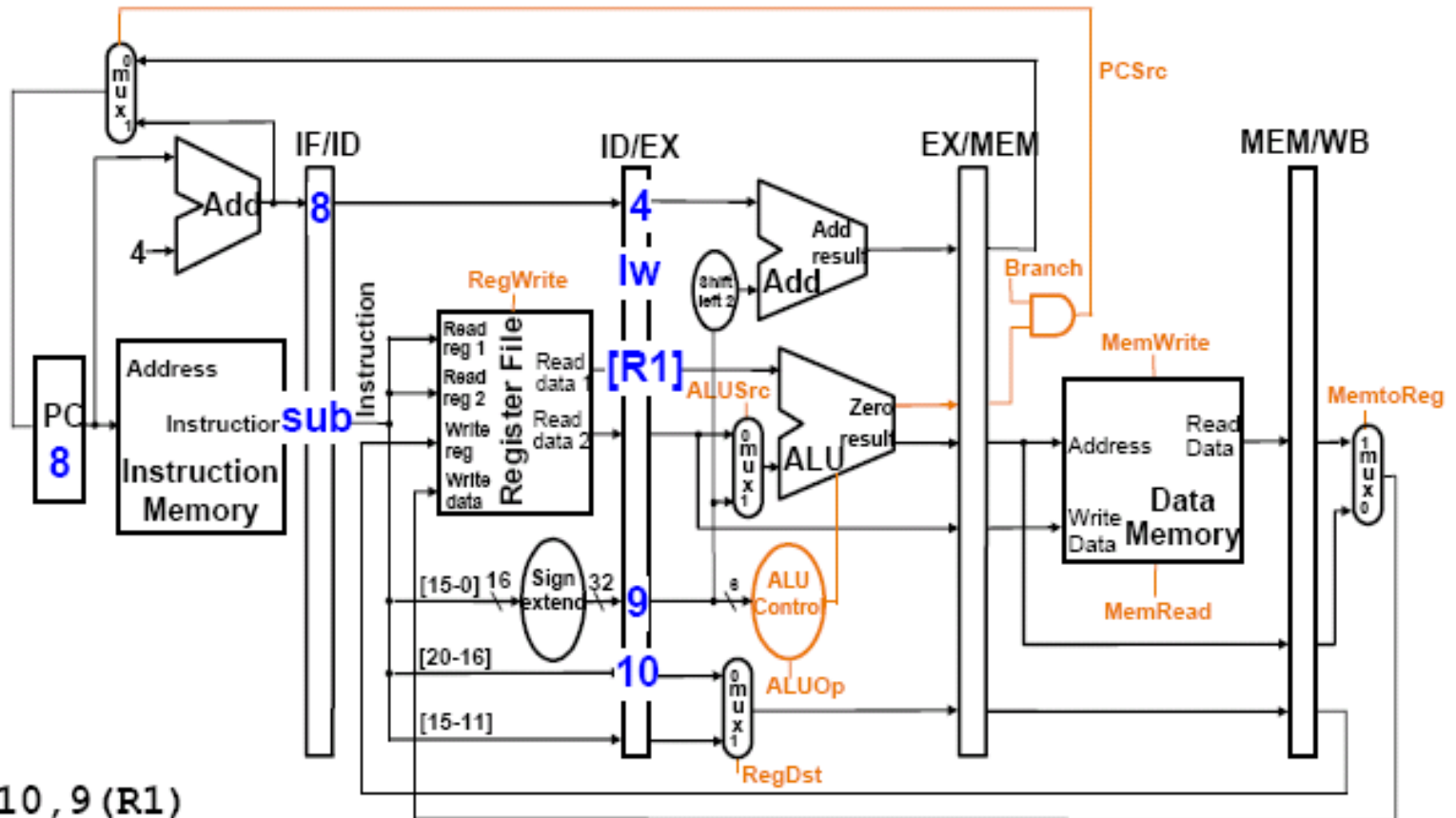


Pipeline Example: cycle 1



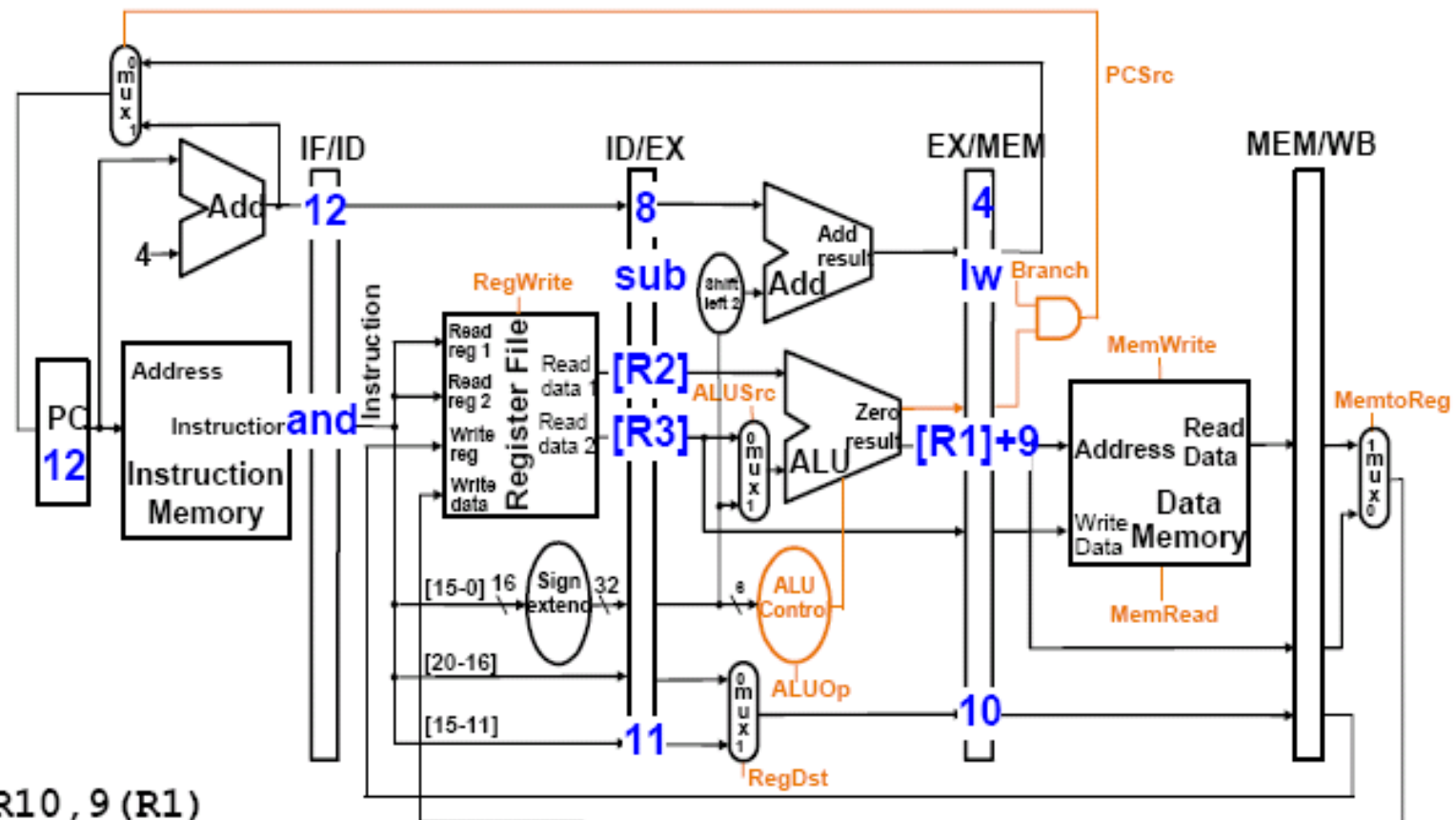
```
lw R10, 9(R1)
sub R11, R2, R3
and R12, R4, R5
or R13, R6, R7
```

