

Over 100x Faster Bootstrapping in Fully Homomorphic Encryption through Memory-centric Optimization with GPUs

Younho Lee
(Dept. Data Science, SeoulTech)

Contents

I. Preliminary

II. Accelerating Full-RNS CKKS (main work)

I. Background & Motivation

II. Approach

I. Overview

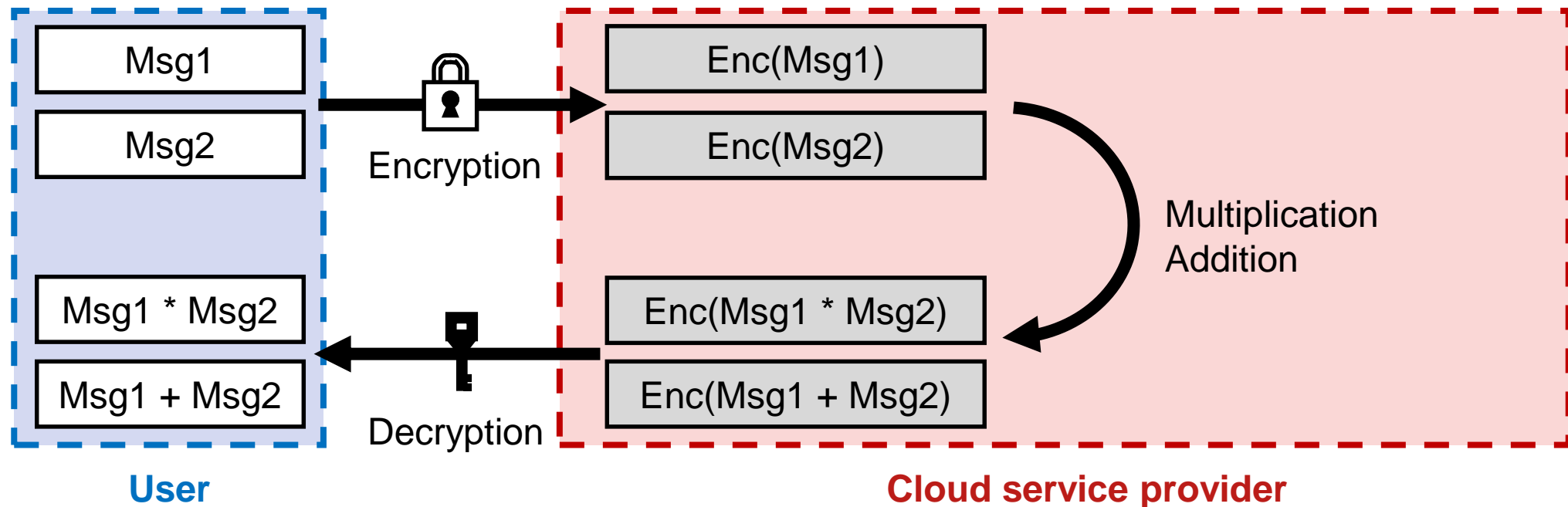
II. Intra-HE-operation Fusions

III. Inter-HE-operation Fusion

III. Summary

Preliminary: Homomorphic encryption (HE) costs a lot

- HE is a cryptographic scheme that enables computation on encrypted messages (*ciphertexts*).
- Users send their data to a cloud service provider without privacy worries.
- Popular HE schemes
 - BFV (integers), TFHE (binary numbers), **CKKS** (fixed-point numbers)



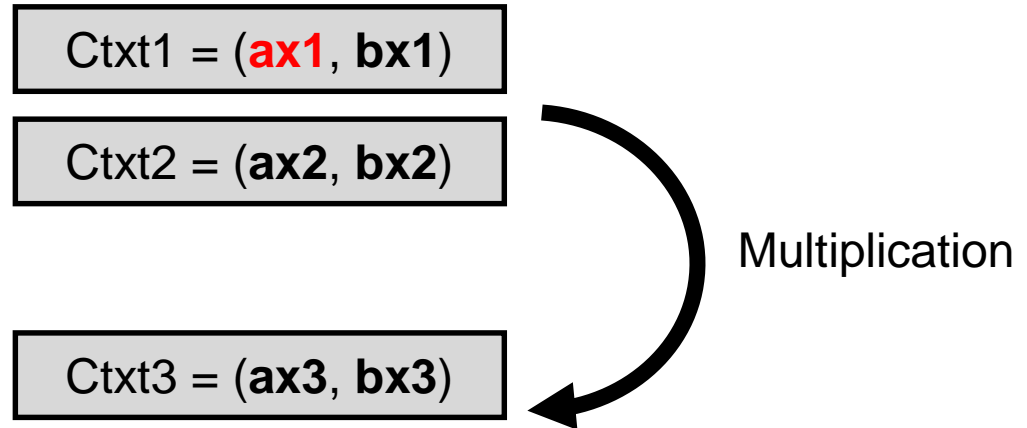
Preliminary: Homomorphic encryption (HE) costs a lot

- Ciphertexts are large.
 - One ciphertext in CKKS = **10-100s of MBs**
- Slow HE operations.
 - (vs. native operations) 100x – 200x in addition
 - (vs. native operations) **1000x – 10000x** in multiplication
- *Bootstrapping* is even heavier.
 - Bootstrapping reduces the noise after HE operations on a ciphertext.
 - = It enables *Fully Homomorphic Encryption* (FHE).
 - It takes **10s of seconds** for a single bootstrapping in CKKS.
 - = (using a single-core of a latest CPU)

Preliminary: Homomorphic encryption (HE) costs a lot

- Each encrypted message is called ciphertext.
- Each ciphertext is represented as a pair of high-degree polynomial in a polynomial ring,

$$\mathbf{ax1} = a_0 + a_1x + a_2x^2 + \dots + a_{65535}x^{65535}$$



Preliminary: Homomorphic encryption (HE) costs a lot

- Each encrypted message is called ciphertext.
- Each ciphertext is represented as a pair of high-degree polynomial in a polynomial ring,
 - ...whose **coefficients** are big-integers.

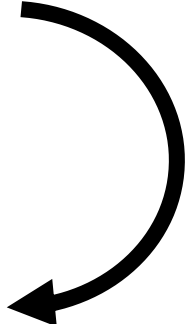
1200-bit
 $\mathbf{ax1} = a_0 + a_1x + a_2x^2 + \dots + a_{65535}x^{65535} \in \mathbb{Z}_q[X]/(X^N + 1)$
 $q = 2^{1200}, N = 65536$

Ctxt1 = (**ax1**, **bx1**)

Ctxt2 = (**ax2**, **bx2**)

Ctxt3 = (**ax3**, **bx3**)

Multiplication



Multiplication between such polynomials is expensive.

Preliminary: Massive parallelism in HE

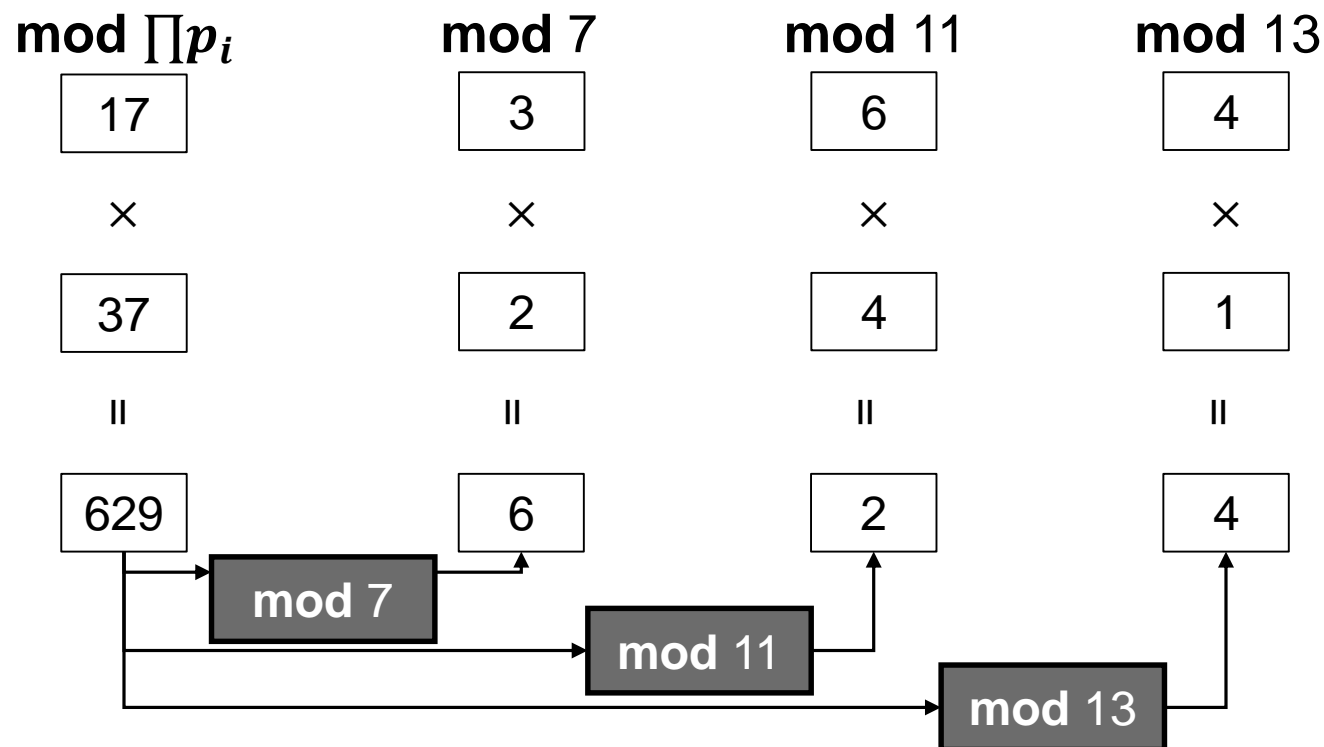
- Fortunately, massive parallelism exists in polynomial multiplication:
 - Chinese Remainder Theorem (CRT) to exploit Residue Number System (RNS)
 - Number Theoretic Transform (NTT)

Preliminary: Chinese Remainder Theorem (CRT) in HE

- CRT reduces the computational complexity of multiplication of two big numbers.

Ex) multiplying 17 and 37 using coprimes $(p_1, p_2, p_3) = (7, 11, 13)$

(CRT moduli)

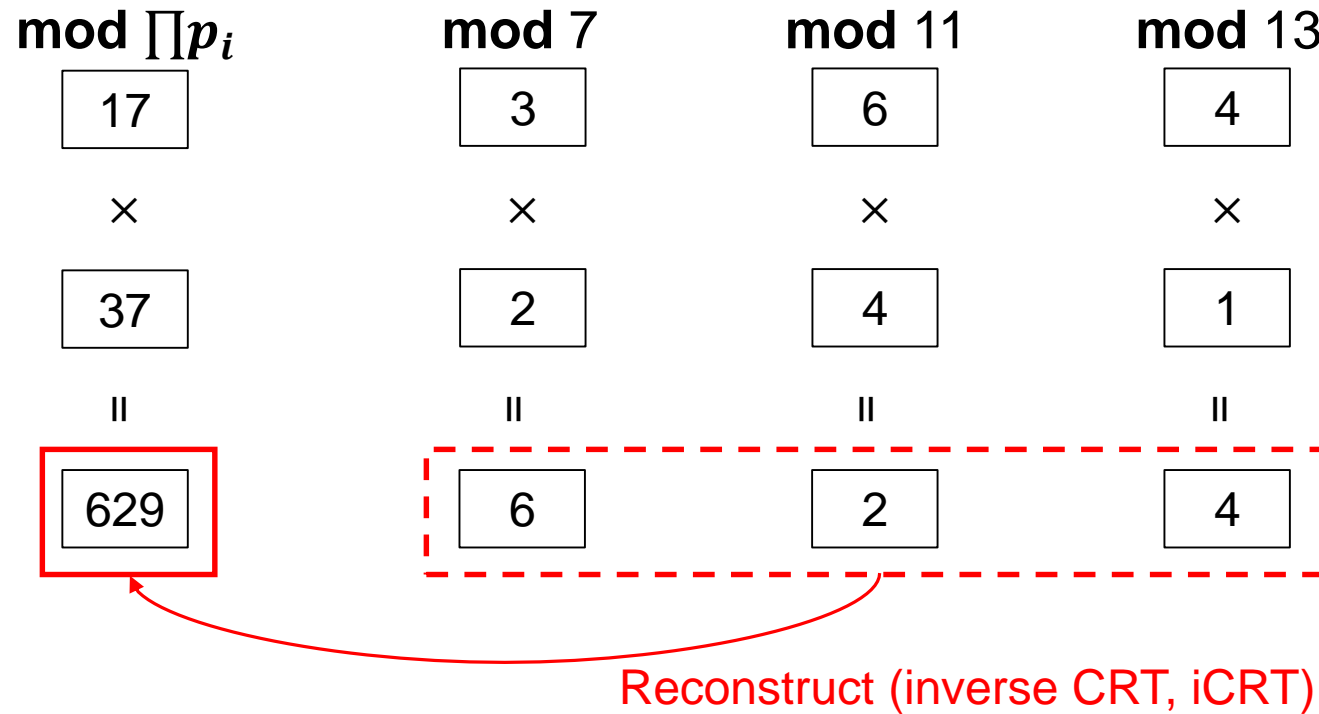


Preliminary: Chinese Remainder Theorem (CRT) in HE

- CRT reduces the computational complexity of multiplication of two big numbers.

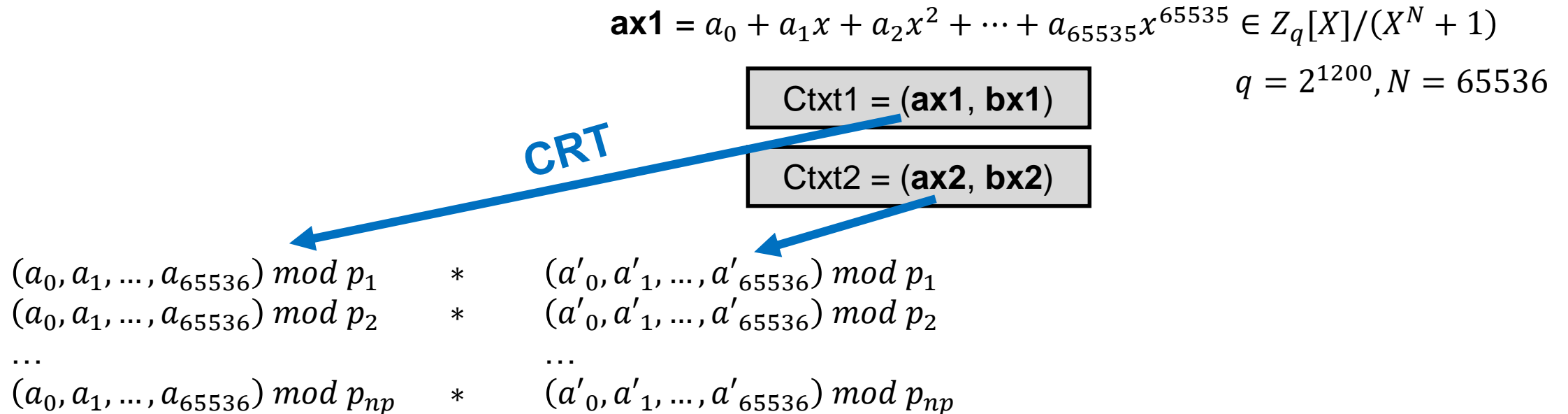
Ex) multiplying 17 and 37 using coprimes $(p_1, p_2, p_3) = (7, 11, 13)$

(CRT moduli)



Preliminary: Number Theoretic Transformation (NTT) in HE

- Each multiplication between big-integer coefficients can be solved by CRT, but still with $O(N^2)$.



