ASP-DAC 2025

Efficient Key Switching Accelerator for Fully Homomorphic Encryption

Seoyoon Jang, Sungjin Park, Dongsuk Jeon Seoul National University Seoul, South Korea

Mobile Multimedia Systems Group

Motivation – Advent of FHE

- Fully Homomorphic Encryption (FHE)
 - The savior of privacy-preserving computation in cloud service
- The main bottleneck operation, Key-Switching (KS)



- Expensive Number Theoretic Transform (NTT) and inverse-NTT operations in KS
 - $\{\hat{a}_i\}_i = NTT(\{a_i\}_i) = \{\Sigma_{j=0}^{N-1} a_j \omega^{ij} \mod q\}_i$
 - At least $O(N \log N)$ computations
 - Irregular memory access pattern



- Expensive Number Theoretic Transform (NTT) and inverse-NTT operations in KS
- Frequent transitions in data access patterns (Element/Ring/Coefficient-wise)

Element-wise



- Expensive Number Theoretic Transform (NTT) and inverse-NTT operations in KS
- Frequent transitions in data access patterns (Element/Ring/Coefficient-wise)



Data access within each polynomial

(NTT/INTT,...)

- Expensive Number Theoretic Transform (NTT) and inverse-NTT operations in KS
- Frequent transitions in data access patterns (Element/Ring/Coefficient-wise)

Element-wise











- Expensive Number Theoretic Transform (NTT) and inverse-NTT operations in KS
- Frequent transitions in data access patterns (Element/Ring/Coefficient-wise)



- Expensive Number Theoretic Transform (NTT) and inverse-NTT operations in KS
- Frequent transitions in data access patterns (Element/Ring/Coefficient-wise)



- Expensive Number Theoretic Transform (NTT) and inverse-NTT operations in KS
- Frequent transitions in data access patterns (Element/Ring/Coefficient-wise)



Parameter selected for KS accelerator

- $(\log N, L, dnum, \alpha) = (17, 35, 9, 4)$
 - *N*: Number of coefficients for each polynomial
 - *L*: Maximum circuit depth level
 - *dnum*: decomposition number
 - $\alpha = \left[(L+1)/dnum \right]$
 - $\alpha > 1 \rightarrow$ computations \downarrow , required *swk* \downarrow
- Constraints for parameters selection
 - Security level $\lambda > 128$
 - Directly related with $N / \log PQ$
 - *L* large enough to guarantee $L_{boot} = 15 \sim 20$ for bootstrapping that enables FHE
 - $\rightarrow N$ providing a sufficient number of <u>lightweight prime</u> moduli for large L

$$q = 2^{l} \pm 2^{sh_0} \pm 2^{sh_1} \pm 2^{sh_2} + 1$$

sparse! \rightarrow enable efficient HW implementation of modular multipliers



- Router transferring instructions and external data to the target core
- LUT for moduli set and modulus-related constants required by each core

Overall Design of KS accelerator Core 1



Core 2

NT

NTT

unit

Ring-wise

MMAC

unit

Conv_{NT1}

Coefficient

-wise

NTT

unit

Ring-wise

- NTT unit
 - NTT/iNTT operation (unified)
 - Process $N_{NTT} = 2^9$ coefficients per each unit ($N = N_{NTT}N_r$)
- Modular-Multiply-and-Accumulate (MMAC) unit
 - Conv operation
- Local distributor: internal router

Proposed Design Techniques for Energy Efficiency



- A. Efficient Twiddle Factor Generator (TFG)
- B. Conflict-free Addressing Scheme for Single-port Memory

III. Bandwidth-efficient Behavior in Core

 $q = 2^{l} \pm 2^{sh_{0}} \pm 2^{sh_{1}} \pm 2^{sh_{2}} + 1 \quad (l > 2 \cdot sh_{1} + 1, sh_{m} > sh_{n}, \forall m > n)$

