



7th Conference on Real Numbers and Computers (RNC 7)

LORIA, Nancy, France, July 10–12, 2006

Pipelined Architecture for Accurate Floating Point Range Reduction

**Francisco J. Jaime, Julio Villalba, Javier Hormigo, Emilio L.
Zapata**

**Dept. Computer Architecture
University of Malga
SPAIN**





Talk Outline

- Introduction
 - Posing the problem
- Double residue modular range reduction
- Pipelined Architecture
- Radix-4 coding
- Summary and conclusions



Introduction

- Range reduction is required for elementary functions evaluation
- Reductions:
 - Additive
 - Subtracting a multiple of a constant
 - Multiplicative
 - Multiplying by the inverse of a power of a constant
- Our algorithm
 - Additive reduction
 - Trigonometric functions



Introduction

- The reduction can be done in software or in hardware
- Our proposal is oriented toward a hardware implementation for application-specific systems
 - Consider both fixed-point & floating point representation
 - Maximum accuracy for floating point for the specific representation



Introduction

- A poor range reduction may lead to catastrophic accuracy problems in the evaluation of trigonometric functions when
 1. the input argument is close to a multiple of the constant (i.e. π)
 2. the input argument is large [1]



Introduction

– Steps in the evaluation of trigonometric functions:

- 1) Reduction of the input argument to the interval $[0, 2\pi)$
- 2) Compute the function inside the interval

Ej.: $\sin(114 \text{ rad})$

1) Reduction $114 \text{ rad} \rightarrow 0.9 \text{ rad}$

2) Computation $\sin(0.9 \text{ rad}) = 0,7850$

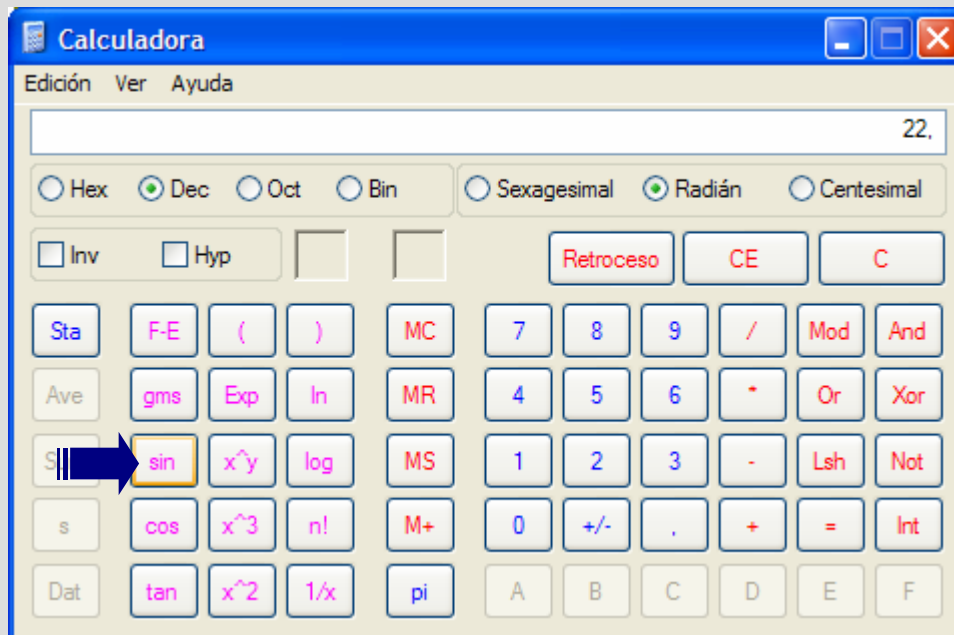


A not accurate reduction (step 1) can lead to inaccurate results



Introduction

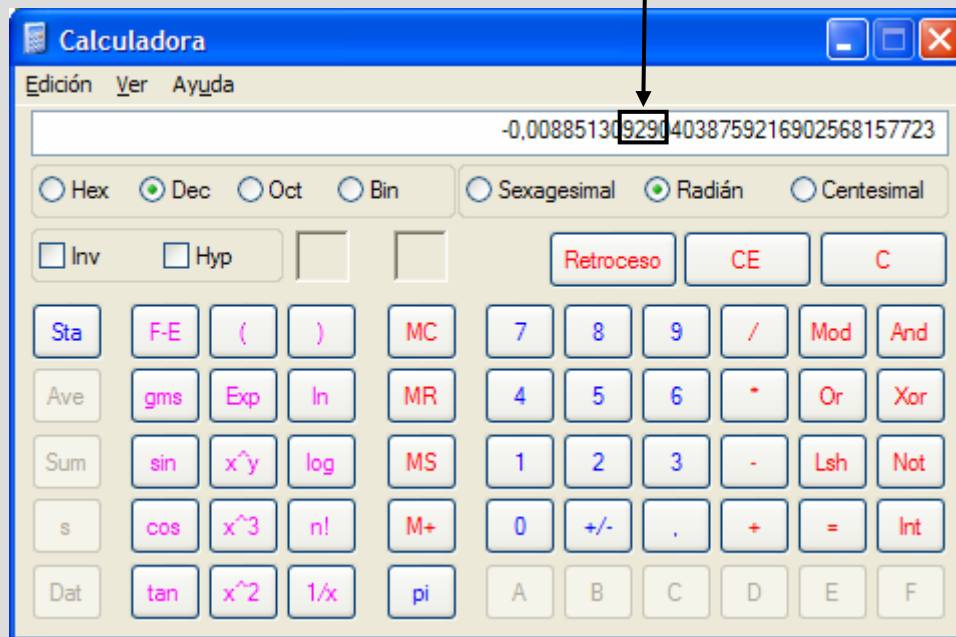
- 1) Example of input argument close a multiple of π
 - Computation of $\sin(22 \text{ rad})$





