

VISA SYNTHESIS: VARIATION-AWARE INSTRUCTION SET ARCHITECTURE SYNTHESIS

Yuko Hara-Azumi^{*} Takuya Azumi[†] Nikil D. Dutt[‡]

^{*}Nara Institute of Science and Technology

[†]Ritsumeikan University [‡]University of California, Irvine

ASP-DAC : Jan. 23rd 2013

Outlines

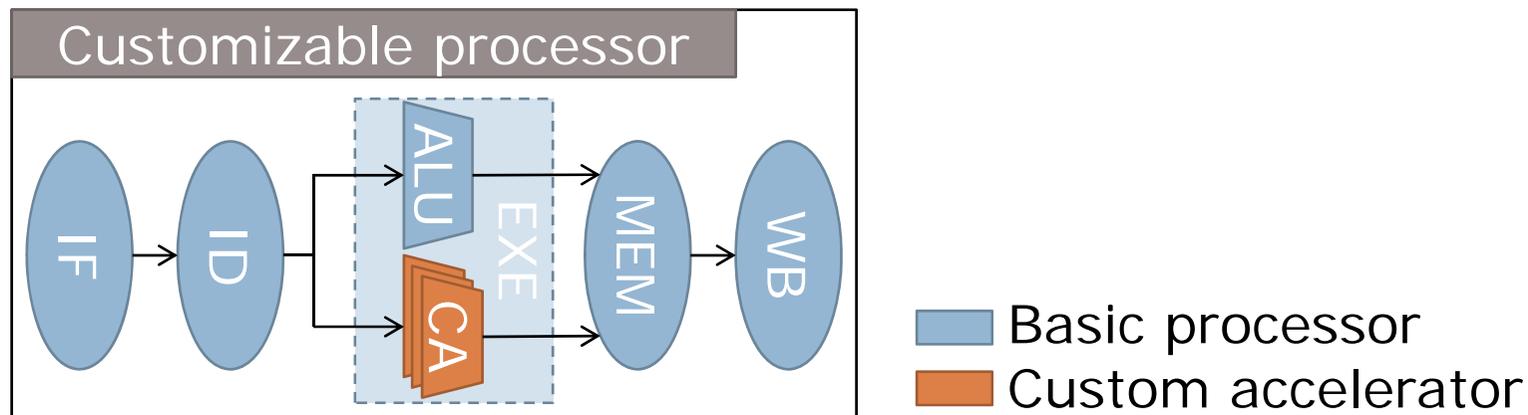
2

- Background
 - ▣ Previous works
- Variation-aware ISA (VISA) synthesis
 - ▣ Razor architecture
 - ▣ Our architecture (**HW**-side approach)
 - ▣ SSTA-based CI selection (**SW**-side approach)
- Experiments
- Conclusions

Instruction-set architecture (ISA) synthesis

3

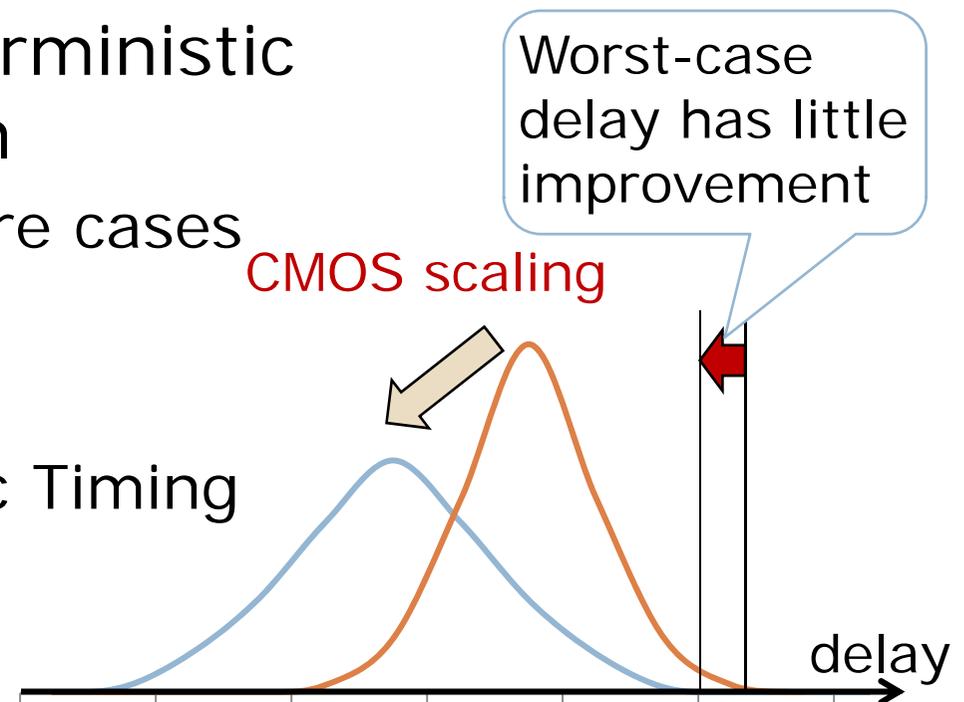
- Embedded processors are widely used in various applications
 - ▣ ISA synthesis: application-specific extension
 - Efficient **speedup** with less cost (**area, power**, etc.)
 - ▣ Custom instruction (CI) selection
 - Critical computation: **CIs** → Custom accelerator (CA)
 - The others: basic instructions (**BIs**) → ALU



