

# Montgomery Modular Multiplication Algorithm for Multi-Core Systems

Junfeng Fan, Kazuo Sakiyama and Ingrid Verbauwhede

Katholieke Universiteit Leuven, ESAT/SCD-COSIC,  
Kasteelpark Arenberg 10  
B-3001 Leuven-Heverlee, Belgium

June 12, 2007

- 1 Outline
  - Outline
- 2 Introduction
  - Montgomery Modular Multiplication Algorithm
  - Implementation Considerations
- 3 Design Platform
  - Multi-Core Systems
  - A Prototype Processor
  - Instruction Set Architecture
- 4 Instruction Scheduling Methods
  - Data Dependency
  - Scheduling Method-I
  - Scheduling Method-II
  - Performance Comparison
- 5 Implementation Results
  - Scalability
  - Performance Comparison
- 6 Future Work

## 1 What is Montgomery Modular Multiplication (MMM) Algorithm?

### The Montgomery Multiplication Algorithm

Given  $n$ -bit modulo  $M$ , integer  $x, y \in \mathbb{Z}_M$ ,  $R = 2^n$

$$\text{Mont}(x, y) = x \cdot y \cdot R^{-1} \bmod M$$

## 1 What is Montgomery Modular Multiplication (MMM) Algorithm?

### The Montgomery Multiplication Algorithm

Given  $n$ -bit modulo  $M$ , integer  $x, y \in \mathbb{Z}_M$ ,  $R = 2^n$

$$\text{Mont}(x, y) = x \cdot y \cdot R^{-1} \bmod M$$

## 2 Why Use Montgomery Modular Multiplication Algorithm?

### Use Normal Multiplication

$$Z = A \cdot B \bmod M$$

- 1  $C = A \cdot B$
- 2  $Z = C - \lfloor \frac{C}{M} \rfloor \cdot M$

### Use MMM

$$Z = A \cdot B \bmod M$$

- 1  $A' = \text{Mont}(A, R^2) = A \cdot R \bmod M$
- 2  $Z = \text{Mont}(A', B) = A \cdot B \bmod M$

- 1 What is Montgomery Modular Multiplication (MMM) Algorithm?

## The Montgomery Multiplication Algorithm

Given  $n$ -bit modulo  $M$ , integer  $x, y \in \mathbb{Z}_M$ ,  $R = 2^n$

$$\text{Mont}(x, y) = x \cdot y \cdot R^{-1} \bmod M$$

- 2 Why Use Montgomery Modular Multiplication Algorithm?

### Use Normal Multiplication

$$Z = A \cdot B \bmod M$$

- 1  $C = A \cdot B$
- 2  $Z = C - \lfloor \frac{C}{M} \rfloor \cdot M$

### Use MMM

$$Z = A \cdot B \bmod M$$

- 1  $A' = \text{Mont}(A, R^2) = A \cdot R \bmod M$
  - 2  $Z = \text{Mont}(A', B) = A \cdot B \bmod M$
- 3 Widely used in RSA, ECC, Diffie-Hellman...

## Radix- $2^w$ Montgomery Modular Multiplication Algorithm

**Input:** integers  $M = (M_{s-1}, \dots, M_0)_r$ ,  $X = (X_{s-1}, \dots, X_0)_r$ ,  
 $Y = (Y_{s-1}, \dots, Y_0)_r$ , where  $0 \leq X, Y < M$ ,  $r = 2^w$ ,  $s = \lceil \frac{n}{w} \rceil$ ,  $R = r^s$  with  
 $\gcd(M, r) = 1$  and  $M' = -M^{-1} \bmod r$ .

**Output:**  $X \cdot Y \cdot R^{-1} \bmod M$

$Z = (Z_{s-1}, \dots, Z_0)_r \leftarrow 0$

**for**  $i = 0$  to  $s - 1$  **do**

$T \leftarrow (Z_0 + X_0 \cdot Y_i) \cdot M' \bmod r$

$Z \leftarrow (Z + X \cdot Y_i + M \cdot T) / r$

**end for**

**if**  $Z > M$  **then**

$Z \leftarrow Z - M$

**end if**

**return**  $Z$

## Hardware Implementations

- 1 Fast, Power efficient
  - 1 special data-path
  - 2 multiple processing elements (PE)
- 2 expensive, fixed functions
  - 1 Cost extra hardware
  - 2 Hard to update