Preliminary: Number Theoretic Transformation (NTT) in HE

- Each multiplication between big-integer coefficients can be solved by CRT, but still with $O(N^2)$.
- NTT, an integer-variant DFT, reduces $O(N^2)$ to $O(N \log N)$.

Point-wise mul, $O(N)$

$$(A_0, A_1, \dots, A_{65536}) \quad \odot \quad (A'_0, A'_1, \dots, A'_{65536})$$



NTT, $O(N \log N)$



$$\begin{array}{lcl} (a_0, a_1, \dots, a_{65536}) \bmod p_1 & * & (a'_0, a'_1, \dots, a'_{65536}) \bmod p_1 \\ (a_0, a_1, \dots, a_{65536}) \bmod p_2 & * & (a'_0, a'_1, \dots, a'_{65536}) \bmod p_2 \\ \dots & & \dots \\ (a_0, a_1, \dots, a_{65536}) \bmod p_{np} & * & (a'_0, a'_1, \dots, a'_{65536}) \bmod p_{np} \end{array}$$

Preliminary: Variants of Cheon-Kim-Kim-Song (CKKS)

- Representation of plaintext and ciphertext in RLWE:
 - Plaintext $m(X) \in R_q = \mathbb{Z}_q[X]/(X^N + 1)$, $q \leq Q (= \text{ciphertext modulus})$
 - Ciphertext $ct(X) = (b(X), a(X)) \in R_q^2$
- Variants of CKKS
 - Binary CKKS [1]
 - = arbitrary q and Q
 - = needs expensive CRT & iCRT conversion.
 - RNS-CKKS [2]
 - = Fixed q and Q
 - = No (i)CRT; exploits *Fast base conversion*.
 - (Improved) RNS-CKKS [3]
 - = [2] + better algorithms (*key-switching* & bootstrapping)

[1] Cheon et al., Homomorphic encryption for arithmetic of approximate numbers, AsiaCrypt'17

[2] Cheon et al., A full RNS variant of approximate homomorphic encryption, SAC'18

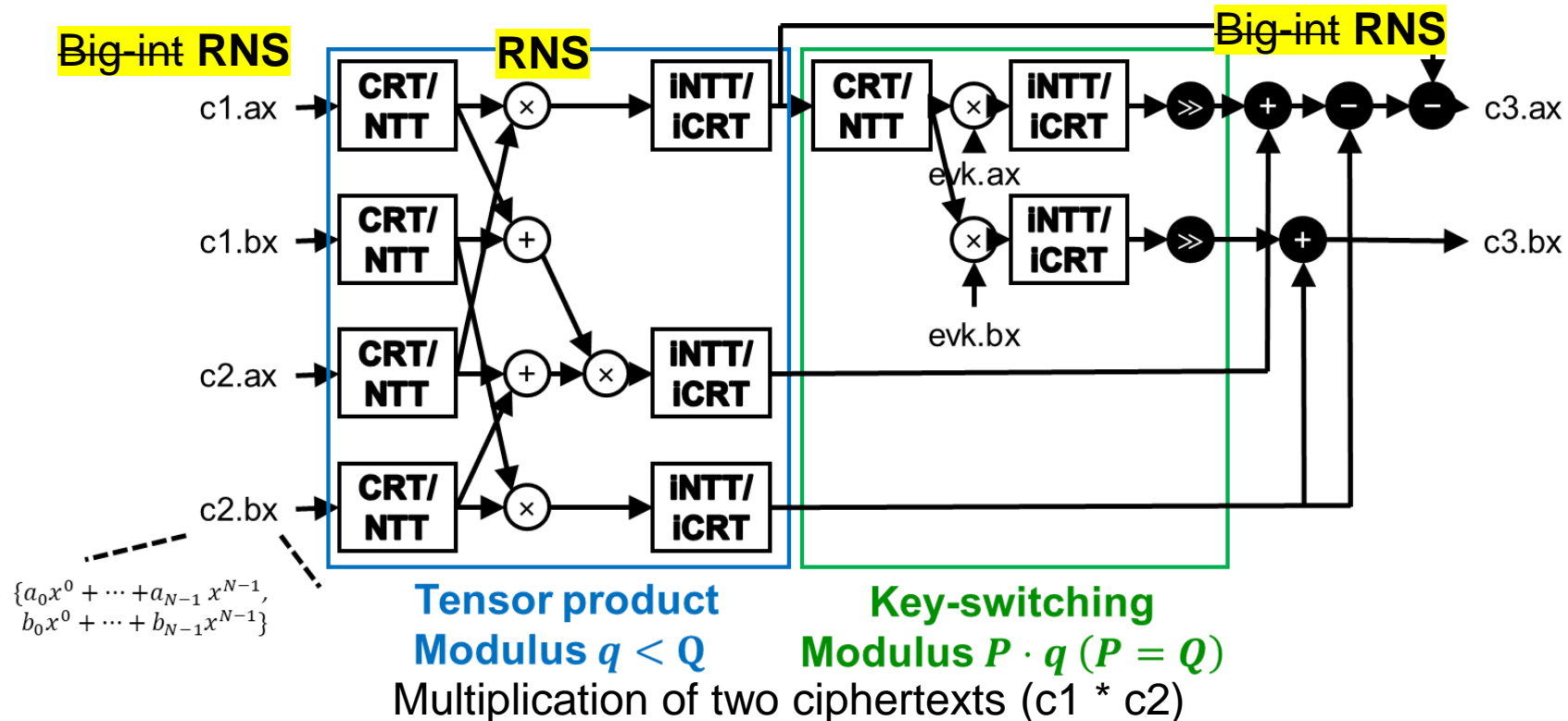
[3] Han et al., Better bootstrapping for approximate homomorphic encryption, CT-RSA'20

Contents

- I. Preliminary
- II. Accelerating Full-RNS CKKS (main work)
 - I. Background & Motivation**
 - II. Approach
 - I. Overview
 - II. Intra-HE-operation Fusions
 - III. Inter-HE-operation Fusion
- III. Summary

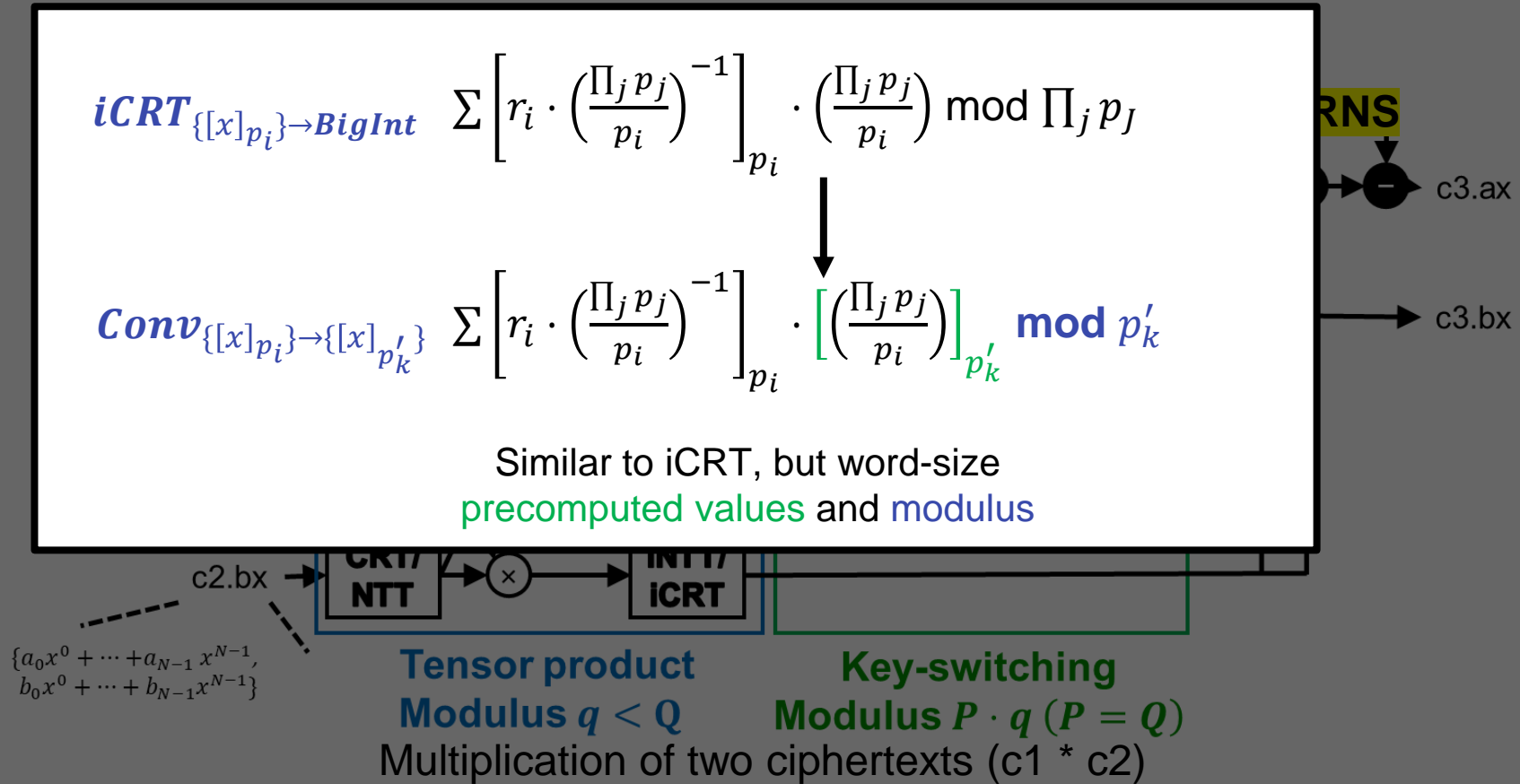
Why Full-RNS CKKS?

- Polynomials are always in RNS domain.
 - no big-int computations



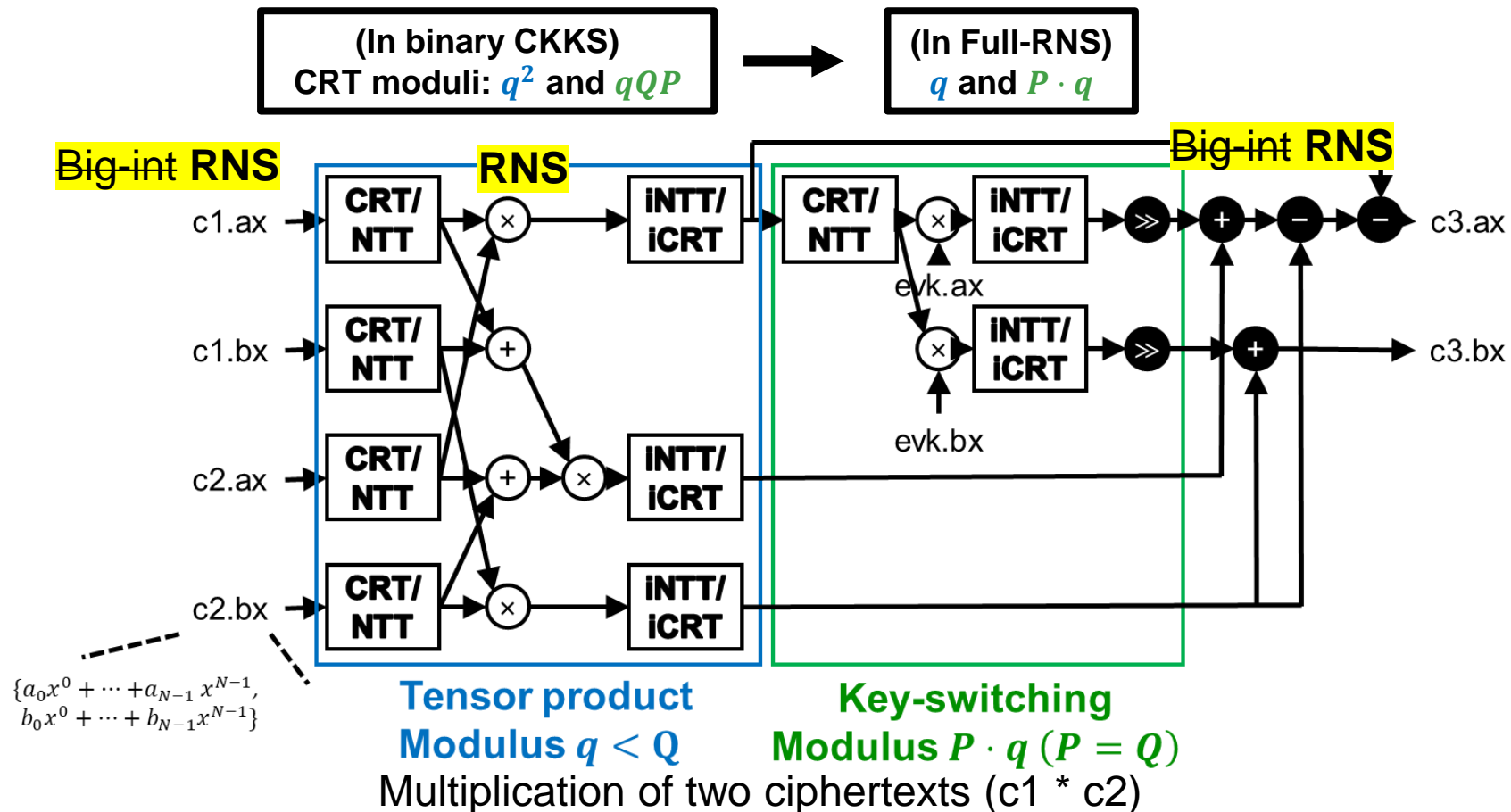
Why Full-RNS CKKS?

- Polynomials are always in RNS domain.
- no big-int computations (iCRT+CRT \rightarrow *Fast base conversion*)



Why Full-RNS CKKS?

- Polynomials are always in RNS domain.
 - no big-int computations (iCRT+CRT \rightarrow *Fast base conversion*)
- Small *CRT modulus*: ciphertext modulus is a product of moduli ($Q = q_0 q_1 \dots q_L$, $q = q_0 q_1 \dots q_\ell$)



Accelerating Full-RNS CKKS

- The first GPU implementation of bootstrapping in the latest Full-RNS CKKS scheme.
- Analyze memory bandwidth bottleneck of the GPU implementation.
- Apply memory-centric optimizations, resulting in sub-second bootstrapping.
 - = a speedup of 257x, compared to single-threaded CPU.