Clock frequency

4

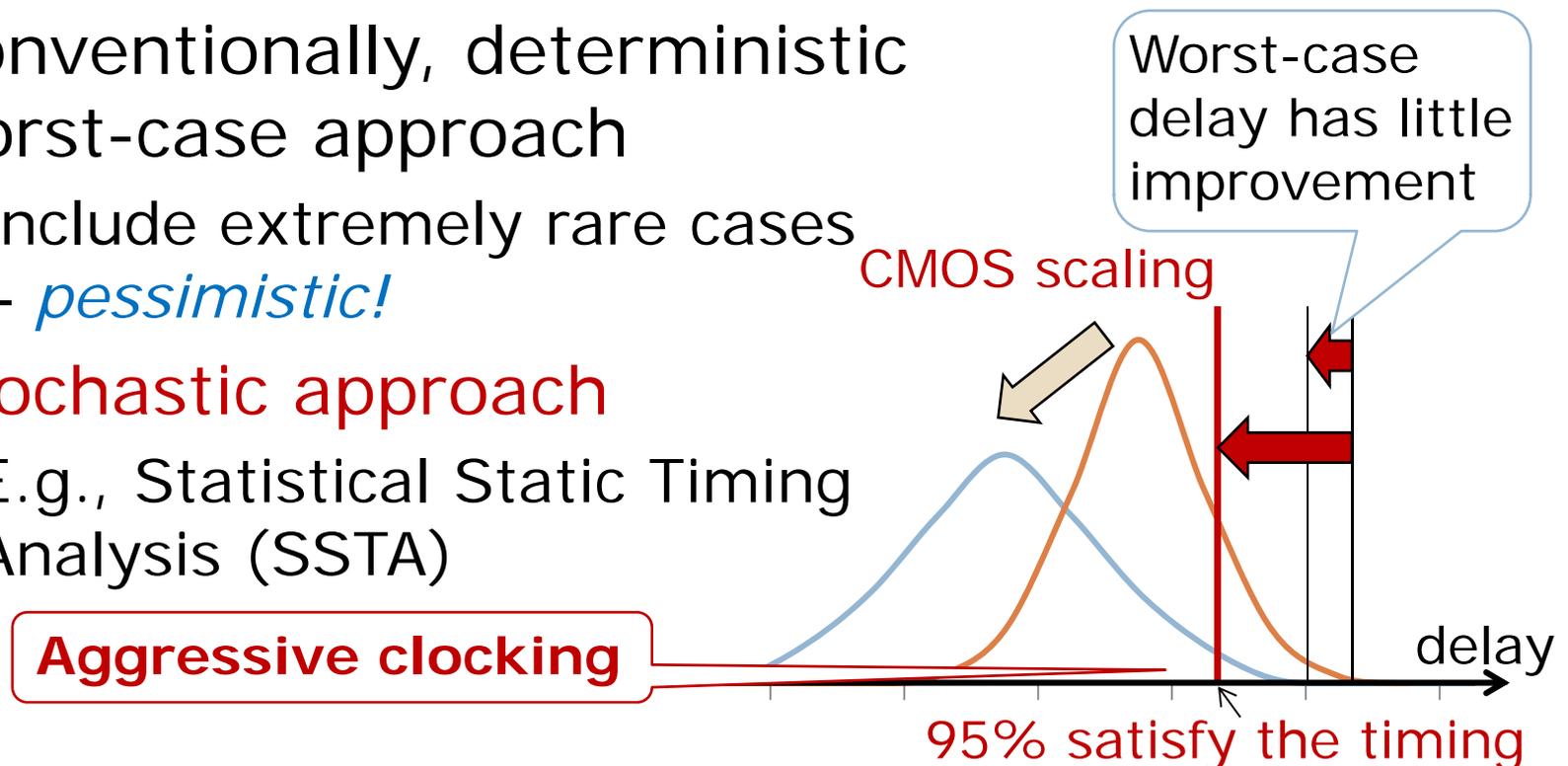
- A lot of challenging issues of CMOS scaling
 - ▣ Cannot expect the frequency scaling
- **Process variation**
 - ▣ Propagation delay varies by environments
 - ▣ Conventionally, deterministic worst-case approach
 - Include extremely rare cases
 - *pessimistic!*
 - ▣ **Stochastic approach**
 - E.g., Statistical Static Timing Analysis (SSTA)



Clock frequency

5

- A lot of challenging issues of CMOS scaling
 - ▣ Cannot expect the frequency scaling
- **Process variation**
 - ▣ Propagation delay varies by environments
 - ▣ Conventionally, deterministic worst-case approach
 - Include extremely rare cases
 - *pessimistic!*
 - ▣ **Stochastic approach**
 - E.g., Statistical Static Timing Analysis (SSTA)



Previous works (1)

6

- SSTA-based ISA synthesis works
 - ▣ [*Kamal'11*]: CI selection with minimum timing yield degradation
 - **Timing yield**: possibility to complete operations for a given target clock
 - **Timing yield degradation**: may be **intolerable** for some applications
 - ▣ [*Kamal'12*]: Maximum speedup with no timing yield degradation
 - An extra cycle to CIs with less-than-1.0 timing yield
 - Static approach only: **extra cycles** even when no timing faults occur – *pessimistic!*
 - ▣ **Timing violation by BIs** are not considered

Previous works (2)

7

- Variation-tolerable works in other fields
 - ▣ Architectures: Razor [*Ernst'03*], etc.
 - Detect and correct timing faults dynamically
 - ▣ High-level synthesis: SSTA-based works ([*Cong'09*], etc.)
 - Use Razor-flipflops for aggressive clocking without timing yield degradation
- HW approach only - *very costly!*
 - ▣ Comprehensive approach from both HW/SW viewpoints is necessary

Outlines

8

- Background
 - Previous works
- Variation-aware ISA (VISA) synthesis
 - Razor architecture
 - Our architecture (**HW**-side approach)
 - SSTA-based CI selection (**SW**-side approach)
- Experiments
- Conclusions

VISA synthesis

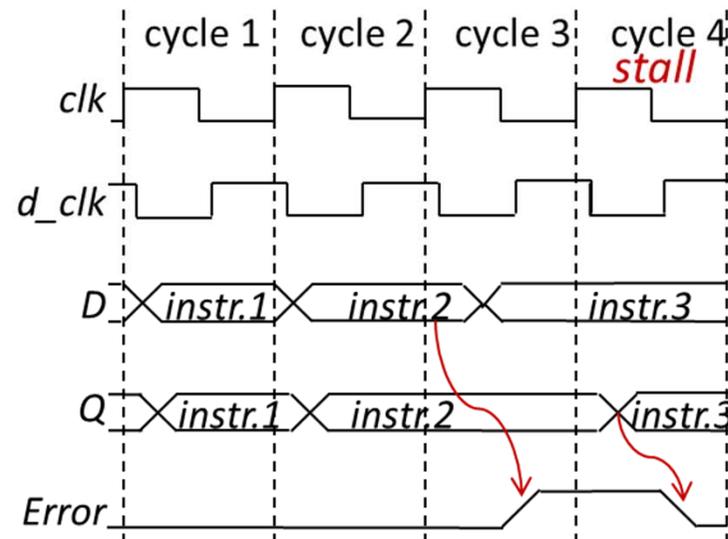
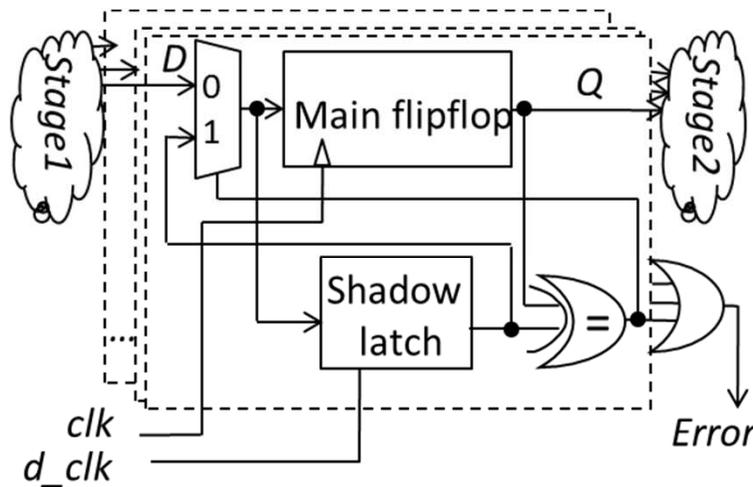
9

- **VISA: Variation-aware ISA synthesis**
 - ▣ Performance improvement by **making effective use of process variation** on both HW and SW
 - **HW**: Dynamic fault detection & correction with minimum performance degradation for aggressive clocking (based on Razor)
 - **SW**: SSTA-based CI selection effectively by exploiting application features
 - ▣ Handle timing violation of **both BIs & CIs**
 - ▣ Applicable to any processors which have at least a mechanism of dynamic timing fault detection

Razor processor architecture (1)