## Software Implementations

- 1 Cheap, flexible
  - 1 Sharing CPU with other applications
  - 2 Easy to modify
- 2 Slow
  - 1 General purpose data-path
  - 2 Normally single core

## Hardware Implementations

- 1 Fast, Power efficient
  - 1 special data-path
  - 2 multiple processing elements (PE)
- 2 expensive, fixed functions
  - 1 Cost extra hardware
  - 2 Hard to update

## Software Implementations

- 1 Cheap, flexible
  - 1 Sharing CPU with other applications
  - 2 Easy to modify
- 2 Slow
  - 1 General purpose data-path
  - 2 Normally single core

The question is:

How about using multi-core systems?

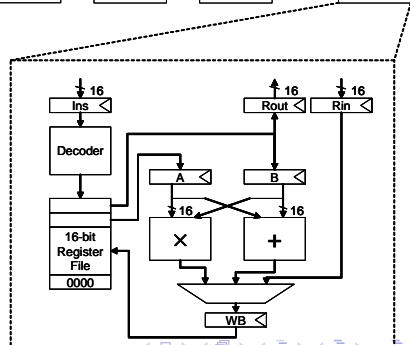
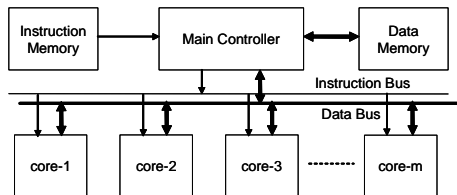


In the real world, a multi-core system can be

- 1 A processor with multiple cores: shared cache
- 2 A system with multiple processors: shared memory

Our prototype processor

- 1 Very Long Instruction Set (VLIW)
- 2 Shared single-port data memory



Opcode 4-bit	Operand 1 4-bit	Operand 2 4-bit	Operand 3 4-bit	Description
<b>Nop</b>				No operation
<b>Load</b>	<b>Ri</b>	<b>#Addr</b>		Load the data from location <b>Addr</b> of the data memory into register <b>Ri</b>
<b>Store</b>	<b>Ri</b>	<b>#Addr</b>		Store the data of register <b>Ri</b> to location <b>Addr</b> or the data memory
<b>Mul</b>	<b>Ri</b>	<b>Rj</b>	<b>Rk</b>	$R(i+1), Ri = Rj \cdot Rk$
<b>Add</b>	<b>Ri</b>	<b>Rj</b>	<b>Rk</b>	$Ca, Ri = Rj + Rk$ , <b>Ca</b> is the carry out and is stored in the status register
<b>Adc</b>	<b>Ri</b>	<b>Rj</b>	<b>Rk</b>	$Ca, Ri = Rj + Rk + Ca$
<b>Sub</b>	<b>Ri</b>	<b>Rj</b>	<b>Rk</b>	$Ri = Rj - Rk$

The question is:

How to map the Montgomery Modular Multiplication to this platform?