Introduction

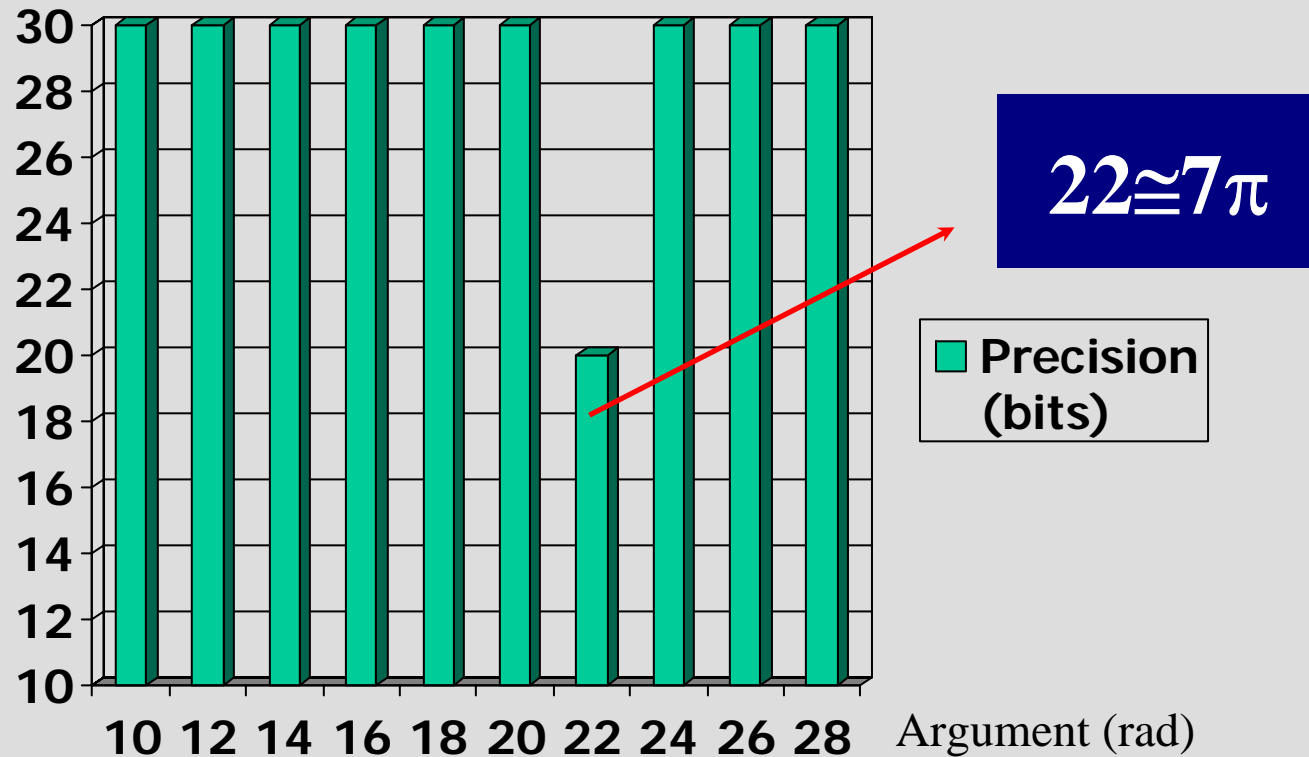
- 1) Example of an input argument close a multiple of π
 - Computation of $\sin(22 \text{ rad})$





Introduction

- 1) Arguments close to a multiple of π
 - Computation of $\sin(22 \text{ rad})$



$$R = 22 - 7\pi \approx 0.0000\ 0001\ 0001\ 1101\ 0011\ 0001\ 1001$$



Introduction

2) Large input argument

INSTRUCTION SET REFERENCE intel®

FSIN—Sine

Opcode	Instruction	Description
D9 FE	FSIN	Replace ST(0) with its sine.

Description

This instruction calculates the sine of the source operand in register ST(0) and stores the result in ST(0). The source operand must be given in radians and must be within the range -2^{63} to $+2^{63}$. The following table shows the results obtained when taking the sine of various classes of numbers, assuming that underflow does not occur.

SRC (ST(0))	DEST (ST(0))
$-\infty$	*

Range of representable numbers for IEEE 754 (single precision): $[-2^{127}, 2^{127}]$



Introduction

- Solutions
 - Software
 - Hardware



Introduction

□ Solutions

– Software

- **Cody and Waite [3]**
 - Software Manual for elementary functions (Prent. Hall 1980)
- **Payne and Hanek algorithm [5]**
 - Radian reduction for trigonometric functions (SIGNUM let. 1983)
- **Brisebarre, Defour, Kornerup, Muller, Revol**
 - A new range reduction algoritm (Trans. Comp. 2005)

Software solutions: good accuracy, slow



Introduction

□ Solutions

– Hardware

- Double Residue Modular range reduction [6]
 - Based on a table with positive and negative residues
 - Solves both problems
 - » Input argument close to multiple of a constant
 - » Large argument
 - Architecture: fast, simple hardware and valid for fixed-point & floating-point when the module is a rational or an irrational number.
 - » IEEE 754 full range
 - Handicap for floating point:
 - » word-serial implementation only

[6] Julio Villalba, Tomas Lang, Mario A. González, Double residue modular range reduction for floating-point hardware implementation, IEEE Transactions On Computers, vol. 55, no. 3, march 2006



Introduction

- Our aim:
 - Design a pipelined architecture supporting the double residue modular range reduction
 - preventing the replication of the table of residues



Double residue modular range reduction



Double residue MRR

- Fundamental of double residue MRR (*fixed-point*)
 - For a number X and a module M , obtain a number R :

$$R = X - k \cdot M \quad 0 \leq R < M$$

- Example: $M=7$, $X=372$

$$R = 372 - k \cdot 7 \quad 0 \leq R < 7$$

- Select $k=53$, thus $R=1$

- We call “residue” to R and “mod” to the operation, such that:

$$R = X \bmod M \quad (372 \bmod 7 = 1)$$



Double residue MRR

- Double residue MRR
 - Define elementary positive residues and elementary negative residues, which are stored in a table:

$$m_i^+ = 2^i \bmod M$$

$$m_i^- = 2^i \bmod M - M$$

Elementary Positive residue

Elementary negative residue

- Example: $M=7$, we define

$$m_6^+ = 2^6 \bmod 7 = 64 \bmod 7 = 1$$

$$m_6^- = 2^6 \bmod 7 - 7 = -6$$

$$m_5^+ = 2^5 \bmod 7 = 32 \bmod 7 = 4$$

$$m_5^- = 2^5 \bmod 7 - 7 = -3$$

$$m_4^+ = 2^4 \bmod 7 = 16 \bmod 7 = 2$$

$$m_4^- = 2^4 \bmod 7 - 7 = -5$$



Double residue MRR

- Double residue MRR

- Algorithm:

- Define an accumulated residue $R(i)$

$$X = x_{n-1} \dots x_1 x_0 . x_{-1} x_{-2} \dots x_{-f}$$

$$R(i+1) = R(i) + m_i^* x_i \quad \text{with}$$

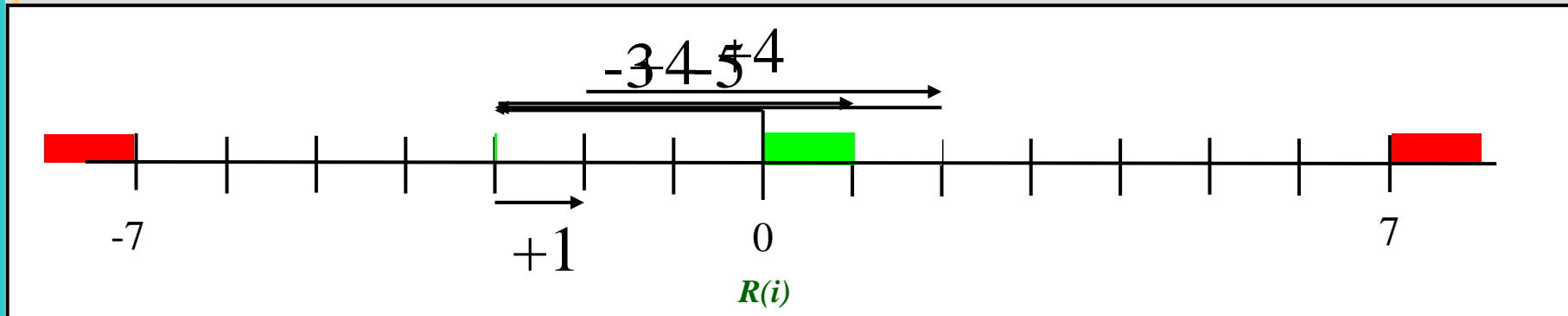
$$m_i^* = \begin{cases} m_i^+ & \text{si } R(i) < 0 \\ m_i^- & \text{si } R(i) \geq 0 \end{cases}$$

