# ATTACKS AGAINST THE CPA-D SECURITY OF EXACT FHE SCHEMES

## Damien Stehlé

May 25, 2024

Talk based on Eprint 2024/127
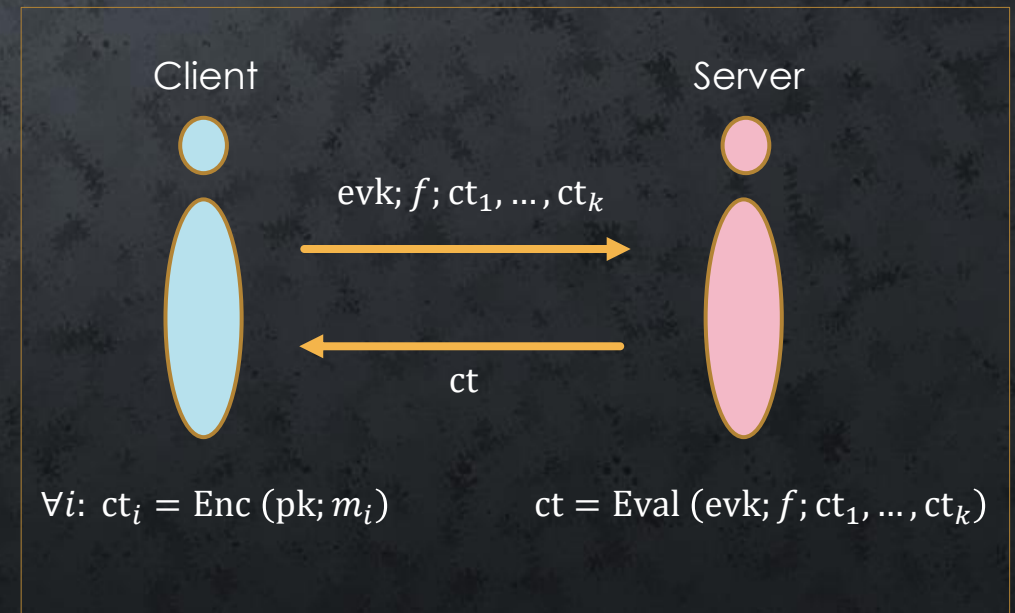Joint work with J. H. Cheon, H. Choe, A. Passelègue & E. Suvanto

HEAAN CRYPTO LAB

# FULLY HOMOMORPHIC ENCRYPTION

An FHE scheme consists of $(\text{KeyGen}, \text{Enc}, \text{Eval}, \text{Dec})$:

- KeyGen $\rightarrow$ $(\text{sk}, \text{pk}, \text{evk})$
- Enc $(\text{pk}; m)$ $\rightarrow$ $\text{ct}$
- Eval $(\text{evk}; f; \text{ct}_1, \dots, \text{ct}_k)$ $\rightarrow$ $\text{ct}$
- Dec $(\text{sk}; \text{ct})$ $\rightarrow$ $m$

$$\forall f, m_1, \dots, m_k :$$

$$\text{Dec}\left(\text{Eval}\left(f;\ \text{Enc}(m_1), \dots, \text{Enc}(m_k)\right)\right) =\ f(m_1, \dots, m_k)$$

Client        Server

$$\text{evk}; f; \text{ct}_1, \dots, \text{ct}_k$$

$$\text{ct}$$

$$\forall i:\ \text{ct}_i = \text{Enc}(\text{pk}; m_i) \qquad \text{ct} = \text{Eval}(\text{evk}; f; \text{ct}_1, \dots, \text{ct}_k)$$

# MAIN FHE SCHEMES

| | Plaintext space | Basic operations | Ctxt format |
|---|---|---|---|
| BFV/BGV (2012) | $\left(\mathbb{F}_{p^k}\right)^{N/k}$ | Add & Mult in // $\mathbb{F}_{p^k}$-automorph. in // Slot rotate | RLWE |
| DM/CGGI (2015) | $\{0,1\}$ | Binary gates | LWE (and RLWE internally) |
| CKKS (2017) | $\mathbb{C}^{N/2}$ | Add & Mult in // Conj in // Slot rotate | RLWE |

