High-precision RNS-CKKS on small word-size architecture

Duhyeong Kim, Intel Labs FHE.org Meetup Jan 11th, 2024

Overview

- Enable high-precision RNS-CKKS on fixed but smaller word-size architectures
 - Single scaling → Composite scaling
- Enable functionally correct CKKS composite scaling in two open-source libraries
 - OpenFHE: C++, enabled by Intel labs
 - Lattigo: Go, enabled by Seoul National University (SNU)
- Demonstrate with secure parameters the **equivalence** between single and composite scaling
 - 7-layer CNN Inference with longitudinal packing in OpenFHE-CKKS with composite scaling
 - 7-layer CNN Inference with multiplexed packing in Lattigo-CKKS with composite scaling
 - Logistic Regression Training in OpenFHE-CKKS with composite scaling

Fully Homomorphic Encryption (FHE)

Any computation on encrypted data "without decryption process"



CKKS: FHE for real-number arithmetic

How can we think of the "approximate" computation in CKKS?

- Imitation of "fixed-point" arithmetic in cleartext version
- Example: computation of $1.584 \times 2.4835 \times 9.5937 \times 8.7264 \times 6.12743$ (≈ 2017.9897)



CKKS: FHE for real-number arithmetic

How can we think of the "approximate" computation in CKKS?

- Imitation of "fixed-point" arithmetic in cleartext version
- Example: computation of $1.584 \times 2.4835 \times 9.5937 \times 8.7264 \times 6.12743$ (≈ 2017.9897)



Scaling Factor in CKKS

- Determine the "initial precision bits" under the decimal point
- CKKS Encoding/Encryption results in



- Larger Δ , start with higher precision
- Smaller Δ , start with lower precision

Scaling Factor in CKKS

• Exponential growth of Scaling Factor

$$\succ (\Delta \cdot m) \cdot (\Delta \cdot m') = \Delta^2 \cdot mm'$$

$$\succ (\Delta^{2^k} \cdot m) \cdot (\Delta^{2^k} \cdot m') = \Delta^{2^{k+1}} \cdot mm'$$

• How to control the growth of scaling factor?

"rescale"

• Rescale(*ct*): *ct* mod $\Delta^{\ell} \mapsto \left[\frac{ct}{\Delta}\right] \mod \Delta^{\ell-1}$ (from the context of "original" CKKS) $\succ (\Delta \cdot m) \cdot (\Delta \cdot m') = \Delta^2 \cdot mm' \xrightarrow{\text{Rescale } (1/\Delta)} \Delta \cdot mm'$

• RNS-CKKS

> An efficient way to implement CKKS w/o big-number arithmetic

 \succ Ctxt moduli $Q_{\ell} = q_0 q_1 \cdots q_{\ell}$ for level ℓ (instead of modulo Δ^{ℓ})

 $\operatorname{RNS}_{Q_{\ell}}(x) \coloneqq (x \mod q_0, x \mod q_1, \dots, x \mod q_{\ell})$

• Rescale modulo Q_ℓ in RNS?

> No efficient way to compute $\mathbf{x} \mapsto \begin{bmatrix} \frac{1}{\Delta} \cdot \mathbf{x} \end{bmatrix}$

> Instead, we can efficiently compute $\mathbf{x} \mapsto \left| \frac{1}{a_{\ell}} \cdot \mathbf{x} \right|$

$$\circ \left\lfloor \frac{1}{q_{\ell}} \cdot x \right\rfloor = q_{\ell}^{-1} \cdot (x - x \mod q_{\ell})$$

 $\circ\,$ Easy to obtain the RNS representation of $x \mod q_\ell$

• $\operatorname{RNS}_{Q_{\ell-1}}(x \mod q_{\ell}) = (x \mod q_{\ell}, x \mod q_{\ell}, \dots, x \mod q_{\ell})$

• Case 1: $\log \Delta < \text{word-size}$

 \succ Set each prime q_{ℓ} to be $\log \Delta$ bits

Perform the "single scaling"

$$\mathbf{x} \mod Q_{\ell} \mapsto \left[\frac{\mathbf{1}}{\boldsymbol{q}_{\ell}} \cdot \boldsymbol{x} \right] \mod Q_{\ell-1}$$