- For each power of 2, select the suitable positive or negative residue
 - The accumulated residue is **bounded by M**
 - » $R(i)$ is always inside convergence bound



Double residue MRR

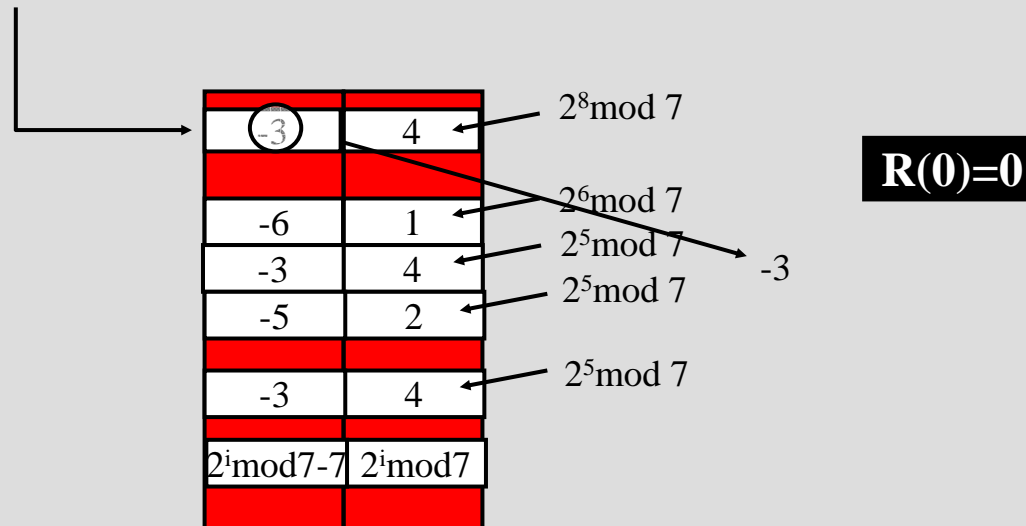
- Double residue MRR: an intuitive idea
 - Computation of $372 \bmod 7 = 1$
 - $372 = 101110100 = 2^8 + 2^6 + 2^5 + 2^4 + 2^2 = 256 + 64 + 32 + 16 + 4$
 $(256 \bmod 7 + 64 \bmod 7 + 32 \bmod 7 + 16 \bmod 7 + 4 \bmod 7) =$
 $\begin{matrix} \searrow (4,-3) & \searrow (1,-6) & \searrow (4,-3) & \searrow (2,-5) & \searrow (4,-3) \\ -3 & +1 & +4 & +(-5) & +4 \end{matrix} = 1$





Double residue MRR

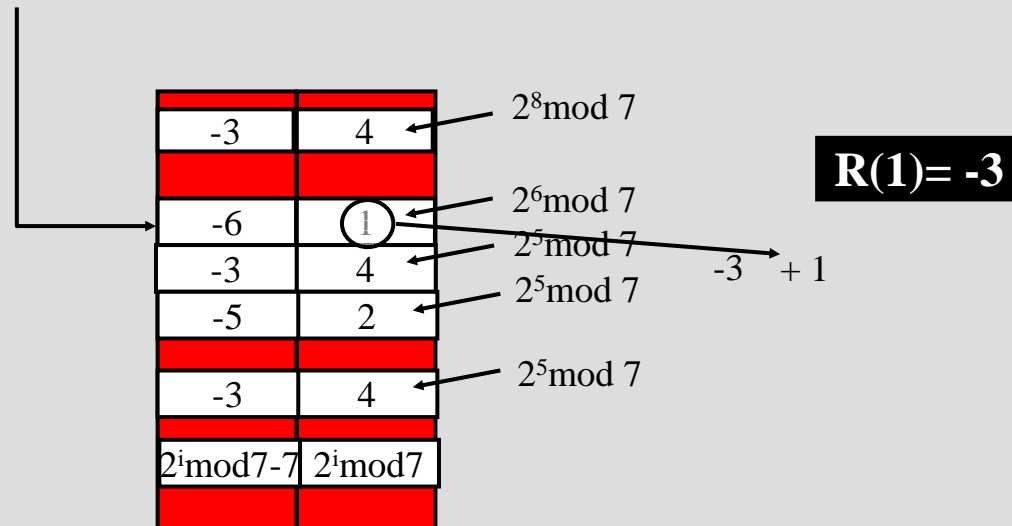
- Double residue MRR: an intuitive idea
 - Computation of $372 \bmod 7 = 1$
 - $372 = 101110100 = 2^8 + 2^6 + 2^5 + 2^4 + 2^2 = 256 + 64 + 32 + 16 + 4$





Double residue MRR

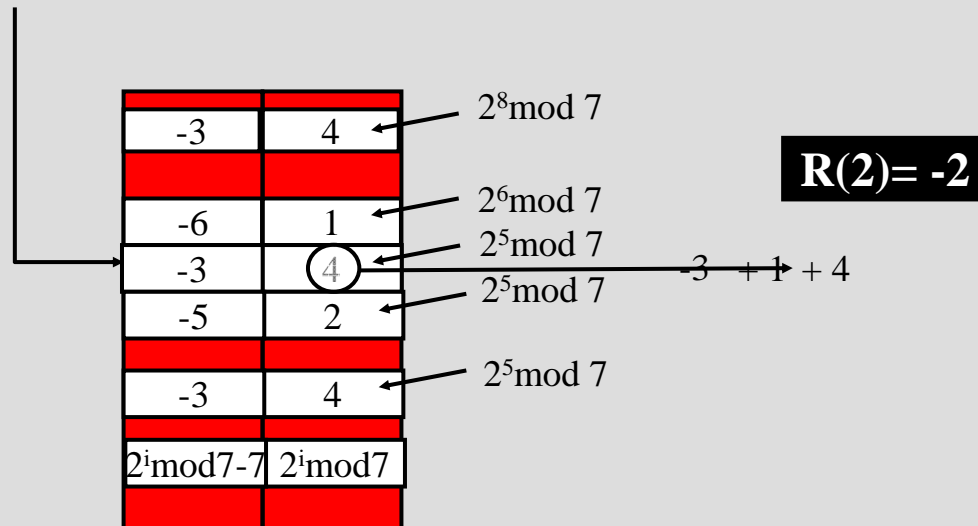
- Double residue MRR: an intuitive idea
 - Computation of $372 \bmod 7 = 1$
 - $372 = 101110100 = 2^8 + 2^6 + 2^5 + 2^4 + 2^2 = 256 + 64 + 32 + 16 + 4$