# MAIN FHE SCHEMES

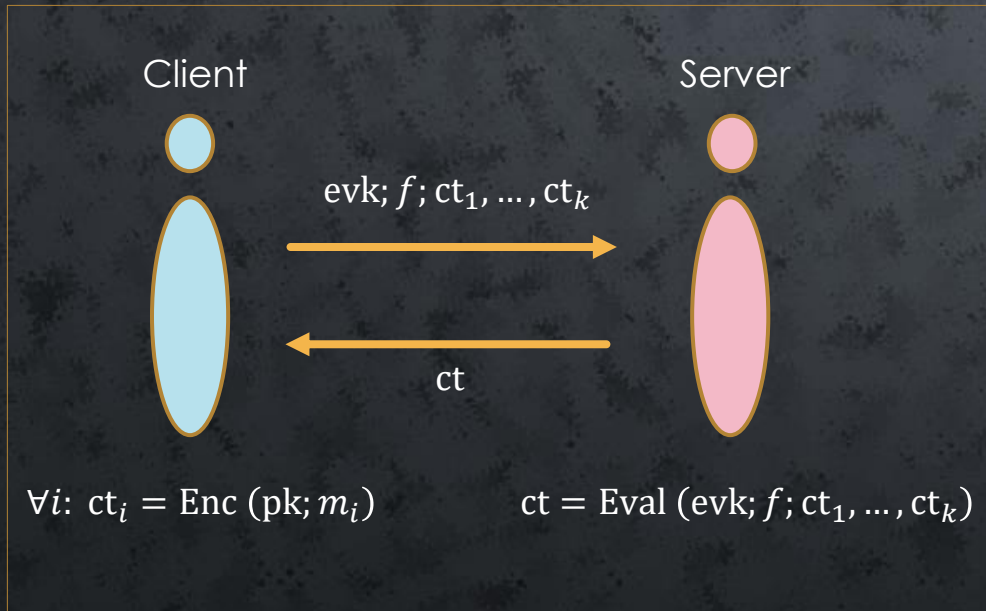| | Plaintext space | Basic operations | Ctxt format | |
|---|---|---|---|---|
| BFV/BGV (2012) | $\left(F_{p^k}\right)^{N/k}$ | Add & Mult in // $F_{p^k}$-automorph. in // Slot rotate | RLWE | **EXACT** |
| DM/CGGI (2015) | $\{0,1\}$ | Binary gates | LWE (and RLWE internally) | |
| CKKS (2017) | $\mathbb{C}^{N/2}$ | Add & Mult in // Conj in // Slot rotate | RLWE | **APPROXIMATE** (there is an exact mode for CKKS, see you on Thursday) |

$$\forall f, m_1, \ldots, m_k : \quad \mathrm{Dec}\left(\mathrm{Eval}\left(f;\ \mathrm{Enc}(m_1), \ldots, \mathrm{Enc}(m_k)\right)\right) \approx f(m_1, \ldots, m_k)$$

# FHE SECURITY

Client

Server

$\text{evk}; f; \text{ct}_1, \dots, \text{ct}_k$

ct

$\forall i: \text{ct}_i = \text{Enc}(\text{pk}; m_i)$

$\text{ct} = \text{Eval}(\text{evk}; f; \text{ct}_1, \dots, \text{ct}_k)$

Eve

**IND-CPA security**

**one cannot distinguish between encryptions of two different plaintexts**

# IND-CPA-D SECURITY

**IND-CPA security**

one cannot distinguish between encryptions of two different plaintexts

**IND-CPA-D security**

Same, but the attacker may ask for decryption of ciphertexts for which it is supposed to know the underlying plaintext

Adversary has $\mathrm{pk}$ and $\mathrm{evk}$

It can make queries:

- $\mathrm{Enc}\ (m)$ $\rightarrow$ $\mathrm{ct}$  // challenger knows the ptxts corresponding to all ctxts
- $\mathrm{ChallEnc}\ (m_0, m_1)$ $\rightarrow$ $\mathrm{ct}$  // challenge ctxts: $m_b$ is encrypted
- $\mathrm{Eval}\ (\mathrm{evk}; f; \mathrm{ct}_1, \dots, \mathrm{ct}_k)$ $\rightarrow$ $\mathrm{ct}$  // for $\mathrm{ct}_1, \dots, \mathrm{ct}_k$ in the databasis
- $\mathrm{Dec}\ (\mathrm{sk}; \mathrm{ct})$ $\rightarrow$ $m$  // for $\mathrm{ct}$ in the databasis
  **if the corresponding plaintext does not depend on $b$**

Adversary guesses $b$

# THE TOPIC OF THIS TALK

B. Li, D. Micciancio: *On the security of homomorphic encryption on approximate numbers*. EUROCRYPT'21

**IND-CPA security**

one cannot distinguish between encryptions of two different plaintexts

**IND-CPA-D security**

Same, but the attacker may ask for decryption of ciphertexts for which it is supposed to know the underlying plaintext
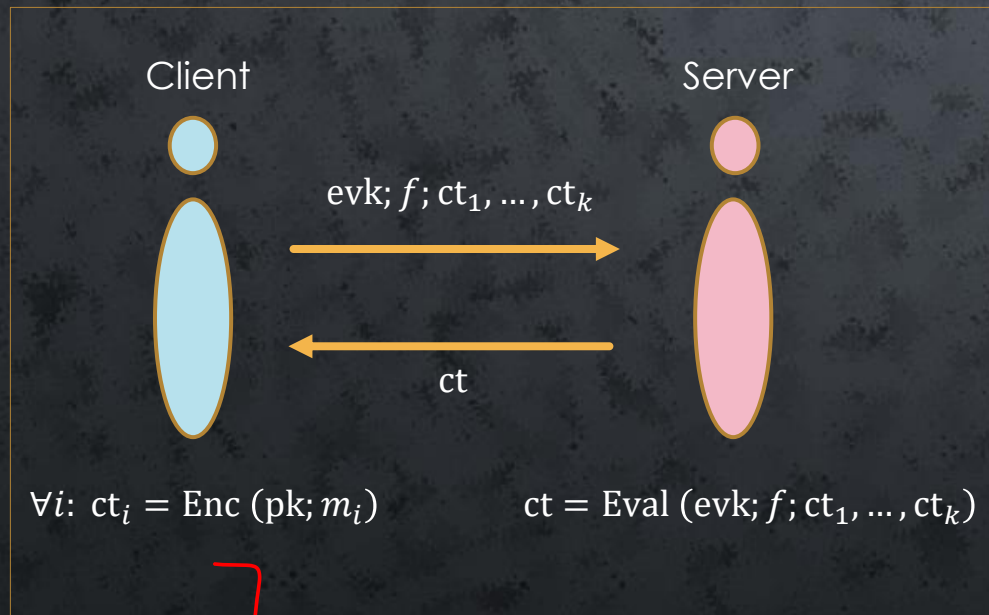
**IND-CPA-D attacks on exact schemes**

BGV / BFV
DM / CGGI
(Exact) CKKS

"when applied to standard (**exact**) encryption schemes, IND-CPA-D is perfectly equivalent to IND-CPA"

