



WPI

Accelerating LTV Based Homomorphic Encryption in Reconfigurable Hardware

Yarkın Doröz, Erdinç Öztürk, Erkey Savaş, Berk Sunar

Worcester Polytechnic Institute
100 Institute Road, Worcester, MA, USA, 01609



VERNAM Group

Security & Privacy @ WPI

Homomorphic Encryption

- What is it: Permits computations on encrypted data
- Partially homomorphic encryption
 - Permits evaluation of **restricted** circuit on encrypted inputs
 - E.g. Goldwasser Micali (XOR), Paillier (Int. Adds), BGN (2-DNF)
- Fully homomorphic encryption (Gentry 2009)
 - Allows **efficient** evaluation of **arbitrary** circuits
 - Gentry's blueprint:
 - Builds on *somewhat homomorphic encryption scheme* (SWHE):
 - Supports homomorphic additions and only few levels of multiplication
 - Add Bootstrapping SWHE -> FHE
- Many new constructions: lattice and integer based
 - Orders of magnitude improvement each year since 2010!
- Efficiency still remains as biggest problem!

Overview of FHE Implementations

Table 1. Overview of specialized FHE Implementations. GH-FHE: Gentry & Halevi's FHE scheme; CMNT-FHE: Coron et al.'s FHE schemes [10, 11] [22]; NTRU based FHE, e.g. [27, 34]

DESIGN	SCHEME	PLATFORM	PERFORMANCE
CPU			
AES [25]	BGV-FHE	2.0 GHz Intel Xeon	5 min / AES block
AES [16]	NTRU-FHE	2.9 GHz Intel Xeon	55 sec / AES block
Full FHE [31]	NTRU-FHE	2.1 GHz Intel Xeon	275 sec / per bootst.
GPU			
NTT mul / reduction [35]	GH-FHE	NVIDIA C250 GPU	0.765 ms
NTT mul [35]	GH-FHE	NVIDIA GTX 690	0.583 ms
AES [14]	NTRU-FHE	NVIDIA GTX 690	7 sec / AES block
FPGA			
NTT transform [37]	GH-FHE	Stratix V FPGA	0.125 ms
NTT modmul / enc. [7]	CMNT-FHE	Xilinx Virtex7 FPGA	13 msec / enc.
ASIC			
NTT modmul [17]	GH-FHE	90nm TSMC	2.09 sec
Full FHE [18]	GH-FHE	90nm TSMC	3.1 sec / decrypt

LTV

- Proposed by Lopez-Alt, Tromer and Vaikuntanathan
- Variant of NTRU Encryption by Stehlé and Seinfeld
- Designed as a multi-user, **leveled** FHE scheme

The construction:

- Operations are performed in $R_q = \mathbb{Z}_q[x]/\langle x^n + 1 \rangle$ where q is the prime coefficient modulus
- Sequence of primes $q_0 > q_1 > \dots > q_d$, a different q_i for each level
- Error distribution X , a truncated discrete Gaussian distribution, for sampling random polynomials

LTV - Key Generation

Private Key:

$$f^{(i)}(x) = 2u^{(i)} + 1$$

Public Key:

$$h^{(i)}(x) = 2g^{(i)}(f^{(i)})^{-1}$$

Evaluation Keys:

$$\zeta_{\tau}^{(i)}(x) = h^{(i)}s^{(i)} + 2e_{\tau}^{(i)} + 2^{\tau}f^{(i-1)}$$

For each $i = 0, 1, \dots, d$

- $\{h^{(i)}, f^{(i)}, \zeta_{\tau}^{(i)}\} \in R_q$
- $\{g^{(i)}, u^{(i)}, s_{\tau}^{(i)}, e_{\tau}^{(i)}\} \in X$
- $\tau = [0, \lfloor \log q_i \rfloor]$

