

Revisiting Key Decomposition Techniques for FHE: Simpler, Faster and More Generic

M. G. Belorgey, S. Carpov, **N. Gama**, S. Guasch, D. Jetchev



SANDBOXAQ



inpher

FHE.org Conference - 24 Mars 2024

① **External products:**

The missing public-secret product.

② **Auxiliary Gadgets and Bivariate representation**

Gadget decomposition of the source of truth

The return of base 2^K representation in large depth.

Part I - External products: Two halves make a whole

External products are omnipresent in FHE

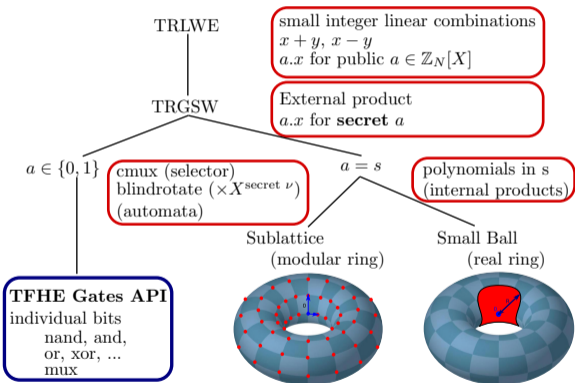
Ring-LWE has never really been about Rings

- Ring-LWE FHE was never really about rings.
- FHE arithmetic flows from **secret linear combinations**, a.k.a. **external products**:
 - with small integer polynomial coefficients
 - over some high precision space (Torus, or big integers)

$$\sum \underbrace{a_i}_{\text{small int}} \cdot \underbrace{c_i}_{\text{High Prec}}$$

External products are omnipresent in FHE

- Covers every "tools":
keyswitches (2009), gadget decomposition, relinearization, products, automorphisms, bootstrappings (2024)
- Yet, proper formalization of external products appeared only from 2016.



External products are omnipresent in FHE

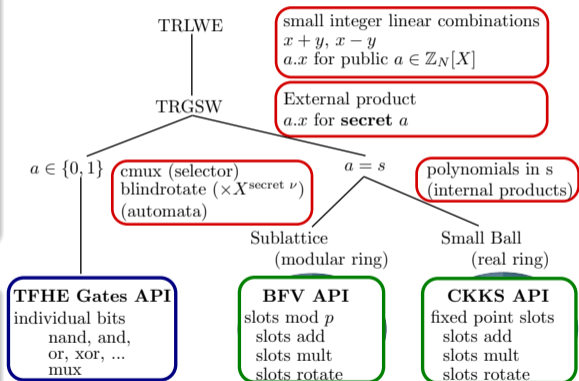
Ring-LWE has never really been about Rings

- Ring-LWE FHE was never really about rings.
- FHE arithmetic flows from **secret linear combinations**, a.k.a. **external products**:
 - with small integer polynomial coefficients
 - over some high precision space (Torus, or big integers)

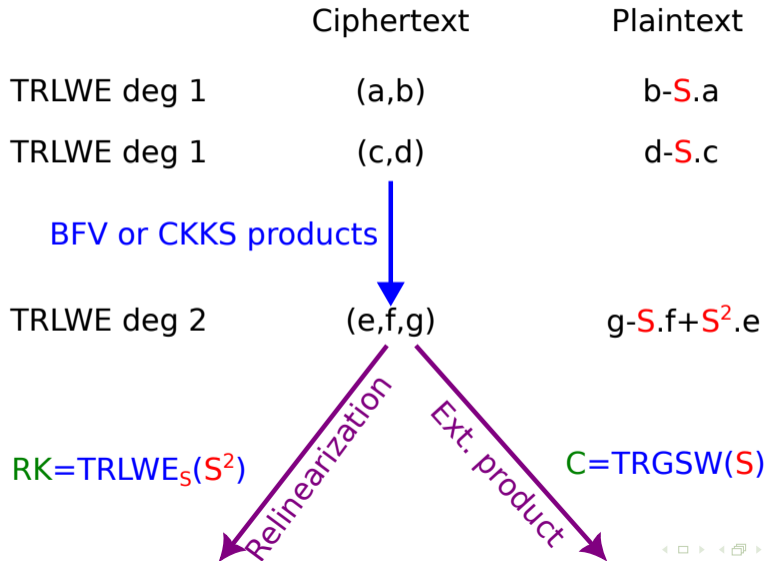
$$\sum \underbrace{a_i}_{\text{small int}} \cdot \underbrace{c_i}_{\text{High Prec}}$$

External products are omnipresent in FHE

- Covers every "tools":
keyswitches (2009), gadget decomposition, relinearization, products, automorphisms, bootstrappings (2024)
- Yet, proper formalization of external products appeared only from 2016.