Background: a recent Full-RNS CKKS [HK20]

- Parameters
 - Multiplicative level L
 - Primes $q_0, \dots, q_L, p_1, \dots, p_k$
 - Modulus $Q = q_0 q_1 \dots q_L, P = p_1 p_2 \dots p_k$
- Plaintext representation:
 - (RNS form) $m(X) \in R_Q$
 - (NTT form) $m = NTT(m(X))$
- Ciphertext representation:
 - (RNS form) $ct = (b(X), a(X)) \in R_Q^2$
 - (NTT form) $ct = (b, a)$

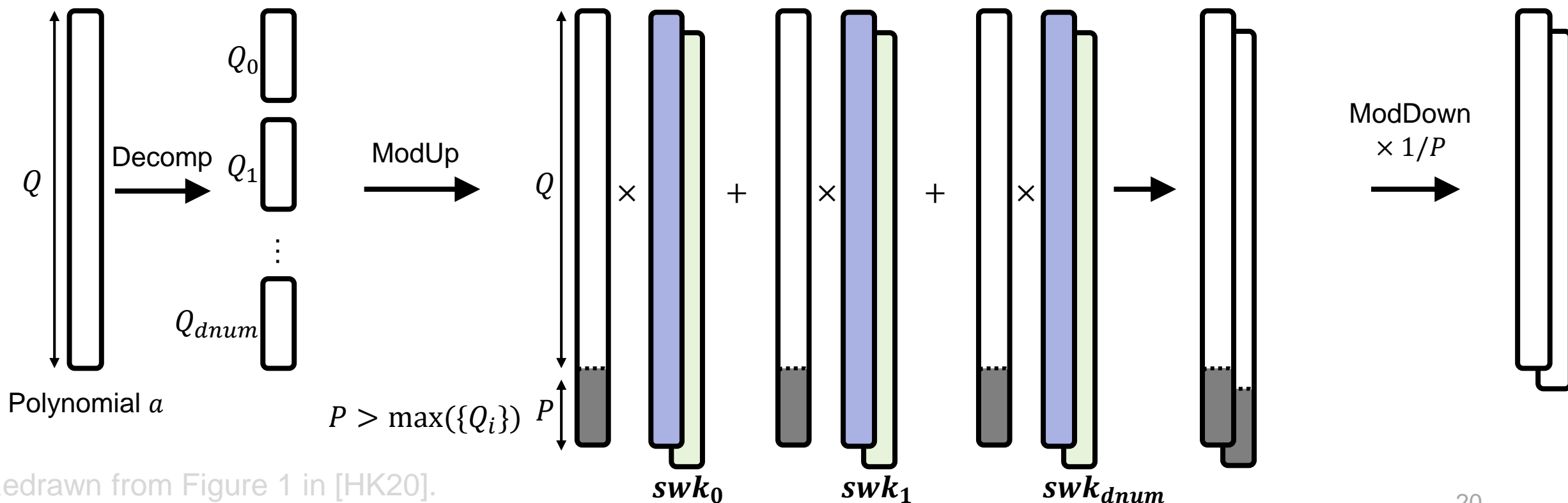
Background: a recent Full-RNS CKKS [HK20]

- Key operations
 - Plaintext multiplication
 $= CMult(ct, m) = (ma, mb)$
 - Ciphertext multiplication
 $= Mult(ct_1, ct_2)$
 1. computes a tensor product $(d_0, d_1, d_2) = (b_2b_1, b_1a_2 + b_2a_1, a_1a_2)$
 2. return $(d_0, d_1) + KeySwitch_{MultKey}(d_2)$
 - Ciphertext rotation by a rotation index r
 $= Rotate(ct, r)$
 1. computes an automorphism $b'(X) = b(X^{5^r}), a'(X) = a(X^{5^r})$
 2. return $(b', 0) + KeySwitch_{RotationKey_r}(a')$
- [HK20] introduces a new key-switching method for CKKS.
 - Generalized key switching (a.k.a. hybrid key-switching)

Background: a recent Full-RNS CKKS [HK20]

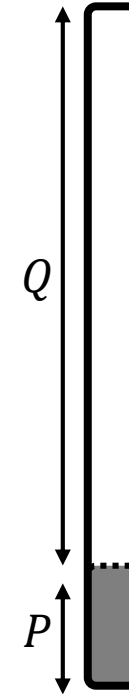
- $KeySwitch_{swk}(a)$

1. (*Decomp*) decompose into up to **dnum** parts, each having α moduli: $Q_0 = q_0 q_1 \cdots q_{\alpha-1}$, $Q_1 = q_\alpha q_{\alpha+1} \cdots$
2. (*ModUp*) raise the modulus of each part to PQ .
3. (*Inner-product*) do inner-product with the key-switching key $swk = (swk_0, swk_1, \dots, swk_{dnum})$.
4. (*ModDown*) reduce the modulus to Q and multiply $1/P$.



Background: Choosing a good *dnum* parameter

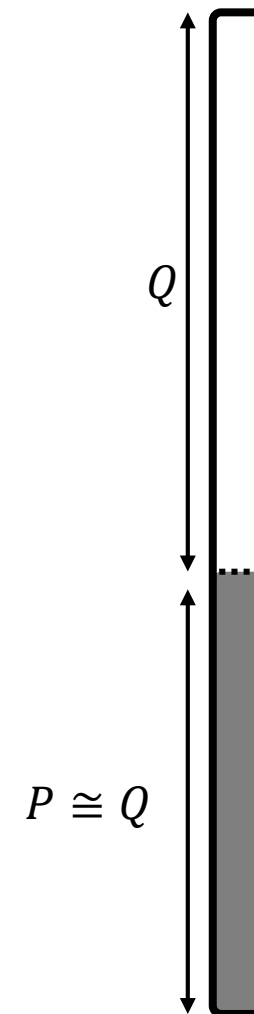
- $\uparrow \log PQ = \downarrow$ security
- The choice of **dnum** affects complexity, key size, and security.



(after ModUp)

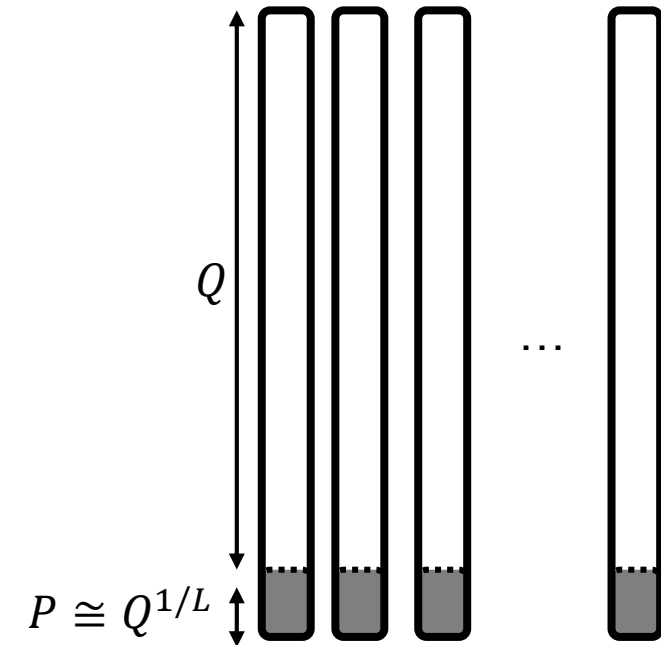
Background: Choosing a good *dnum* parameter

- $\uparrow \log PQ = \downarrow$ security
- The choice of **dnum** affects complexity, key size, and security.
- Given Q ,
 - **dnum** = 1 (minimum)
 - $\uparrow P$
 - \downarrow security
 - \downarrow size of keys



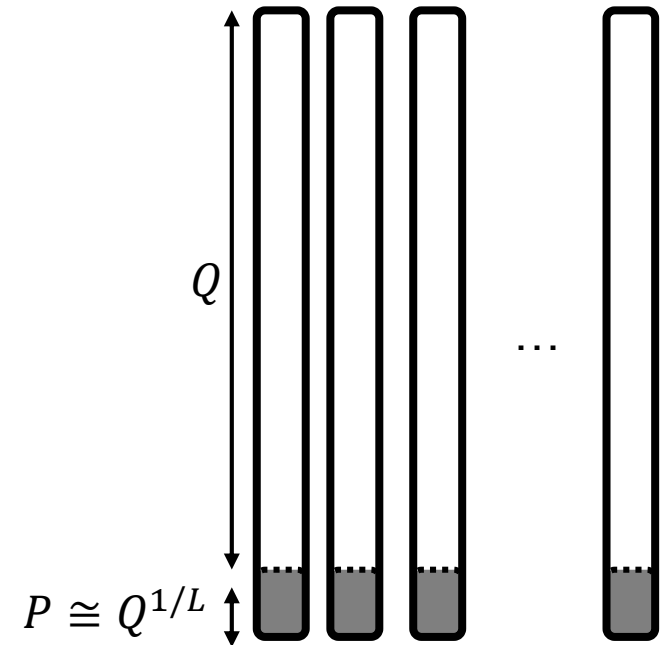
Background: Choosing a good *dnum* parameter

- $\uparrow \log PQ = \downarrow$ security
- The choice of **dnum** affects complexity, key size, and security.
- Given Q ,
 - **dnum** = 1 (minimum)
 - $\uparrow P$
 - \downarrow security
 - \downarrow size of keys
 - **dnum** = L (maximum)
 - $\downarrow P$
 - \uparrow security
 - \uparrow size of keys



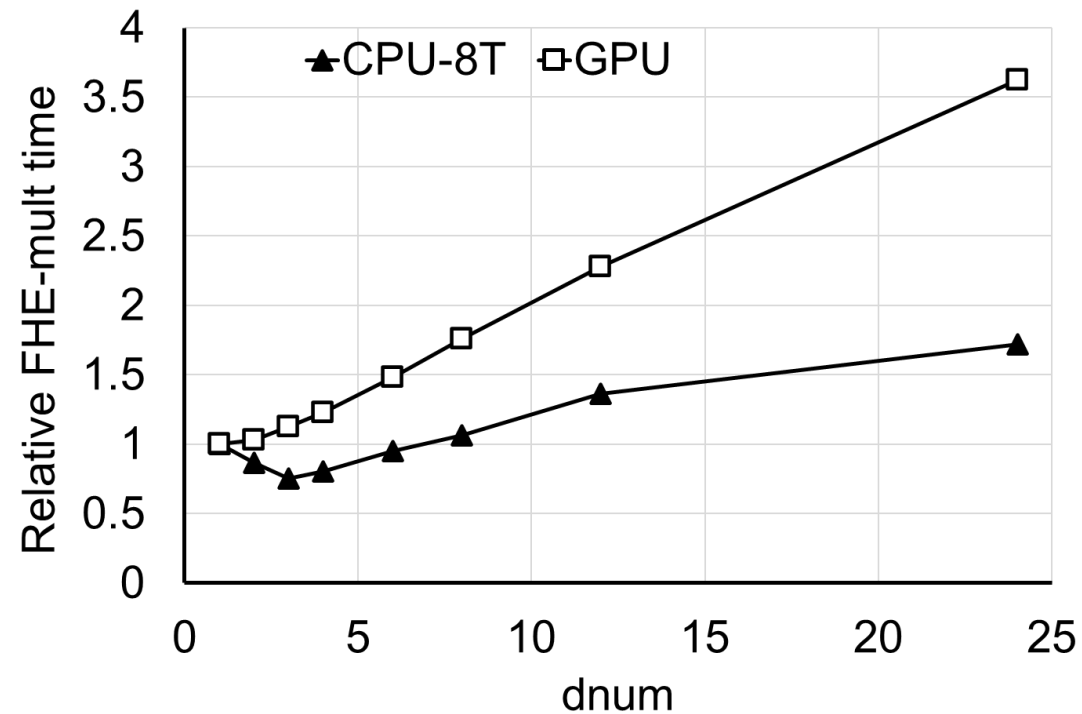
Background: Choosing a good *dnum* parameter

- $\uparrow \log PQ = \downarrow$ security
- The choice of **dnum** affects complexity, key size, and security.
- Given Q ,
 - **dnum** = 1 (minimum)
 - $\uparrow P$
 - \downarrow security
 - \downarrow size of keys
 - **dnum** = L (maximum)
 - $\downarrow P$
 - \uparrow security
 - \uparrow size of keys
 - Computational complexity becomes minimal somewhere between them.
 - = [HK20] chooses that value.



Motivation: GPU implementations are more memory-bottlenecked

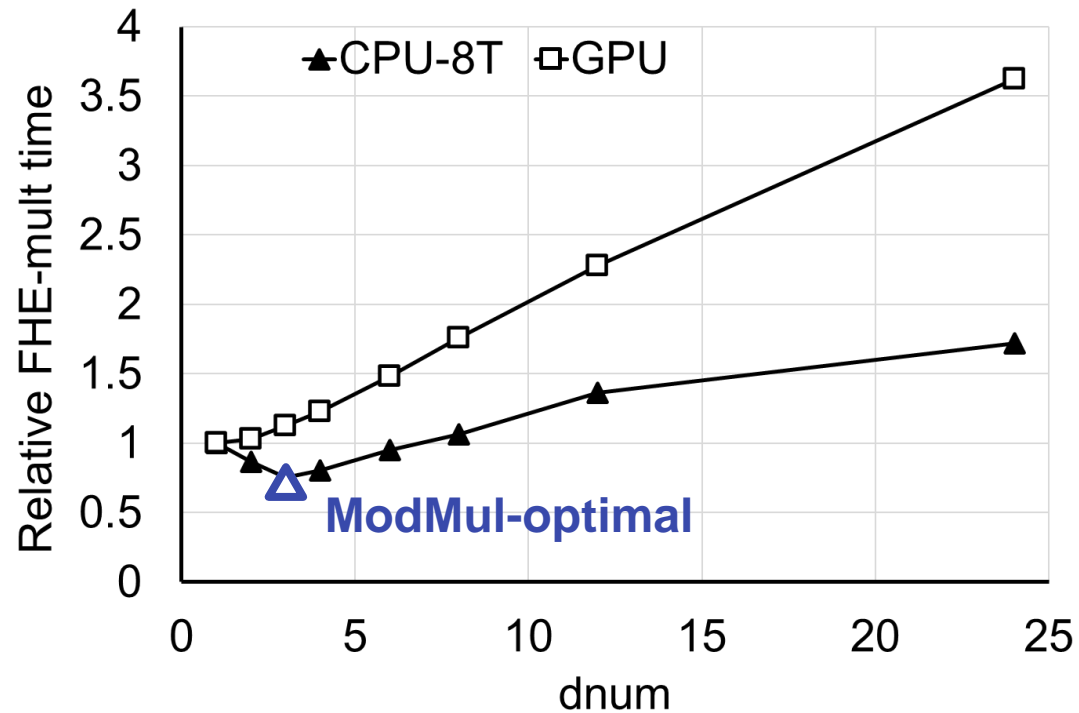
- Last-level cache size of a server class
 - CPU: ~256MB
 - GPU: ~6MB
- = Hardly accommodates ciphertexts.