Pipeline Example: cycle 2



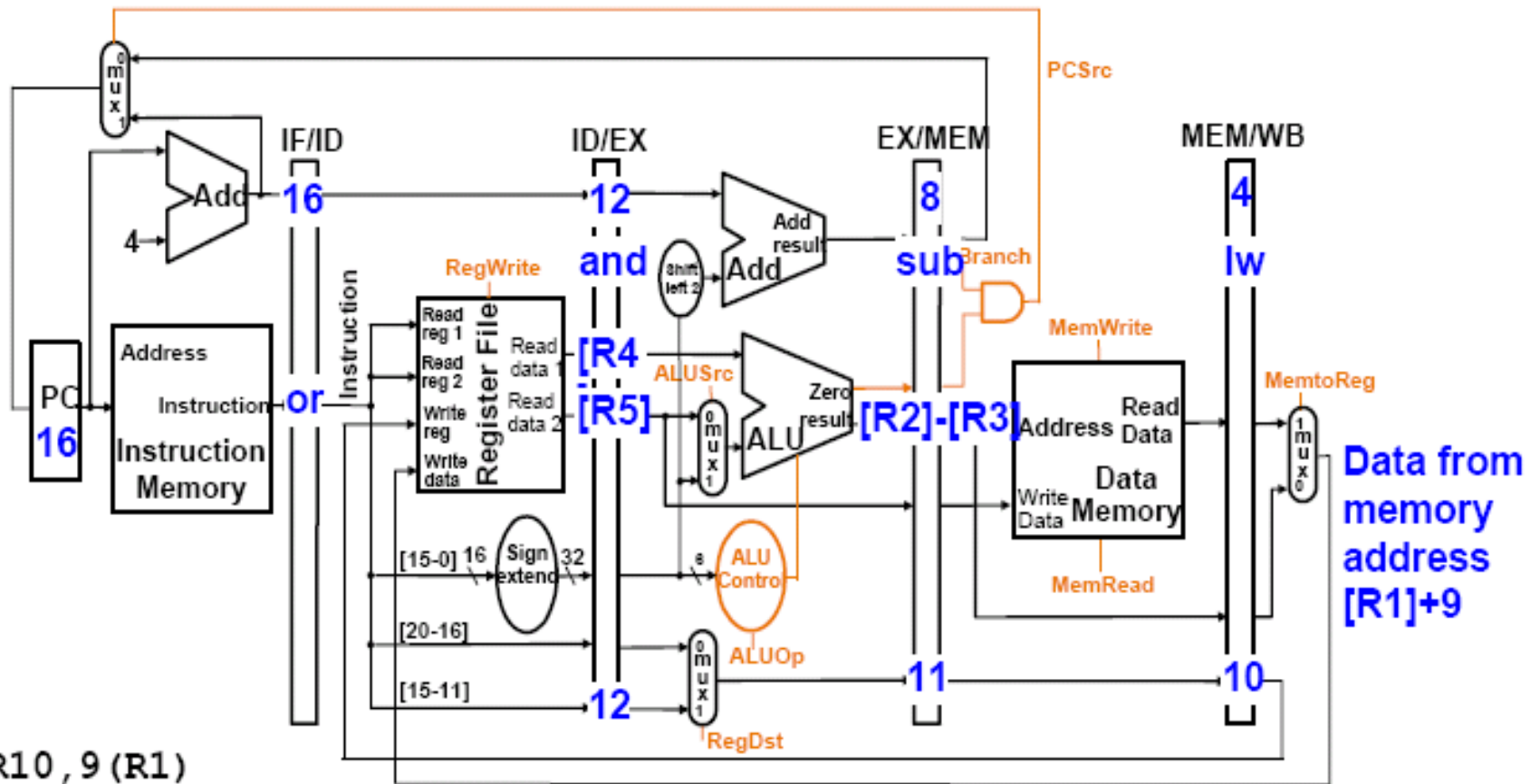
```
lw    R10,9(R1)
sub   R11,R2, R3
and   R12,R4, R5
or    R13,R6, R7
```

Pipeline Example: cycle 3



```
lw    R10, 9(R1)
sub   R11, R2, R3
and   R12, R4, R5
or    R13, R6, R7
```

Pipeline Example: cycle 4



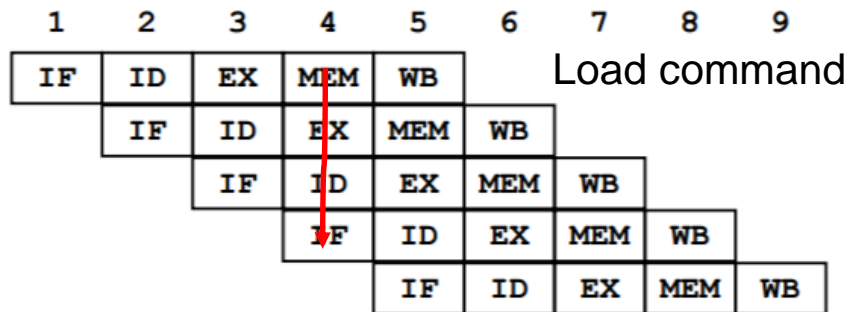
```
lw  R10, 9(R1)
sub R11, R2, R3
and R12, R4, R5
or  R13, R6, R7
```

סיכונים ב- Pipeline - תקציר

- **Structural hazards** - גישה לזיכרון ע"י שני סגמנטים.
- **חציית אוגרים** – גישה לאוגרים לקריאה וכתיבה.
- **Data hazards** - תלות בין הנתונים המתבצעים.
 - **Hardware interlock** – מעגל המשהה את הפקודה
 - **Operand forwarding** – ניתוב תוצאה במסלול חלופי .
 - **Delayed load** - הקומפיילר מסדר את הפקודות מחדש.
- **Control Hazards** - בקרת זרימה-הסתעפות.

Structural Hazards

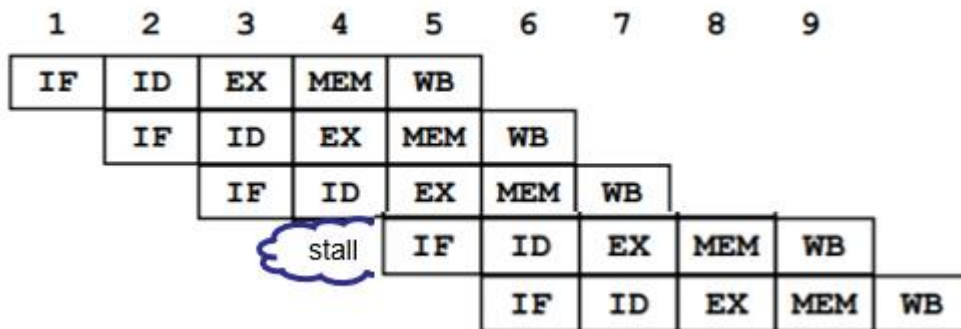
– גישה לזיכרון ע"י שתי בפקודות, בשלבי ביצוע שונים.



פתרונות

– הכנסת STALL

– הפרדת זיכרונות



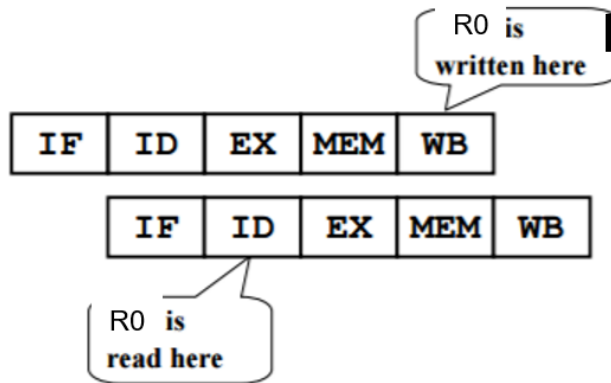
What if Instruction 1 is also a LOAD?