LTV – Encrypt/Decrypt

Encryption:

$$b \in \{0,1\}, \quad \{s, e\} \in X$$

$$c^{(i)} = h^{(i)}s + 2e + b$$

Decryption:

$$\mu = \lfloor c^{(i)} f^{(i)} \rfloor_{q_i}$$

$$m = \mu \pmod{2}$$

LTV - Evaluation

$$c_1^{(i)} = \text{Encrypt}(b_1) \quad c_2^{(i)} = \text{Encrypt}(b_2)$$

$$\text{Decrypt}(c_1^{(i)} + c_2^{(i)}) = b_1 + b_2$$

$$\text{Decrypt}(c_1^{(i)} \cdot c_2^{(i)}) = b_1 \cdot b_2$$

$$\begin{aligned} & (h^{(i)}s_1 + 2e_1 + b_1) + (h^{(i)}s_2 + 2e_2 + b_2) \\ & (h^{(i)}s_1 + 2e_1 + b_1) \cdot (h^{(i)}s_2 + 2e_2 + b_2) \end{aligned}$$

$$\begin{aligned} & h^{(i)}\check{s} + 2\check{e} + (b_1 + b_2) \\ & h^{(i)^2}\check{s} + 2\check{E} + (b_1 \cdot b_2) \end{aligned}$$

LTV - Evaluation

Relinearization:

$$\check{c}^{(i-1)}(x) = \sum_{\tau} 2^{\tau} \check{c}_{\tau}^{(i-1)}(x)$$
$$\check{c}^{(i)}(x) = \sum_{\tau} \zeta_{\tau}^{(i)}(x) \check{c}_{\tau}^{(i-1)}(x)$$

Modulus Switch:

- Noise Reduction

$$\check{c}^{(i)}(x) = \left\lfloor \frac{q_i}{q_{i-1}} \check{c}^{(i)}(x) \right\rfloor_2$$

Specialization of LTV

- Operations are performed in $R_{q_i} = \mathbb{Z}_{q_i}[x]/\langle \varphi_m(x) \rangle$
 - $\varphi_m(x) = m^{th}$ cyclotomic polynomial
 - Degree of $\varphi_m(x)$, n is $\phi(m)$
 - $\varphi_m(x)$ is factorized over F_2 into equal degree polynomials $F_i(x)$
 - Using CRT, we batch $S = n/t$ messages, t is the smallest integer that satisfies $m|(2^t - 1)$
- Select number of primes p_i with size of $\log p$ – noise cutting size
 - $q_i = \prod_{i \leq j} p_j$
- Evaluation keys are promoted to the next level via
$$\zeta_\tau^{(i)}(x) = \zeta_\tau^{(0)}(x) \pmod{q_i}$$

Arithmetic Operations

- Large polynomial multiplications and relinearization operations are costly
 - convert input polynomials using CRT
 - Same degree, smaller word-size coefficients
- Pairwise polynomial product
 - computed using Number Theoretical Transform (NTT) multiplication
- Resulting polynomial is recovered from the partial products using inverse CRT
 - postpone ICRT until switching a level

Polynomial Multiplication

- The classical multiplication techniques are too costly:
 - Schoolbook algorithm have quadratic complexity $O(N^2)$
 - Karatsuba algorithm reduce the complexity to $O(N^{\log_2 3})$
- Very large parameters/operands:
 - Need asymptotically better algorithms!
- We use Schönhage–Strassen algorithm to compute multiplications:
 - Number theoretic transform (NTT) has $O(N \log N \log \log N)$ complexity
 - Fast Fourier Transform is needed
 - We use Cooley-Tukey decomposition

Polynomial Multiplication

$$A(x) = \sum_{i=0}^{2N-1} A_i x^i \quad A_i \in \mathbb{Z}_p$$

$$\check{A}_i = \sum_{j=0}^{2N-1} A_j \cdot w^{ij} \pmod{p} \quad w^{2N} = 1 \pmod{p}$$

$$\check{C}_i = \check{A}_i \cdot \check{E}_i \pmod{p} \text{ for all } i \in 2N$$

$$C_i = (2N)^{-1} \cdot \sum_{j=0}^{2N-1} \check{C}_j \cdot w^{-ij} \pmod{p}$$

Polynomial Multiplication

- Cooley-Tukey

- SW: Recursive

- $\check{A}_i = \sum_{j=0}^{2N-1} A_j \cdot w^{ij} \pmod{p}$
 $= \sum_{j=0}^{N-1} A_{2j} \cdot w^{i(2j)}$
 $+ w^i \sum_{j=0}^{N-1} A_{2j+1} \cdot w^{i(2j)} \pmod{p}$