External products vs. relinearization



We just need half of the TRGSW

Relinearization needs only half of the TRGSW material

- TRGSW was designed for **secret** \times **secret** products
- Relinearization uses **secret** \times **public** product

Half External Product (i.e. secret-public)

Plaintext:	$\mathbb{Z}_N[X]$	\cdot	$\mathbb{T}_N[X]$	\rightarrow	$\mathbb{T}_N[X]$
Ciphertext:	$\text{HTRGSW}(A)$	\triangle	b	\rightarrow	$\text{TRLWE}(A \cdot b)$
Noise:	ϵ			\rightarrow	$\approx N\epsilon$

External product from [CGG16]

Plaintext:	$\mathbb{Z}_N[X]$	\cdot	$\mathbb{T}_N[X]$	\rightarrow	$\mathbb{T}_N[X]$
Ciphertext:	$\text{TRGSW}(A)$	\square	$\text{TRLWE}(b)$	\rightarrow	$\text{TRLWE}(A \cdot b)$
Noise:	ϵ		α	\rightarrow	$\approx \ A\ _2 \cdot \alpha$

Application of half TRGSW products

Relinearization application

- Relinearization:

$$(f, g) + \text{HTRGSW}(S^2) \triangleq e$$

Original external product

- Original TRGSW:

$$\underbrace{\text{TRGSW}(A)}_C = \underbrace{(\text{HTRGSW}(A))}_D, \underbrace{(\text{HTRGSW}(SA))}_E$$

- 1 full = 2 halves:

$$C \square (a, b) = E \triangleq a + D \triangleq b$$

Half TRGSW: It is not just syntax!

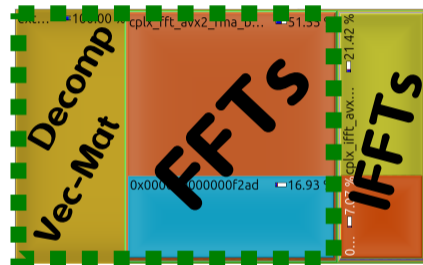
$$C \boxplus (a, b) = E \boxtimes a + D \boxtimes b$$

Approximate decomposition: Faster FHE

- Since b needs less precision than a , $D \boxtimes b$ is faster than $E \boxtimes a$
- expect 1 less FFT per external product.

Practical improvements

- 8 FFTs \rightarrow 7 FFTs in the original TFHE lib.
- Possibility to improve also the Circuit bootstrapping (impacts TFHE-rs).



17% faster

overall: 15% faster

Table 2. Performance comparison of gate bootstrapping with a n2-standard GCP instance with 64GB of RAM and a 12-th Gen i7-1260p laptop with 64GB of RAM. All the benchmarks are single core.

Library	Instruction set	n2-standard	12-Gen i7-1260p
TFHE-lib, spqlios-fma	AVX2	22.4ms	10.4ms
TFHE-rs, TFHE_LIB_PARAMETERS	AVX2	18.2ms	8.6ms
	AVX512	14.4ms	<i>not supported</i>
TFHE-rs, DEFAULT_PARAMETERS	AVX2	14.4ms	7.6ms
	AVX512	13.7ms	<i>not supported</i>
Our work, halfTRGSW	AVX2	11.2ms	5.3ms

Part II - Auxiliary Gadgets: The return of base- 2^K

Idash 2018 - a pivoting moment for FHE

- Chimera: combine TFHE's blind-rotate with BFV-style arithmetic
- CKKS: first really efficient Full-RNS POC

External Product Formula

Ciphertext layer: $\text{prod} = \text{GadgetDec}(a, b) \times \text{Fixed_GSW_key}$