Double residue MRR

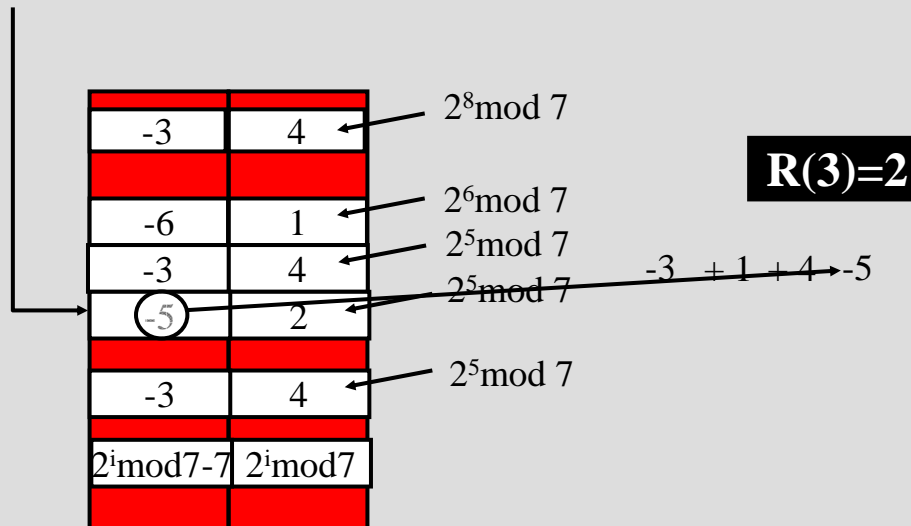
- Double residue MRR: an intuitive idea
 - Computation of $372 \bmod 7 = 1$
 - $372 = 101110100 = 2^8 + 2^6 + 2^5 + 2^4 + 2^2 = 256 + 64 + 32 + 16 + 4$





Double residue MRR

- Double residue MRR: an intuitive idea
 - Computation of $372 \bmod 7 = 1$
 - $372 = 101110100 = 2^8 + 2^6 + 2^5 + 2^4 + 2^2 = 256 + 64 + 32 + 16 + 4$



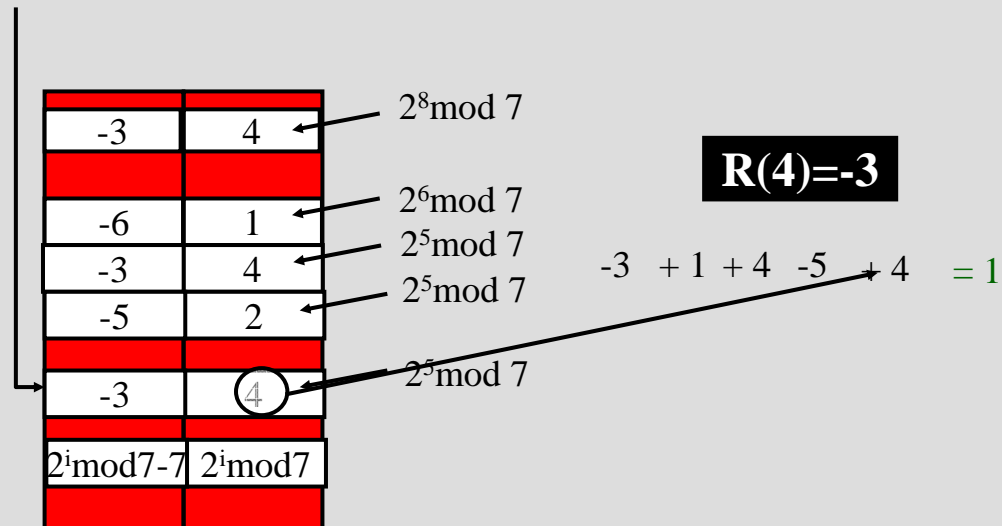


Double residue MRR

- Double residue MRR: an intuitive idea

- Computation of $372 \bmod 7 = 1$

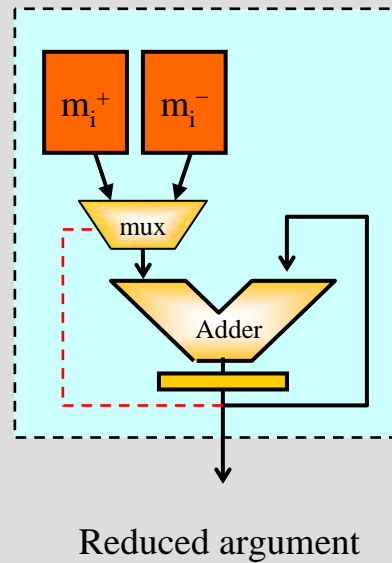
- $372 = 101110100 = 2^8 + 2^6 + 2^5 + 2^4 + 2^2 = 256 + 64 + 32 + 16 + 4$





Double residue MRR

- Word-serial Architecture:





Double residue MRR

- Implementation in carry-save arithmetic
 - $R(i) = S(i) + C(i)$
 - **Estimation** of the sign of $R(i)$ from the 3- MSBs of $R(i)$:
 - Define a **new selection function** of the elementary residues:

$$m_i^* = \begin{cases} m_i^+ & \text{if } R(i) < 0 \\ m_i^- & \text{if } R(i) \geq 0 \end{cases} \quad \longrightarrow \quad m_i^* = \begin{cases} m_i^+ & \text{if } \bar{R}(i) < 0 \\ m_i^- & \text{if } \bar{R}(i) \geq 0 \end{cases}$$