- HW: Non-Recursive

ALGORITHM 1: Iterative Version of Number Theoretic Transformation

input : $A(x) = A_0 + A_1x + \dots + A_{N-1}x^{N-1}$, $N = 2^n$, and w
output: $\mathcal{A}(x) = \mathcal{A}_0 + \mathcal{A}_1x + \dots + \mathcal{A}_{N-1}x^{N-1}$

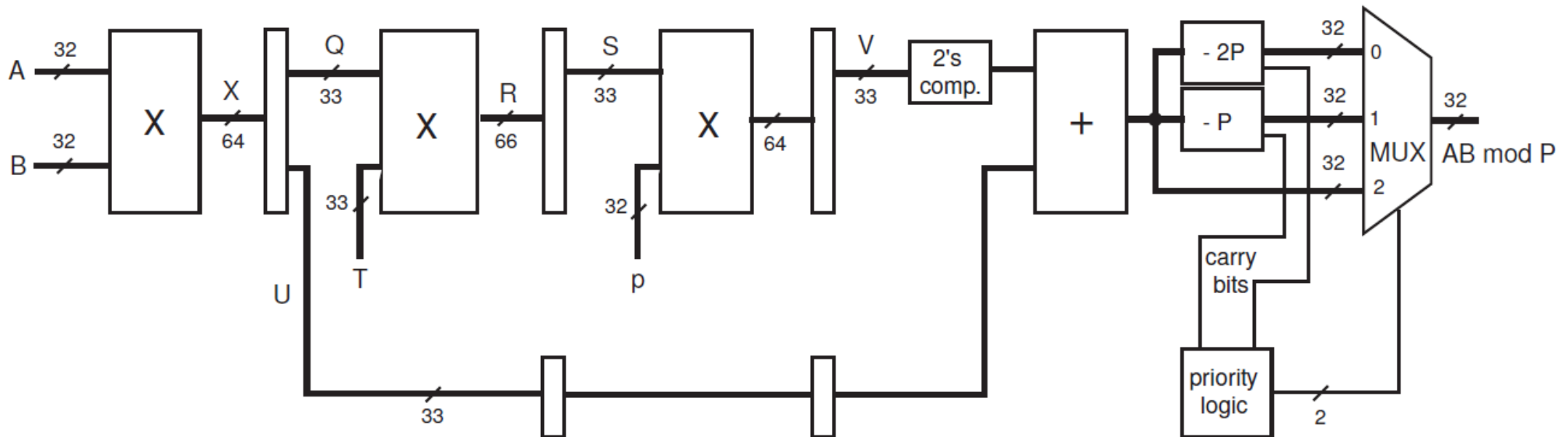
```
1 for  $i = N$  to  $2N - 1$  do
  |  $A_i = 0$ ;
  end
2  $(\mathcal{A}_0, \mathcal{A}_1, \dots, \mathcal{A}_{2N-1}) \leftarrow \text{Permutation}(A_0, A_1, \dots, A_{2N-1})$ ;
3 for  $M = 2$  to  $2N$  do
4   for  $j = 0$  to  $2N - 1$  do
5     for  $i = 0$  to  $\frac{M}{2} - 1$  do
6        $x \leftarrow i \times \frac{2N}{M}$ ;
7        $\mathcal{I} \leftarrow j + i$ ;
8        $\mathcal{J} \leftarrow j + i + \frac{M}{2}$ ;
9        $\mathcal{A}[\mathcal{I}] \leftarrow \mathcal{A}[\mathcal{I}] + w^{x \bmod 2N} \times \mathcal{A}[\mathcal{J}] \pmod{p}$ ;
10       $\mathcal{A}[\mathcal{J}] \leftarrow \mathcal{A}[\mathcal{J}] - w^{x \bmod 2N} \times \mathcal{A}[\mathcal{I}] \pmod{p}$ ;
         $i \leftarrow i + 1$ ;
      end
       $j \leftarrow j + M$ ;
    end
  end
   $M \leftarrow M \times 2$ ;
end
```