Arithmetic layer: $\text{vector}\langle \mathbb{T}_N[X] \rangle_{\text{Huge}} = \sum_{i=1}^{\text{level } \ell} \mathbb{Z}_N[X]_{\text{small}} \cdot \text{vector}\langle \mathbb{T}_N[X] \rangle_{\text{Huge}}$

$$\text{vector}_{\text{Huge}} \langle \mathbb{T}_N[X] \rangle = \sum_{i=1}^{\text{level } \ell} \underset{\text{small}}{\mathbb{Z}_N[X]} \cdot \underset{\text{Huge}}{\text{vector} \langle \mathbb{T}_N[X] \rangle}$$

$$\underset{\text{small}}{\mathbb{Z}_N[X]} \cdot \underset{\text{Huge}}{\mathbb{T}_N[X]}$$

Two big blockers

- FFT products: does not take advantage of the small coefficients.
- High precision base- 2^K : carry propagation is prohibitive!

Main resolution: Give up, abort. (<2020) – Separation of concerns

- **Small depth FHE**: Use Approx-decomposition with FP.
- **Large depth FHE**: Use CRT + RNS.

$$\text{vector}_{\text{Huge}} \langle \mathbb{T}_N[X] \rangle = \sum_{i=1}^{\text{level } \ell} \mathbb{Z}_N[X]_{\text{small}} \cdot \text{vector}_{\text{Huge}} \langle \mathbb{T}_N[X] \rangle$$

$$\mathbb{Z}_N[X]_{\text{small}} \cdot \mathbb{T}_N[X]_{\text{Huge}}$$

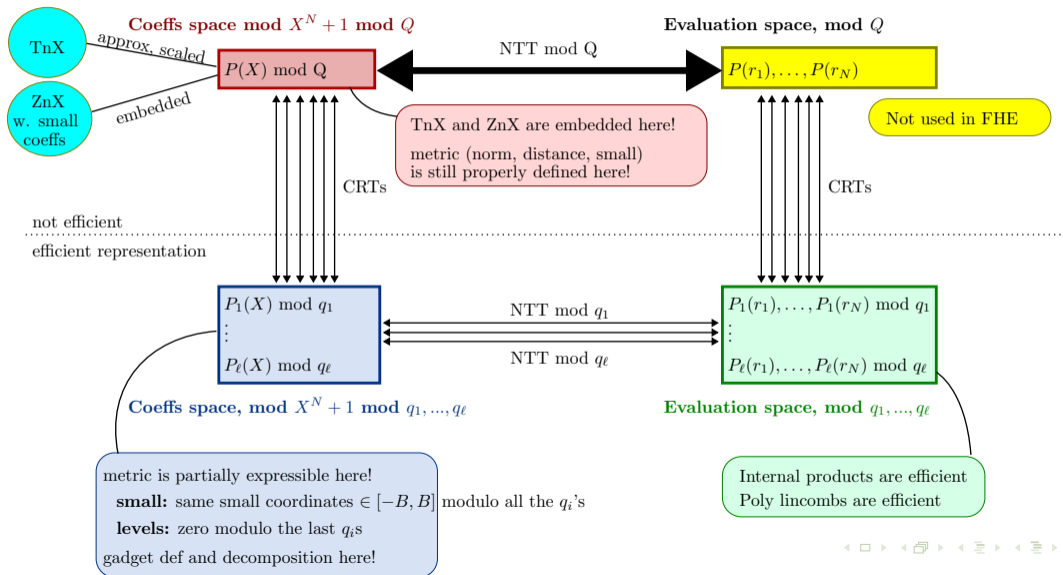
Two big blockers

- FFT products: does not take advantage of the small coefficients.
- High precision base- 2^K : carry propagation is prohibitive!

Main resolution: Give up, abort. (<2020) – Separation of concerns

- **Small depth FHE:** Use Approx-decomposition with FP.
- **Large depth FHE:** Use CRT + RNS.

Full-RNS CRT-CKKS since 2018



[KLSS23]: Breaking news!! Full-RNS was not optimal!

External Product Formula (CRT)

$$\text{prod} = \text{GadgetDec}(a, b) \times \text{Fixed_GSW_key}$$

$$\begin{bmatrix} \mathbb{Z}_N[X] \pmod{q_1} \\ \vdots \\ \mathbb{Z}_N[X] \pmod{q_\ell} \end{bmatrix} = \sum_{i=1}^{\text{level } \ell} \underset{\text{small}}{\mathbb{Z}_N[X]} \cdot \begin{bmatrix} \mathbb{Z}_N[X] \pmod{q_1} \\ \vdots \\ \mathbb{Z}_N[X] \pmod{q_\ell} \end{bmatrix}$$

Observation by [KLSS23]

- Although the gadget decomposition is small, we still compute ℓ NTT's per polynomials (so ℓ^2 NTT's in total)
- This is the same as if the gadget decomposition was Huge.