CKKS shouldn't be singled out

# HOW RELEVANT IS IND-CPA-D SECURITY?

Client

Server

$evk; f; ct_1, \dots, ct_k$

ct

Eve

$\forall i: ct_i = Enc(pk; m_i)$

$ct = Eval(evk; f; ct_1, \dots, ct_k)$
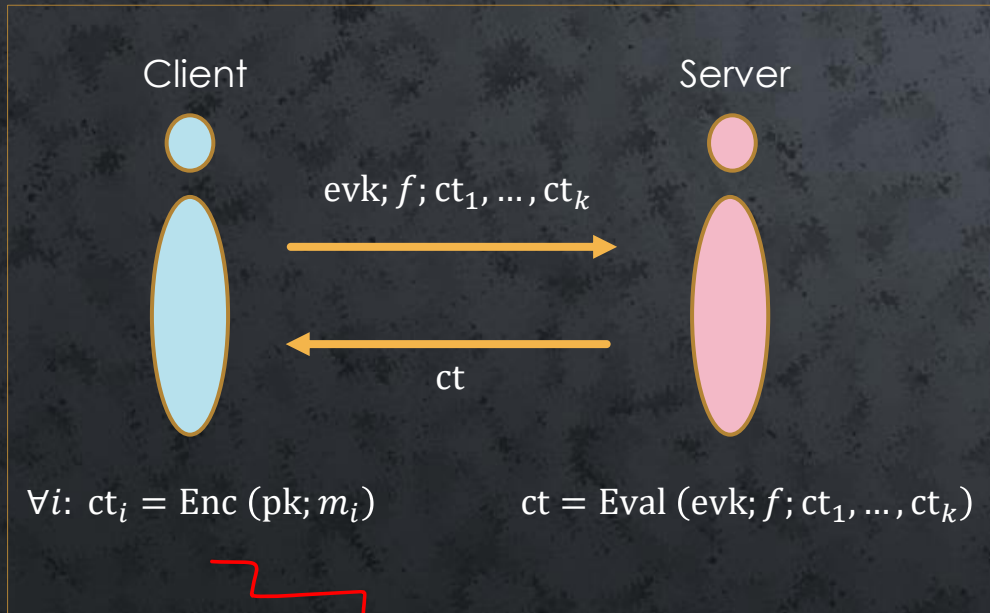
$m = Dec(sk; ct)$

**IND-CPA-D security**

**Same, but the attacker may ask for decryption of ciphertexts for which it is supposed to know the underlying plaintext**

If the computation is **confidential**, why would the client make the output of a confidential computation **public**?

# HOW RELEVANT IS IND-CPA-D SECURITY?

Client

Server

$$\text{evk}; f; \text{ct}_1, \dots, \text{ct}_k$$

ct

$\forall i: \text{ct}_i = \text{Enc}(\text{pk}; m_i)$

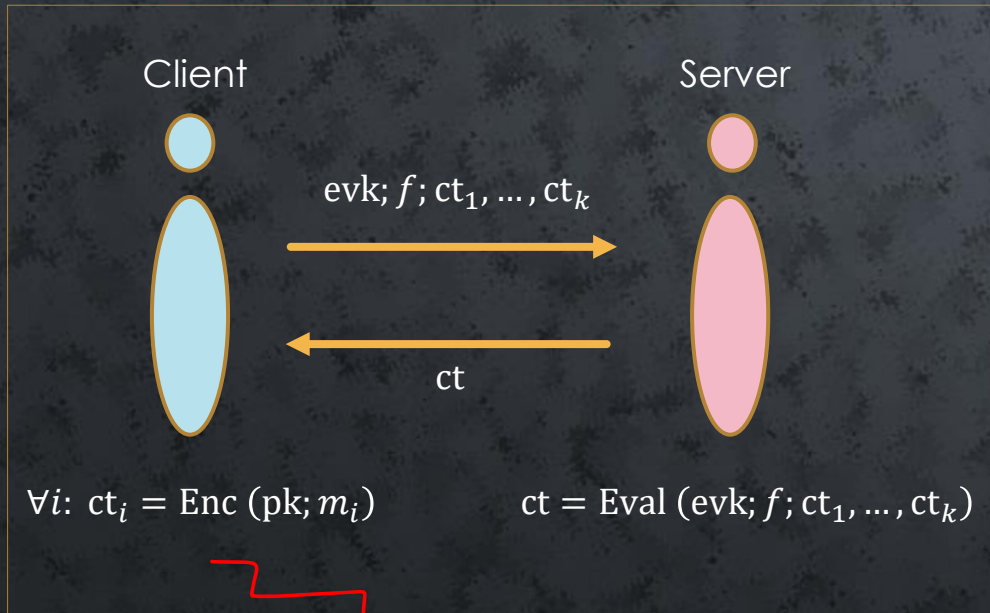$\text{ct} = \text{Eval}(\text{evk}; f; \text{ct}_1, \dots, \text{ct}_k)$

"Dec (sk; ct) is weird, restart!"