Architecture – 32-bit Modular Multiplier

- The homomorphic AES algorithm requires 31-bit cutting size
 - small CRT primes need to be 31-bits.
- We designed our hardware for 32-bit multiplications.
 - We require four 16x16-bit multiplications
 - 1 DSP unit can compute in 4 clock cycles
- We used Barrett's algorithm for modular reduction
 - 33x33-bit multiplications
 - 6 clock cycle latency

Architecture – 32-bit Modular Multiplier

- Using 3 multiplier and pipelining we compute modular multiplication in 19 clock cycles
 - Throughput: 4 clock cycles



Architecture - $2^{15} \times 2^{15}$ Polynomial Multiplier

- K modular multipliers
 - $3 \times K$ DSP units
 - power of two
- K Block RAMs
 - To feed multipliers with digits
 - Store the results
 - Block RAM indices generation
 - degree size (M), NTT level (m) and modular multiplier number (K)
 - We make sure to access in 4 clock cycles:
 - powers of primitive root w
 - polynomial coefficients

Architecture - 2¹⁵x2¹⁵ Polynomial Multiplier

ALGORITHM 4: Parallel access to polynomial coefficients

```

input  :  $A(x) = A_0 + A_1x + \dots + A_{2N-1}x^{2N-1}$ ,  $n$ ,  $m$ , and  $\kappa < n$ 
output:  $B_i[j]$ 