dnum vs. multiplication time for a fixed Q

Motivation : GPU implementations are more memory-bottlenecked

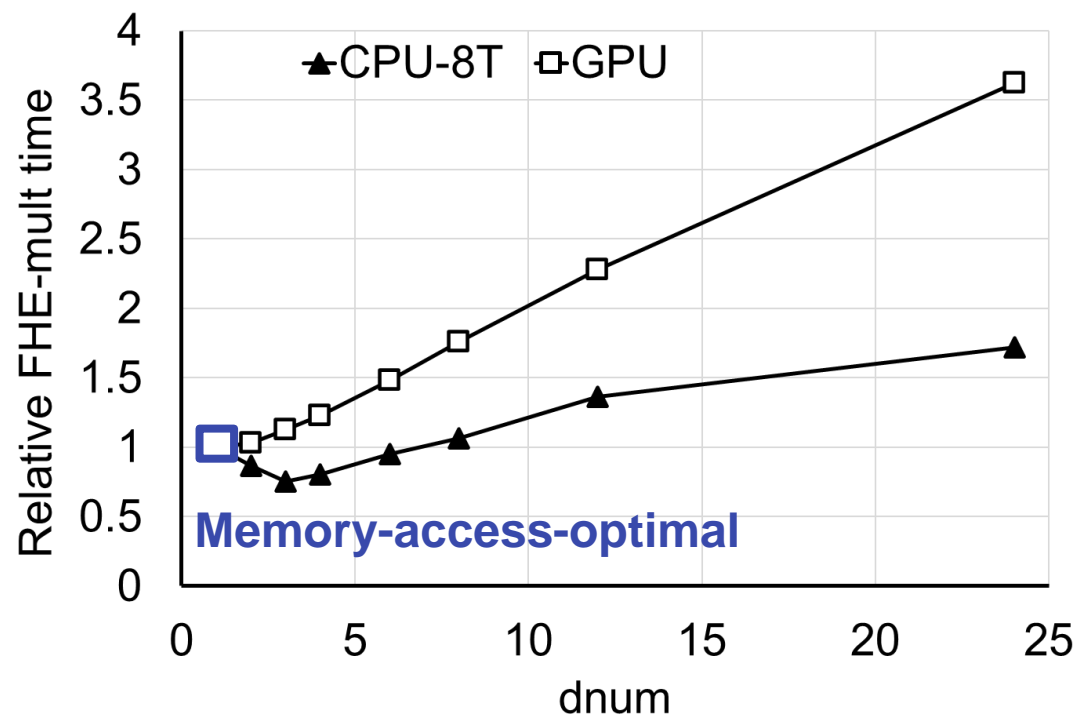
- Last-level cache size of a server class
 - CPU: ~256MB
 - GPU: ~6MB
- = Hardly accommodates ciphertexts.



dnum vs. multiplication time for a fixed Q

Motivation : GPU implementations are more memory-bottlenecked

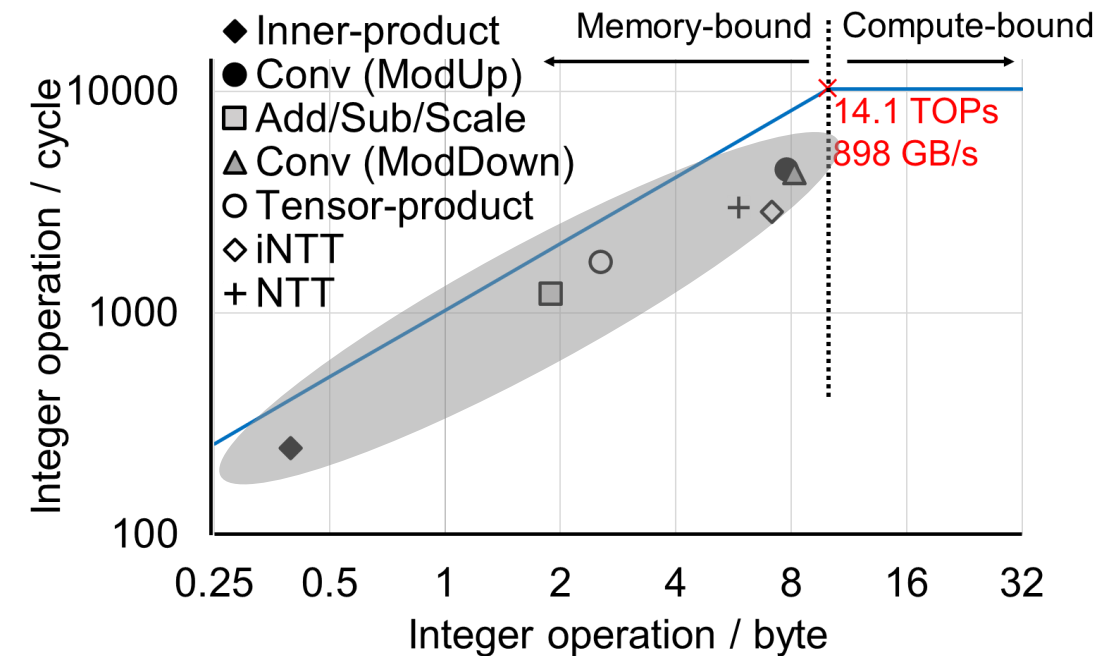
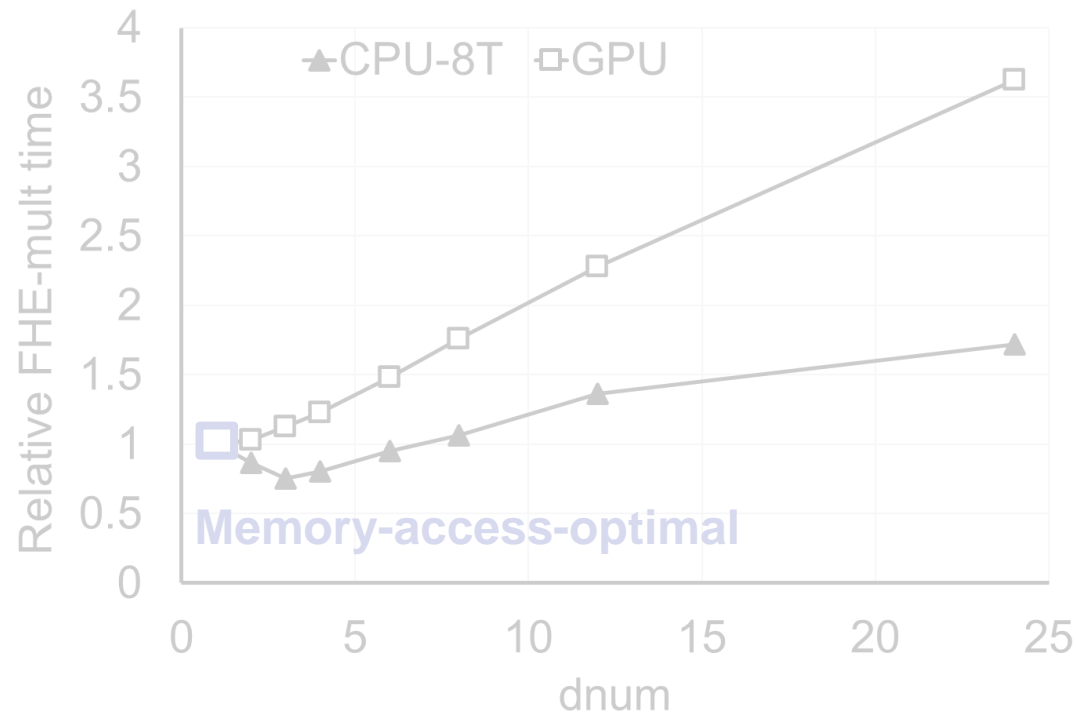
- Last-level cache size of a server class
 - CPU: ~256MB
 - GPU: ~6MB
- = Hardly accommodates ciphertexts.



dnum vs. multiplication time for a fixed Q

Motivation: GPU implementations are more memory-bottlenecked

- Last-level cache size of a server class
 - CPU: ~256MB
 - GPU: ~6MB
- = Hardly accommodates ciphertexts.



Roofline plot of a multiplication

HE operations are memory-bound, especially on GPUs.

Contents

I. Preliminary

II. Accelerating Full-RNS CKKS (main work)

I. Background & Motivation

II. Approach

I. Overview

II. Intra-HE-operation Fusions

III. Inter-HE-operation Fusion

III. Summary

Overview: Brief introduction to the programming model in GPUs

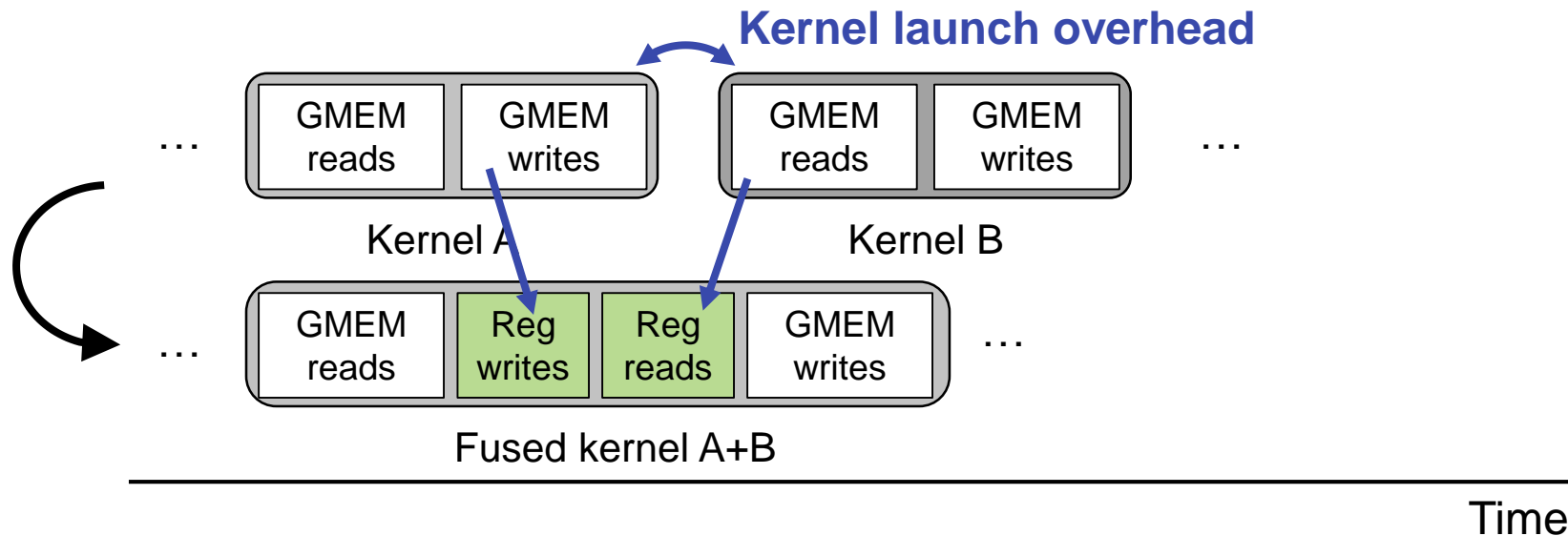
- A GPU has many *Streaming Multiprocessors (SM)*.
 - Each SM manages and executes threads in groups of 32 parallel threads.
- A function executed in a GPU is called *Kernel*.
 - Many GPU threads run the same kernel in parallel in SMs.
- Each kernel is configured with
 - the number of threads running on the kernel
 - the amount of the *shared memory* it will use.
- Shared memory is a user-configurable scratch-pad memory.

Overview: Baseline GPU implementation

- Based on prior works:
 - NTT-friendly degree-first data layout [CLP17, HS14, CHK+19]
 - Executing $(L \times N)$ RNS operations in a kernel [BPA+19]
 - Fast NTT/iNTT implementations [KJPA20]
 - ...

Overview: Kernel fusion - memory-centric optimization for GPUs

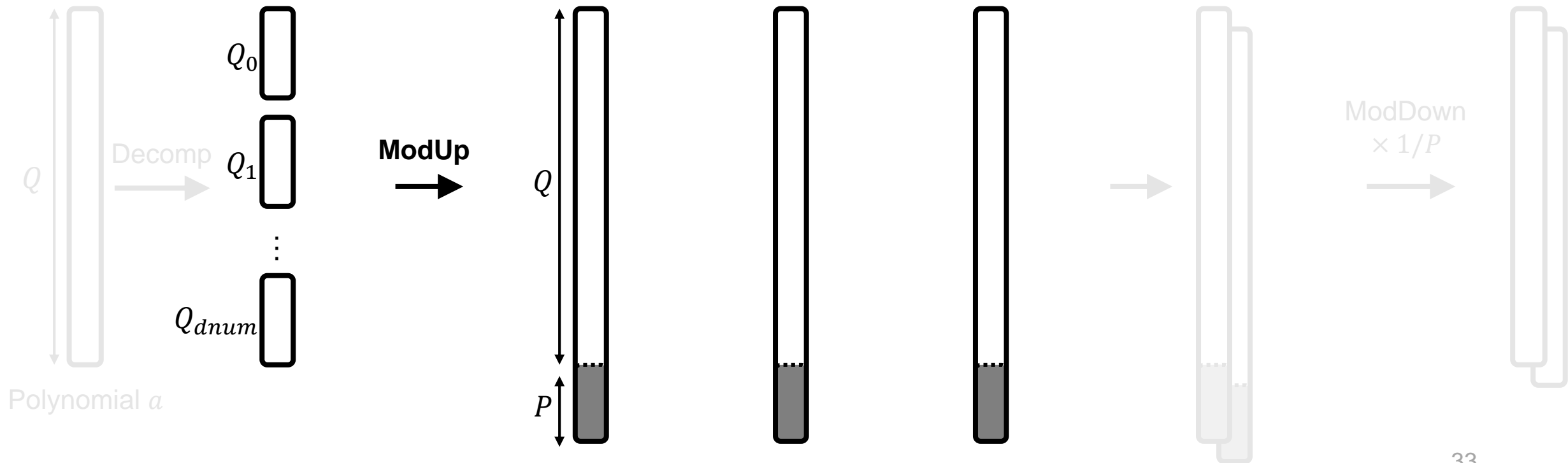
- Kernel fusion (operation fusion)
 - A common technique that fuses multiple GPU kernels into a single kernel
- Kernel fusion can
 - save global memory (DRAM) reads and writes by reusing the data in register file/shared memory.
 - = Good for the kernels with low OP/B (memory-bound).
 - reduce *kernel launch overhead* by merging small kernels
- We found many kernel fusion opportunities in **intra- / inter-HE-operation** manner.



Kernel fusion reduces global memory (GMEM) accesses.