10

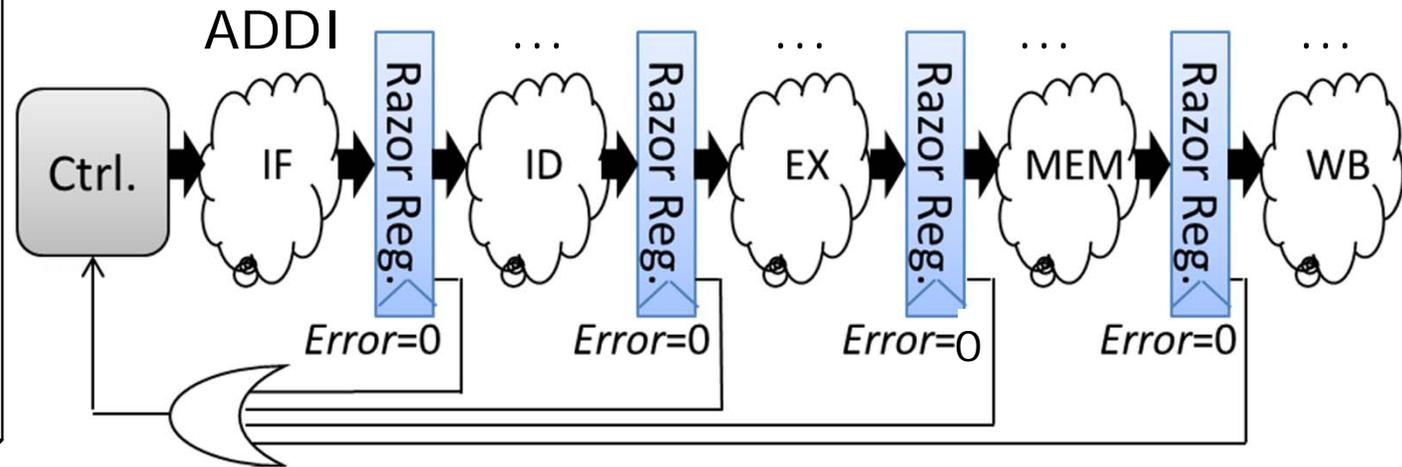
- Razor flipflop: Main flipflop + shadow latch (delayed clock)
 - ▣ **Fault detection**: compares the results
 - ▣ **Fault correction**: simply copies the correct result to the main flipflop during 1-cycle stall



Razor processor architecture (2)

11

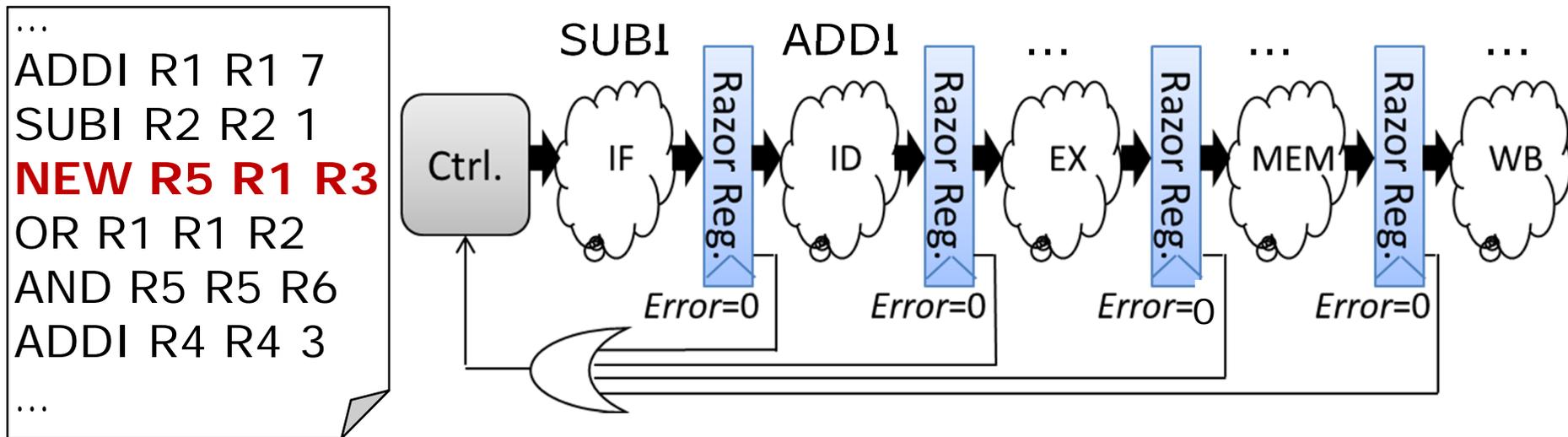
...
 ADDI R1 R1 7
 SUBI R2 R2 1
NEW R5 R1 R3
 OR R1 R1 R2
 AND R5 R5 R6
 ADDI R4 R4 3
 ...



cycle stage	t
IF	ADDI
ID	...
EX	...
MEM	...
WB	...

Razor processor architecture (2)

12



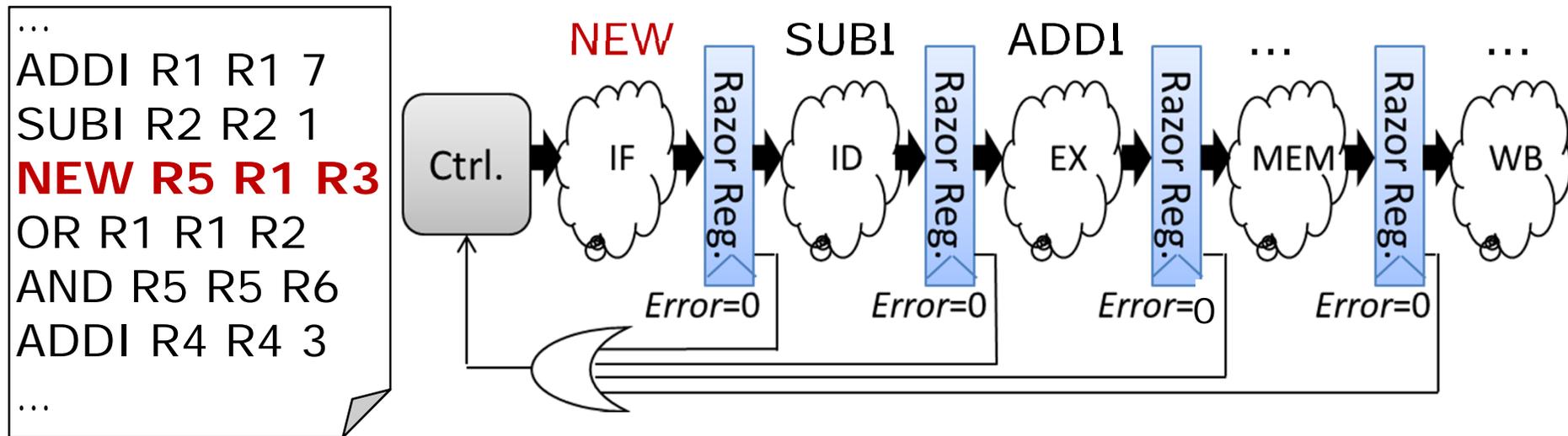
```

...
ADDI R1 R1 7
SUBI R2 R2 1
NEW R5 R1 R3
OR R1 R1 R2
AND R5 R5 R6
ADDI R4 R4 3
...
    
```

cycle stage \	t	t+1
IF	ADDI	SUBI
ID	...	ADDI
EX
MEM
WB

Razor processor architecture (2)