- Set each product of q_ℓ 's to be $\log \Delta$ bits
- Perform the "composite scaling"

• Case 1: $\log \Delta < \text{word-size}$

 \succ Set each prime q_{ℓ} to be $\log \Delta$ bits

Perform the "single scaling"

$$\mathbf{x} \mod Q_{\ell} \mapsto \left[\frac{\mathbf{1}}{\boldsymbol{q}_{\ell}} \cdot \boldsymbol{x} \right] \mod Q_{\ell-1}$$

- Set each product of q_ℓ 's to be $\log \Delta$ bits
- Perform the "composite scaling" (degree = 2)

$$\mathbf{x} \mod Q_{\ell} \mapsto \left[\frac{\mathbf{1}}{\boldsymbol{q}_{\ell} \boldsymbol{q}_{\ell-1}} \cdot \boldsymbol{x} \right] \mod Q_{\ell-2}$$

• Case 1: $\log \Delta < \text{word-size}$

 \succ Set each prime q_{ℓ} to be $\log \Delta$ bits

Perform the "single scaling"

$$\mathbf{x} \mod Q_{\ell} \mapsto \left[\frac{\mathbf{1}}{\boldsymbol{q}_{\ell}} \cdot \boldsymbol{x} \right] \mod Q_{\ell-1}$$

- Set each product of q_ℓ 's to be $\log \Delta$ bits
- Perform the "composite scaling" (degree = 3)

$$\mathbf{x} \mod Q_{\ell} \mapsto \left[\frac{1}{\boldsymbol{q}_{\ell} \boldsymbol{q}_{\ell-1} \boldsymbol{q}_{\ell-2}} \cdot \boldsymbol{x} \right] \mod Q_{\ell-3}$$

• Case 1: $\log \Delta < \text{word-size}$

 \succ Set each prime q_{ℓ} to be $\log \Delta$ bits

Perform the "single scaling"

$$\mathbf{x} \mod Q_{\ell} \mapsto \left[\frac{\mathbf{1}}{\boldsymbol{q}_{\ell}} \cdot \boldsymbol{x} \right] \mod Q_{\ell-1}$$

- Set each product of q_ℓ 's to be $\log \Delta$ bits
- Perform the "composite scaling" (degree = t)

$$\mathbf{x} \mod Q_{\ell} \mapsto \left[\frac{1}{\boldsymbol{q}_{\ell} \cdots \boldsymbol{q}_{\ell-t+1}} \cdot \boldsymbol{x} \right] \mod Q_{\ell-t}$$

• Examples

(double-prime) (double-prime) (triple-prime)

Precision Issue due to Rescale

- Original CKKS: NO precision issue
 - \succ Scaling factor is **ALWAYS** preserved as Δ
- RNS-CKKS: **YES** precision issue
 - Scaling factor is **NOT** be preserved as Δ \circ Division by q_{ℓ} 's, instead of Δ $\circ \Delta^2/q_{\ell} \neq \Delta$
 - Critical Impact to Homomorphic Addition
 Enc($\Delta \cdot m$) + Enc($\Delta' \cdot m'$) = Enc($\Delta \cdot (m + \Delta' / \Delta \cdot m')$) \neq Enc($\Delta \cdot (m + m')$)
 The ratio Δ' / Δ (≠ 1) directly harms the precision

Precision Issue due to Rescale

• Solution 1: Choose the primes properly

 \succ To keep the scaling factors (not equal but) very close to Δ

Single Scaling

 \circ Requirement: $q_\ell \simeq \Delta$ (proposed in original RNS-CKKS)

$$\circ \ \Delta^2/q_\ell \ \simeq \Delta$$

- Composite Scaling
 - Requirement: $q_ℓ q_{ℓ-1} ≃ Δ$

$$\circ \Delta^2/q_\ell q_{\ell-1} \simeq \Delta$$

- Precision (Single Scaling v.s. Composite Scaling)
 - **NO Difference** in Mult + Relin + Rescale
 - **Closeness** of q_ℓ (resp. $q_\ell q_{\ell-1}$) and ∆ affects the **Add Precision**