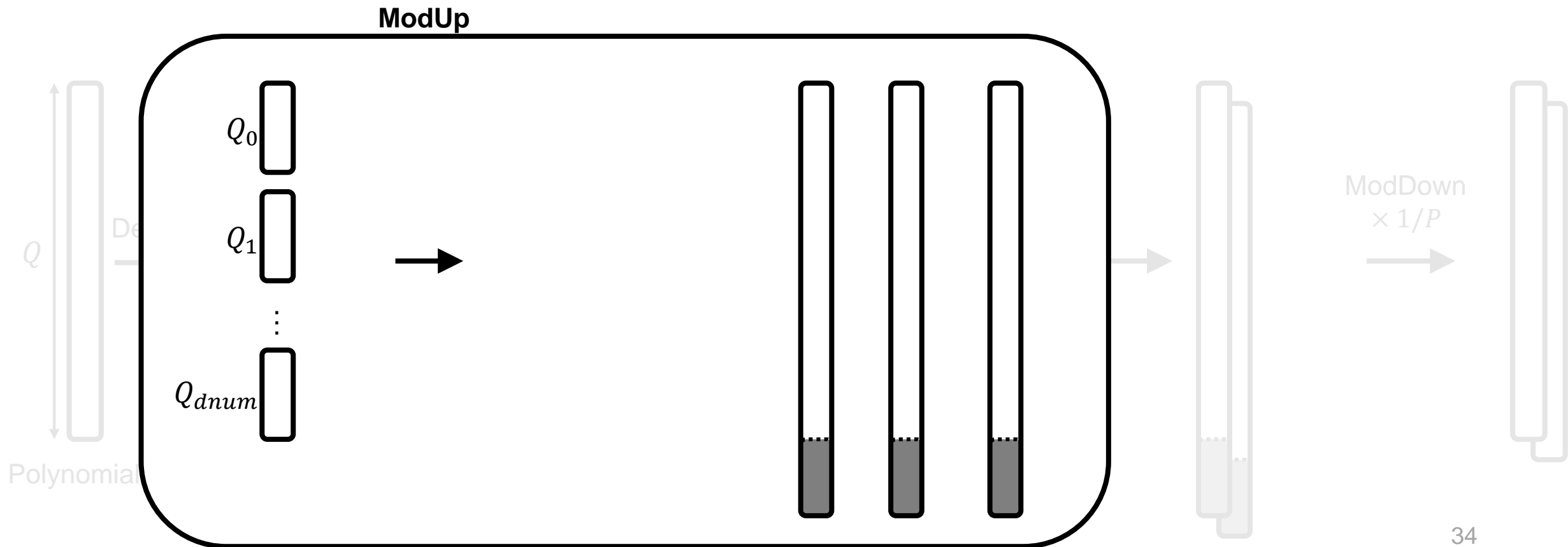
Intra-HE-operation Fusions

- ModUp Fusion (MF)*



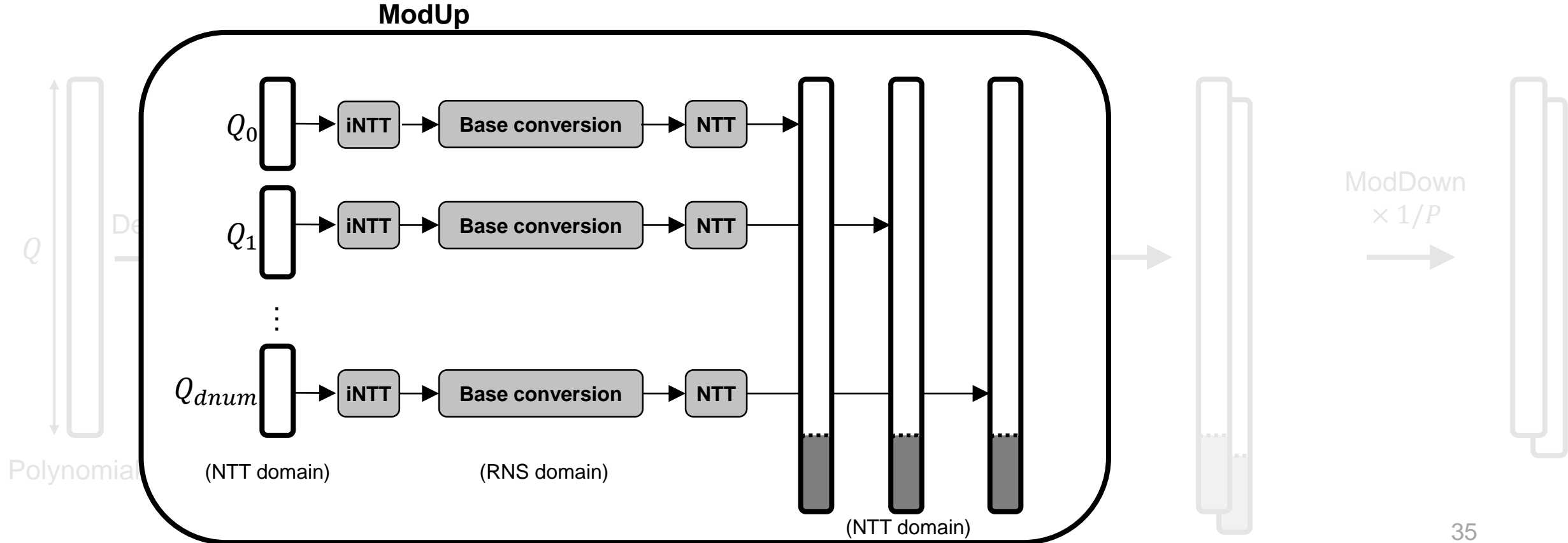
Intra-HE-operation Fusions

- *ModUp Fusion (MF)*



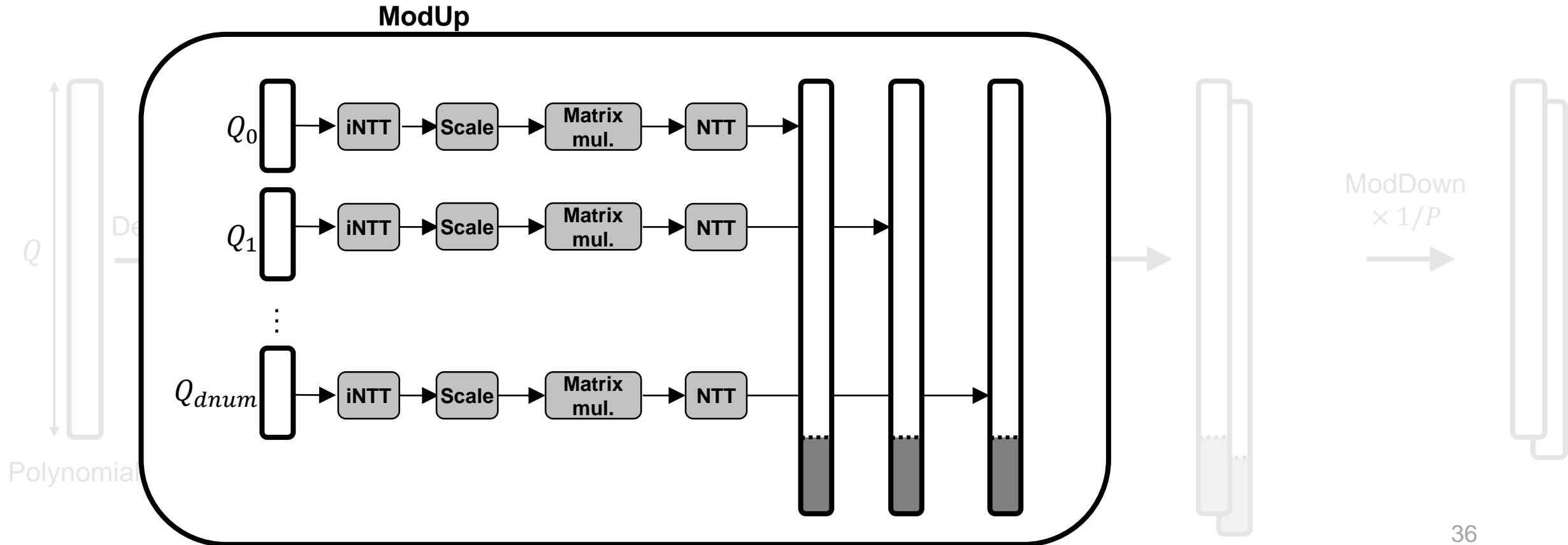
Intra-HE-operation Fusions

- *ModUp Fusion (MF)*



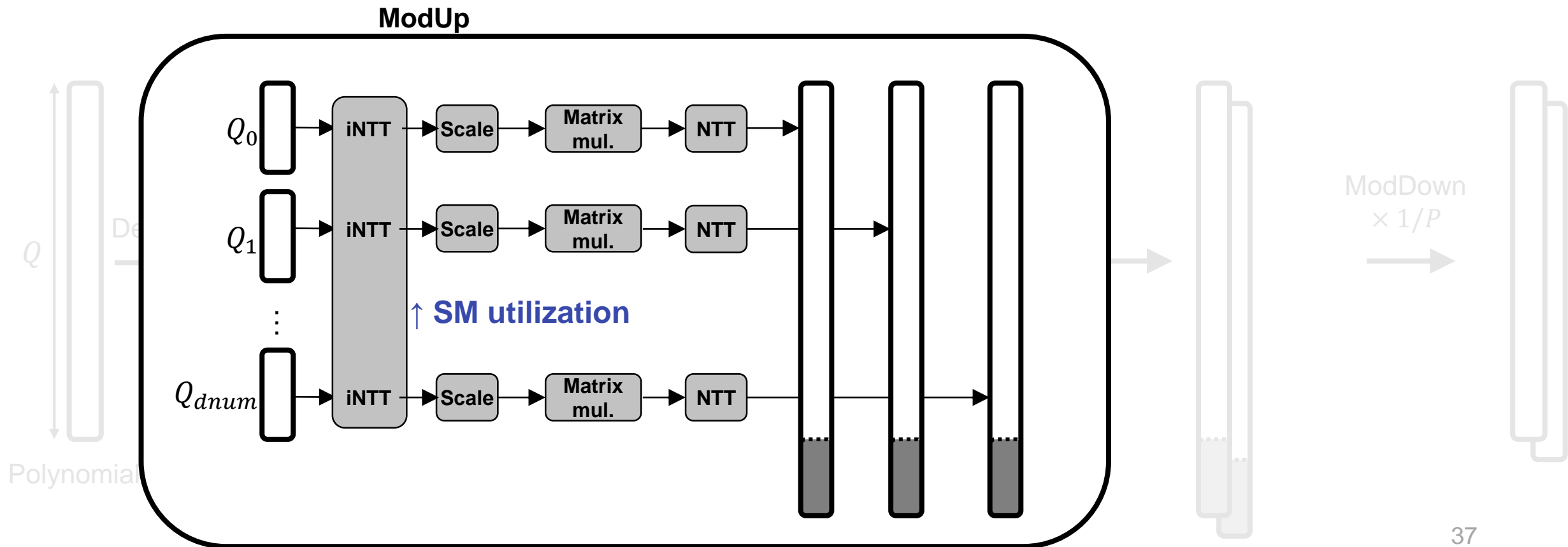
Intra-HE-operation Fusions

- *ModUp Fusion (MF)*
 1. fuses 'small' iNTT kernels into a large one.



Intra-HE-operation Fusions

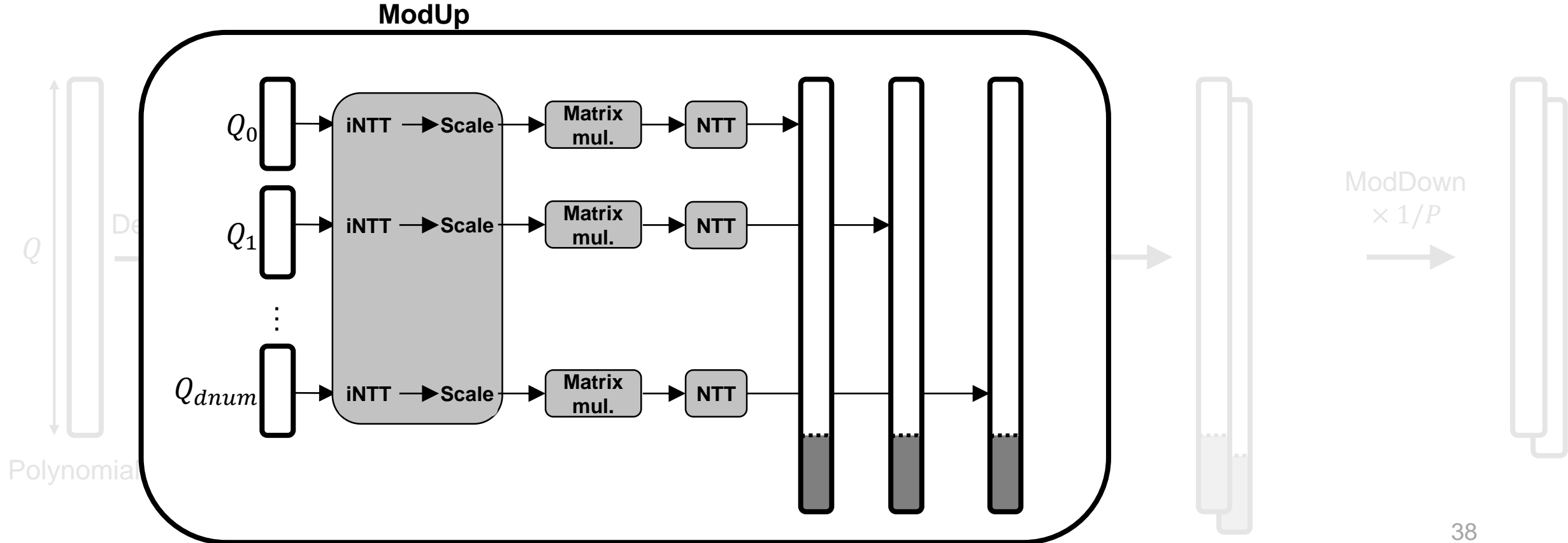
- *ModUp Fusion (MF)*
 1. fuses 'small' iNTT kernels into a large one.



Intra-HE-operation Fusions

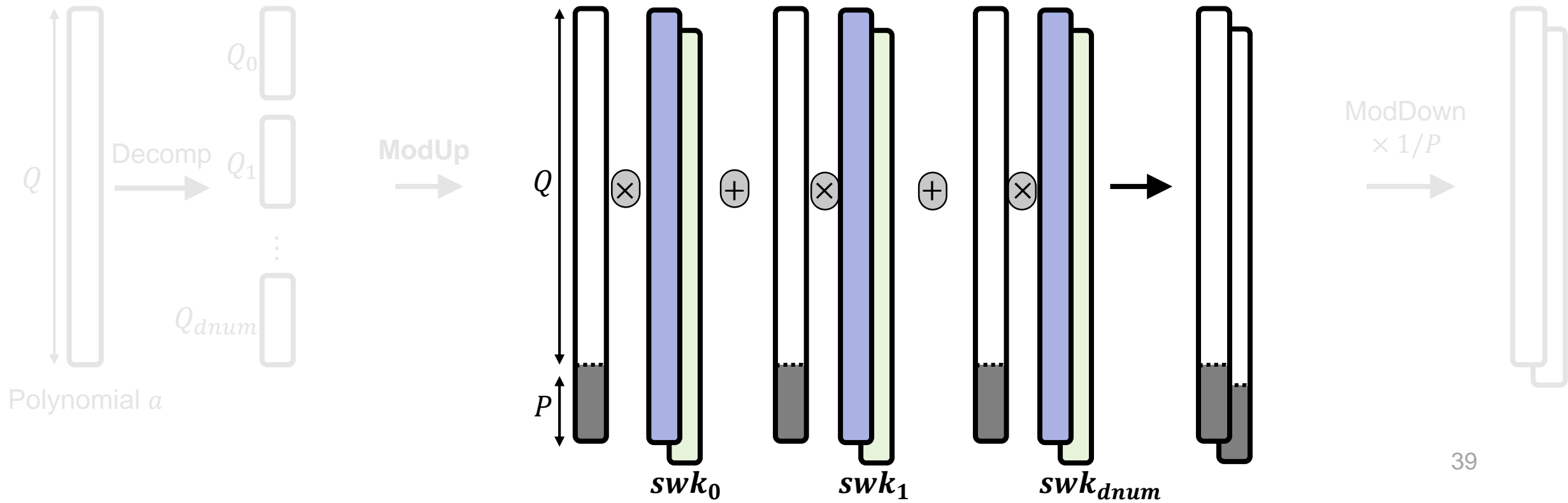
- *ModUp Fusion (MF)*

1. fuses 'small' iNTT kernels into a large one.
2. fuses the scaling kernels with iNTTs.



