# Montgomery Modular Multiplication via Single-Base Residue Number Systems

31<sup>st</sup> IEEE International Symposium on Computer Arithmetic

Zabihollah Ahmadpour Ghassem Jaberipur Jeong-A Lee

# ARITH 2024



Málaga, Spain, June 10-12, 2024.

#### Outline

- RNS-MMM implementation problems
- Solutions
- RNS-MMM: Conventional Base extension algorithm
- New single-base algorithm coupled with free- $\delta$  RNS
- Definitive derivation of the reduction factors
- Advantages of SB-MMM vs. BE-MMM
- Conclusion

# RNS-MMM for Crypto algorithms (e.g., RSA)

- Typical Key-length:
  - $n = \mathbf{k} \times \mathbf{r} = 2^i (9 \le i \le 13)$

k = # of co-prime moduli,

r = Residue channel bit-width



Source: AI generated image (DALLE)

• Ultra-fast long key-length RNS-MMM via free- $\delta$  scheme:

 $r = 16, k \in \{32, 64, 128, 256, 512\}$ 

 # of dynamically changeable k-moduli pair sets for k ∈ {32,64,128,256,512}: {2<sup>505</sup>, 2<sup>887</sup>, 2<sup>1302</sup>, 2<sup>2727</sup>, 2<sup>3833</sup>}

# MMM and RNS-MMM

Assumptions:	Algorithm								
$\Gamma > N$ ,	MMM	RNS-MMM	Problem	Solution					
$CGD(\Gamma, N) = 1,$	$W = X \times Y$	$w_i =  x_i \times y_i _{m_i}$							
$\Gamma_i =  \Gamma _{m_i},$	$Q = \left  W \times \tilde{N} \right _{\Gamma}$	$q_i = \left  \left  W \times \widetilde{N} \right _{\Gamma} \right _{m_i}$	$ W \times \widetilde{N} _{\Gamma}$	$ \begin{array}{c} \overline{\boldsymbol{\Gamma}} = \boldsymbol{M} \Longrightarrow \\ q_i = \left  \boldsymbol{w}_i \times \widetilde{\boldsymbol{N}_i} \right _{\boldsymbol{m}_i} \end{array} $					
$\left \Gamma_{i} \times \Gamma_{i}^{-1}\right _{N} = 1,$	$U = Q \times N$	$u_i = \left  q_i  N _{m_i} \right _{m_i}$							
$\widetilde{N} =  (-N)^{-1} _{\mathrm{D}}$	V = W + U	$v_i =  w_i + u_i _{m_i}$							
$\widetilde{N}_i =  \widetilde{N} $	$Z = V/\Gamma$	$z_i =  v_i \Gamma_i^{-1} _{m_i}$	$ \Gamma_i^{-1} _{m_i} =  M^{-1} _{m_i}^*$	Base Extension					
	* $ M^{-1} _{m_i}$ does not exist								

#### Half a century BE development and usage



# **BE-MMM:** Architecture and formulas



# **Motivation**

- Problem: RNS-unfriendly division operation
- Current Solution: Base extension
  - BE-MMM: Two non-overlapped k-size co-prime (2<sup>r</sup> δ) moduli sets: Dynamic range (DR) of each set: 2<sup>kr</sup>
  - BE-MMM drawbacks

✓ Requires double of the essential number of co-prime moduli to support the working DR

- Proposed Solution: SB-MMM: No base extension ⇒
  - $X^2$  DR (i.e.,  $2^{2kr}$ ) size of the same r
  - or same DR with half bit-width residue channels (i.e., r/2)
  - Example : *r* = 24

 $\{2^{170}, 2^{294}, 2^{506}, 2^{835}, 2^{1309}\}, \text{ for } k \in \{22, 43, 86, 171, 342\}$ 

## Proposed algorithm (SB-MMM)

Inputs:  $m_i, x_i, y_i, 0 \le i \le k$ , Outputs:  $\widehat{z'_i}^+, 0 \le i \le k$ 1) For i = 0 to k do par  $\{w_i = |x_i y_i|_{m_i}; \xi_{x_i} = |x_i M_i^{-1}|_{m_i}; \xi_{y_i} = |y_i M_i^{-1}|_{m_i}\};$ 2) For i = 0 to k do par  $\begin{cases} \xi_{\omega_i} = \left| w_i \left| \widehat{N} M_i^{-1} \right|_{m_i} \right|_{m_i}; \\ u_i = \left| w_i \left| M_i^{-1} \left( \left| \frac{M_i}{m_i} \right| M_i^{-1} + \left| \frac{N}{m_i} \right| \widetilde{N} \right) \right|_{m_i} \right|_{m_i} \end{cases};$ 3) For i = 0 to k do par  $\left\{ s_i = \left| \left| \xi_{x_i} \xi_{y_i} | M_i |_{m_i} m_i^{-1} + \xi_{\omega_i} | N |_{m_i} m_i^{-1} \right|_{2^{2r}} \right|_{m_i} \right\};$ 4) For i = 0 to k do par  $\begin{cases} p_{x_i} = \left(\sum_{j=0, j\neq i}^{k} \xi_{x_j} m_j^{-1}\right) - \left|\gamma_X^e + \frac{1}{2}\right|; \\ p_{y_i} = \left(\sum_{j=0, j\neq i}^{k} \xi_{y_j} m_j^{-1}\right) - \left|\gamma_Y^e + \frac{1}{2}\right|; \\ p_{\omega_i} = \left(\sum_{j=0, j\neq i}^{k} \xi_{\omega_j} m_j^{-1}\right) - \left|\gamma_{\Omega}^e + \frac{1}{2}\right|; \end{cases};$ 5) For i = 0 to k do par  $\widehat{z'_i}^+ = |N + x_i| p_{y_i}|_{m_i} + y_i |p_{x_i}|_{m_i} + |N|_{m_i} |p_{\omega_i}|_{m_i} + u_i + s_i |_{m_i}$ ;