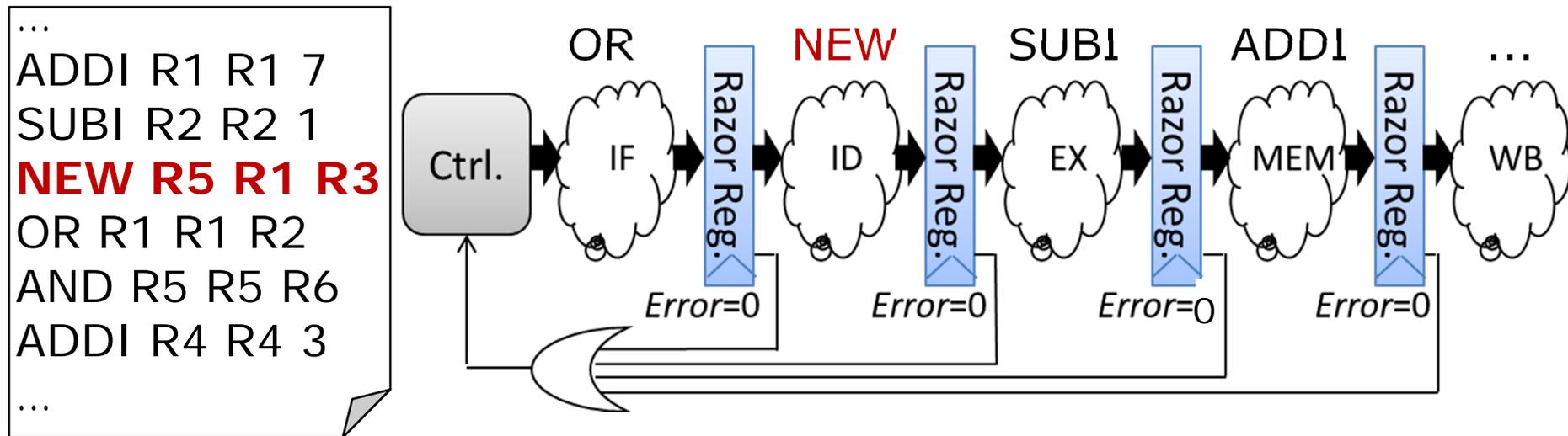
13



cycle stage \	t	t+1	t+2
IF	ADDI	SUBI	NEW
ID	...	ADDI	SUBI
EX	ADDI
MEM
WB

Razor processor architecture (2)

14



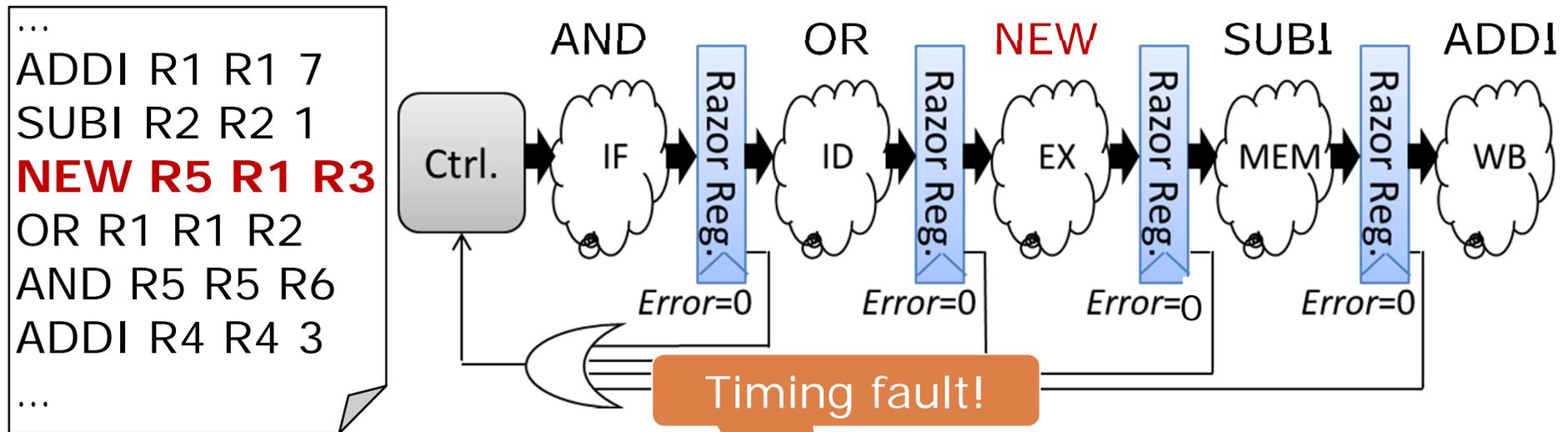
```

...
ADDI R1 R1 7
SUBI R2 R2 1
NEW R5 R1 R3
OR R1 R1 R2
AND R5 R5 R6
ADDI R4 R4 3
...
    
```

cycle stage \	t	t+1	t+2	t+3
IF	ADDI	SUBI	NEW	OR
ID	...	ADDI	SUBI	NEW
EX	ADDI	SUBI
MEM	ADDI
WB

Razor processor architecture (2)

15



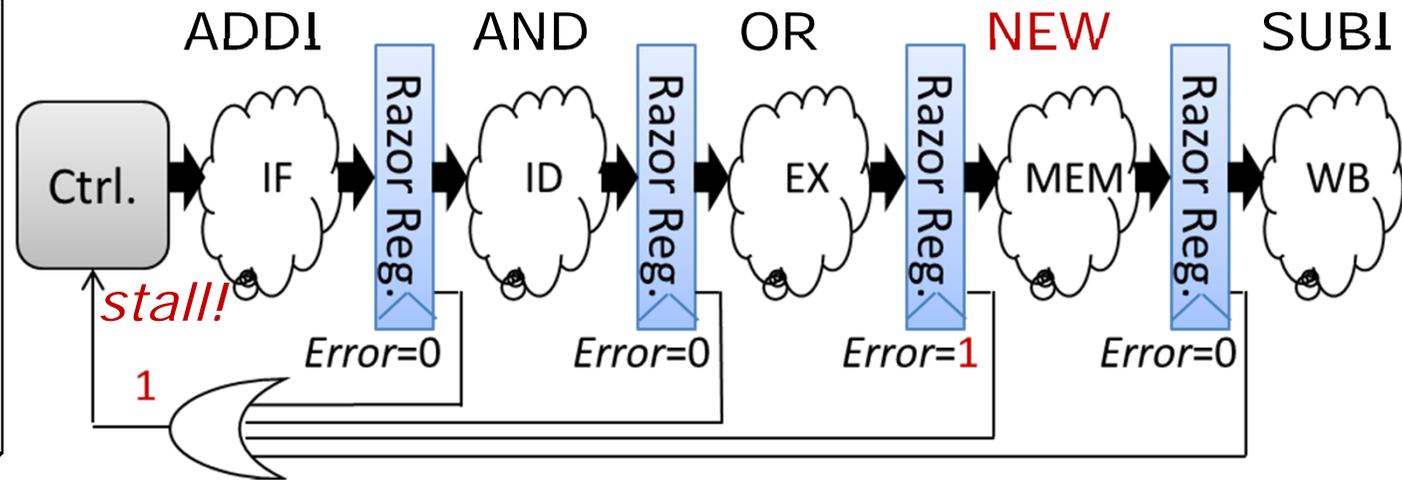
...
 ADDI R1 R1 7
 SUBI R2 R2 1
NEW R5 R1 R3
 OR R1 R1 R2
 AND R5 R5 R6
 ADDI R4 R4 3
 ...