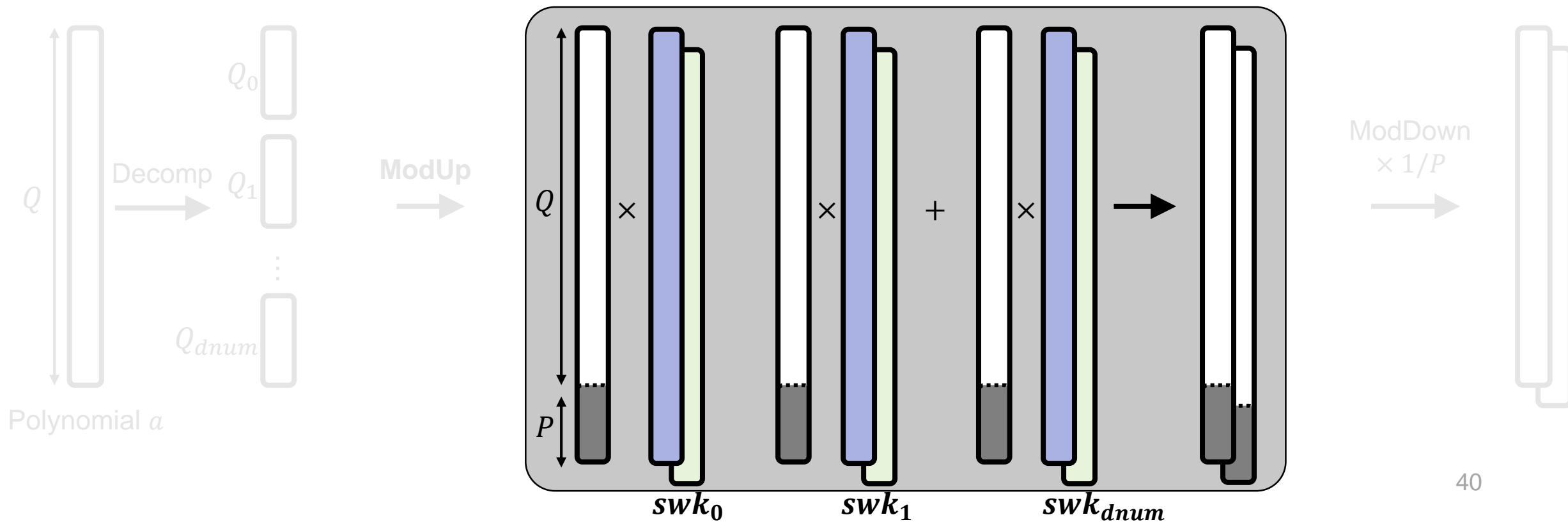
Intra-HE-operation Fusions

- *Inner-product Fusion (IF)*



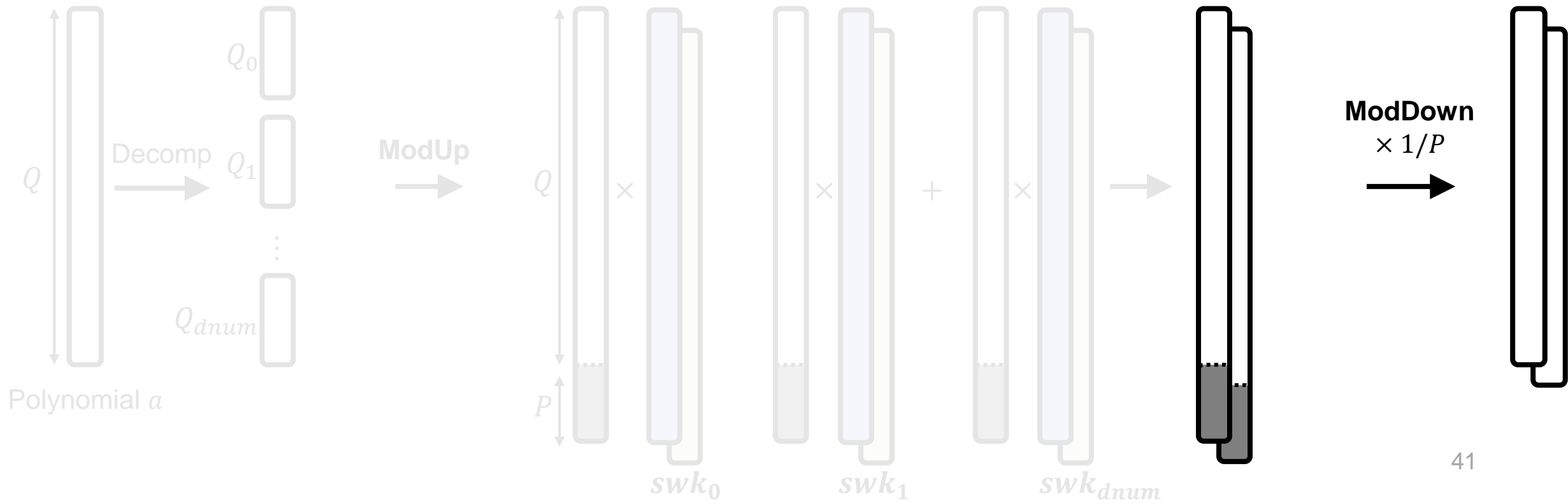
Intra-HE-operation Fusions

- *Inner-product Fusion (IF)*
 - fuses multiplications and additions into one kernel.