$\bar{R}(i)$ **Truncated value of $R(i)$ using the 3 MSBs of $S(i), C(i)$**



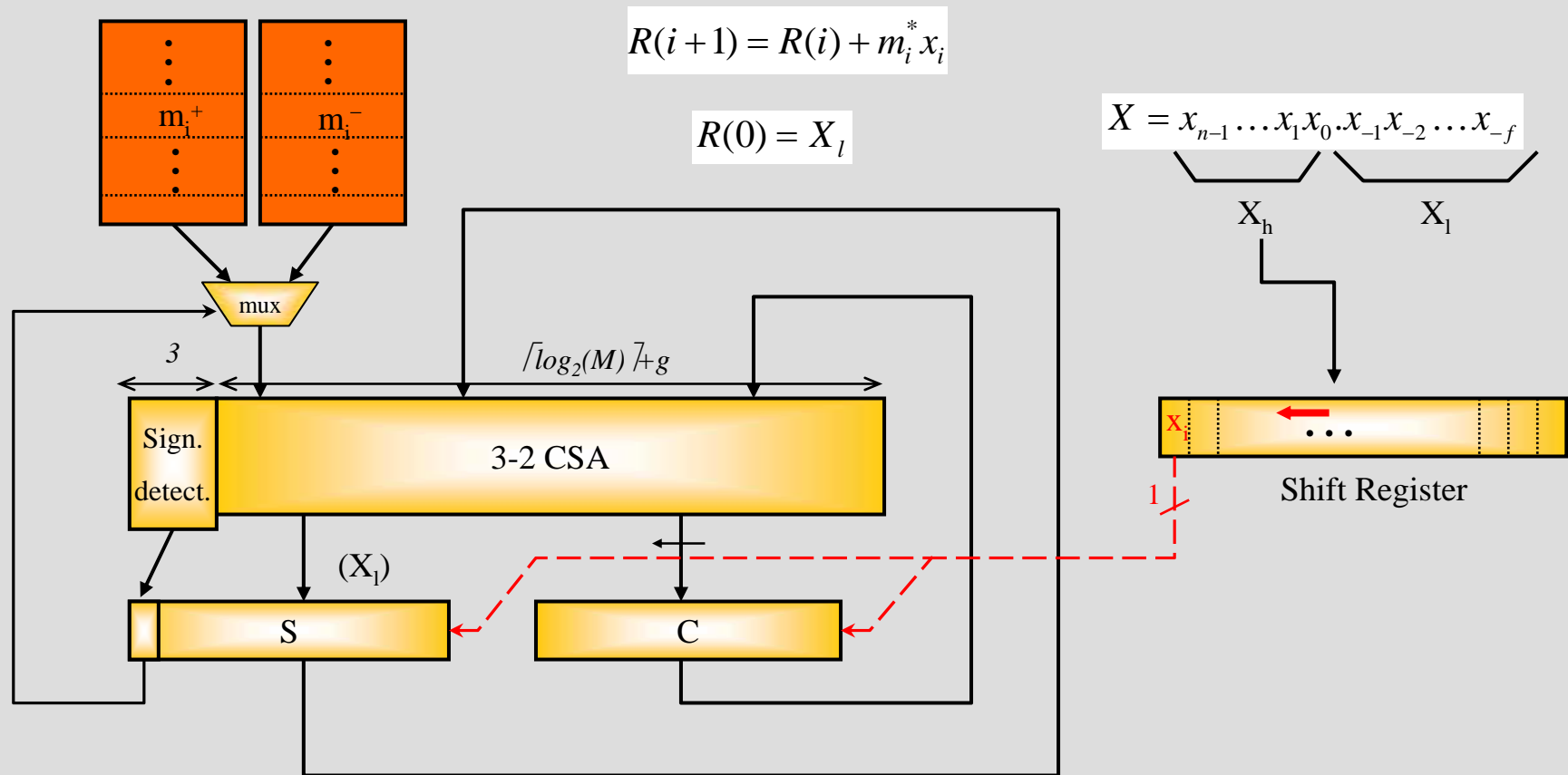
Double residue MRR

- Implementation in carry-save arithmetic
 - The iterations are self-correcting
 - The deviation produced by a wrong estimation of the sign of $R(i)$ is corrected by the next iteration
 - A final correction iteration is required to ensure the convergence



Double residue MRR

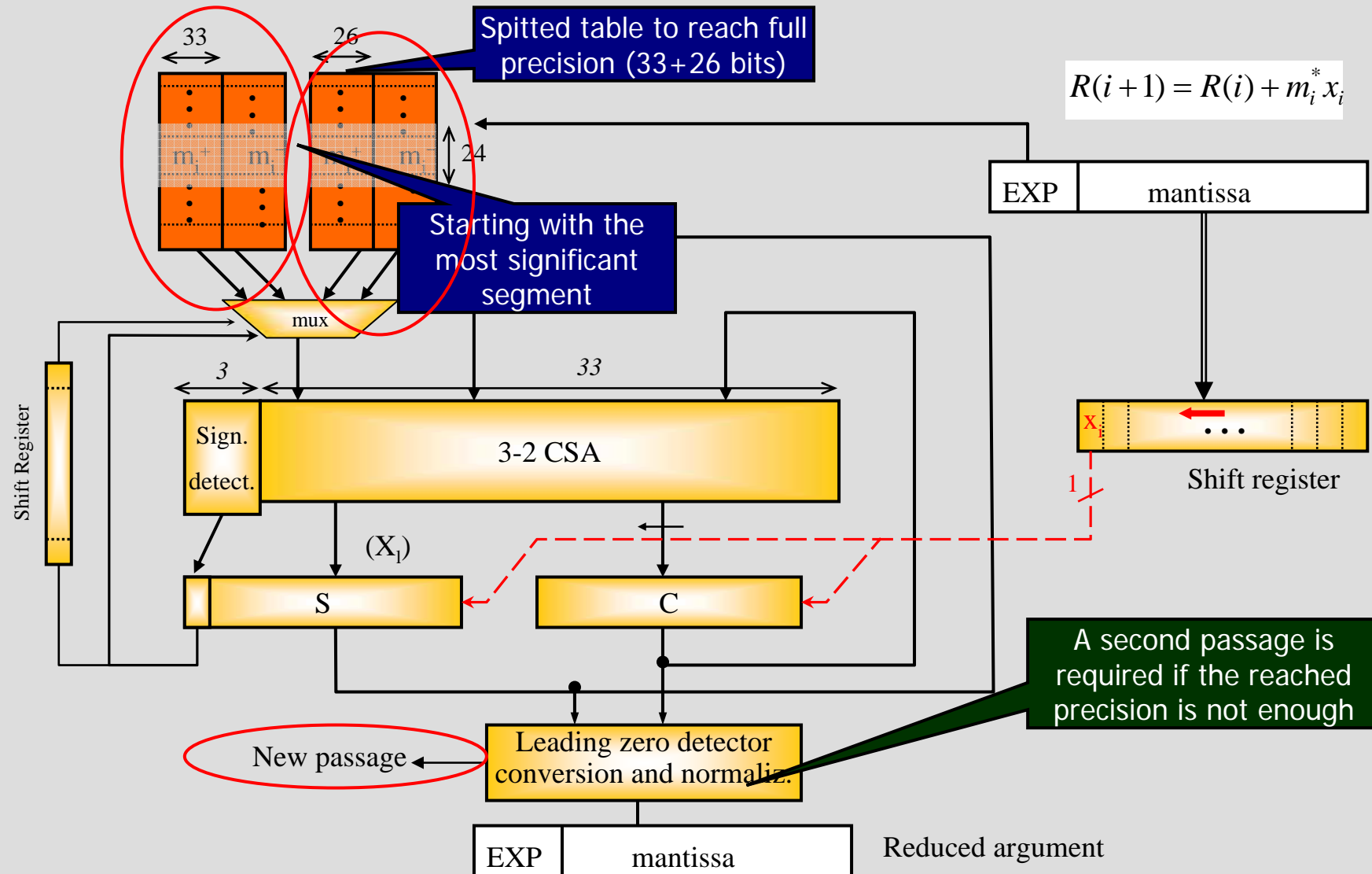
- Architecture in carry-save arithmetic (fixed-point)