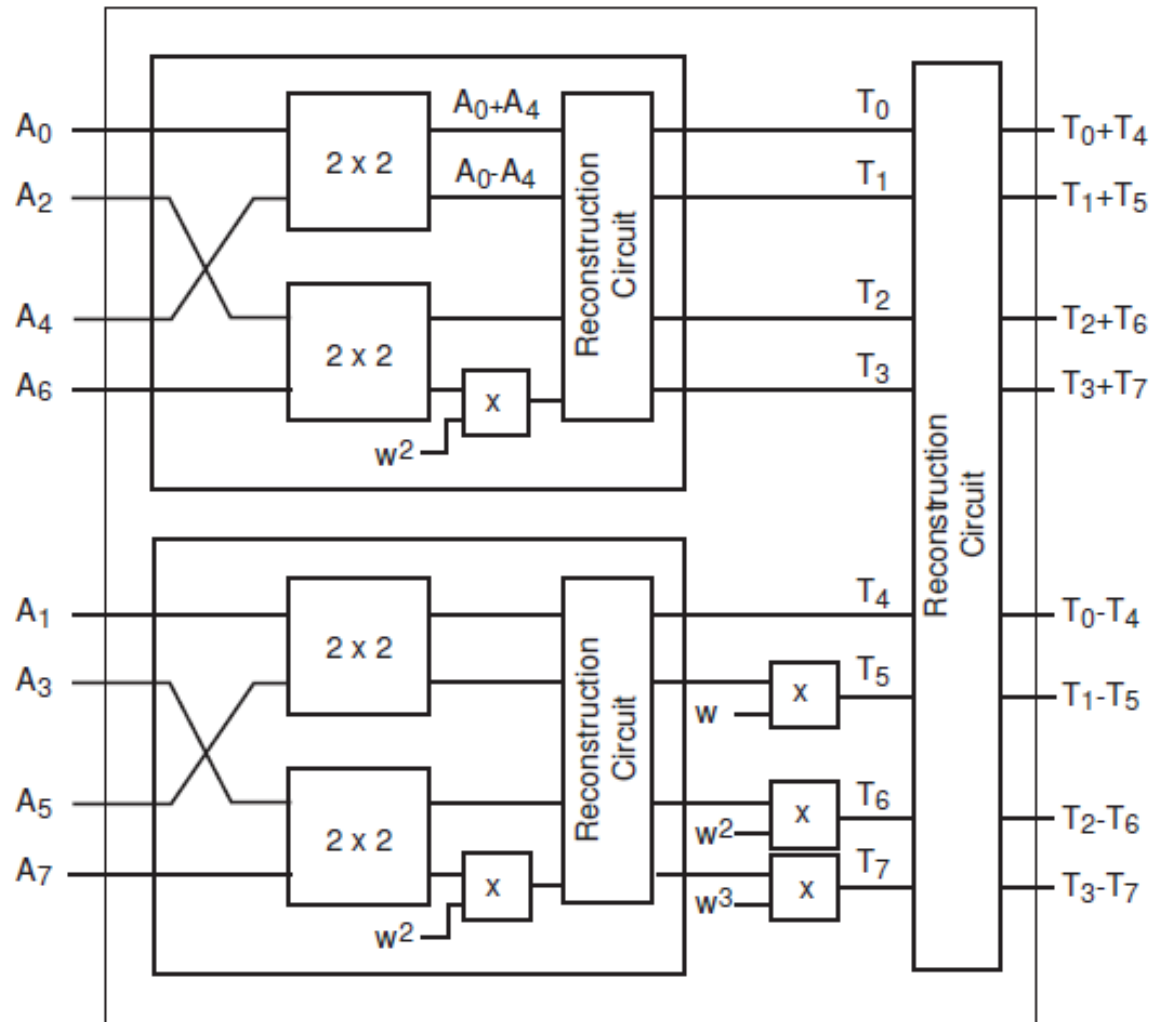
1  $mCnt \leftarrow 2^{m-1} - 1$ ;          /* number of multiplications in a block */
2  $bSize \leftarrow 2^m$ ;             /* size of a block */
3  $BRAMCnt \leftarrow 2^{\kappa-2}$ ;       /* number of BRAMs */
4 if  $bSize \leq 2^{\kappa-2}$  then
    for  $t = 0$  to  $1024$  do
        for  $i = 0$  to  $BRAMCnt$  do in parallel
            for  $j = i + bSize - mCnt$  to  $i + bSize$  do
                for  $k = 0$  to  $3$  do
5                 Access  $BRAM_j[t + 2k]$ ;
6                 Access  $BRAM_j[t + 2k + 1]$ ;
                     $k \leftarrow k + 1$ ;
                end
                 $j \leftarrow j + 1$ ;
            end
             $i \leftarrow i + bSize$ ;
        end
         $t \leftarrow t + 8$ ;
    end
end
7 else
    for  $i = 0$  to  $BRAMCnt$  do in parallel
        for  $j = 0$  to  $1024$  do
            for  $k = 2^{m-\kappa+1}$  to  $2^{m-\kappa+2}$  do
8                 Access  $BRAM_i[k + j]$ ;
                     $k \leftarrow k + 1$ ;
            end
             $j \leftarrow j + 2^{m-\kappa+2}$ ;
        end
         $i \leftarrow i + 1$ ;
    end
end
end