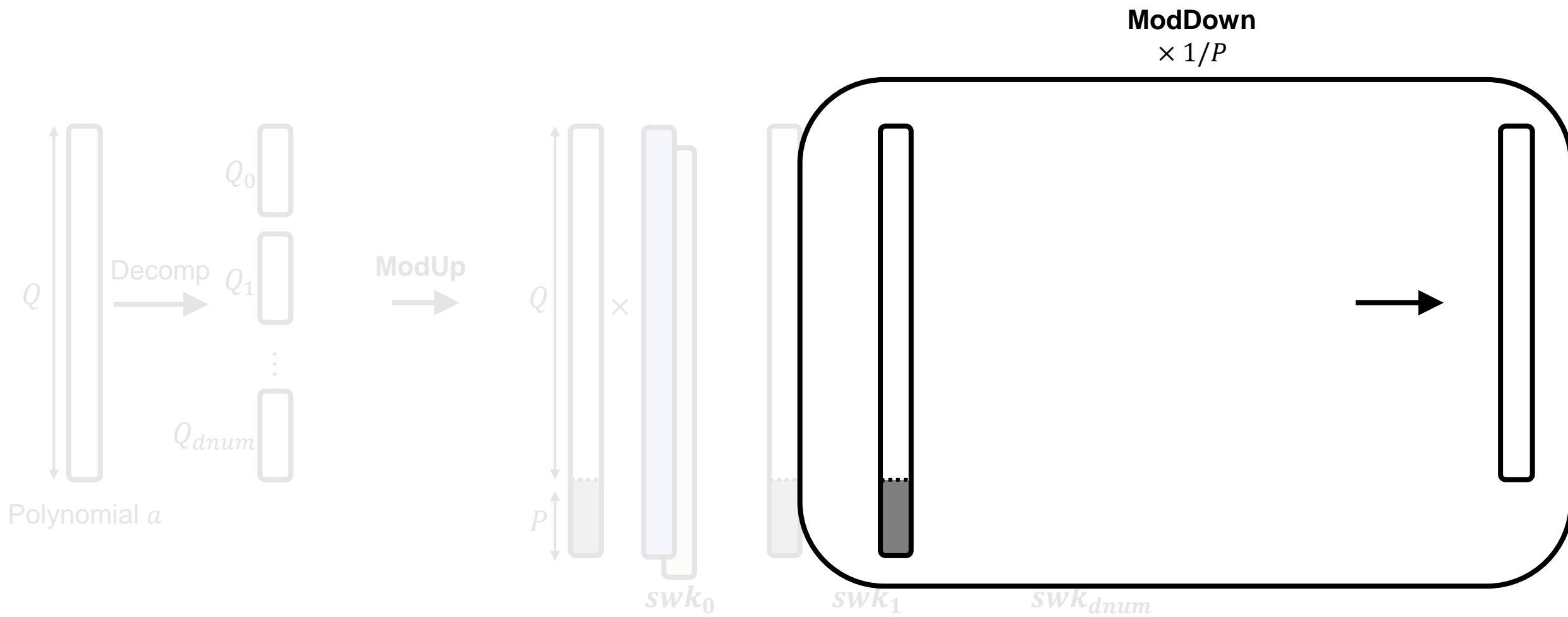
Intra-HE-operation Fusions

- *ModDown Fusion (MDF)*



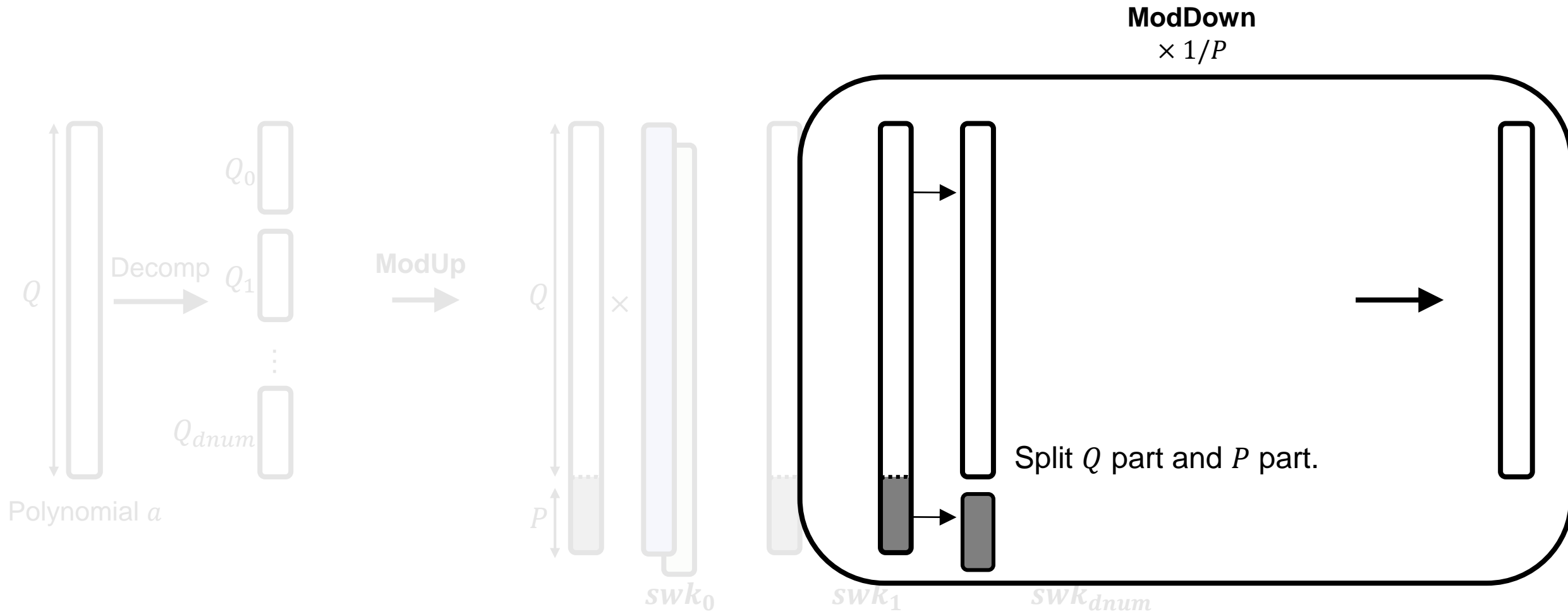
Intra-HE-operation Fusions

- *ModDown Fusion (MDF)*



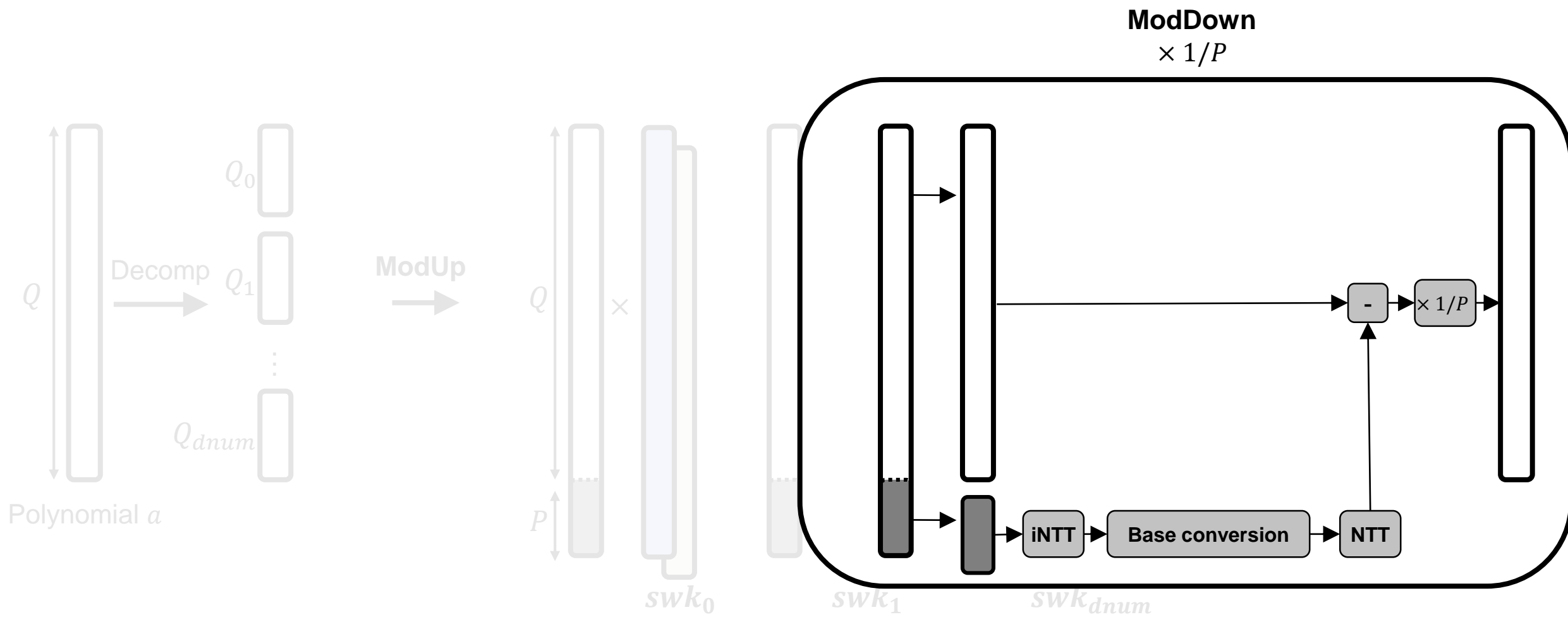
Intra-HE-operation Fusions

- *ModDown Fusion (MDF)*



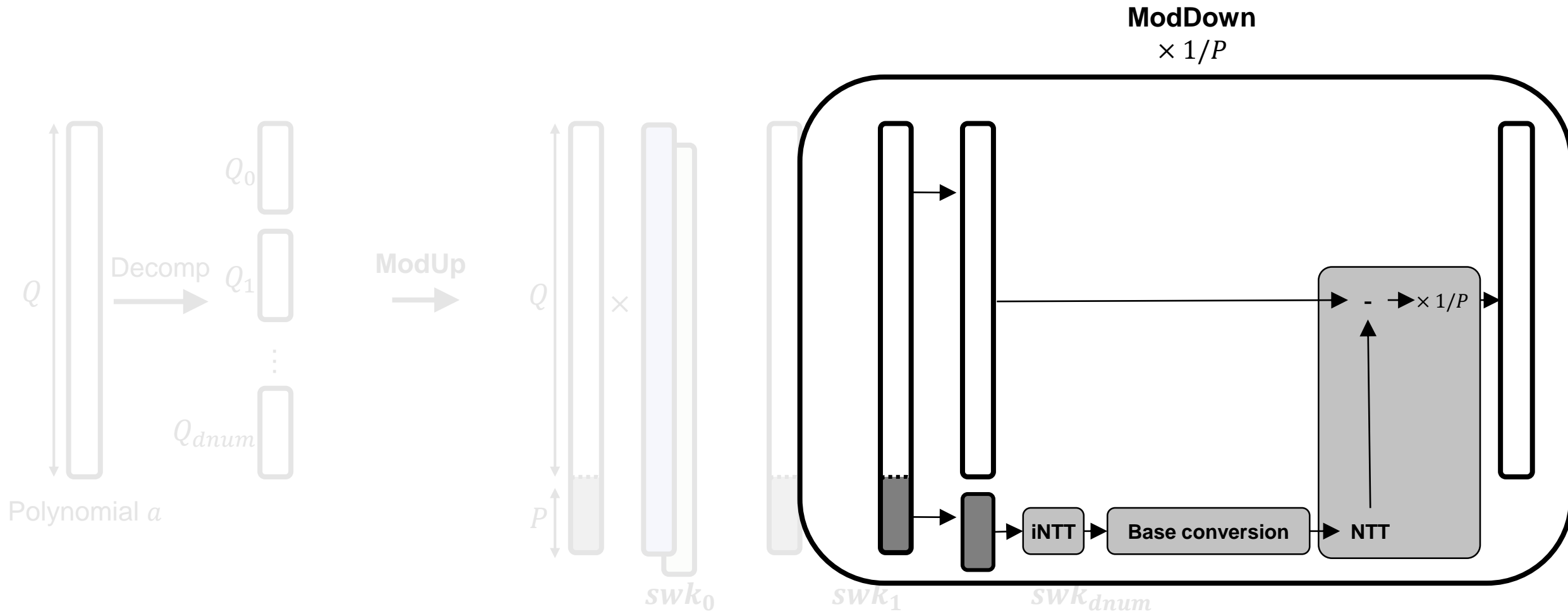
Intra-HE-operation Fusions

- ModDown Fusion (MDF)*



Intra-HE-operation Fusions

- *ModDown Fusion (MDF)*
 - fuses the element-wise, memory-bound kernels with NTTs.



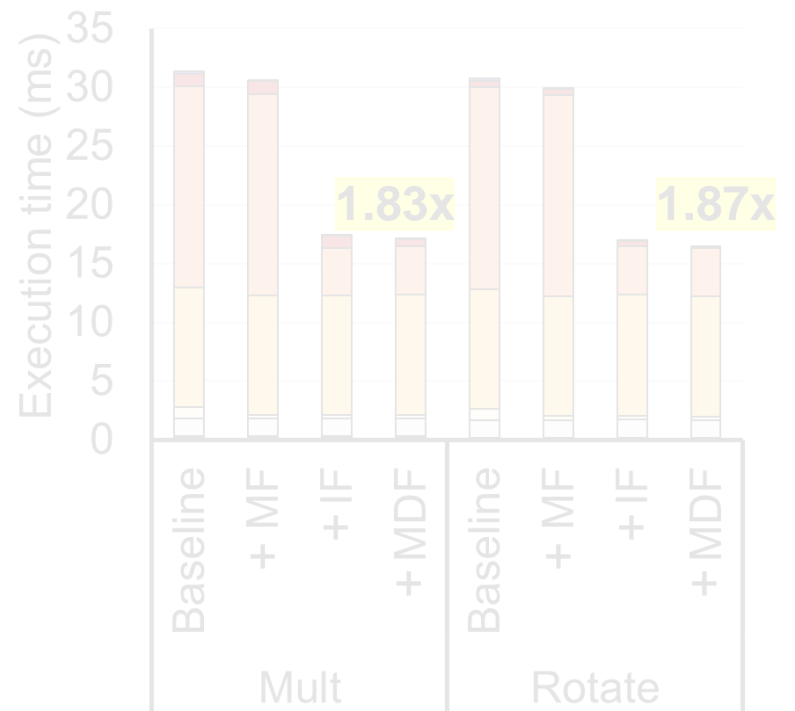
Experiment results

- Latency of Mult and Rotate after the optimizations.

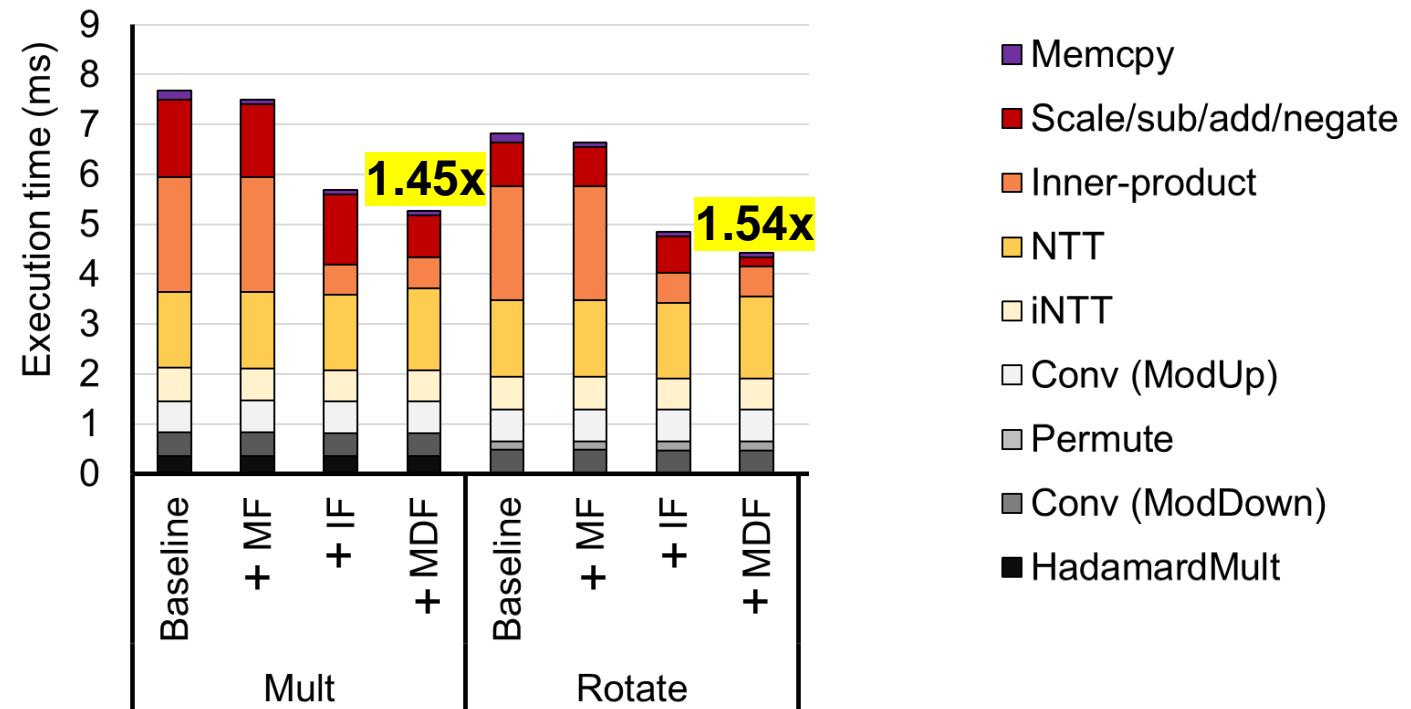


Experiment results

- Latency of Mult and Rotate after the optimizations.
- **7.02x** (7.96 vs. 55.88 ms) faster Mult over a prior GPU impl. [BHM+20] which uses **dnum = L**
 - (same security & level)



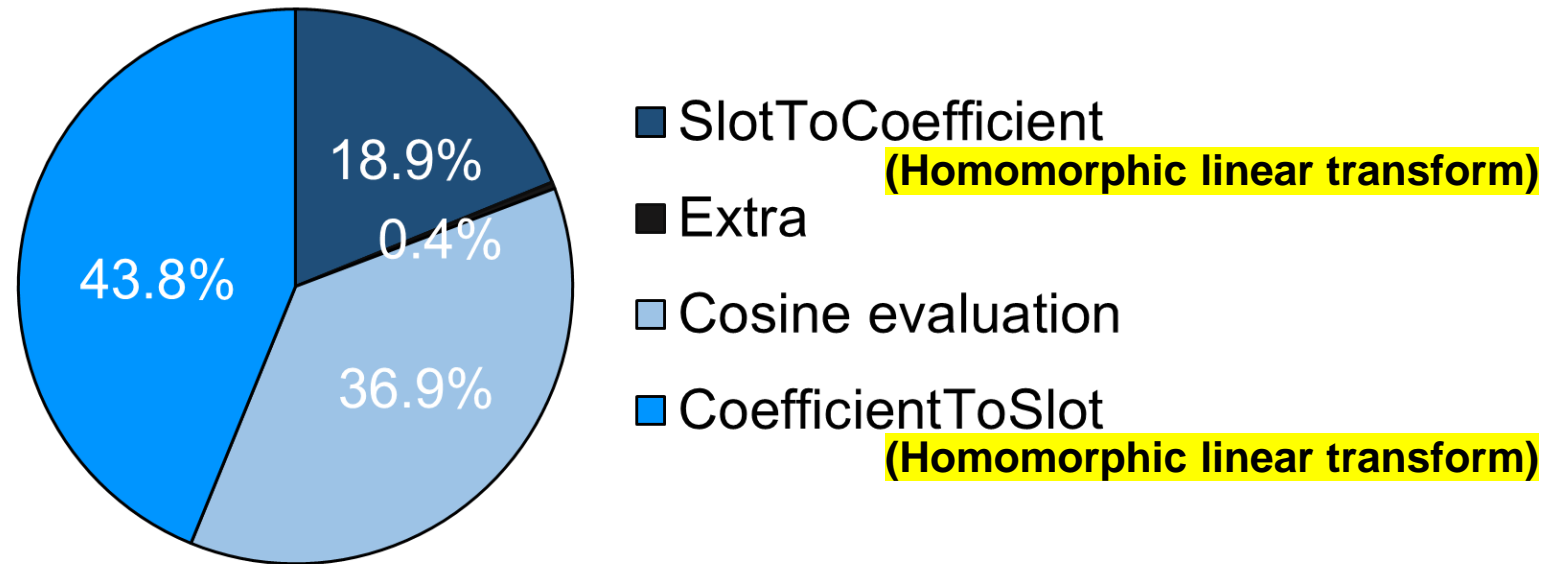
Latency with $dnum = L$ (maximum)



Latency with $dnum = 3$ [HK20]

Bootstrapping in CKKS

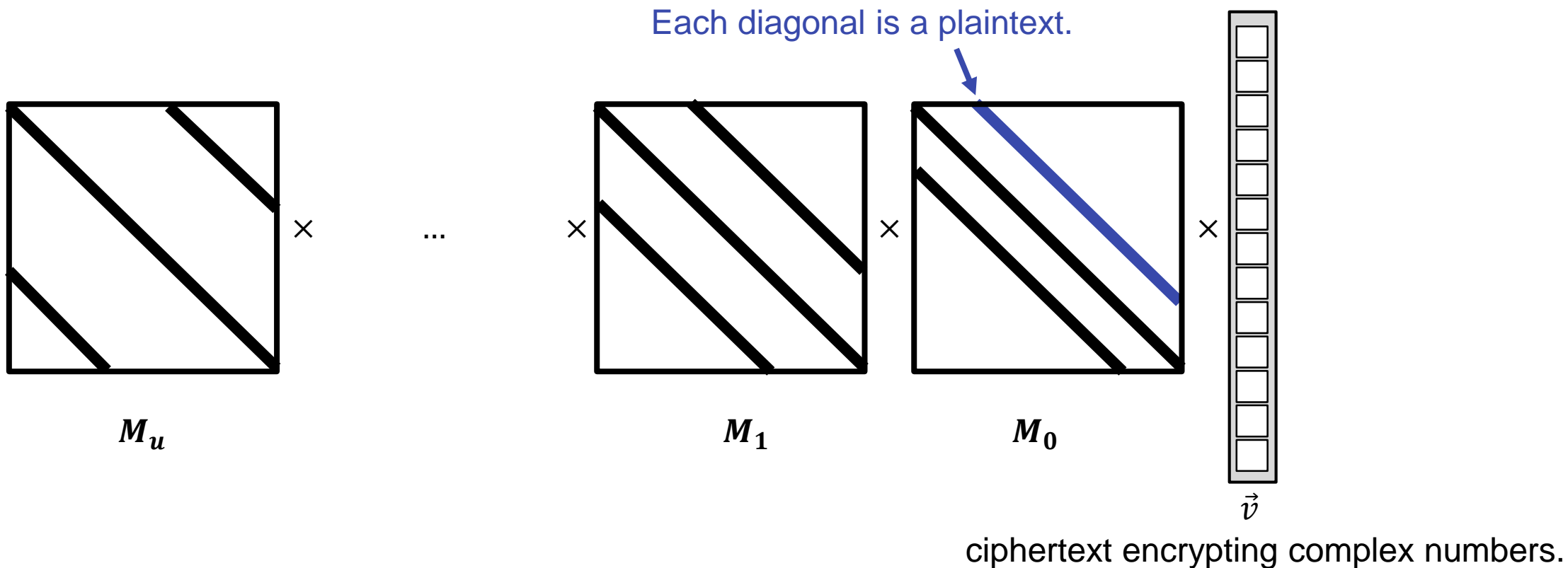
- Bootstrapping is an HE circuit, consisting of Mult, CMult, Rotate, ...
- Breakdown of a bootstrapping latency
 - Homomorphic linear transformations take up of 62.7%.



Bootstrapping latency breakdown on a GPU.

Homomorphic linear transformation in bootstrapping

- Represented as a series of matrix-vector multiplications, where
 - the matrix is sparse-diagonal matrices whose diagonals are plaintexts.
 - the vector is a ciphertext.
- How do we compute this?