Double residue MRR

- Architecture in carry-save arithmetic (IEEE 754 floating-point, sing. prec. $M=2\pi$) [6]



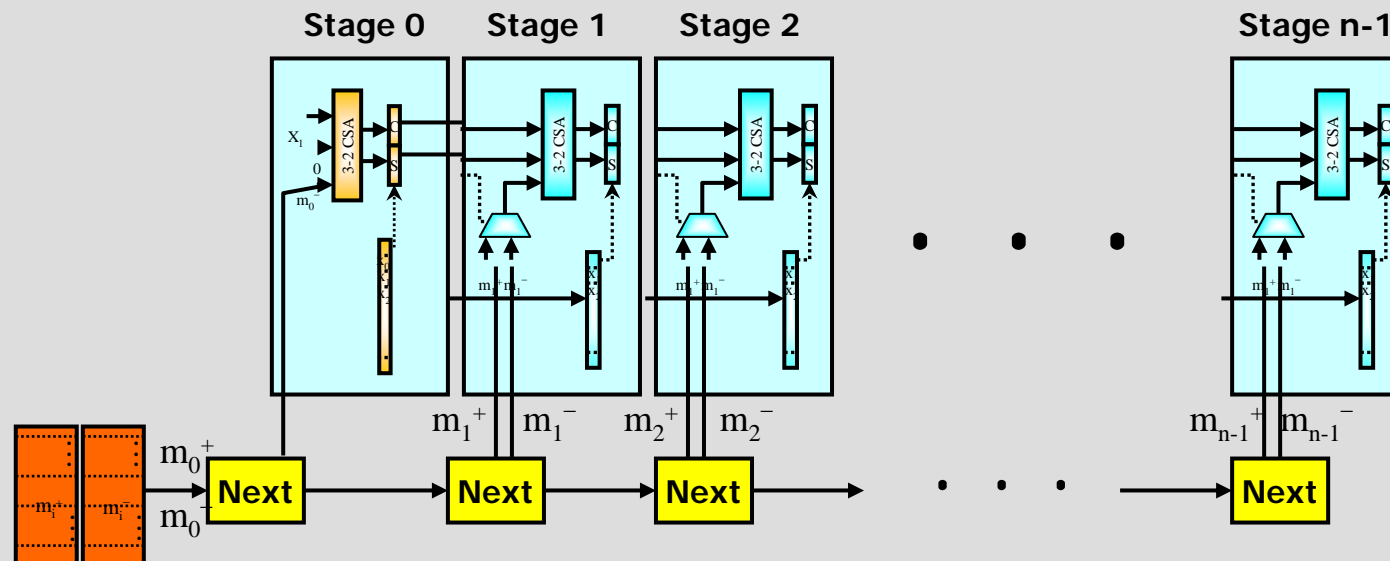


Pipelined architecture



Pipelined architecture

- Our proposal:
 - Prevents the replication of the table by generation of the elementary residue for the next stage





Pipelined architecture

- Preventing the replication:
 - Define a elementary residue generator as:

$$m_{i+1}^+ = \begin{cases} 2m_i^+ & \text{if } 2m_i^+ < M \\ 2m_i^+ - M & \text{if } 2m_i^+ \geq M \end{cases}$$

$$m_{i+1}^- = \begin{cases} 2m_i^- + M & \text{if } 2m_i^- < M \\ 2m_i^- & \text{if } 2m_i^- \geq M \end{cases}$$

- A single table of elementary residues is required
- The elementary residue for the stage i is obtained from the elementary residue of the previous stage $i-1$



Pipelined architecture

- Reducing the hardware cost of the elementary residue generator

$$\begin{array}{l} \text{if } 2m_i^+ < M \text{ then} \\ \quad \left\{ \begin{array}{l} m_{i+1}^+ = 2m_i^+ \\ m_{i+1}^- = 2m_i^- + M \end{array} \right. \quad \text{FLAG}_i = 0 \\ \\ \text{else} \\ \quad \left\{ \begin{array}{l} m_{i+1}^+ = 2m_i^+ - M \\ m_{i+1}^- = 2m_i^- \end{array} \right. \quad \text{FLAG}_i = 1 \end{array}$$

- A single comparison is required
- The evolution of the generation of elementary residues can be precomputed and stored (Flags)
 - » **Comparison is prevented**
 - » **Only one addition is required in each stage**