Precision Issue due to Rescale

Solution 2: Exact Scaling

- Differences v.s. Solution 1
 - \circ Scaling factor Δ_i for each level i
 - $\circ \Delta_i$'s are **NOT** required to be **very close to** Δ
 - \circ Adjust the ciphertext scaling factors to Δ_i before Add and Mult
 - As a result, we "always" add two ciphertexts with "same" scaling factors
- > **Precision** (Single Scaling v.s. Composite Scaling)
 - O NO Difference in Mult + Relin + Rescale
 - NO Difference in Add
- > We implemented 32-bit RNS-CKKS in OpenFHE and Lattigo with **Solution 2**
 - "FLEXIBLEAUTO" mode in OpenFHE
 - Bootstrapping enabled in both libraries

Theoretical Analysis on Precision

Rescale
$$(ct)$$
: $ct \mod Q_i \mapsto \left\lfloor \frac{1}{q_i} \cdot ct \right\rfloor \mod Q_{i-1}$ (single scaling)
Rescale $^{(t)}(ct)$: $ct \mod Q_i \mapsto \left\lfloor \frac{1}{q_i q_{i-1} \cdots q_{i-t+1}} \cdot ct \right\rfloor \mod Q_{i-t}$ (composite scaling)

Theorem. Let B_{rs} , $B_{comp-rs}$ be the upper bounds of the error induced by $\text{Rescale}(\cdot)$ and $\text{Rescale}^{(t)}(\cdot)$, respectively. Then, it holds that

$$B_{comp-rs} \leq \left(\frac{1}{q_i q_{i-1} \cdots q_{i-t+1}} + \frac{1}{q_i q_{i-1} \cdots q_{i-t+2}} + \cdots + \frac{1}{q_i} + 1\right) B_{rs} \approx \left(\frac{1}{q_i} + 1\right) B_{rs}$$
Hence, composite scaling results in less than $\log\left(\frac{1}{q_i} + 1\right) \approx \frac{3.322}{q_i}$ bit **precision loss**, which is **negligible**, compared to single scaling

compared to single scaling.

Experimental Results

7-layer CNN Inference (CIFAR-10)

- Implementation in OpenFHE with longitudinal packing
 - Unit tests with Same Precision

Unit tests	Precision bits 64-bit single scaling	Precision bits 32-bit composite scaling	======== Parameters ====================================
Fully connected	39	39	 Same for both cases Primes 58-bit primes for 64-bit case (29, 29)-bit primes for 32-bit case > Double-prime scaling > 58 = 29 + 29
ReLU	40	40	
Mean pool	41	41	
Convolution	39	39	
Bootstrapping	12	12	Security
Bootstrapping	12	12	• Same for both cases

- Implementation in Lattigo with multiplexed packing
 - > The end-to-end CNN Inference results match up to 5 digits after the decimal point
 - > 14 consecutive bootstrapping (2 per layer, before and after ReLU)

Experimental Results

Logistic Regression Training

- Reference code: <u>https://github.com/openfheorg/openfhe-logreg-training-examples</u>
- 1 bootstrapping per epoch





Wrap-up

- **Result:** Enable high-precision RNS-CKKS on small word-size architectures without multi-precision arithmetic
 - Use of small word-size: GPU, FPGA, Embedded devices, etc.
 - Arbitrary precision for bootstrapping combined with Meta-BTS
- Limitation: Choice of scaling factor
 - Lower bound exists on each prime (NTT condition)
 - E.g., $\Delta = 2^{40} \rightarrow$ two 20-bit primes for double-prime scaling
 - How many 20-bit "NTT-friendly" primes exist for the dimension $N = 2^{16}$?
 - Several small intervals that are not usable as scaling factor
- Implementation: Not public yet but planning for open-sourcing compositescaling variant of OpenFHE-CKKS



https://eprint.iacr.org/2023/1462