cycle stage \	t	t+1	t+2	t+3	t+4
IF	ADDI	SUBI	NEW	OR	AND
ID	...	ADDI	SUBI	NEW	OR
EX	ADDI	SUBI	NEW 
MEM	ADDI	SUBI
WB	ADDI

Razor processor architecture (2)

```

...
ADDI R1 R1 7
SUBI R2 R2 1
NEW R5 R1 R3
OR R1 R1 R2
AND R5 R5 R6
ADDI R4 R4 3
...
    
```



cycle stage \	t	t+1	t+2	t+3	t+4	t+5
IF	ADDI	SUBI	NEW	OR	AND	ADDI
ID	...	ADDI	SUBI	NEW	OR	AND
EX	ADDI	SUBI	NEW*	OR
MEM	ADDI	SUBI	NEW*
WB	ADDI	SUBI

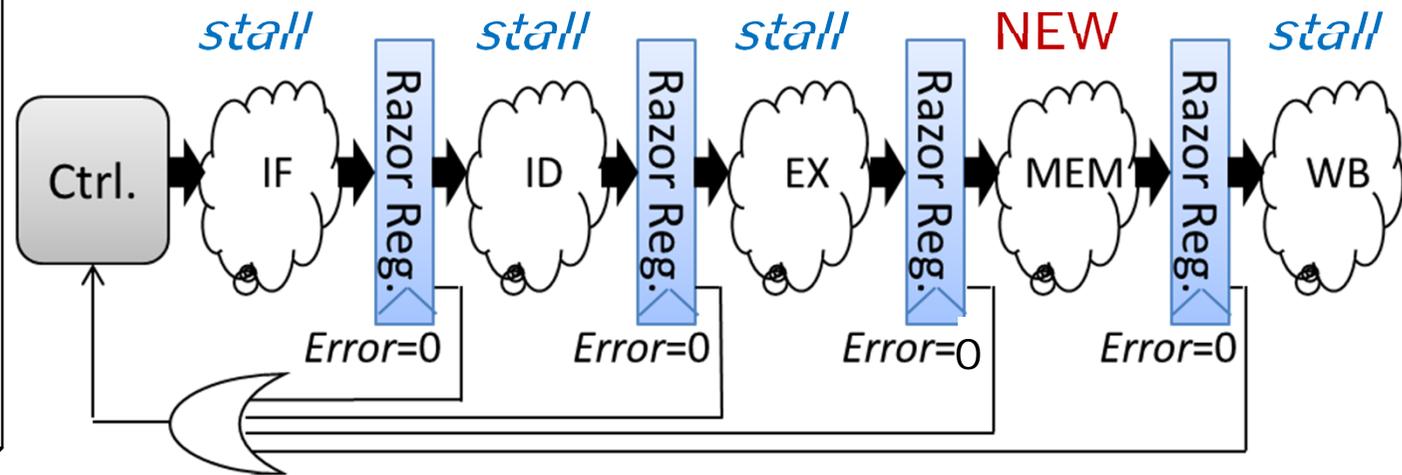
The shadow latch detects the timing fault and issues the Error signal

Razor processor architecture (2)

17

```

...
ADDI R1 R1 7
SUBI R2 R2 1
NEW R5 R1 R3
OR R1 R1 R2
AND R5 R5 R6
ADDI R4 R4 3
...
    
```

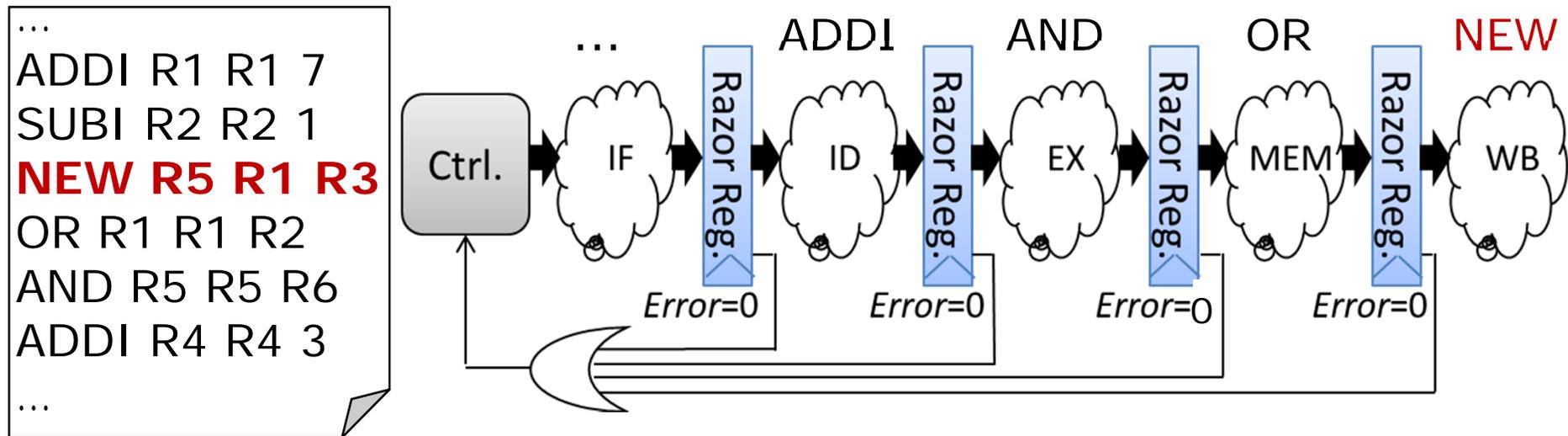


cycle stage \	t	t+1	t+2	t+3	t+4	t+5	t+6
IF	ADDI	SUBI	NEW	OR	AND	ADDI	<i>stall</i>
ID	...	ADDI	SUBI	NEW	OR	AND	<i>stall</i>
EX	ADDI	SUBI	NEW*	OR	<i>stall</i>
MEM	ADDI	SUBI	NEW*	NEW
WB	ADDI	SUBI	<i>stall</i>

The correct result of MULT is provided to the MEM stage

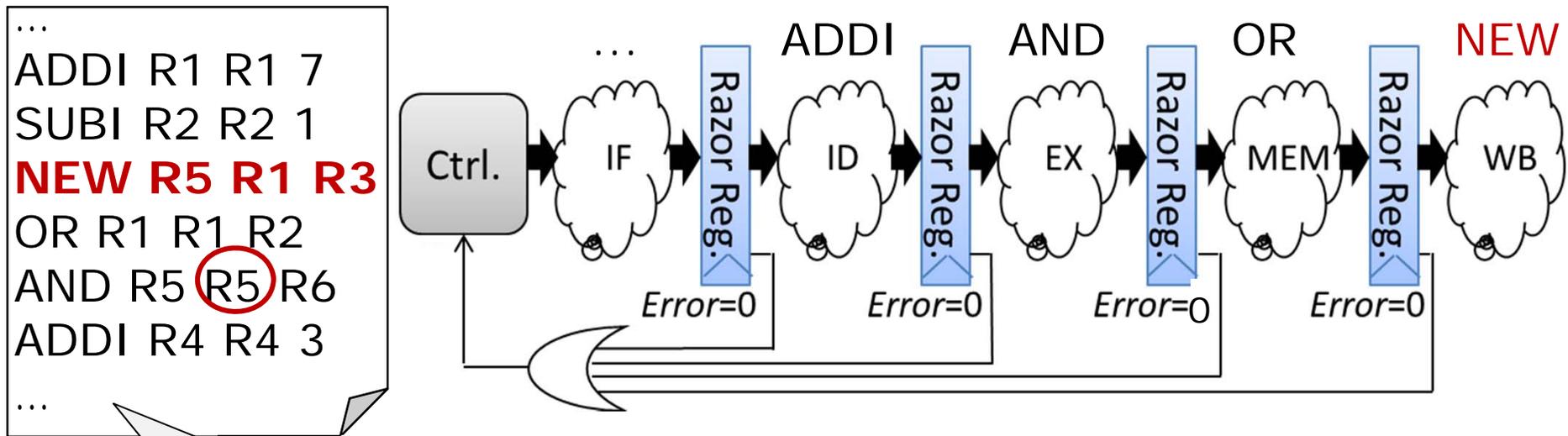
Razor processor architecture (2)