Baby-step Giant-step (BSGS) algorithm

- A core algorithm for homomorphic linear transformation.
 - Used in matrix-vector multiplications in SlotToCoefficient/CoefficientToSlot
 - Setting $\ell, k \cong \sqrt{n}$, the # of rotations becomes $O(n) \rightarrow O(\sqrt{n})$.

$$\begin{aligned}
 \mathbf{M} \cdot \vec{v} &= \sum_{i=0}^n \text{diag}_i(\mathbf{M}) \odot \text{rot}_i(\vec{v}) \\
 &= \sum_{i=0}^{\ell} \sum_{j=0}^k \text{diag}_{ki+j}(\mathbf{M}) \odot \text{rot}_{ki+j}(\vec{v}) \\
 &= \sum_{i=0}^{\ell} \text{rot}_{ki} \left(\sum_{j=0}^k \underbrace{\text{rot}_{-ki}(\text{diag}_{ki+j}(\mathbf{M}))}_{\text{Precomputed plaintext}} \odot \text{rot}_j(\vec{v}) \right)
 \end{aligned}$$

The i-th diagonal element

Vector rotated by i

Ciphertext

Inter-HE-operation Fusion: Mult-and-add batching in BSGS

- Naïve multiplication-and-add requires multiple memory accesses on a temporal ciphertext.
 - Many reads and writes on ct .

$$\begin{aligned}
 ct &= m_1 \times ct_1 \\
 ct &= m_2 \times ct_2 + ct \\
 ct &= m_3 \times ct_1 + ct \\
 &\dots
 \end{aligned}$$

$$\begin{aligned}
 M \cdot \vec{v} &= \sum_{i=0}^n \text{diag}_i(M) \odot \text{rot}_i(\vec{v}) \\
 &= \sum_{i=0}^{\ell} \sum_{j=0}^k \text{diag}_{ki+j}(M) \odot \text{rot}_{ki+j}(\vec{v}) \\
 &= \sum_{i=0}^{\ell} \text{rot}_{ki} \left(\sum_{j=0}^k \underbrace{\text{rot}_{-ki}(\text{diag}_{ki+j}(M))}_{\text{Precomputed plaintext}} \odot \text{rot}_j(\vec{v}) \right)
 \end{aligned}$$

The i-th diagonal element

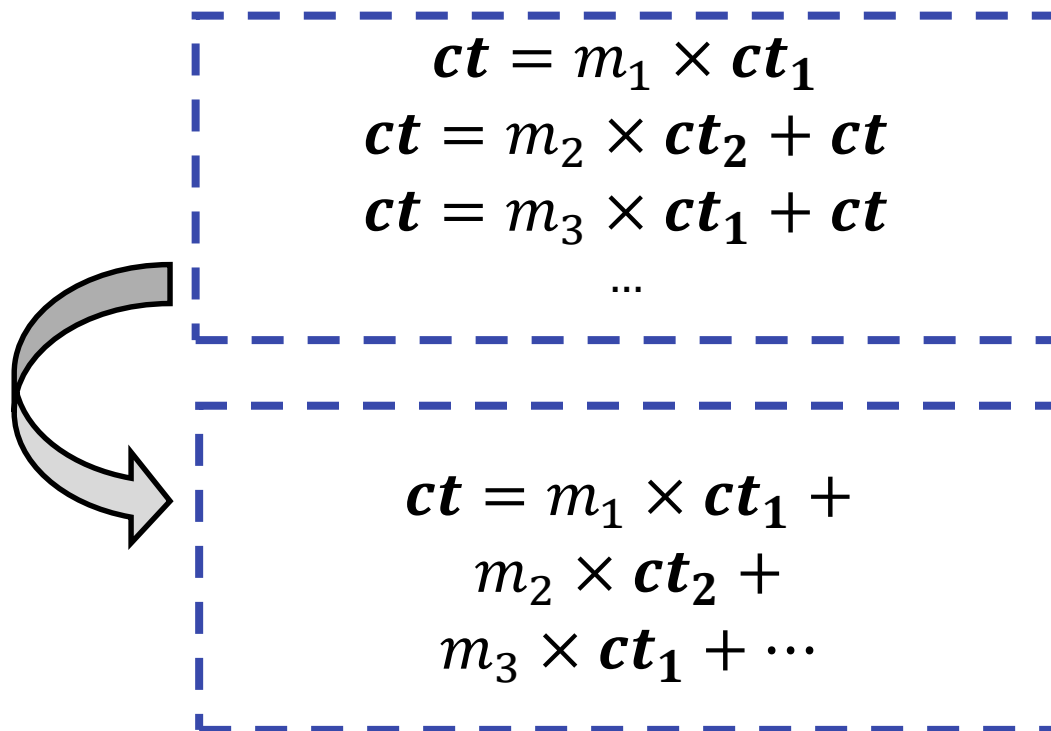
Vector rotated by i

Ciphertext

Precomputed plaintext

Inter-HE-operation Fusion: Mult-and-add batching in BSGS

- Naïve multiplication-and-add requires multiple memory accesses on a temporal ciphertext.
 - Many reads and writes on ct .
- Fusing multiple mul-and-add operations
 - removes most of read and writes on ct .



$$\begin{aligned}
 M \cdot \vec{v} &= \sum_{i=0}^n \text{diag}_i(M) \odot \text{rot}_i(\vec{v}) \\
 &= \sum_{i=0}^{\ell} \sum_{j=0}^k \text{diag}_{ki+j}(M) \odot \text{rot}_{ki+j}(\vec{v}) \\
 &= \sum_{i=0}^{\ell} \text{rot}_{ki} \left(\sum_{j=0}^k \underbrace{\text{rot}_{-ki}(\text{diag}_{ki+j}(M))}_{\text{Precomputed plaintext}} \odot \text{rot}_j(\vec{v}) \right)
 \end{aligned}$$

The i-th diagonal element

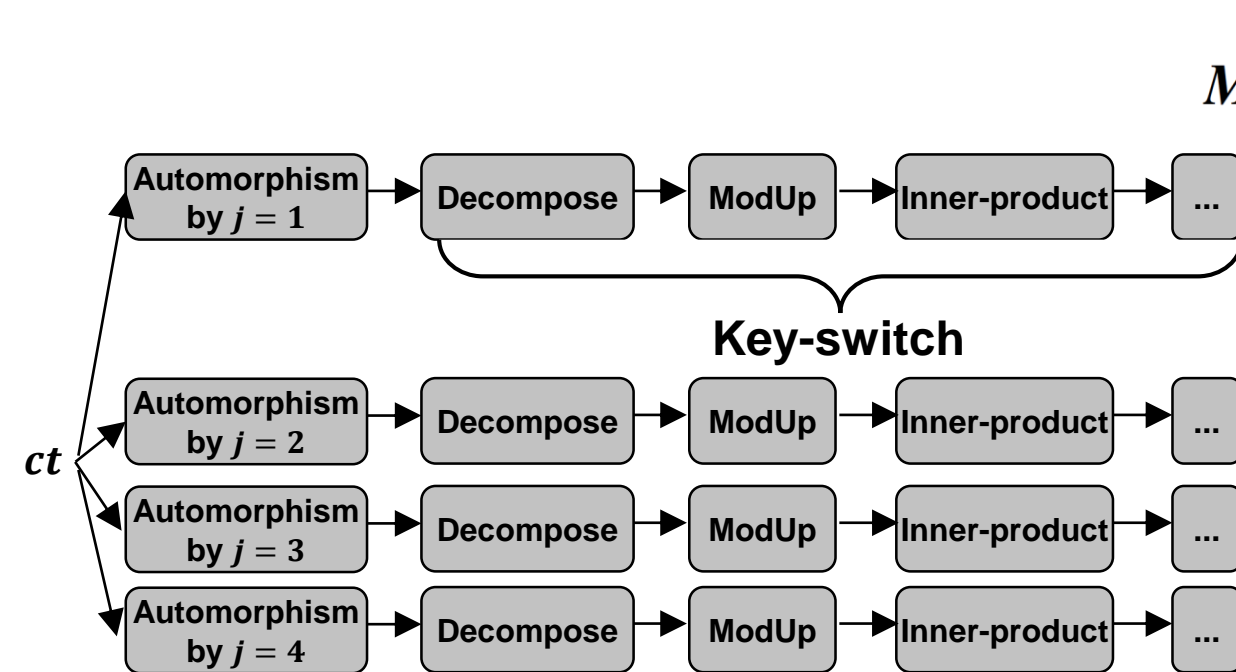
Vector rotated by i

Ciphertext

Precomputed plaintext

Hoisting [HS18]

- We also evaluated an optimization in [HS18]
- When rotating a single ciphertext multiple times, we save some ModUp computations.
 - Precomputes ModUp first, then permutes by each rotation index.



$$\mathbf{M} \cdot \vec{v} = \sum_{i=0}^n \text{diag}_i(\mathbf{M}) \odot \text{rot}_i(\vec{v})$$

$$= \sum_{i=0}^{\ell} \sum_{j=0}^k \text{diag}_{ki+j}(\mathbf{M}) \odot \text{rot}_{ki+j}(\vec{v})$$

$$= \sum_{i=0}^{\ell} \text{rot}_{ki} \left(\sum_{j=0}^k \text{rot}_{-ki}(\text{diag}_{ki+j}(\mathbf{M})) \odot \boxed{\text{rot}_j(\vec{v})} \right)$$

Performs rotations with different indices.

Experiment results

- Bootstrapping latency

	Latency (ms)				
Parameter set (N, L, dnum, λ)	Baseline	Intra-HE fusion	Mult-and-add Batching	Hoisting [HS18]	Speedup (vs. single-thread CPU)
$(2^{16}, 34, 5, 106)$	428.94	377.78	351.09	328.25	242x
$(2^{17}, 29, 3, 173)$	719.87	623.92	568.2	526.96	257x

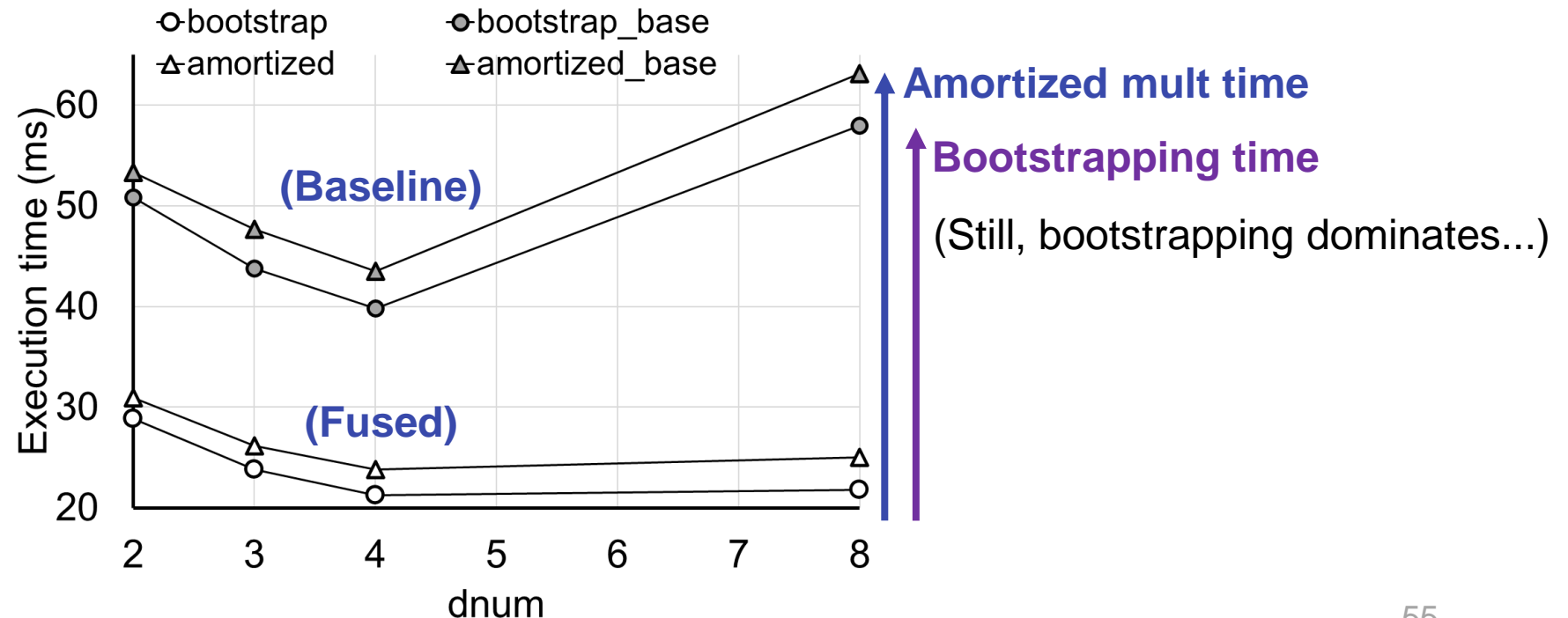
- End-to-end evaluation: training a binary classification (logistic regression) model
 - Compared to 8-threaded CPU case, GPU implementation exhibits a speedup of 40x in total.

Amortized Mult-time

- We suggest a new metric for performance considering bootstrapping:

$$\text{Amortized mult time} = \frac{\text{Bootstrapping time} + \text{multiplication times}}{\text{\# of multiplications after bootstrapping}}$$

- Our optimizations largely amortize the memory bottleneck.



Summary

- Based on our performance analysis, we accelerated binary CKKS (HEAAN) with
 - AVX-512 & GPU implementation with microarchitecture-aware optimizations.
- Further, we accelerated the recent Full-RNS CKKS using GPU, including
 - the first implementation of bootstrapping
 - analysis on the severe memory-bandwidth bottlenecks
 - applying memory-centric optimizations, giving
 - = 257x of speedup (vs. single-threaded CPU),
 - = application-level speedups (40x vs. 8-threaded CPU in logistic regression)

Reference

- [HK20] Han et al., Better Bootstrapping for Approximate Homomorphic Encryption, CT-RSA, 2020
- [BHM+20] Ahmad Al Badawi, Louie Hoang, Chan Fook Mun, Kim Laine, and Khin Mi Mi Aung. Privft: Private and fast text classification with homomorphic encryption, IEEE Access, 2020
- [CLP17] Hao Chen, Kim Laine, and Rachel Player. Simple encrypted arithmetic library - SEAL v2.1. In FC 2017 Workshops, volume 10323 of LNCS, pages 3–18. Springer, Heidelberg, April 2017.
- [HS14] Shai Halevi and Victor Shoup. Algorithms in HElib. In CRYPTO 2014, Part I, volume 8616 of LNCS, pages 554–571. Springer, Heidelberg, August 2014.
- [CHK+19] Jung Hee Cheon, Kyoohyung Han, Andrey Kim, Miran Kim, and Yongsoo Song. A full RNS variant of approximate homomorphic encryption. In SAC 2018, volume 11349 of LNCS, pages 347–368. Springer, Heidelberg, August 2019.
- [BPA+19] Ahmad Al Badawi, Yuriy Polyakov, Khin Mi Mi Aung, Bharadwaj Veeravalli, and Kurt Rohloff. Implementation and performance evaluation of RNS variants of the BFV homomorphic encryption scheme. IEEE Transactions on Emerging Topics in Computing, 2019.
- [KJPA20] Sangpyo Kim, Wonkyung Jung, Jaiyoung Park, and Jung Ho Ahn. Accelerating number theoretic transformations for bootstrappable homomorphic encryption on gpus. In IEEE International Symposium on Workload Characterization, IISWC 2020, Beijing, China, October 27-30, 2020, pages 264–275. IEEE, 2020.
- [HS18] Shai Halevi and Victor Shoup. Faster homomorphic linear transformations in HElib. In CRYPTO 2018, Part I, volume 10991 of LNCS, pages 93–120. Springer, Heidelberg, August 2018.

감사합니다