Data Hazards (1)

- תלות בין הנתונים המתבצעים

$R0 \leftarrow R1 + R2$

$R4 \leftarrow R3 + R0$

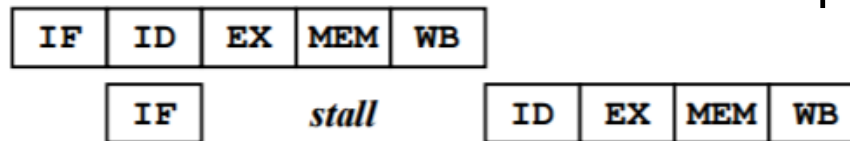


- **Hardware interlock**

– הכנסת Stalls עד שהפקודה זמינה

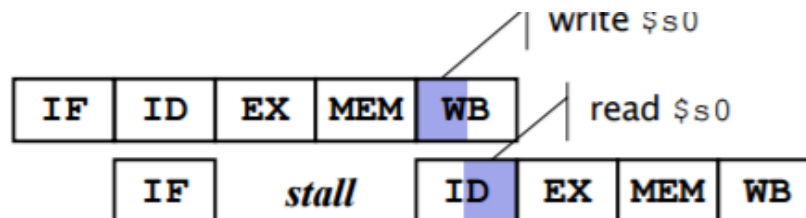
$R0 \leftarrow R1 + R2$

$R4 \leftarrow R3 + R0$



- **חציית אוגרים**

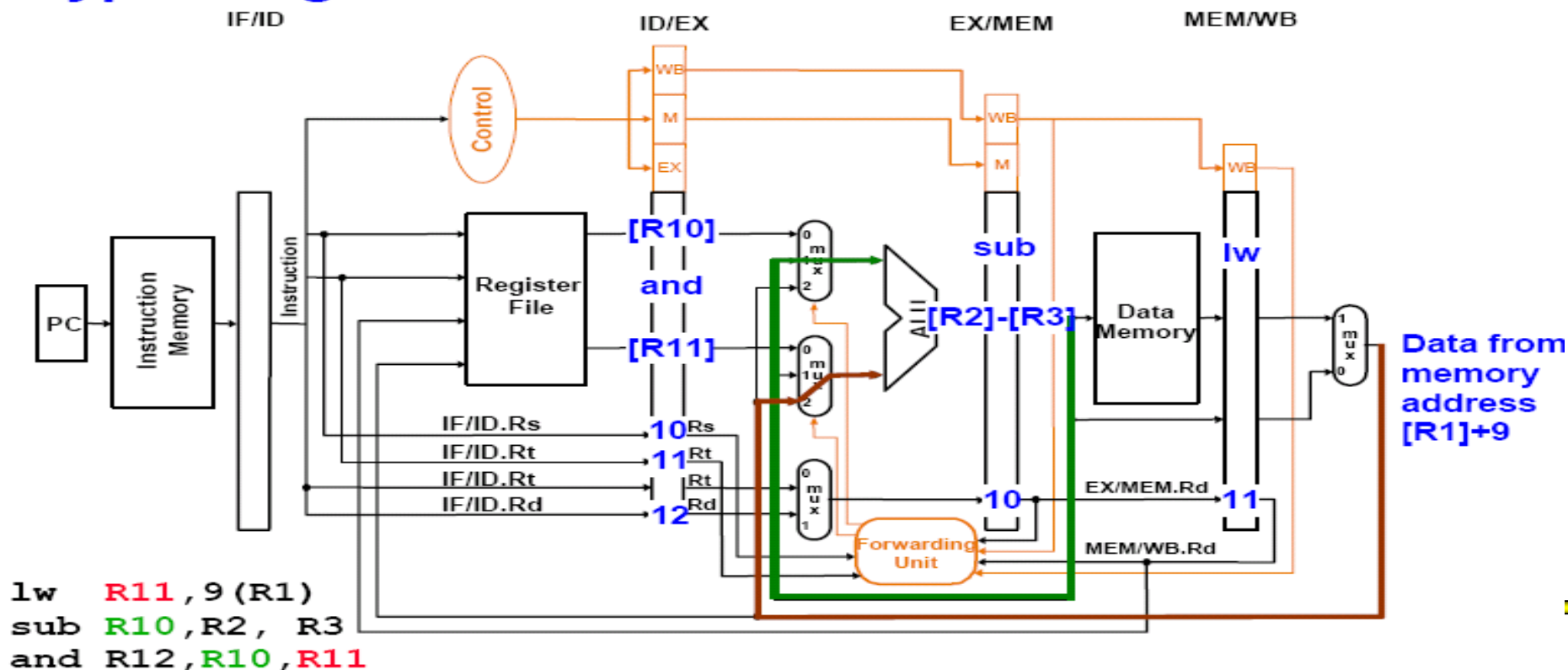
– כתוב לרגיסטר בחצי שיעון ראשון, קרא בחצי שיעון שני



Data Hazards – Operand Forwarding

- חומרה שמגלה תלות בין פקודות, ומנתבת תוצאה במסלול חלופי ישירות לסגמנט המתאים.

Forwarding Hardware Example: Bypassing From EX to Src1 and From WB to Src2

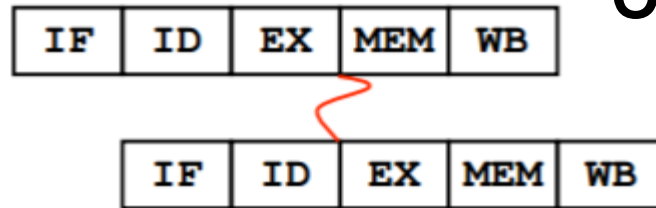


דוגמא

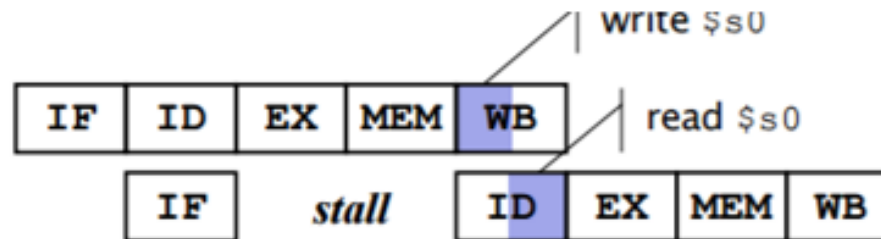
Operand Forwarding

$R0 \leftarrow R1 + R2$

$R4 \leftarrow R3 + R0$



ביחס ל HW Interlock



Data Hazard- Delayed Load

- הקומפייילר מסדר את הפקודות מחדש במידת הצורך מכניס NOPs

דוגמא - Delayed Load+ Forwarding

$a = b + c;$
 $d = e - f;$

סידור פקודות מחדש ע"י המהדר

Unscheduled code:

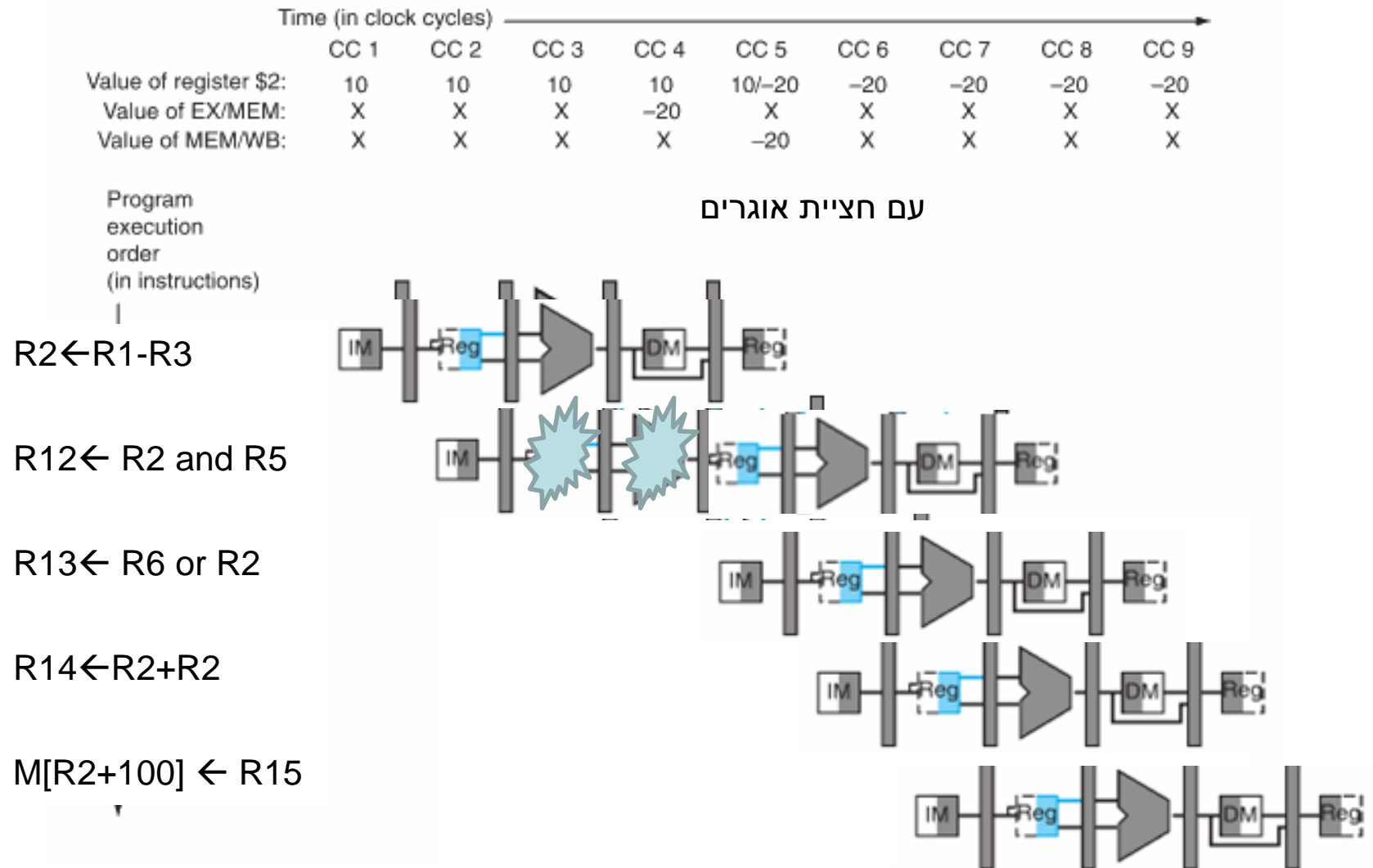
```
ld    r2, b(0)
ld    r3, c(0)
→ add r1, r2, r3
→ st  r1, a(0)
ld    r5, e(0)
ld    r6, f
→ sub r4, r5, r6
→ st  r4, d
```

Scheduled Code:

```
ld    r2, b(0)
ld    r3, c(0)
ld    r5, e(0)
add   r1, r2, r3
ld    r6, f
st    r1, a(0)
sub   r4, r5, r6
nop   (or stalls by the HW)
st    r4, d(0)
```

Delayed Load - A load requiring that the following instruction not use its result