18



cycle stage \	t	t+1	t+2	t+3	t+4	t+5	t+6	t+7
IF	ADDI	SUBI	NEW	OR	AND	ADDI	<i>stall</i>	...
ID	...	ADDI	SUBI	NEW	OR	AND	<i>stall</i>	ADDI
EX	ADDI	SUBI	NEW*	OR	<i>stall</i>	AND
MEM	ADDI	SUBI	NEW*	NEW	OR
WB	ADDI	SUBI	<i>stall</i>	NEW

Razor processor architecture (2)



...
 ADDI R1 R1 7
 SUBI R2 R2 1
NEW R5 R1 R3
 OR R1 R1 R2
 AND R5 **R5** R6
 ADDI R4 R4 3
 ...

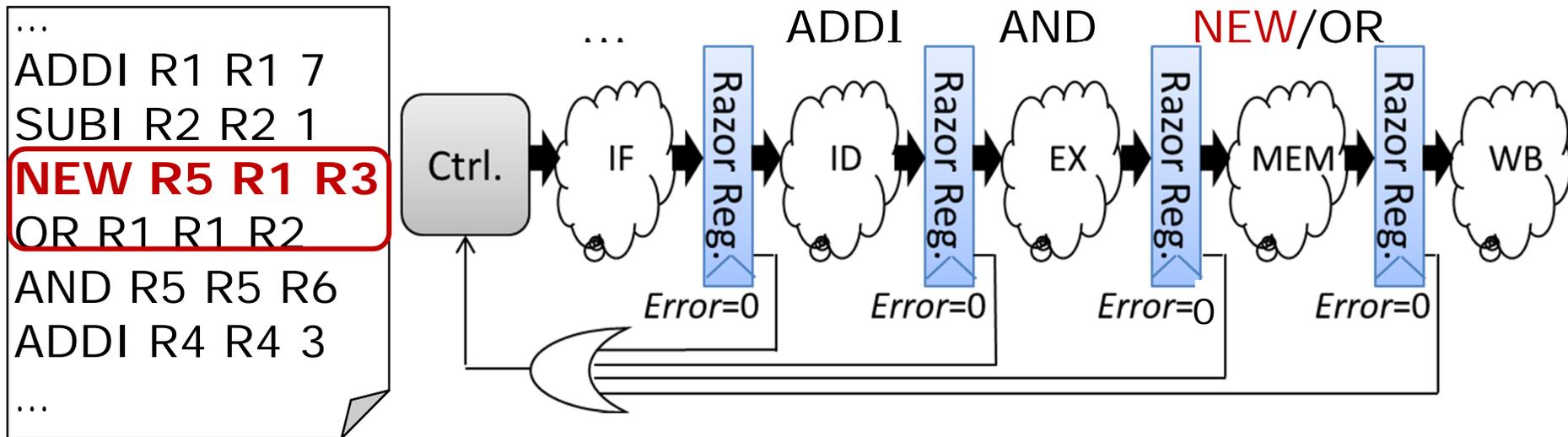
Instruction	t	t+1	t+2	t+3	t+4	t+5	t+6	t+7
...
OR	MEM							
AND	MEM							
ADDI	MEM							
SUBI	MEM							
NEW*	MEM							
OR	MEM							
AND	MEM							
NEW	MEM							
OR	MEM							
ADDI	MEM							
SUBI	MEM							
NEW	MEM							

OR has no dependency from & share no arithmetic unit with **NEW** → **NEW** can be executed *with no stall*

Actually this stall is **not** necessary!

Razor processor architecture (2)

20



cycle stage \	t	t+1	t+2	t+3	t+4	t+5	t+6	t+7
IF	ADDI	SUBI	NEW	OR	AND	ADDI
ID	...	ADDI	SUBI	NEW	OR	AND	ADDI	...
EX	ADDI	SUBI	NEW*	OR	AND	ADDI
MEM	ADDI	SUBI	NEW*	NEW/OR	AND
WB	ADDI	SUBI	-	NEW/OR

CI selection

22

- SSTA + application features
 - ▣ Speedup effects of individual instances of CIs

CI (NEW)

2 instances

Speedup = 10 cycles/exec.

Timing yield = 0.9

- ▣ Obtain the speedup of CIs considering
 - Penalty for register file accesses
 - **Neighboring instructions of its instances (=application features)**

NEW₁:
8 iterations

```
...  
SUBI R2 R2 1  
NEW R5 R1 R3  
OR R1 R1 R2  
...
```

NEW₂:
20 iterations

```
...  
ADDI R4 R4 8  
NEW R6 R4 R2  
AND R4 R4 R6  
...
```

Need a stall

CI selection

- SSTA + application features
 - ▣ Speedup effects of **individual instances** of CIs

CI (NEW)

2 instances

Speedup = 10 cycles/exec.

Timing yield = 0.9

NEW₁:
8 iterations

NEW₂:
20 iterations

```

...
SUBI R2 R2 1
NEW R5 R1 R3
OR R1 R1 R2
...
ADDI R4 R4 8
NEW R6 R4 R2
AND R4 R4 R6
...
    
```

Need a stall

Benefit by NEW

$$\begin{aligned}
 &= \text{average speedup by NEW}_1 \\
 &+ \text{average speedup by NEW}_2 \\
 &= \{0.9 \times 10 + 0.1 \times (10 - 0)\} \times 8 \\
 &+ \{0.9 \times 10 + 0.1 \times (10 - \mathbf{1})\} \times 20 \\
 &= 80 + 198 \\
 &= 278
 \end{aligned}$$

↑
Penalty (stall)

CI selection: constraint

24

- Only CIs which **always finish by the setup time of the shadow latch** are selectable
 - ▣ Const.: $y_i(c_i T + d) = 1.0$
 - $y_i(t)$: timing yield of instruction i at time t
 - c_i : minimum latency of instruction i in the EX stage
 - T : target clock period
 - d : delayed setup time of the shadow latch
 - ▣ All BIs must hold the Const.
 - ▣ CIs are pruned by the Const.
 - CIs are also pruned by an *area constraint*

Outlines

25

- Background
 - ▣ Previous works
- Variation-aware ISA (VISA) synthesis
 - ▣ Razor architecture
 - ▣ Our architecture (**HW**-side approach)
 - ▣ SSTA-based CI selection (**SW**-side approach)
- Experiments
- Conclusions