The big idea (simplified scoop!)

- Forget/Delay the annoying mod (q_1, \dots, q_ℓ) , do the computation over \mathbb{Z} !!!

[KLSS23]: Breaking news!! Full-RNS was not optimal!

External Product Formula (CRT)

$$\text{prod} = \text{GadgetDec}(a, b) \times \text{Fixed_GSW_key}$$

$$\left(\begin{bmatrix} \mathbb{Z}_N[X] \text{---mod } q_1 \\ \vdots \\ \mathbb{Z}_N[X] \text{---mod } q_\ell \end{bmatrix} \right) = \sum_{i=1}^{\text{level } \ell} \mathbb{Z}_N[X]_{\text{small}} \cdot \left(\begin{bmatrix} \mathbb{Z}_N[X] \text{---mod } q_1 \\ \vdots \\ \mathbb{Z}_N[X] \text{---mod } q_\ell \end{bmatrix} \right) \begin{matrix} \text{mod } q_1 \\ \vdots \\ \text{mod } q_\ell \end{matrix}$$

Observation by [KLSS23]

- Although the gadget decomposition is small, we still compute ℓ NTT's per polynomials (so ℓ^2 NTT's in total)
- This is the same as if the gadget decomposition was Huge.

The big idea (simplified scoop!)

- Forget/Delay the annoying mod (q_1, \dots, q_ℓ) , do the computation over $\mathbb{Z}!!!$

[BCGGJ23] Extends the scope of [KLSS23])

- If reductions mod (q_1, \dots, q_ℓ) can be delayed and amortized.
- Carry propagation too!!!

When we consider base- 2^k limbs:

$$\mathbb{T}_N[X] = \underbrace{\mathbb{Z}_N[X]}_{\text{Huge}} \cdot \frac{1}{2^K} + \underbrace{\mathbb{Z}_N[X]}_{\text{small}} \cdot \frac{1}{2^{2K}} + \dots + \underbrace{\mathbb{Z}_N[X]}_{\text{small}} \cdot \frac{1}{2^{\ell K}}$$

Linear combinations can be done per limb:

$$\begin{bmatrix} \underbrace{\mathbb{Z}_N[X]}_{\text{bounded}} \\ \underbrace{\mathbb{Z}_N[X]}_{\text{bounded}} \\ \vdots \\ \underbrace{\mathbb{Z}_N[X]}_{\text{bounded}} \end{bmatrix} = \sum_{i=1}^{\text{level}\ell} \underbrace{\mathbb{Z}_N[X]}_{\text{small}} \cdot \begin{bmatrix} \underbrace{\mathbb{Z}_N[X]}_{\text{small}} \\ \underbrace{\mathbb{Z}_N[X]}_{\text{small}} \\ \vdots \\ \underbrace{\mathbb{Z}_N[X]}_{\text{small}} \end{bmatrix}$$

and carry propagation is delayed until the very end! (small DFTs, bounded integers!)

Instead of:

$$\mathbb{T}_N[X] = \underbrace{\mathbb{Z}_N[X]}_{\text{Huge}} \cdot \frac{1}{2^K} + \underbrace{\mathbb{Z}_N[X]}_{\text{small}} \cdot \frac{1}{2^{2K}} + \cdots + \underbrace{\mathbb{Z}_N[X]}_{\text{small}} \cdot \frac{1}{2^{\ell K}}$$

Write:

$$\mathbb{T}_N[X] = \underbrace{\mathbb{Z}_N[X]}_{\text{Huge}} \cdot Y + \underbrace{\mathbb{Z}_N[X]}_{\text{small}} \cdot Y^2 + \cdots + \underbrace{\mathbb{Z}_N[X]}_{\text{small}} \cdot Y^\ell$$

Bivariate representation




- Carry propagation decoupling formalized by the presence of a variable Y . – hence "Bivariate" representation
- Multiplication of $\mathbb{Z}(X, Y)$ do make sense (genuine morphism between lifts over $\mathbb{R}[X]$) – leveraged in BFV and CKKS internal products.
- Fast arithmetic possible product via DFT over X and over Y !!!