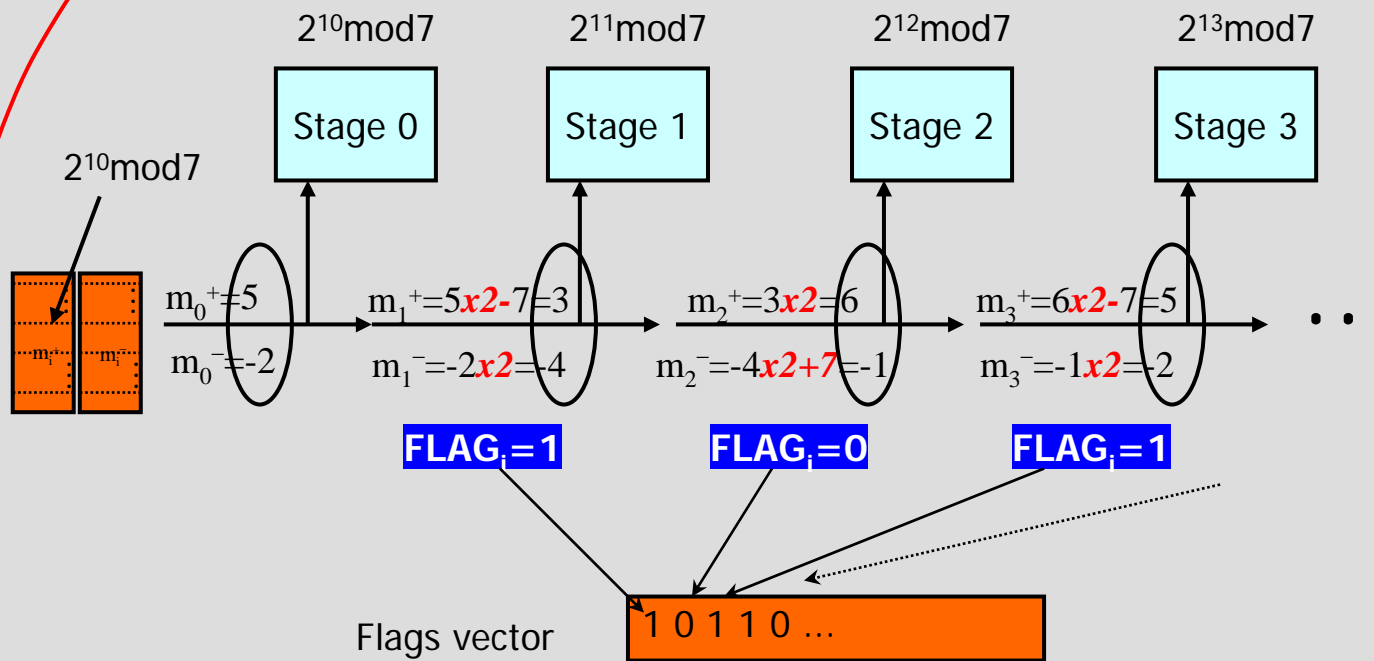
Pipelined architecture

– Example (M=7)

X=1.101001.. 2¹⁰

$$\text{if } 2m_i^+ < M \text{ then } \begin{cases} m_{i+1}^+ = 2m_i^+ \\ m_{i+1}^- = 2m_i^- - M \end{cases} \text{ FLAG}_i = 0$$

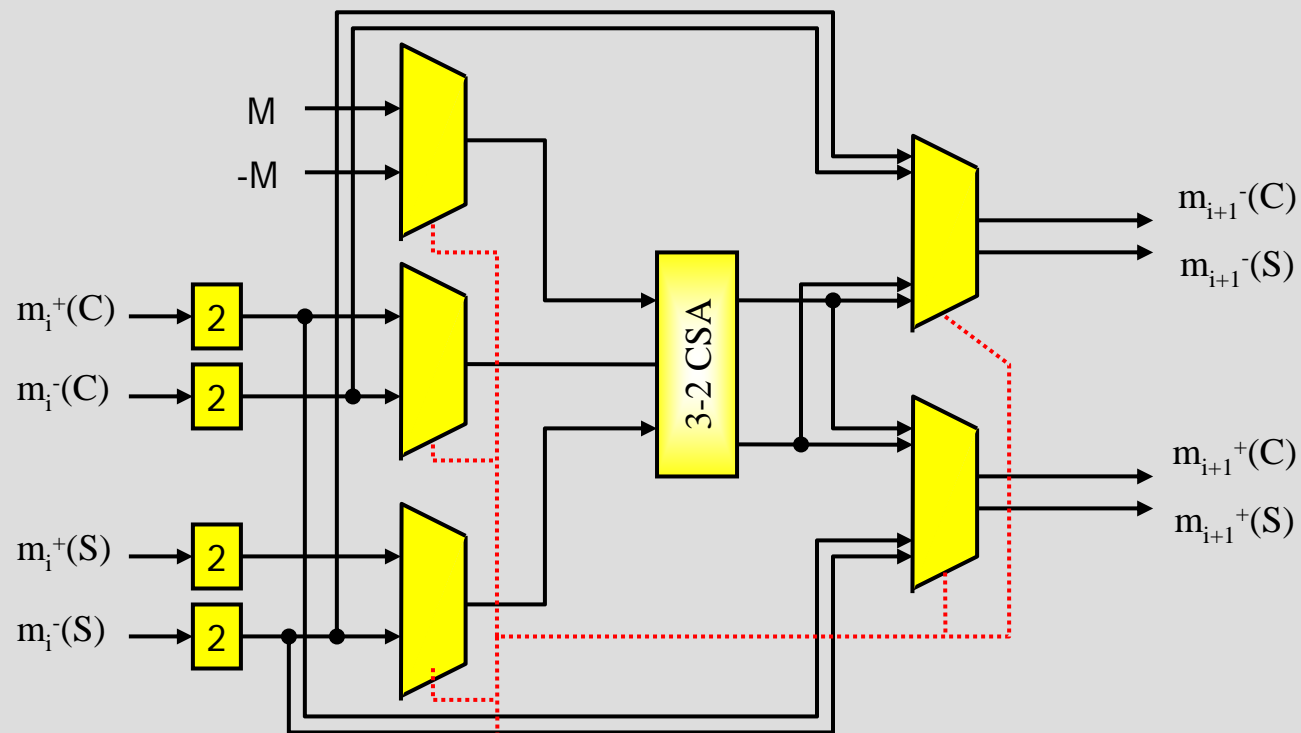
$$\text{else } \begin{cases} m_{i+1}^+ = 2m_i^+ - M \\ m_{i+1}^- = 2m_i^- \end{cases} \text{ FLAG}_i = 1$$





Pipelined architecture

- Carry-save elementary residue generator (module **Next**)



FLAG

FLAG_i=0

$$\begin{cases} m_{i+1}^+ = 2m_i^+ \\ m_{i+1}^- = 2m_i^- + M \end{cases}$$

FLAG_i=1

$$\begin{cases} m_{i+1}^+ = 2m_i^+ - M \\ m_{i+1}^- = 2m_i^- \end{cases}$$



Pipelined architecture

- Error propagation

- Computing the accumulative addition $R(i)$

$$R(i+1) = R(i) + m_i^* x_i$$

- Computing elementary residues m_i

$$m_{i+1}^* = 2m_i^* - M$$

Aim: precision 1 ULP using guard bits



Pipelined architecture

- Error propagation
 - Computing the accumulative addition $R(i)$

$$R(i+1) = R(i) + m_i^* x_i$$

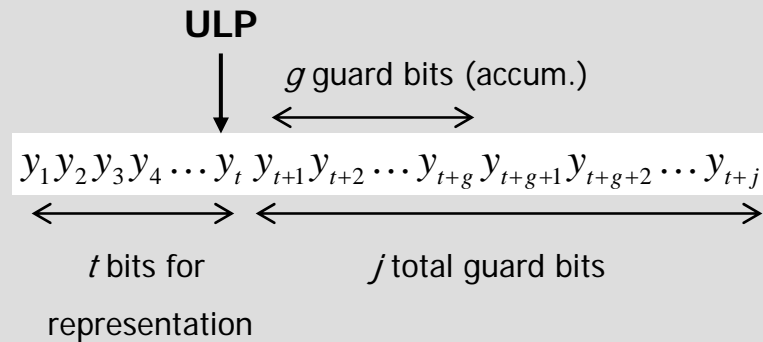
- *n iterations:*
 - » *guard bits:* $g = \lceil \log_2 \sqrt{n} \rceil$



Pipelined architecture

- Error propagation
 - Computing elementary residues m_i

$$m_{i+1}^* = 2m_i^* - M$$



Initial error ε :