### SB-MMM: Architecture and formulas



#### **SB-MMM:** Detailed Architecture



#### Definitive derivation of the reduction factors

Main architecture

The cyclic register file





Initial value of MAC =  $2^{r-1}$  for supporting  $|\gamma + \frac{1}{2}|$ 

Definitive derivation via an additional modulo  $m_0 = 8$ that allows  $\Gamma = M > 6N$ 

#### Three scenarios for using SB-MMM

1) Speedup via reducing the bit-widths to half, thus doubling the number of moduli, for the same working DR.

2) More dynamism in moduli set selection to reduce the probability of successful side-channel attacks.

 Keeping the same bit-widths, but reducing the total number of moduli to half, for faster CRT-like operations. **Example**: Key-length = 1024,  $r = 64 \implies k = 16 \implies$ BE-MMM delay:  $2 \times 16M_{64}$  delay. Total # of moduli: 32.

 $r = 32 \implies k = 32 \implies$  **SB-MMM** delay:  $32M_{32}$  delay  $(M_{32\approx} \text{ delay} \approx 80\% M_{64} \text{ delay in parallel PPR and PPA}).$ Total # of moduli: 32. **Example**: key-length n = 1024, BE-MMM: r = 32: Total # of free- $\delta$  co-prime moduli of the form  $2^{32} - \delta$ : 384000 Total # of moduli-set selection:  $\binom{384000}{32} \times \binom{383968}{32} \approx 2^{2374}$ . SB-MMM: r = 64. Total # of free- $\delta$  co-prime moduli of the form  $2^{64} - \delta > 2^{51}$ Total # of moduli-set selection:  $> \binom{2^{51}}{16} \approx 2^{816}$ **Example**: key-length n = 1024, BE-MMM:  $r = 32 \implies k = 32 \implies$  Total # of moduli:  $2 \times 32 = 64 \implies$ Delay:  $2 \times 32M_{32}$ SB-MMM:  $r = 32 \implies k = 32 \implies$  Total # of moduli:  $32 \implies$ Delay:  $32M_{32}$ 

#### # of multiplications: SB-MMM versus [5]

Steps of the SB-MMM	# of multiplications					
algorithm	CDP	Total				
1	3	$3 \times 3 = 9$				
2	3	3 + 3 = 6				
3	2	$3 + 1 = 4$ $3 \times 3k$				
4	k					
5	3	5				
Grand Total	k + 11	k(9k + 24)				
Ref. [5] (Unpipelined)	6k + 15	h(ch + 1E)				
Ref. [5] (Pipelined)	2k + 22	κ(οκ + 15)				

The CDP of un-pipelined SB-MMM vs. similar design of [5] is more than 83% shorter and 50% shorter than that of its pipelined version

Cost of the latter advantage is 50% increase In the total # of multiplications vs. that of [5]

# Synthesis result: SB-MMM versus [5]

Design	Size of the moduli pool	Per res chan	sidue nnel							MMM Delay			
		Area ( <b>mm</b> ²)	Power ( <i>mW</i> )	r l	k	# of clock cycle	Clock time (ns)	Of successful attacks	n	(ns)	Ratio	$\begin{array}{c} \mathbf{AT} \\ \begin{pmatrix} ms \times \\ mm^2 \end{pmatrix} \end{array}$	$\begin{array}{c} \mathbf{PDP} \\ \binom{ms \times}{mW} \end{array}$
SB-MMM	1981	0 11 2		24	43	75	0.76	2 <sup>-294</sup>	1024	57	1.21	0.274	0.163
		0.112	00.0	24	86	118		2 <sup>-506</sup>	2048	90	1.45	0.866	0.514
	384000 0	0 125	60.0	20	32	64	0.79	2 <sup>-369</sup>	1024	51	1.08	0.221	0.102
		0.135	02.3	32	64	96		2 <sup>-678</sup>	2048	76	1.22	0.660	0.303
	> 2 <sup>51</sup> 0.	0.264	107.0	64	16	48	0.97	$< 2^{-816}$	1024	47	1	0.274	0.096
		0.304	127.9	04	32	64		$< 2^{-1631}$	2048	62	1	0.723	0.253
[5]	251	0.079	17.4	32	32	80	1.86	2 <sup>-134</sup>	1024	149	3.17	0.380	0.083

#### Selected results: SB-MMM versus [5]

- 62%, 66%, and 69% delay reductions, for channel widths r = 24, 32, and 64, respectively, versus r = 32 of [5]
- 28%, 42%, and 28% reduced area-time (AT) product, for r = 24, 32, and 64, respectively
- Reduction in the probability of successful side channel attacks by a scale of 2<sup>160</sup>

# Conclusion

- New RNS-MMM with three parallel CRT-like operations instead of two consecutive similar operations of the BE-MMM
- Definitive parallel derivation of three reduction factors
- 60% speedup, and 20% less area cost of proposed algorithm for key-length n = 1024, and channel width r = 24 in comparison to the best previous work [5], with r = 32
- Significant reduction in the probability of successful side channel attacks