Wins

- External products (half and full) are fastest on base- 2^K (optimal)
- Compatible with automorphisms and internal products (optimal)
- Prefix property: modulus rescaling is free (optimal)
- CKKS Noise levels: continuous! no gap, no artificial rescaling (optimal)
- Parametrization: every parameters flows from the noise level, no mysterious additional moduli (optimal)

Minor Setback

- CKKS and BFV: dominant terms improve, negligible terms are a bit larger. (good enough!)

		Cyclotomic Arithmetic (Backend)	
		NTT	FFT
Big-Number (Frontend)	CRT	RNS-BFV RNS-CKKS (large depth)	 Also Valid!
	Bivariate (Base-2 ^K)	 Also Valid!	TFHE (low depth)  large depth!

Frontend vs. Backend

- The big number mechanism and the cyclotomic mechanism must be decoupled.
- All combinations are as efficient (up to a factor 2)
- Switching between representations is easy! (i.e. external product to switch between representations)

Aritmetic over vector/matrices of small integer polynomials. – 3 key operations

- 1 **dft/idft**: via an NTT or FFT backend (whichever is faster)
- 2 **vmp_prepare/apply**: vector \times preprocessed matrix
- 3 **cnv_prepare/apply**: (precomputed) vector \times vector convolution

It is omnipotent!

- Can all **CRT and bivariate frontends at any depth!**
 - BlindRotate (CGGI bootstrapping in 6ms).
 - CKKS and BFV products (depth 30 in 0.3s)
 - Keyswitches and Automorphisms (depth 30 in 0.2s)
- Easier and sufficient **target for hardware developers!**

Aritmetic over vector/matrices of small integer polynomials. – 3 key operations

- 1 **dft/idft**: via an NTT or FFT backend (whichever is faster)
- 2 **vmp_prepare/apply**: vector \times preprocessed matrix
- 3 **cnv_prepare/apply**: (precomputed) vector \times vector convolution

It is omnipotent!

- Can all **CRT and bivariate frontends at any depth!**
 - BlindRotate (CGGI bootstrapping in 6ms).
 - CKKS and BFV products (depth 30 in 0.3s)
 - Keyswitches and Automorphisms (depth 30 in 0.2s)
- Easier and sufficient **target for hardware developers!**

Some benchmarks

Operation	Keyswitch		Automorphism	CKKS product
Size	$N = 64k$ L=1729	$N = 32k$ L=865	$N = 64k$ L=1729	$N = 64k$ L=1729
Hardware	n2-standard VM Xeon(R) CPU @ 2.8GHz, 64GB RAM			
- Full-RNS (best r)	3.111s	0.359s	3.279s	3.311s
- [26] (best r)	0.965s	0.161s	1.134s	1.155s
- ours: biv + fft-f64 ($K = 19$)	0.589s	0.086s	0.602s	0.862s
- ours: biv + ntt120 ($K = 52$)	0.541s	0.073s	0.547s	0.777s
Hardware	Laptop with Intel Core i7-1260P @ 4.7GHz, 64GB RAM			
- Full-RNS (best r)	1.598s	0.192s	1.796s	1.759s
- [26] (best r)	0.521s	0.085s	0.578s	0.598s
- ours: biv + fft-f64 ($K = 19$)	0.228s	0.027s	0.233s	0.335s
- ours: biv + ntt120 ($K = 52$)	0.218s	0.029s	0.221s	0.314s

Conclusions and Key Takeways

Key Takeway for hardware developers

- FFT or NTT are both valid and will provide efficient FHE in all levels: just pick one.
- Aim for native hardware support of the largest possible word size.
64-bit arithmetic is good, 128-bit would be much better!

Key Takeway for FHE libraries developers

- Use gadget decompositions as primary software ↔ hardware API.
- Do not expose parameters that the end-user won't be able to set!

Key Takeway for FHE compilers/transpilers developers

- CRT frontend is best for large dimension matrix or matrix/vector algebra over ciphertexts.
External products can be amortized.
- Bivariate frontend is best when internal products are rare (e.g. only external products and automorphisms), or chained in sequence.
CKKS noise propagation at any plaintext precision is easier. (no clamping to the nearest integer "level")

Thank you!