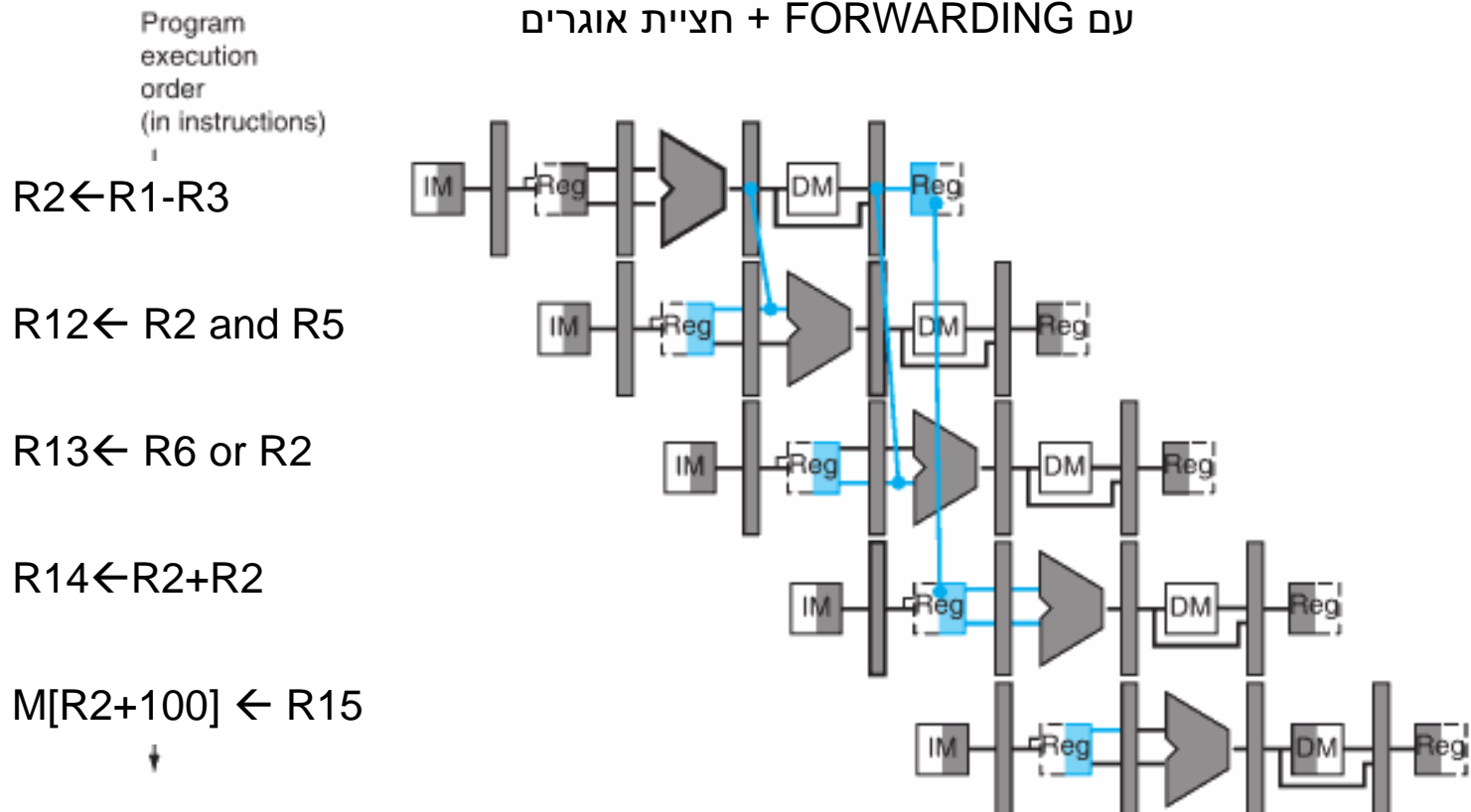
Data hazards – HW interlock



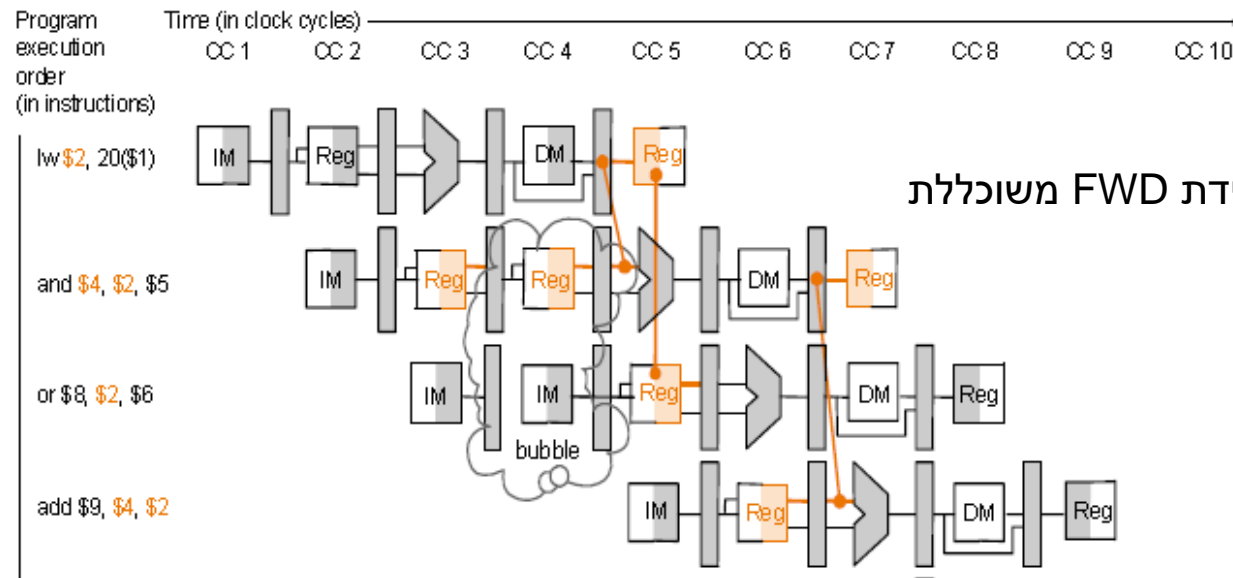
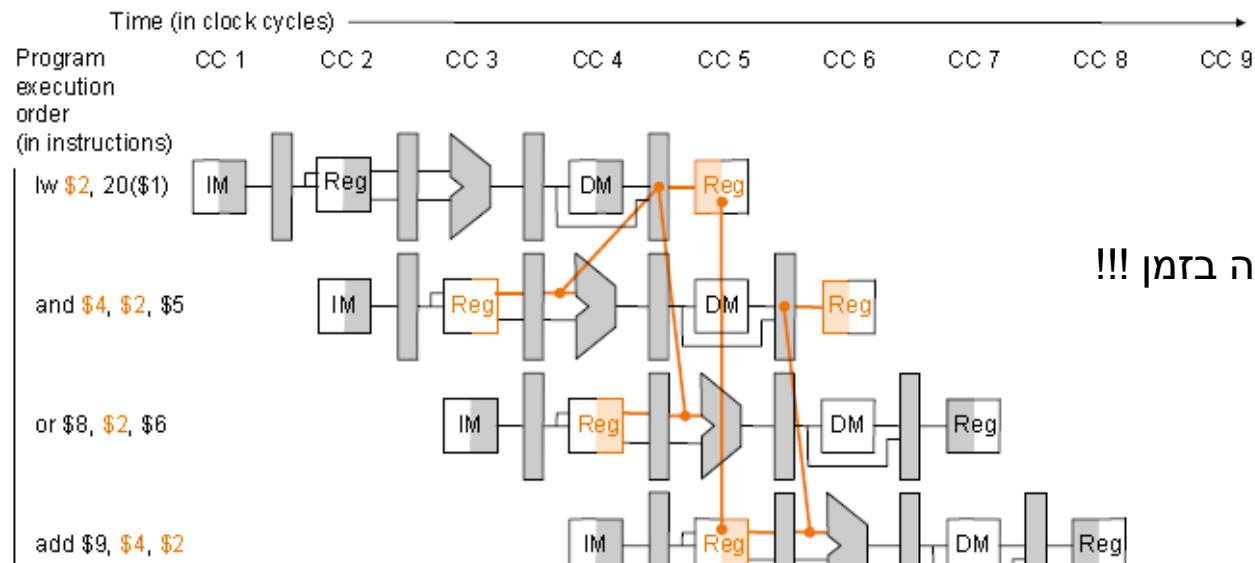
Data hazards - Forwarding

	Time (in clock cycles) →								
	CC 1	CC 2	CC 3	CC 4	CC 5	CC 6	CC 7	CC 8	CC 9
Value of register \$2:	10	10	10	10	10/-20	-20	-20	-20	-20
Value of EX/MEM:	X	X	X	-20	X	X	X	X	X
Value of MEM/WB:	X	X	X	X	-20	X	X	X	X

עם FORWARDING + חציית אוגרים



Data Hazards that cannot solved by FWD



Control Hazards - בקרת זרימה-הסתעפות.

– Branch Prediction - ניסיון לנחש את תוצאת ההסתעפות – והבאת הפקודות המתאימות בהתאם לניחוש

– Branch target Buffer (BTB) (ניהול הזיכרון אסוציאטיבי)

– Delayed Branch - המהדר מסדר את הפקודות מחדש, מכניס פקודות שימושיות – כך שהצינור יישאר מלא כאשר ישנה הסתעפות

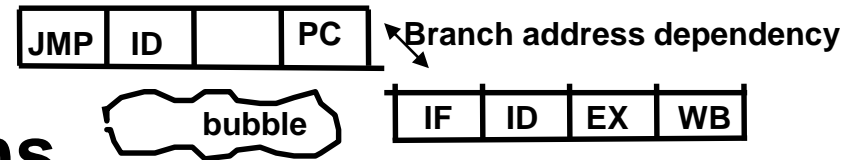
– PRE-FETCH ביצוע *FETCH* לכתובת הבאה ולכתובת הסתעפות.

– Loop Buffer (High Speed Register file) – שמירת הלולאה כולה כולל ההסתעפות באוגרים הלולאה כולה מתבצעת ללא גישה לזיכרון.

Control Hazard on Branches (1)

- **Static Option 1: Stall**

- Stall pipe when branch is encountered until resolved



- **Stall impact: assumptions**

- CPI = 1
 - 20% of instructions are branches
 - Stall 3 cycles on every taken branch

- **$\text{CPI}_{\text{new}} = 1 + 0.2 \times 3 = 1.6$**

- $(\text{CPI}_{\text{new}} = \text{CPI}_{\text{Ideal}} + \text{avg. stall cycles} / \text{instr.})$

- **We lose 60% of the performance**

Control Hazard on Branches (2)

- **Static Option 2: Predict Not Taken**
 - **Execute instructions from the fall-through (not-taken), path**
 - As if there is no branch
 - If the branch is not-taken (~50%), no penalty is paid
 - **If branch actually taken**
 - Flush the fall-through path instructions before they change the machine state (memory / registers).
 - Fetch the instructions from the correct (taken) path
 - **Assuming ~50% branches not taken on average**
 - $\text{CPI}_{\text{new}} = 1 + (0.2 \times 0.5) \times 3 = 1.3$

Branch Prediction

- חיזוי דינמי נעשה בהתאם להיסטוריה.