• *q and *T in Barrett modular multiplication

 \rightarrow replaced with shift-adders

•
$$T = \lfloor 2^{2l}/q \rfloor = 2^l + 2^{sh_0} + 2^{sh_1} + 2^{sh_2} - 1$$

Benefits of sparsity

$$q = 2^{l} \pm 2^{sh_{0}} \pm 2^{sh_{1}} \pm 2^{sh_{2}} + 1$$
$$(l > 2 \cdot sh_{0} + 1, sh_{m} > sh_{n}, \forall m < n$$

* q and * T in Barrett modular multiplication
→ replaced with shift-adders
T = [2^{2l}/q] = 2^l + 2^{sh} + 2^{sh} + 2^{sh} + 2^{sh} + 2^{sh} - 1

Benefits of sparsity

Using
$$l = 59$$

Take the most advantage out of this inherent sparsity!

Table 1: Moduli Set.

sh_0	sh_1	sh_2	Index	sh_0	sh_1	sh_2
28	-	-	21	27	-21	-18
25	19	-	22	27	-23	18
26	-21	-	23	-27	23	19
-29	20	-	24	27	25	22
29	23	-	25	27	26	-19
-22	19	18	26	28	20	19
-23	21	20	27	28	-22	-18
23	22	-19	28	-28	22	20
24	-19	-18	29	-28	23	20
-24	21	18	30	-28	24	19
24	23	-18	31	28	24	23
¹ 41 moduli available						
•••		,ffiai			22	-18
, -	7 31				23	-18
<i>b</i> 00	tstra	ppin	ng in F	HE	23	-20
25	24	18	36	29	24	18
-26	22	19	37	29	-25	-23
26	-22	-21	38	29	26	23
26	23	-19	39	29	27	18
-26	24	19	40	29	28	26
26	25	20	-	-	-	-
	sh_0 28 25 26 -29 29 -22 -23 23 24 -24 24 24 41 booo 25 -26 26 26 26 26 26 26		$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	sh_0 sh_1 sh_2 Index28212519-2226-21-23-2920-242923-25-22191826-232120272322-192824-19-1829-242118302423-1831 41 moduli availabSufficient forbootstrapping in F 25241836-2622193726-22-21382623-1939-26241940262520-	sh_0 sh_1 sh_2 Index sh_0 2821272519-222726-21-23-27-2920-24272923-2527-2219182628-23212027282322-1928-2824-19-1829-2824211830-282423-183128 Atl moduli availableAtl moduli availableAtl moduli available25 24183629-262219372926-22-2138292623-193929-2624194029262520	sh_0 sh_1 sh_2 Index sh_0 sh_1 282127-212519-2227-2326-21-23-2723-2920-2427252923-252726-221918262820-2321202728-222322-1928-2823-24211830-28242423-183128242423-18312824252418362923252418362924-2622193729-2526-22-213829262623-19392927-262419402928262520



(VxT)≫(I+1)





 $q = 2^{l} \pm 2^{sh_{0}} \pm 2^{sh_{1}} \pm 2^{sh_{2}} + 1$ $T = \lfloor 2^{2l}/q \rfloor = 2^{l} \mp 2^{sh_{0}} \mp 2^{sh_{1}} \mp 2^{sh_{2}} - 1$

Simplified computation

- shift-adding
- removing one multiplication (V*T)



$$q = 2^{l} \pm 2^{sh_{0}} \pm 2^{sh_{1}} \pm 2^{sh_{2}} + 1$$
$$T = \lfloor 2^{2l}/q \rfloor = 2^{l} \mp 2^{sh_{0}} \mp 2^{sh_{1}} \mp 2^{sh_{2}} - 1$$

Simplified computation

- shift-adding
- removing one multiplication



<u>Area (Power)</u> vs. Non-Sparse: **47.8 (46.0) %**↓ vs. Sparse [1]: **24.6 (22.5) %**↓

[1] Kim et al., "Fpga-based accelerators of fully pipelined modular multipliers for homomorphic encryption," ReConFig, 2019.