Experimental setup

26

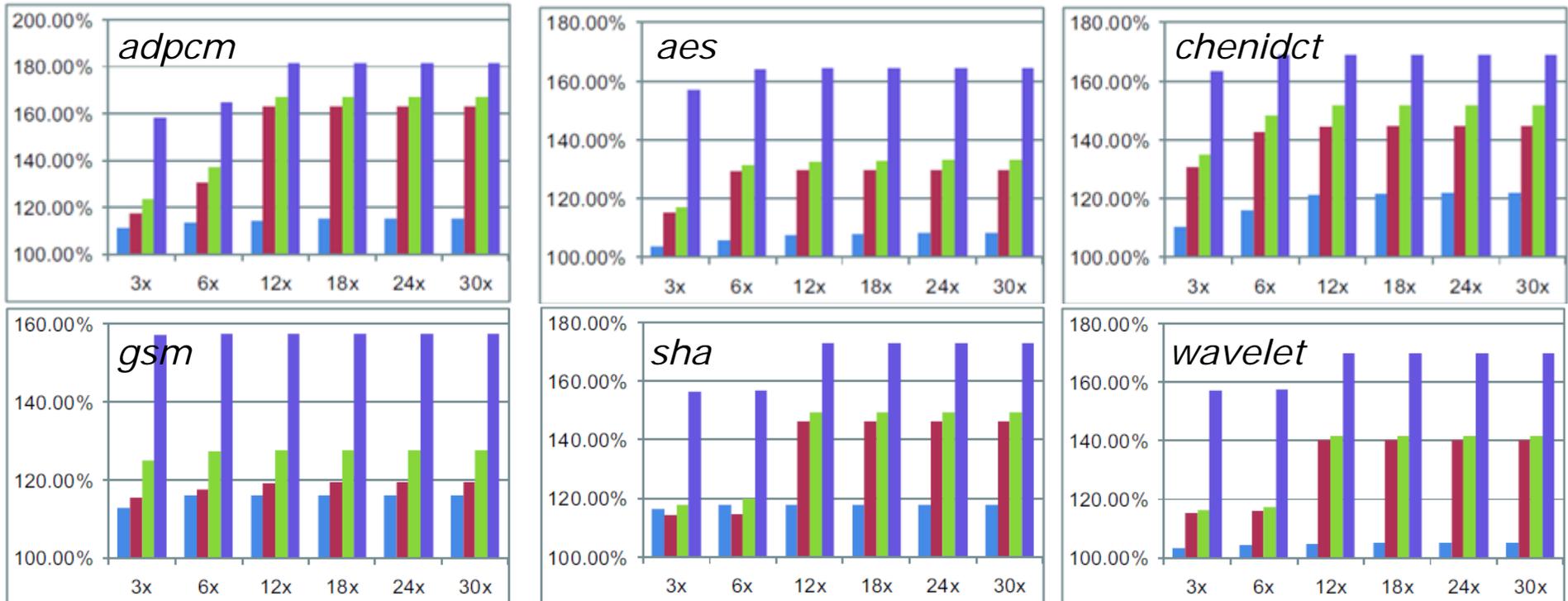
- Benchmarks: *adpcm*, *aes*, *chenidct*, *gsm*, and *sha*, and *wavelet*
- Target device: 90nm technology
 1. T_0 : $y_i(c_i T_0) = 1.0$ for all BIs and CIs
 2. T_1 : $y_i(c_i T_1) = 1.0$ for all BIs
 3. T_2 : $y_i(c_i T_2) < 1.0$ & $y_i(c_i T_2 + d) = 1.0$ for some BIs
- Simulator: SimpleScalar
 - ▣ MIPS (PISA)
- Comparative methods
 1. A deterministic worst-case method (**DW**): only for T_0
 2. An existing SSTA-based method (**ES**) [Kamal'12]: an extra cycle (stall) is always given to CIs with less-than-1.0 timing yield – only for T_1
 3. Our proposed method (**VISA**): compensation by both HW and SW – for T_1 and T_2

All performed greedily

Experimental results: speedup

27

DW ES VISA(T₁) VISA(T₂)



x-axis: Area introduced for CIs (#x ALU's area)

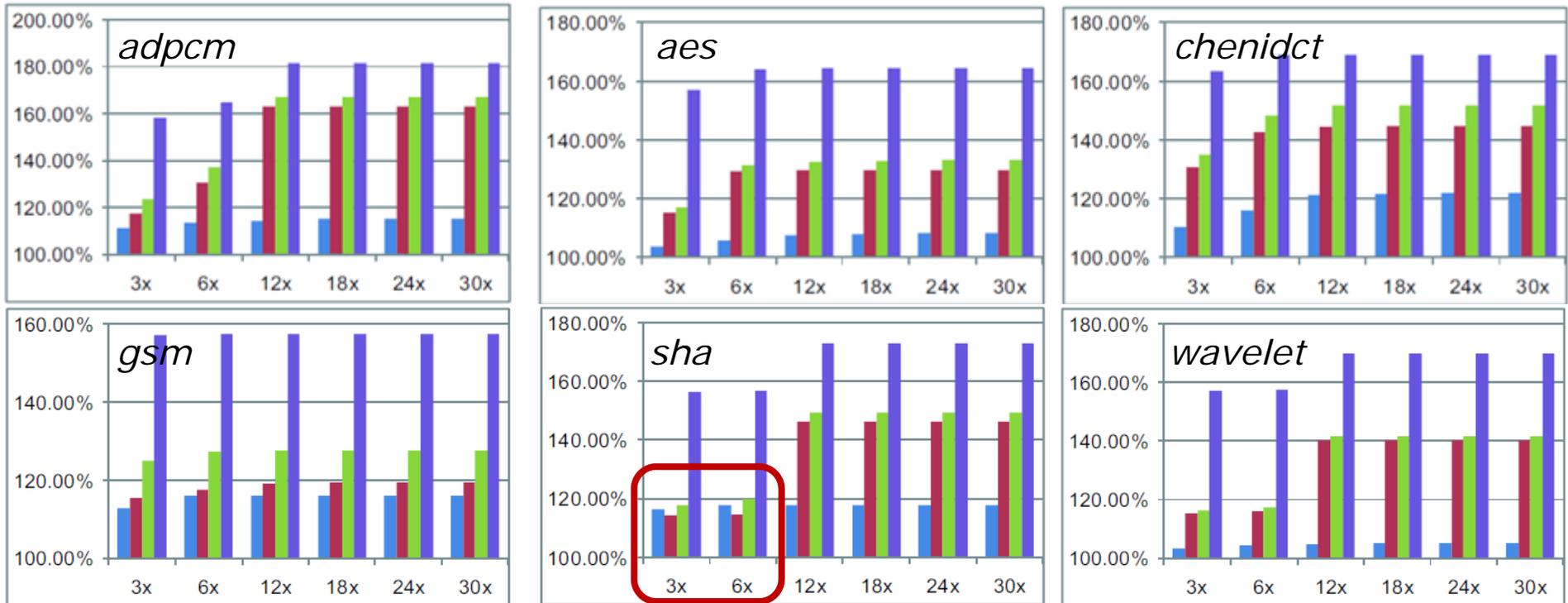
y-axis: Execution time improvement

- DM
 - ▣ Not very large speedup
 - ▣ Similarly with DM, deterministic approaches quickly face the clock wall

Experimental results: speedup

28

DW ES VISA(T₁) VISA(T₂)



x-axis: Area introduced for CIs (#x ALU's area)

y-axis: Execution time improvement

- ES
 - Larger speedup than DM, but still pessimistic in that ES *always* gives a stall for CIs with less-than-1.0 timing yield
 - For *sha* with 3x and 6x, even less speedup than DM

Experimental results: stalls

29

No. of instr. with a stall / No. of instr. with less-than-1.0 timing yield

Method	<i>adpcm</i>	<i>aes</i>	<i>chenidct</i>	<i>gsm</i>	<i>sha</i>	<i>wavelet</i>
ES (T_1)	2/2	2/2	18/18	7/7	1/1	3/3
VISA(T_1)	0/70	0/112	6/78	4/95	0/17	1/22
VISA(T_2)	497/875	618/1155	417/633	547/907	390/689	366/655