- טכניקה 1 – סיבית חיזוי בהתאם להסתעפות האחרונה (Taken/Not Taken).

- נניח מצב התחלתי דגל חיזוי Not Taken =

- בכניסה ללולאה Miss, דגל חיזוי Taken =

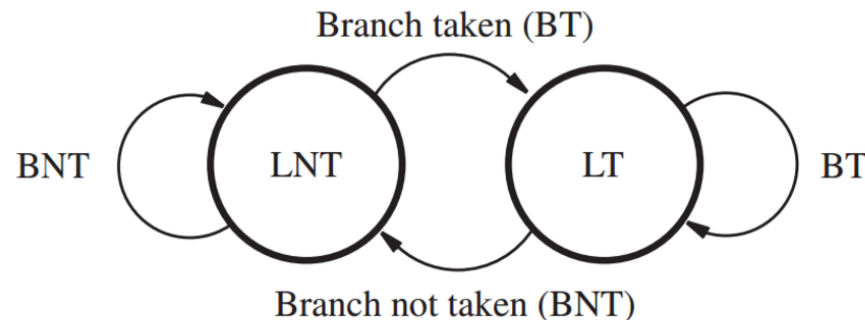
- תשע פעמים חיזוי נכון (HIT)

- בסוף הלולאה MISS דגל חיזוי Not Taken =

- פגיעה ב- 8 מתוך 10

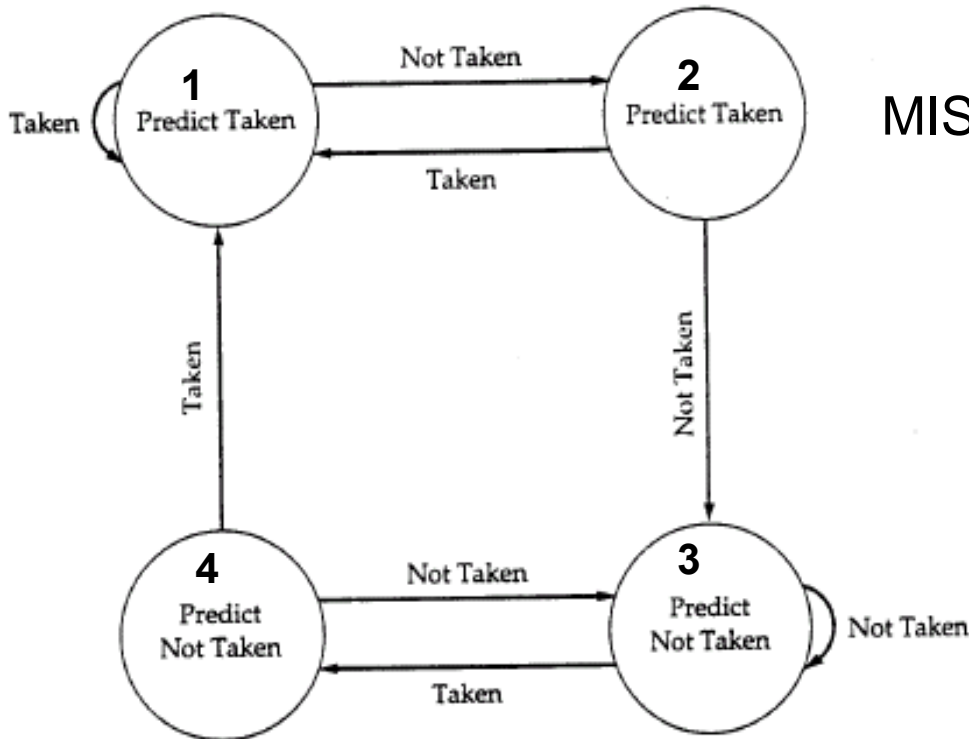
- בכניסה הבאה ללולאה MISS - דגל חיזוי Not Taken =

LOP	LDA	NUM
	ISZ	NUM
	BUN	LOP
	STA	STM
....		
NUM	-10	
STM	0	



Branch Prediction

- טכניקה 2 משתמשת בשתי סיביות – נועדה למנוע את ה-MISS בכניסה השנייה
- כיוון שנדרש MISS פעמיים
- בכניסה ללולאה פעם נוספת:
- התוכנית תהיה במצב 2 ולכן לא יהיה MISS



```
LOP LDA NUM
      ISZ NUM
      BUN LOP
      STA STM
```

....