 $NTT(\{a_i\}_i) = \left\{ \sum_{j=0}^{N-1} a_j \boldsymbol{\omega}^{ij} \mod q \right\}_i$

- N twiddle factors required for each NTT/iNTT on N coefficients
 - \rightarrow Too much overhead in data loading latency and memory area!
- Twiddle Factor Generator
 - Saving memory area for the twiddle factors
 - Generate twiddle factors during NTT/iNTT operations







- Reduction of twiddle factor memory area $O(N) \rightarrow O(\log N)$
- Additional pre-processing stage for further reduction on seed elements

Table 2: TFG comparison.

	[2]	[6]	[17]	Ours			
Modular	\sqrt{N}	N – 2	2k - 1	$2k \pm 1$			
multipliers		$I\mathbf{v} = \mathbf{Z}$	$2\kappa - 1$	$2\kappa \pm 1$			
Area	1.62-	23.89-	0.01~				
reduction*	17.13×	5988.84×	0.91^				
Seed elements	22.63	1	18/	18			
$(\log N = 9)$	22.05	1	104	10			
Peduction*	1.26-	0.03-	5.65-				
Reduction	10.65×	0.06 ×	10.22 ×				
*Reduction for $\log N = 9 - 17$							

[2] Kim et al., "Ark: Fully homomorphic encryption accelerator with run-time data generation and inter-operation key reuse," MICRO, 2022.

[6] Geelen et al., "Basalisc: Flexible asynchronous hardware accelerator for fully homomorphic encryption," preprint, arXiv, 2022.

[17] Kim et al., "Hardware architecture of a number theoretic transform for a bootstrappable rns-based homomorphic encryption scheme," FCCM, 2020.







Algorithm 1 Conflict-free address generation scheme

Note \ll_r is a left-rotate operator. $GR = BR \circ Gray \circ BR$ where BR and Gray are a bit-reverse and gray function, respectively. **Input** *s*, *addr*, *bn* **Output** *BN*, *ADDR*

1: for s = 0; s < 6; s = s + 2 do for addr = 0; addr < 16; addr = addr + 1 do 2: for bn = 0; bn < 4; bn = bn + 1 do 3: $tmp \leftarrow \{BR(addr), BR(bn)\}$ 4: $tmp \leftarrow$ 5: $\{GR(tmp[5:4]), GR(tmp[3:2]), GR(tmp[1:0])\}$ orig_addr \leftarrow tmp \ll_r s 6: $BN \leftarrow ^{orig} addr$ 7: $ADDR \leftarrow orig \ addr[5:1]$ 8: end for 9: end for 10: 11: **end for**

Goal: No conflict!

$$BN(t) \oplus BN(t + 16) = 1, \forall t$$

Network pipelining stage (n_{pp})



Algorithm 1 Conflict-free address generation scheme

Note \ll_r is a left-rotate operator. $GR = BR \circ Gray \circ BR$ where BR and Gray are a bit-reverse and gray function, respectively. **Input** *s*, *addr*, *bn* **Output** *BN*, *ADDR*

1: for s = 0; s < 6; s = s + 2 do for addr = 0; addr < 16; addr = addr + 1 do 2: for bn = 0; bn < 4; bn = bn + 1 do 3: $tmp \leftarrow \{BR(addr), BR(bn)\}$ 4: 5: $tmp \leftarrow$ $\{GR(tmp[5:4]), GR(tmp[3:2]), GR(tmp[1:0])\}$ orig_addr \leftarrow tmp \ll_r s 6: $BN \leftarrow ^{orig} addr$ 7: $ADDR \leftarrow orig_addr[5:1]$ 8: end for 9: end for 10: 11: **end for**

Goal: No conflict!

 $BN(t) \oplus BN(t + 16) = 1, \forall t$ Network pipelining stage (n_{pp})

Then, should satisfy:

 $^{addr[3:2](t)} \oplus ^{addr[3:2](t+16)}$ = $^{addr[3:2](t)} \oplus ^{(addr[3:2](t)+1)} = 1$

However, this is not satisfied for all t...