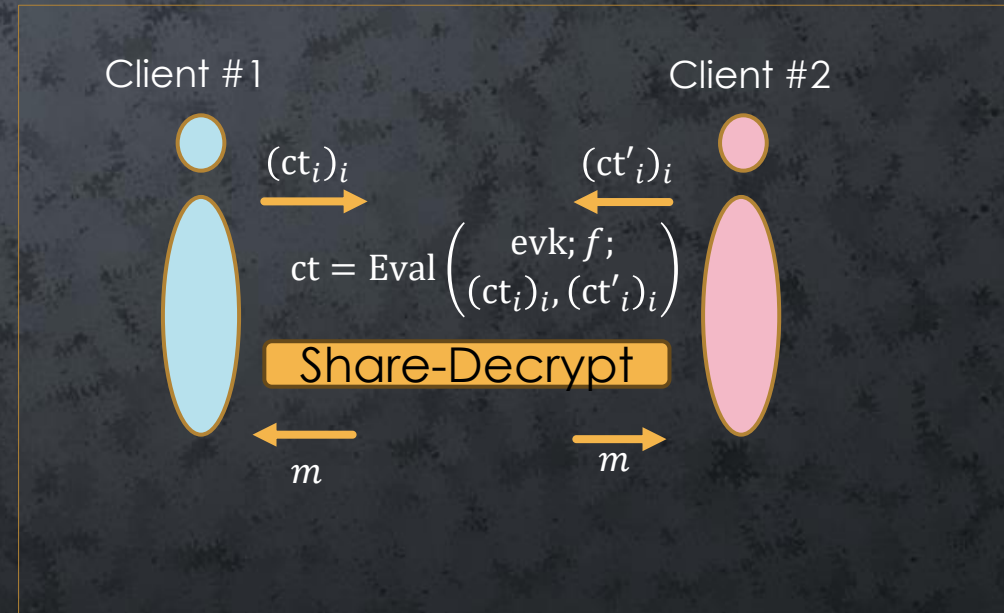**Weak variant of CVA security**

If the result is weird,
the client could ask to redo the computation

# HOW RELEVANT IS IND-CPA-D SECURITY?

Client    Server

$\text{evk}; f; \text{ct}_1, \ldots, \text{ct}_k$

ct

$\forall i: \text{ct}_i = \text{Enc}(\text{pk}; m_i)$    $\text{ct} = \text{Eval}(\text{evk}; f; \text{ct}_1, \ldots, \text{ct}_k)$

"Dec (sk; ct) is weird, restart!"

Client #1    Client #2

$(\text{ct}_i)_i$    $(\text{ct}'_i)_i$

$\text{ct} = \text{Eval}\begin{pmatrix} \text{evk}; f; \\ (\text{ct}_i)_i, (\text{ct}'_i)_i \end{pmatrix}$

Share-Decrypt

$m$    $m$

**Weak variant of CVA security**

If the result is weird,
the client could ask to redo the computation

**Threshold FHE**

sk  is shared across several clients
they collaborate to decrypt
and they all get to know the result

# ROADMAP

1- Motivation

**2- Attacks against CKKS**

3- IND-CPA-D versus IND-CPA for exact schemes
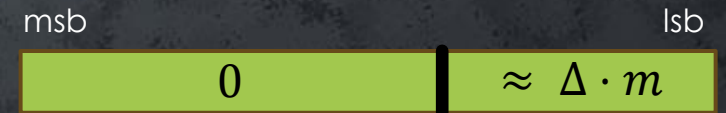
4- An attack against BFV/BGV addition

5- Attacks against bootstrapping algorithms

6- Concluding remarks

# REMINDERS ON CKKS

**Plaintext space**:   vectors of $\mathbb{C}^{N/2}$    (up to some precision)

- add in //
- multiply in //

msb                                                                                    lsb

| $0$ | $\approx \Delta \cdot m$ |
|---|---|

A **ciphertext** is of the form $(a, b) \in R_q^2$   s.t.    $a \cdot s + b \approx \Delta \cdot m$