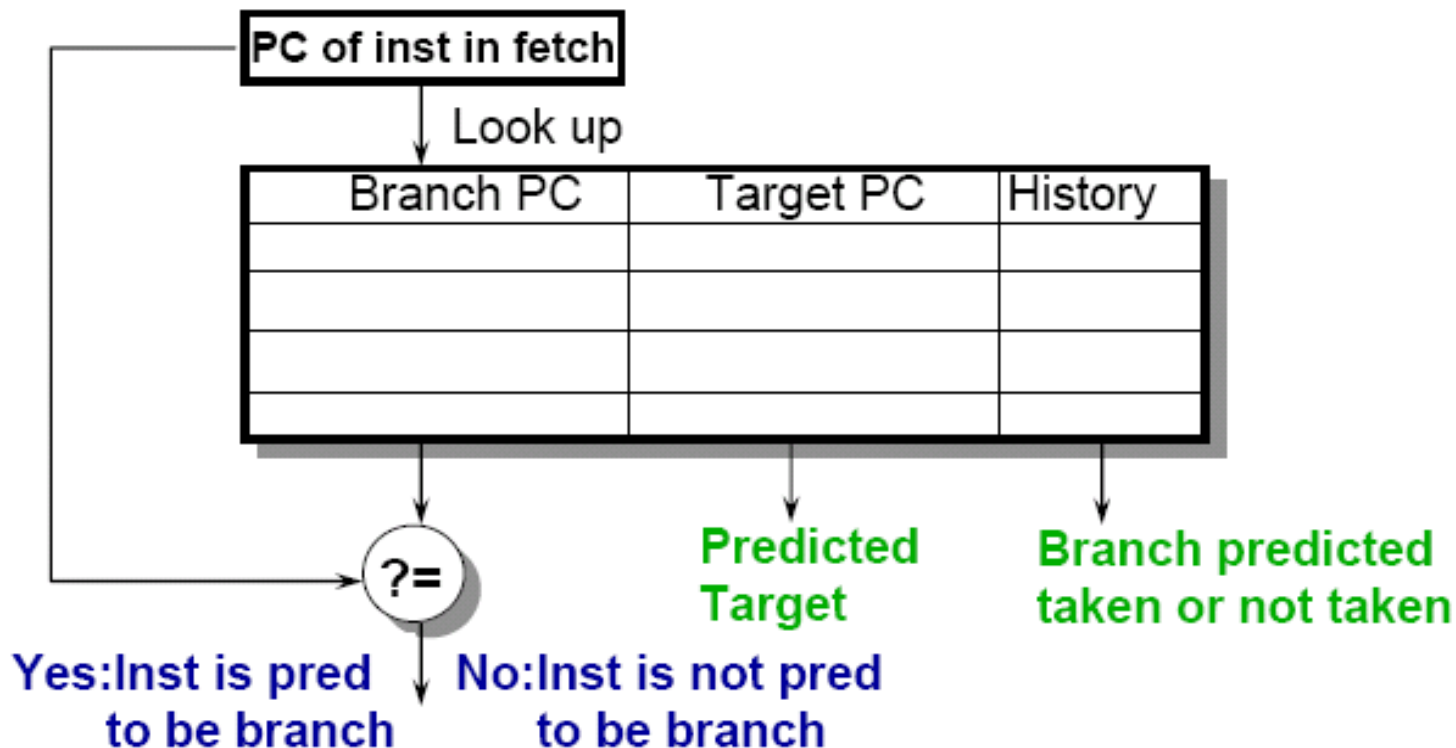
```
NUM -10
```

```
STM 0
```

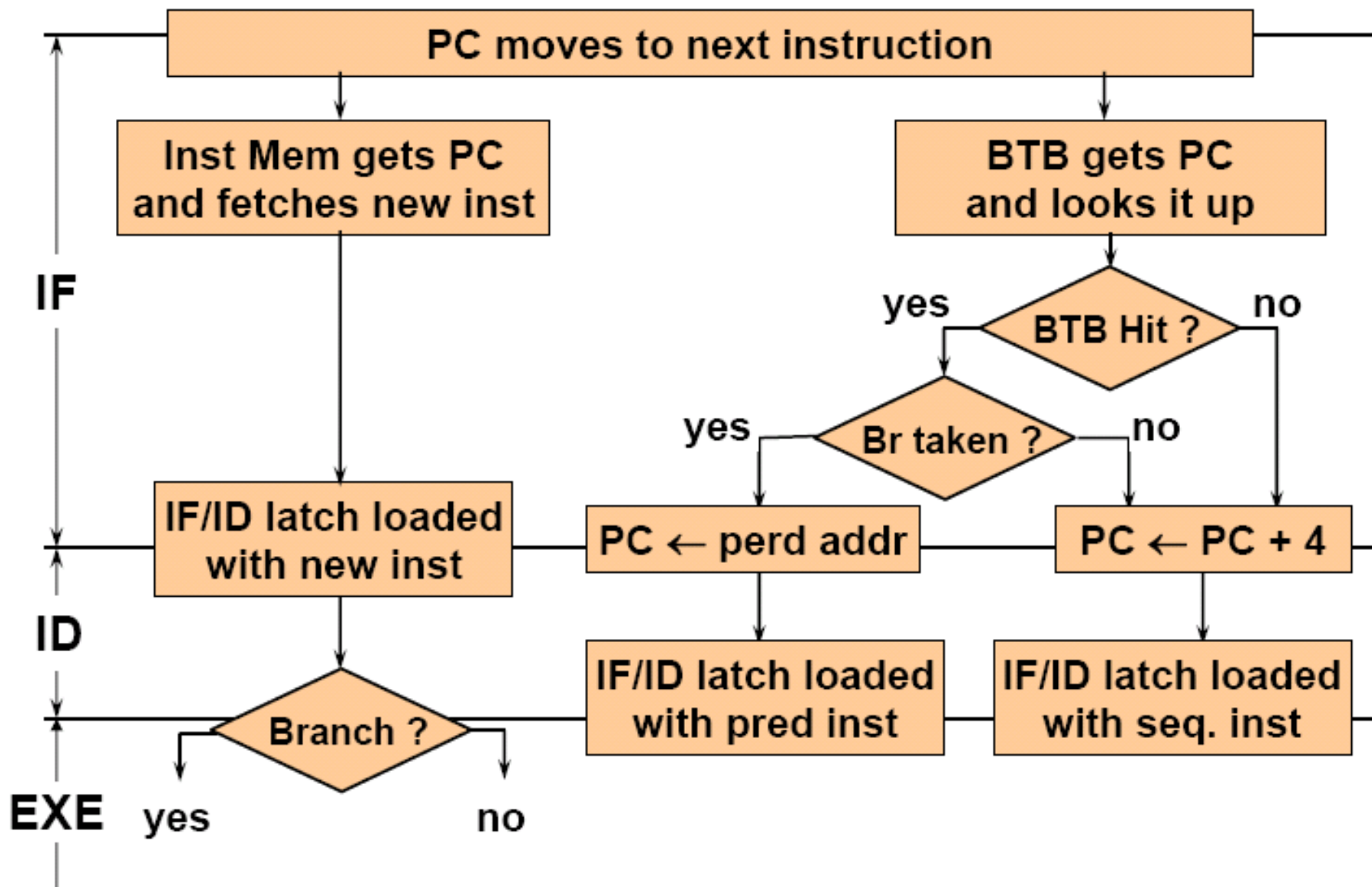
- טכניקה נוספת – תלויה גם בביצוע של ההסתעפות האחרונה שהתבצעה.

Dynamic Branch Prediction

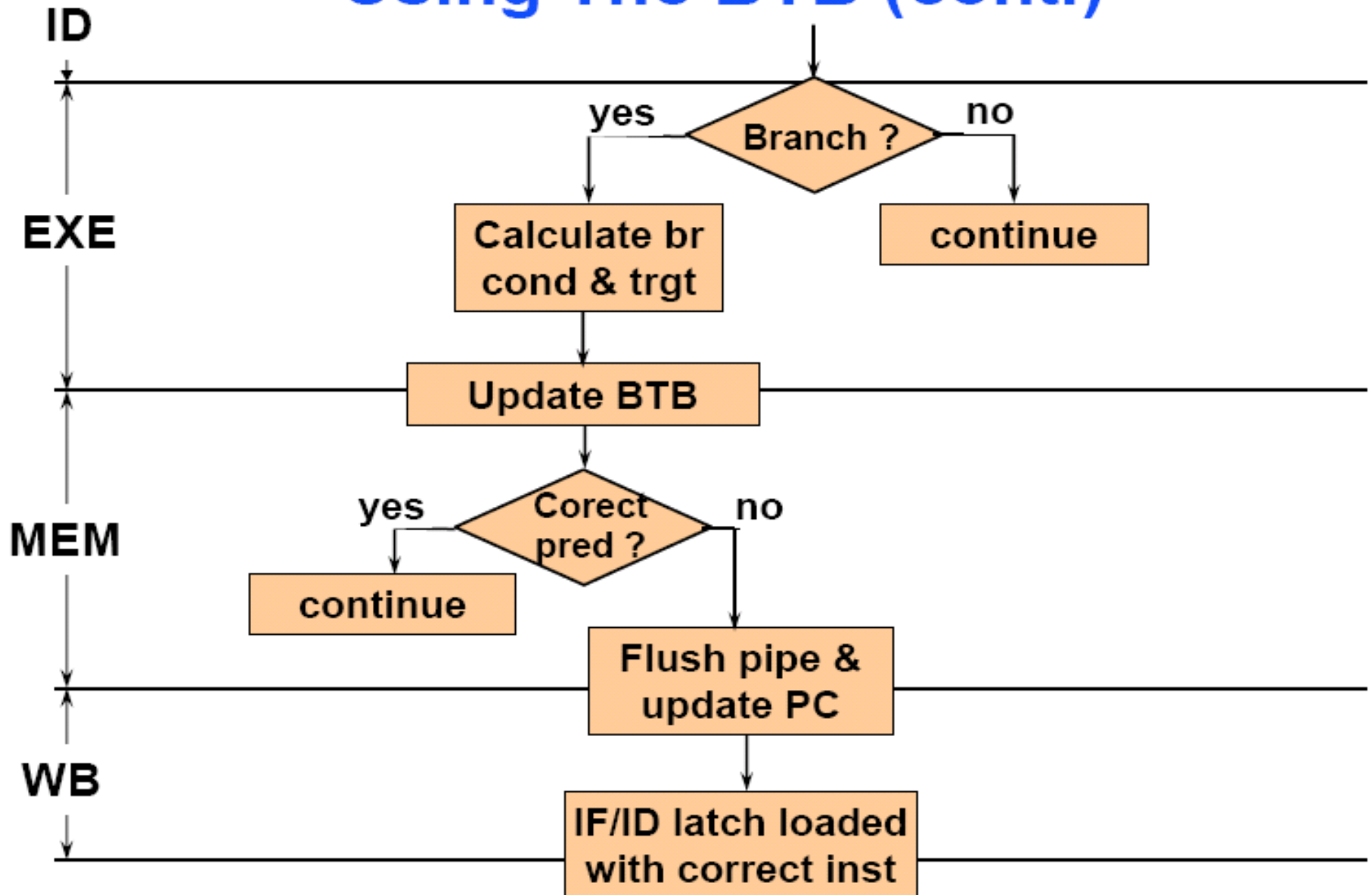
- ❑ Add a **Branch Target Buffer (BTB)** the predicts (at fetch)
 - Instruction is a branch
 - Branch taken / not-taken
 - Taken branch target



Using The BTB



Using The BTB (cont.)