Algorithm 1 Conflict-free address generation scheme

Note \ll_r is a left-rotate operator. $GR = BR \circ Gray \circ BR$ where *BR* and *Gray* are a bit-reverse and gray function, respectively. **Input** *s*, *addr*, *bn* **Output** *BN*, *ADDR*

1: for s = 0; s < 6; s = s + 2 do for addr = 0; addr < 16; addr = addr + 1 do 2: for bn = 0; bn < 4; bn = bn + 1 do 3: $tmp \leftarrow \{BR(addr), BR(bn)\}$ 4: $tmp \leftarrow$ 5: $\{GR(tmp[5:4]), GR(tmp[3:2]), GR(tmp[1:0])\}$ orig_addr \leftarrow tmp \ll_r s 6: $BN \leftarrow ^{orig} addr$ 7: $ADDR \leftarrow orig_addr[5:1]$ 8: end for 9: end for 10: 11: **end for**

Goal: No conflict!

 $BN(t) \oplus BN(t + 16) = 1, \forall t$ Network pipelining stage (n_{pp})

 $^{addr[3:2](t)} \oplus ^{addr[3:2](t+16)}$ = $^{addr[3:2](t)} \oplus ^{(addr[3:2](t)+1)} = 1$

Applying Bit-reverse + Gray

 $^{GR}\left(BR\left(addr[3:2](t)\right)\right) \oplus ^{GR}\left(BR\left(addr[3:2](t+16)\right)\right)$ $= ^{Gray}\left(addr[3:2](t)\right) \oplus ^{Gray}\left(addr[3:2](t)+1\right) = 1$

Now, this is satisfied for all t!

Algorithm 1 Conflict-free address generation scheme

Note \ll_r is a left-rotate operator. $GR = BR \circ Gray \circ BR$ where BRand *Gray* are a bit-reverse and gray function, respectively. **Input** *s*, *addr*, *bn*

Output *BN*, *ADDR*

1:	for $s = 0$; $s < 6$; $s = s + 2$ do
2:	for $addr = 0$; $addr < 16$; $addr = addr + 1$ do
3:	for $bn = 0$; $bn < 4$; $bn = bn + 1$ do
4:	$tmp \leftarrow \{BR(addr), BR(bn)\}$
5:	$tmp \leftarrow$
	$\{GR(tmp[5:4]), GR(tmp[3:2]), GR(tmp[1:0])\}$
6:	$orig_addr \leftarrow tmp \ll_r s$
7:	$BN \leftarrow ^{\circ} orig_addr$
8:	$ADDR \leftarrow orig_addr[5:1]$
9:	end for
10:	end for
11:	end for

Goal: No conflict!

 $BN(t) \oplus BN(t+16) = 1, \forall t$

Memory access conflict most likely occur during stage transition

Example of Stage 0 \rightarrow Stage 1 in INTT s = 0, addr = 12, ..., 15 s = 2, addr = 0, ..., 3