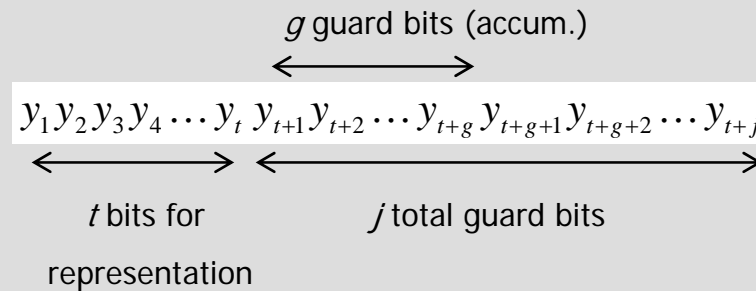
$$\varepsilon < 2^{-t-j}$$



Pipelined architecture

- Error propagation

- Computing elementary residues m_j



Initial error ε : $\varepsilon < 2^{-t-j}$

Recurrence:

$$\begin{cases} a_0 = \varepsilon \\ a_k = (2a_{k-1} + 1)\varepsilon \end{cases} \Rightarrow a_k = (2^{k+1} - 1)\varepsilon$$

} $j > g + \log_2(2^{k+1} - 1)$



Pipelined architecture

- Error propagation conclusion:

- Computing the accumulative addition $R(i)$

$$R(i+1) = R(i) + m_i^* x_i$$

Guard bits:

$$g = \lceil \log_2(n+1) \rceil$$

- Computing elementary residues m_i

$$m_{i+1}^* = 2m_i^* - M$$

Guard bits:

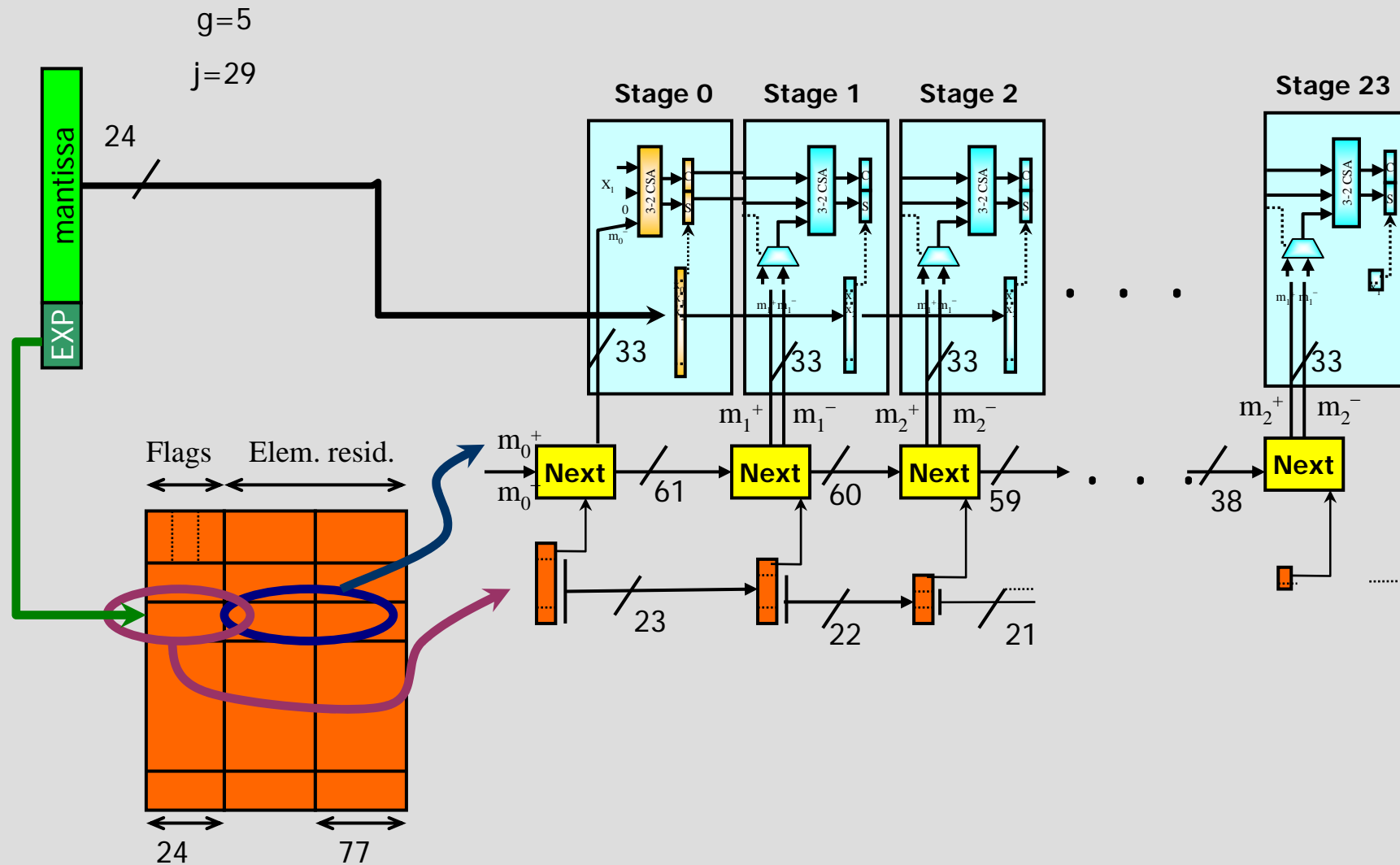
$$j = n + \lceil \log_2(n+1) \rceil$$

- The number of guard bits required for m_i decreases by one in every iteration
 - Hardware reduction in every stage



Pipelined architecture

- Architecture for IEEE-754 single precision ($M=2\pi$, $n=24$)





Future works

- Radix-4 coding $X = x_1x_2x_3\dots x_i\dots x_n$, $x_i \in \{-2, -1, 0, 1, 2\}$
 - Reduces the number of stages by half

$$R(i+1) = R(i) + \sigma_i m_i^*$$

- Selection function more complex
 - Increases the stage delay

$$\sigma_i = \text{sign}(x_i)$$

$$m_i^* = \begin{cases} m_{i+1}^+ & \text{if } R(i) < 0 \ \& \ |x_i| = 2 \\ m_i^+ & \text{if } R(i) < 0 \ \& \ |x_i| = 1 \\ 0 & \text{if } |x_i| = 0 \\ m_i^- & \text{if } R(i) \geq 0 \ \& \ |x_i| = 1 \\ m_{i+1}^- & \text{if } R(i) \geq 0 \ \& \ |x_i| = 2 \end{cases}$$



Summary and Conclusion

- Novel pipeline architecture for range reduction
 - Carry-save arithmetic
 - Based on double elementary residues
 - Prevent replication of tables of the word-serial version
 - obtaining next elementary residue from the previous one
 - Elementary residue generator
 - Precision 1 ULP
- Valid for floating point reduction
 - IEEE 754 compliant
- Radix-4
 - Reduces the number of stages by half



THANK YOU !