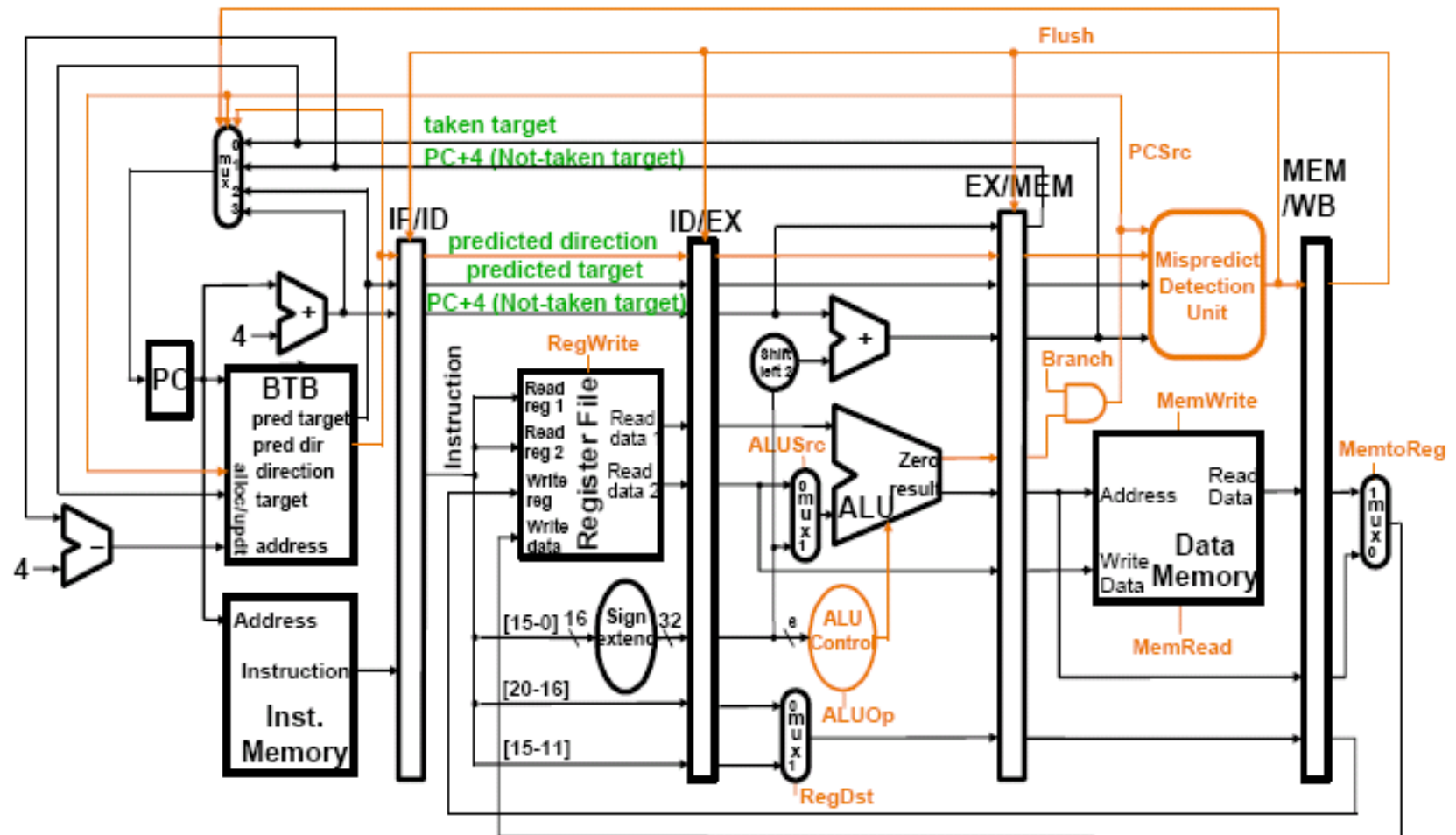
BTB

- **Allocation**
 - Allocate instructions identified as branches (after decode)
 - Both conditional and unconditional branches are allocated
 - Not taken branches need not be allocated
 - BTB miss implicitly predicts not-taken
- **Prediction**
 - BTB lookup is done parallel to IC lookup
 - BTB provides
 - Indication that the instruction is a branch (BTB hits)
 - Branch predicted target
 - Branch predicted direction
 - Branch predicted type (e.g., conditional, unconditional)
- **Update (when branch outcome is known)**
 - Branch target
 - Branch history (taken / not-taken)

BTB (cont)

- **Wrong prediction**
 - Predict not-taken, actual taken
 - Predict taken, actual not-taken.
- **In case of wrong prediction – flush the pipeline**
 - Reset latches (same as making all instructions to be NOPs)
 - Select the PC source to be from the correct path
 - Need get the fall-through with the branch
 - Start fetching instruction from correct path
- **Assuming P% correct prediction rate**
 - 20% of instructions are branches
- **$\text{CPI}_{\text{new}} = 1 + (0.2 \times (1-P)) \times 3$**
 - For example, if $P=0.7$
- **$\text{CPI}_{\text{new}} = 1 + (0.2 \times 0.3) \times 3 = 1.18$**

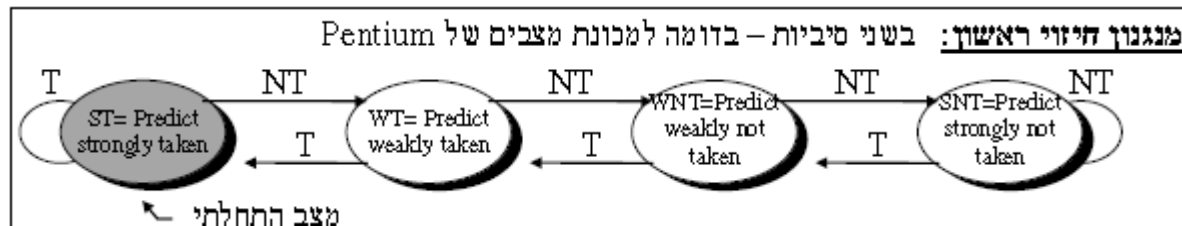
Adding a BTB to the Pipeline



נתון מעבד , אשר משתמש במנגנון BTB לחיזוי התנהגות של פקודות branch.

מריצים את קטע הקוד המייצג הבא על גבי המחשב:

hex	כתובת הפקודה	הפקודה	משמעות
1000		<u>Mov</u> R1 #0	$R1 \leftarrow 0$
1004		<u>Mov</u> R2 #2	$R2 \leftarrow 2$
1008	loop1:	<u>Inc</u> R1	$R1 \leftarrow R1+1$
100C	loop2:	<u>Blt</u> R1 R2 loop1	If ($R1 < R2$) goto loop1
1010		<u>Nop</u>	Do Nothing
1014		<u>Inc</u> R2	$R2 \leftarrow R2+1$
1018		<u>Nop</u>	Do Nothing
101C		<u>Blt</u> R2 #4 loop2	If ($R2 < 4$) goto loop2
1020		<u>Halt</u>	End of program



להניח מצב התחלתי NT
בדוגמה הנ"ל