Algorithm 1 Conflict-free address generation scheme

Memory access conflict-free

 $BN(t) \oplus BN(t+16) = 1, \forall t$

	Example of Stage $0 \rightarrow $ Stage 1 in INTT									
Note \ll_r is a left-rotate operator. $GR = BR \circ Gray \circ BR$ where BF_1		ipie		מַמ						
and <i>Gray</i> are a bit-reverse and gray function, respectively.	Time	ad	dr	bn	BN	Time	dr	bn	ad	BN
Innuts addr hn	$t_{stg,0} - 16$	00	01	00	1	$t_{stg,0}$	00	00	00	0
$\mathbf{O}_{\mathbf{v}} = \mathbf{P} \mathbf{V} + \mathbf{P} \mathbf{V}$	$t_{stg,0} - 15$	00	01	10	0	$t_{stg,0} + 1$	00	10	00	1
Output BN, ADDR	$t_{stg,0} - 14$	00	01	11	1	$t_{stg,0}+2$	00	11	00	0
1: for $s = 0$; $s < 6$; $s = s + 2$ do	$t_{stg,0} - 13$	00	01	01	0	$t_{stg,0} + 3$	00	01	00	1
2: for $addr = 0$; $addr < 16$; $addr = addr + 1$ do	$t_{stg,0} - 12$	10	01	00	0	$t_{stg,0} + 4$	00	00	10	1
3: for $bn = 0$: $bn < 4$: $bn = bn + 1$ do	$t_{stg,0} - 11$	10	01	10	1	$t_{stg,0} + 5$	00	10	10	0
$4: \qquad tmp \leftarrow \{BR(addr), BR(bn)\} \in \mathbb{R}$	$t_{stg,0} - 10$	10	01	11	0	$t_{stg,0} + 6$	00	11	10	1
$\{ad, dr, bn\}$	$t_{stg,0} - 9$	10	01	01	1	$t_{stg,0} + 7$	00	01	10	0
$5: \qquad \lim p \leftarrow \qquad \qquad$	$t_{stg,0} - 8$	11	01	00	1	$t_{stg,0} + 8$	00	00	11	0
$\{GR(tmp[5:4]), GR(tmp[3:2]), GR(tmp[1:0])\}$	$t_{stg,0} - 7$	11	01	10	0	$t_{stg,0} + 9$	00	10	11	1
6: $orig_addr \leftarrow tmp \ll_r s$	$t_{stg,0} - 6$	11	01	11	1	$t_{stg,0} + 10$	00	11	11	0
7: $BN \leftarrow ^{orig} addr$ Same throughput	$t_{stg,0} - 5$	11	01	01	0	$t_{stg,0} + 11$	00	01	11	1
8: $ADDR \leftarrow orig addr[5:1]$ (Power)	$t_{stg,0} - 4$	01	01	00	0	$t_{stg,0} + 12$	00	00	01	1
end for	$t_{stg,0} - 3$	01	01	10	1	$t_{stg,0} + 13$	00	10	01	0
$67.87 (44.39) \% \downarrow$	$t_{stg,0} - 2$	01	01	11	0	$t_{stg,0} + 14$	00	11	01	1
	$t_{stg,0} - 1$	01	01	01	1	$t_{stg,0} + 15$	00	01	01	0
11: end for		Stage	0				Stage	1		

s = 0, addr = 12, ..., 15 s = 2, addr = 0, ..., 3

III. Bandwidth-efficient Behavior in Core



Frequent data access pattern transition between NTT and MMAC unit

→ Expensive external memory access!

III. Bandwidth-efficient Behavior in Core



Phase	Operation	Data Transfer					
1	IN0	NTTs \leftrightarrow Mem					
	IN1	$Mem \rightarrow NTTs \rightarrow Tmp$					
2	Conv	$\mathbf{MMAC} \leftrightarrow \mathbf{Tmp}$					
	N0	$\underline{Tmp} \to NTTs \to Mem$					
3	N1	$Mem \rightarrow NTTs \rightarrow Tmp$					
5	MAC	Tmp, Mem \rightarrow MMAC \rightarrow Mem					
Δ	N1	$Mem \rightarrow NTTs \rightarrow Tmp$					
+	SM	$\frac{\text{Tmp}}{\text{Mem}}, \text{Mem} \rightarrow \text{MMAC} \rightarrow \text{Mem}$					

Dataflow in KS modified for better data utilization, using dedicated buffers (Tmp) between NTT and MMAC units