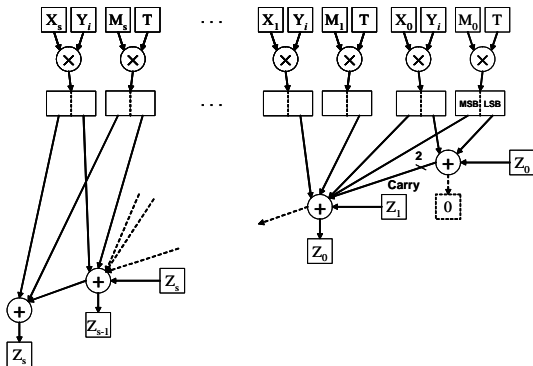
## Data dependency in one loop

for  $i = 0$  to  $s - 1$  do

$$T \leftarrow (Z_0 + X_0 \cdot Y_i) \cdot M' \bmod r$$

$$Z \leftarrow (Z + X \cdot Y_i + M \cdot T) / r$$

end for

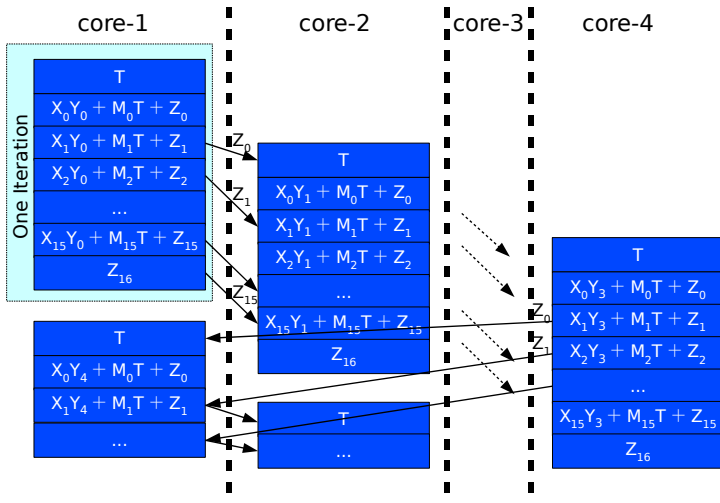


## Basic considerations

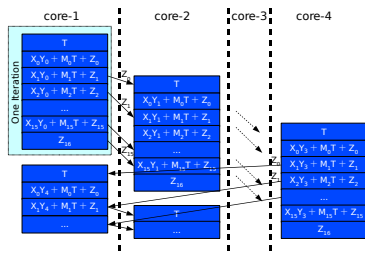
- 1 Number of **Mul** and **Add** are almost constant
- 2 Data transfers are expensive
- 3 Carry should be used in the local core

## We propose

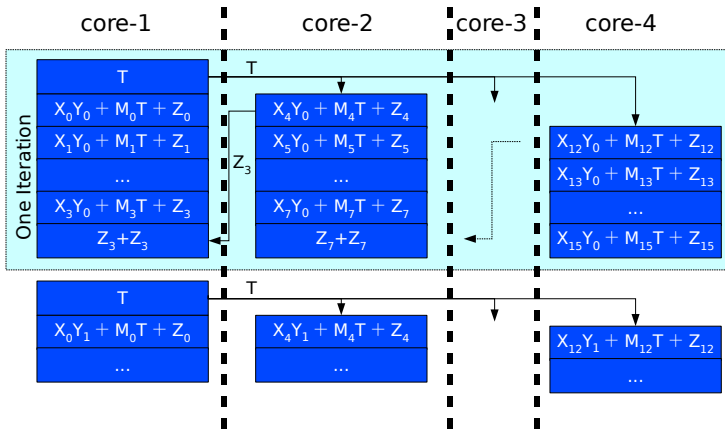
- 1 Instruction scheduling method-I: Each core performs one iteration
- 2 Instruction scheduling method-II: Multiple cores perform one iteration



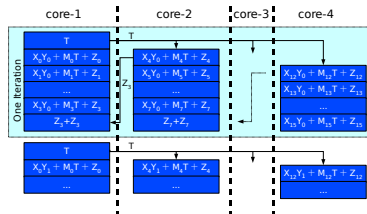
- 1 Carry is always used in the local core
- 2 Data transfers cause a heavy overhead
  - 1 Suppose  $Z$  has  $s$  words, one multiplication requires  $s(s - 1)$  data transfers
  - 2 For example, when performing 256-bit MMM, **240** data transfers are needed
- 3  $X_{s-1}, \dots, X_0$  and  $M_{s-1}, \dots, M_0$  are loaded to each core in each iteration







- 1 Carry is always used in the local core
- 2 Less data transfers are required
  - 1 Suppose  $Z$  has  $s$  words and a  $p$ -core system is used, one multiplication requires  $3ps - 2s$  data transfers
  - 2 For example, when performing 256-bit MMM on a 4-core system, **96** data transfers are needed
- 3 Only  $\lceil \frac{s}{p} \rceil$  words of  $X_{s-1}, \dots, X_0$  and  $M_{s-1}, \dots, M_0$  are loaded to each core in each iteration



Compared to the method-I, the method-II has two major advantages.

- ① Operands and intermediate data are distributed in the register file of each core, thus less registers are required in each core.
- ② Less data transfers reduce memory accesses, as a result, a single-port data memory can support more cores before becoming the bottleneck.

**Table:** Number of memory accesses required for one Montgomery multiplication for various Register File size ( $S_{rf}$ ).