```

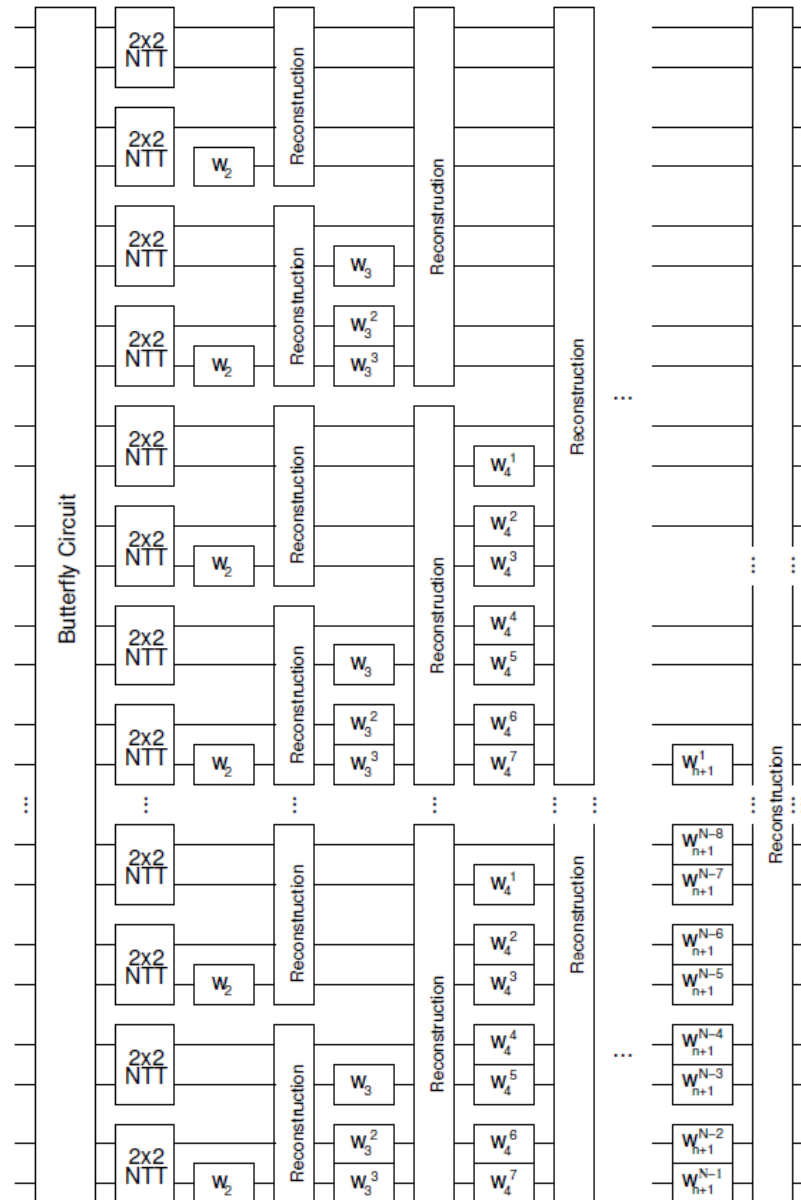
- Recall Our Recursive Algorithm:

$$\begin{aligned}
 \check{A}_i &= \sum_{j=0}^{2N-1} A_j \cdot w^{ij} \pmod{p} \\
 &= \sum_{j=0}^{N-1} A_{2j} \cdot w^{i(2j)} \\
 &\quad + w^i \sum_{j=0}^{N-1} A_{2j+1} \cdot w^{i(2j)} \pmod{p}
 \end{aligned}$$

Architecture – 2-point to 8-point NTT



Architecture – 64K-point NTT



Implementation Results

- We used Virtex-7 XC7VX690T FPGA
- The design reach 250 MHz clock
- The polynomial degree $n = 32,768$ and $\log p = 32$
- We choose the number of modular multiplier as $K=64$

	Total	Used	Used (%)
Slice LUTs	433,200	219,192	50.59
Slice Registers	866,400	90,789	10.47
RAMB36E1	1470	193	13.12
DSP48E1	3600	768	21.33

Timing for Primitive Operations

- We include PCIe timing (8 Gbit/sec per lane)
- CRT is done on the CPU

	Timings (msec)		Timings (msec)
CRT	89	Modulus Switch	89
Multiplication	9.51	Relinearization	526
NTT conversions	6.25	CRT conversions	89
PCIe cost	3.26	NTT conversions	331
Modular Reduction	19	PCIe cost	52

Comparison

- AES is a 40 level circuit
 - 2,880 Relinearizations
 - 5,760 Modular Multiplications
 - 6,080 Modular Switching
- Total AES takes 15 minutes (2048 message slots)
 - Amortized time is 439 msec

	Mul Speedup (msec)		Relin Speedup (sec)		AES Speedup (sec)	
CPU [16]	970	1×	103	1×	55	1×
GPU [14]	340	2.8×	8.97	11.5×	7.3	7.5×
CPU (16-bit)	970	1×	6.5	16×	12.6	4.4×
FPGA (ours)	9.5	102×	0.53	195×	0.44	125×

Thank you!

Questions?