III. Bandwidth-efficient Behavior in Core

		Na	ïve	Ou	urs
Part	Usage	In	In Out		Out
	INTT	$N \cdot 2 \cdot (L+1)$	$N \cdot 2 \cdot (L+1)$	$N \cdot 2 \cdot (L+1)$	$N \cdot (L+1)$
1	Conv	$N\cdot lpha \cdot dnum$	$N \cdot dnum$	0	0
	NTT	$N \cdot 2 \cdot dnum$	$N \cdot 2 \cdot (L + 1) \cdot dnum$	$N \cdot dnum$	$N \cdot dnum \cdot (L+1)$
	(NTT)	$egin{array}{c} N \cdot dnum \cdot \ (L+1+lpha) \end{array}$	-	0	-
2	Key	$2N \cdot dnum \cdot (L+1+lpha)$	-	$2N \cdot dnum \cdot (L+1+lpha)$	-
	Cxt	0	$2N \cdot (L + 1 + lpha)$	0	2N(L+1+lpha)
	INTT	$2N\cdot 2\cdot lpha$	$2N\cdot 2\cdot lpha$	$2N\cdot 2\cdot lpha$	$2N\cdot lpha$
3	Conv	$2N\cdot lpha$	2N	0	0
	NTT	$2N \cdot 2$	$\frac{2N\cdot 2\cdot}{(L+1)}$	2N	2N(L+1)
4	(NTT)	2N(L+1)	-	0	-
4	Cxt	2N(L+1)	2N(L+1)	2N(L+1)	2N(L+1)

Phase	Operation	Data Transfer					
1	IN0	$NTTs \leftrightarrow Mem$					
	IN1	$Mem \rightarrow NTTs \rightarrow Tmp$					
2	Conv	$MMAC\leftrightarrowTmp$					
	N0	$\frac{\text{Tmp}}{\text{Tmp}} \rightarrow \text{NTTs} \rightarrow \text{Mem}$					
3	N1	$Mem \rightarrow NTTs \rightarrow Tmp$					
5	MAC	Tmp, Mem \rightarrow MMAC \rightarrow Mem					
	N1	$Mem \rightarrow NTTs \rightarrow Tmp$					
-	SM	Tmp, Mem \rightarrow MMAC \rightarrow Mem					

Dataflow in KS modified for better data utilization, using dedicated buffers (Tmp) between NTT and MMAC units → External memory access 38.7 % ↓

Chip Implementation



Area (mm ²)	23.48 (5.70×4.12)
Memory voltage (V)	0.9
Memory size (kB)	240
Core voltage (V)	0.8
Core frequency (MHz)	200
Avg. power (mW)	54.7
Throughput (KS/s)	10.58



GD (LD): Global/Local Distributor

Comparison with Prior Works

	HOST'19 [22]	CICC'18 [23]	TVLSI'22 [24]	ISSCC'19 [25]	Ours
Technology	65nm	40nm	65nm	40nm	28nm
$(N, \lceil \log q \rceil)$	(512, 14)	(512, 14)	(512, 14)	(256, 13)	(512, 60)
Energy efficiency (Gbit/W/s)	30.308	106.009	51.358	74.989	431.577
Improvement	14.24×	4.07 ×	8.40 ×	5.755×	
Area efficiency (bit/um ²)	635.846	6371.329	4126.984	1237.580	61033.820
Improvement	95.99×	9.58 ×	14.79 ×	49.32 ×	

Table 3: Comparisons with Prior ASIC NTT Accelerators.

Table 4: Comparison of Energy Efficiency for KS Operation.

	Platform	$\log N$	E_{norm} (J)	Improvement
HEAX [9]	FPGA	14	0.0082	12.66 ×
BASALIC [6]	ASIC	16	0.0784	30.32 ×
100x [3]	GPU	17	0.3129	60.52 ×
Lattigo [16]	CPU	17	67.7411	13104.87 ×
Ours	ASIC	17	0.0052	

Conclusion

- Designed a dedicated accelerator for KS that requires frequent transitions in data access patterns incurring redundant expensive external memory accesses
- Proposed design techniques on various levels for high energy efficiency modular multiplier, NTT unit, and data access behavior in core, thus full-stack optimization
- As a result, the design shows significant improvement in performance in energy efficiency compared with prior FHE implementations.
- Although designed specifically for KS, it remains highly applicable since KS operations dominate power, time, and bandwidth across the entire computation.
- Techniques can be also applied to other parameter sets modular multiplier as long as same moduli set used and $L + \alpha + 1 \le 41$, NTT unit down to $N = 2^{16}$ (the least amount to support bootstrapping reported in literature)