Processor type	$S_{rf}$	$N_{load-opr}$	$N_{load-tr}$	$N_{store-tr}$	$N_{total}$
Single-core	$S_{rf} > 3s$	$3s$	0	0	$3s$
	$2s < S_{rf} \leq 3s$	$s^2 + 2s$	0	0	$s^2 + 2s$
	$s < S_{rf} \leq 2s$	$2s^2 + s$	0	0	$2s^2 + s$
	$S_{rf} \leq s$	$2s^2 + s$	$s(s-1)^*$	$s^2^*$	$4s^2$
Multi-core Method-I	$S_{rf} > 2s$	$2ps + s$	$s(s-1)$	$s^2$	$2s^2 + 2ps$
	$s < S_{rf} \leq 2s$	$s^2 + ps + s$	$s(s-1)$	$s^2$	$3s^2 + ps$
	$S_{rf} \leq s$	$2s^2 + s$	$s(s-1)$	$s^2$	$4s^2$
Multi-core Method-II	$S_{rf} > \frac{3s}{p}$	$2s + ps$	$2(p-1)s$	$ps$	$5ps$
	$\frac{2s}{p} < S_{rf} \leq \frac{3s}{p}$	$s^2 + ps + s$	$2(p-1)s$	$ps$	$s^2 + 4ps - s$
	$\frac{s}{p} < S_{rf} \leq \frac{2s}{p}$	$2s^2 + ps$	$2(p-1)s$	$ps$	$2s^2 + 4ps - 2s$
	$S_{rf} \leq \frac{s}{p}$	$2s^2 + s$	$s^2 + (2p-3)s^*$	$s^2 + s^*$	$4s^2 + 2ps - s$

\*Including store and load operations caused by calculating intermediate data.

**Figure:** Number of data memory accesses for various operand bit-length.  
 ( $w = 16$ ,  $S_{rf} = 16$ ).

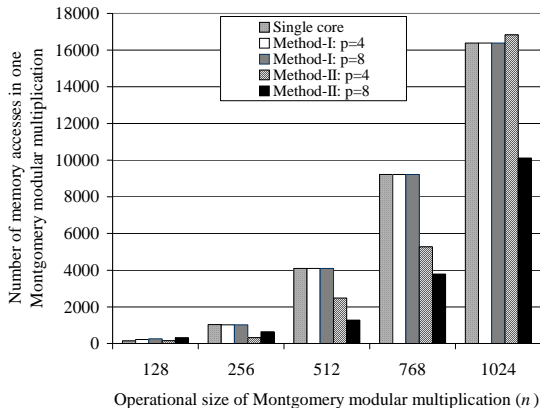
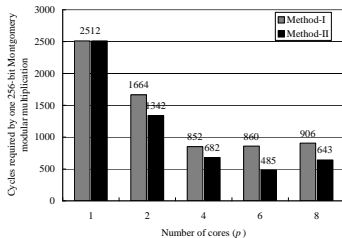


Figure: Performance of 256-bit Montgomery modular multiplication on a multi-core system. ( $n = 256$ ,  $w = 16$ ,  $S_{rf} = 16$ ).



The performance of 256-bit MMM can be improved by a factor of **1.87** and **3.68** when using 2-core and 4-core systems, respectively.[Method-II]

Table: Performance comparison of modular multiplication.

Reference	Description	Platform	Area (Slices)	Freq. (MHz)	256-bit time( $\mu$ s)	1024-bit time( $\mu$ s)
This work (method-I)	4-cores/4 16x16 mults	Xilinx	2029	125	6.8	131.0
	4-cores/4 32x32 mults	XC2VP30	3173	93	2.6	44.0
This work (method-II)	4-cores/4 16x16 mults	Xilinx	2029	125	5.5	134.7
	4-cores/4 32x32 mults	XC2VP30	3173	93	2.2	33.0
Tenca <i>et al.</i>	Software	ARM	-	80	43	570
Itoh <i>et al.</i>	Software	DSP(TMS320C6201)	-	200	2.68 <sup>‡</sup>	—
Brown <i>et al.</i>	Software	Pentium II	-	400	1.57 <sup>§</sup>	—
Kelley <i>et al.</i>	4-PEs/8 16x16 mults	XC2V2000	360*	135	0.68	8.3
Mentens	130 16x16 mults	XC2VP30	7244	64	0.31	1.07

\* Author's estimation from the original paper.

‡ 239-bit Montgomery modular multiplication.

§ Using fixed modulo for fast reduction.

- ① Hardware implementations
  - ① Use specific data-path
  - ② Use specific Register Files
- ② Software implementations
  - ① VLIW DSP
  - ② Intel quad-core processors





# Thanks!