□ ES

- 100% of CIs with less-than-1.0 timing yield *always* take a stall for T_1

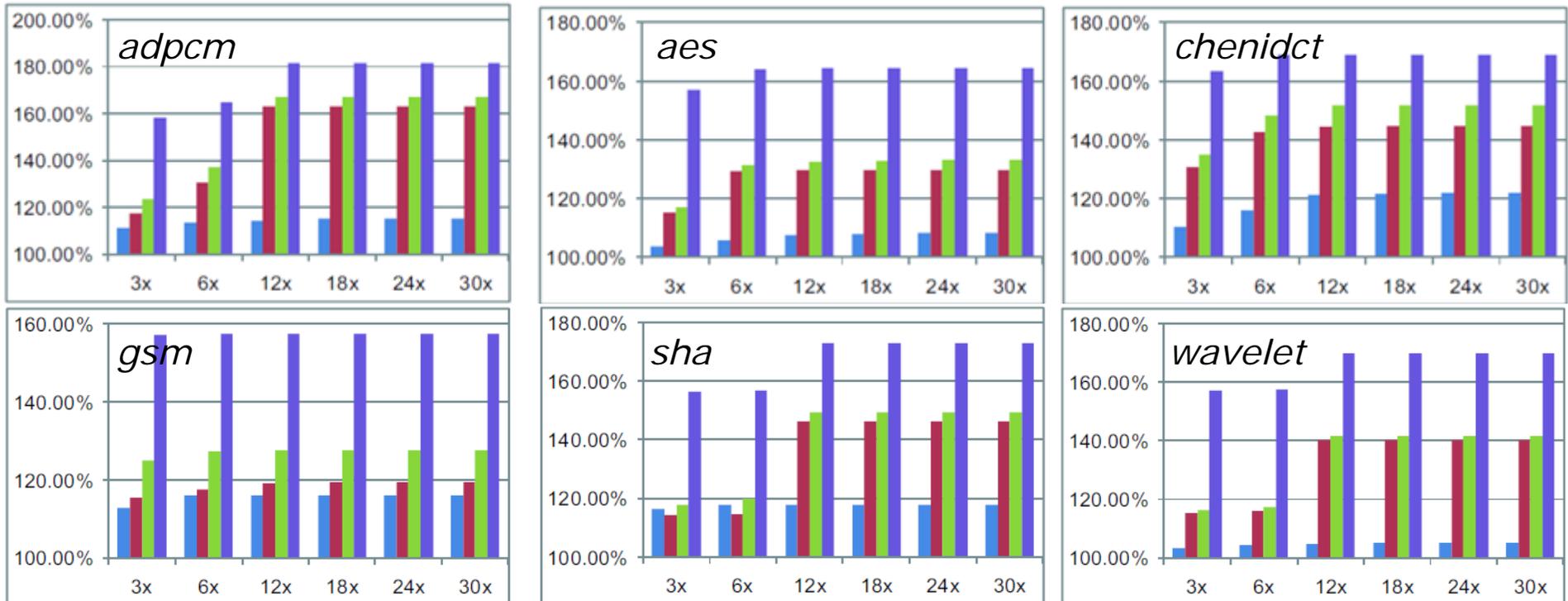
□ VISA

- Up to 8% of such CIs *may* take a stall for T_1
→ Effectively remove stalls
- Can select *more effective CIs* by considering application features

Experimental results: speedup

30

DW ES VISA(T₁) VISA(T₂)



x-axis: Area introduced for CIs (#x ALU's area)

y-axis: Execution time improvement

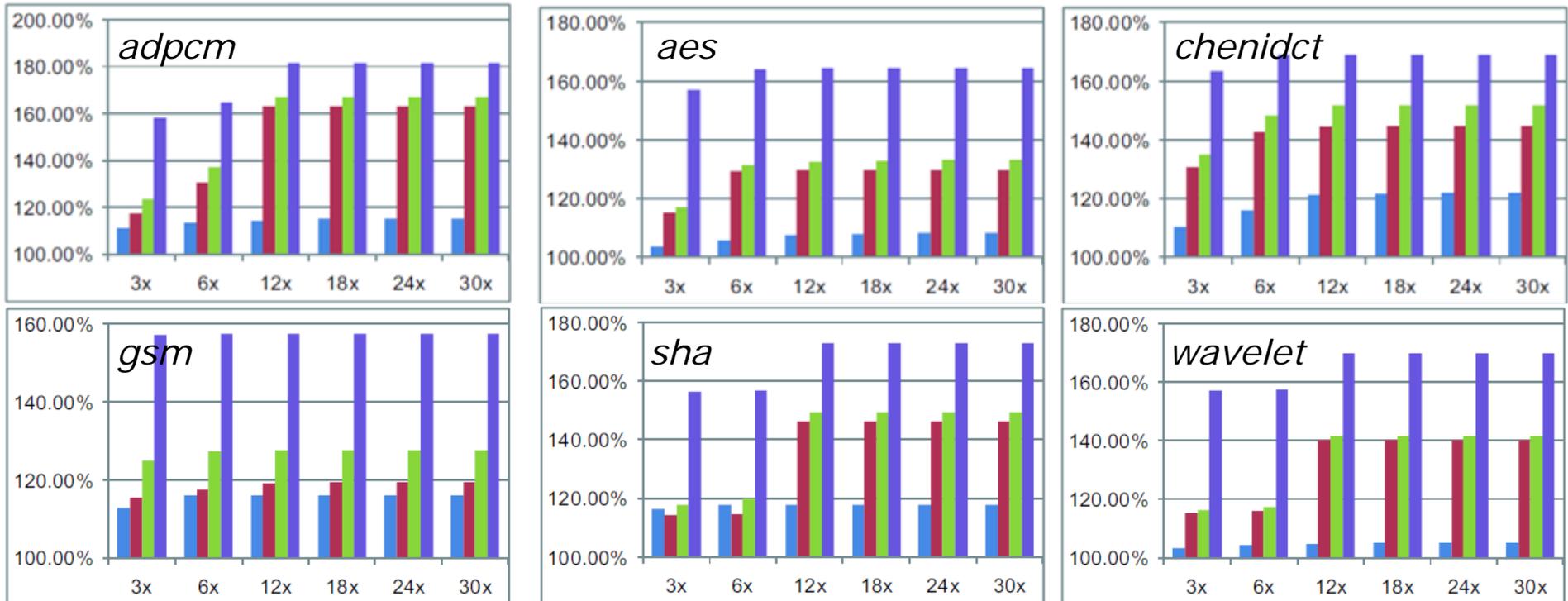
□ VISA

- More speedup than ES for T₁
 - Stall removal & effective CI selection
- Outperform DM and ES by up to **61.3%** and **13.0%**, respectively

Experimental results: speedup

31

DW ES VISA(T₁) VISA(T₂)



x-axis: Area introduced for CIs (#x ALU's area)

y-axis: Execution time improvement

□ VISA

- More effective for more aggressive clocking
- Outperform DM and ES by up to **78.0%** and **49.4%**, respectively, for T₂

Conclusion

32

- **VISA**: novel **V**ariation-aware **I**SA synthesis
 - ▣ Make effective use of process variation comprehensively from both HW and SW
 - HW: **Dynamic fault detection & correction** with minimum performance degradation
 - SW: **SSTA-based CI selection** considering application features
 - ▣ Substantially improves performance compared with existing methods
 - **More effective for more aggressive clocking**
 - Up to **78.0%** and **49.4%** performance improvement over two existing approaches