- $s \in R_q$ is the secret key
- $m$ is the  (encoded)  plaintext

- $\Delta$ is the scaling factor (precision)
- $R_q = \mathbb{Z}_q[x] / x^N + 1$

To **decrypt**:    $(a, b) \mapsto (a \cdot s + b \bmod q) / \Delta$

# THE LI-MICCIANCIO ATTACK

To decrypt: $(a, b) \mapsto (a \cdot s + b \mod q) / \Delta$

Encrypt 0 and decrypt it:

=> We know $(a, b)$ and $a \cdot s + b \mod q$

=> This reveals $s$

**Key recovery**

# A COUNTERMEASURE

Noise flooding: $\quad (a, b) \mapsto (a \cdot s + b \mod q) / \Delta + e$

1- Bound the contributions of all errors (due to encryption and evaluation), for all possible inputs

2- Add to the decrypted value a noise $e$ that is $\geq 2^{\lambda/2}$ larger

**Security**

The output is simulatable from the knowledge of the expected ptxt

# NECESSITY OF LARGE FLOODING

Noise flooding: $(a, b) \mapsto (a \cdot s + b \mod q) / \Delta \; + \; e$

If the noise is smaller, then there is an attack

$$f: \quad x_1, \ldots, x_{2k} \quad \mapsto \quad \begin{array}{c} x_1^2 + \cdots + x_k^2 \\ -x_{k+1}^2 - \cdots - x_{2k}^2 \end{array}$$

$(0, \ldots, 0)$ and $(1, \ldots, 1)$ give the same result
But the noise for $(1, \ldots, 1)$ is larger

(multiplication noise grows with plaintext)

If the flooding is too small, we can distinguish

**Distinguishing attack**

# ROADMAP

## Passive Security

- IND-CPA security is typically sufficient to achieve passive security (for data privacy) for **exact** FHE schemes, including BGV, BFV, DM, and CGGI

- IND-CPA security is not sufficient for **approximate** FHE schemes
  - Li and Micciancio showed that CKKS is not secure if access to a decryption oracle is provided, i.e., when the decryption result is shared with parties that do not have the secret key [LM21]
  - They proposed a new definition IND-CPA$^D$ that provides access to encryption, evaluation, and decryption oracles

**(Borrowed from a talk by Y. Polyakov, given at NIST)**

1- Motivation

2- Attacks against CKKS

**3- IND-CPA-D versus CPA-D for exact schemes**

4- An attack against BFV/BGV addition

5- Attacks against bootstrapping algorithms

6- Concluding remarks

# CPA / CPA-D

**Assume the scheme is exact**

The decryption queries do not help the adversary:

For any valid decryption query    (i.e., the corresponding ptxt does not depend on the challenge $b$),
the adversary already knows the underlying ptxt

# CPA / CPA-D

**Assume the scheme is exact**

The decryption queries do not help the adversary:

For any valid decryption query (i.e., the corresponding ptxt does not depend on the challenge $b$),
the adversary already knows the underlying ptxt

> **Caveat**
> **The above requires perfect correctness**

Let $p_{\mathrm{err}}$ be the maximum over all $f, m_1, \dots, m_k$ of the probability that

$$\mathrm{Dec}\left(\ \mathrm{Eval}\left(\ f;\ \mathrm{Enc}(m_1), \dots, \mathrm{Enc}(\ m_k)\right)\right) \neq\ f(m_1, \dots, m_k)$$

The equivalency still holds if $p_{\mathrm{err}}$ **is extremely small**

# (SEMI-)GENERIC ATTACK FOR INCORRECT SCHEMES

Let $p_{\mathrm{err}}$ be the maximum over all $f, m_1, \ldots, m_k$ of the probability that

$$\mathrm{Dec}\left(\mathrm{Eval}\left(f; \mathrm{Enc}(m_1), \ldots, \mathrm{Enc}(m_k)\right)\right) \neq f(m_1, \ldots, m_k)$$

Assume that the adversary knows $f, m_1, \ldots, m_k, m'_1, \ldots, m'_k$ s.t.
- $f, m_1, \ldots, m_k$ reaches $p_{\mathrm{err}}$
- $f, m'_1, \ldots, m'_k$ has much lower decryption error
- $f(m_1, \ldots, m_k) = f(m'_1, \ldots, m'_k)$

Then:
- request encryptions of $m_1, \ldots, m_k$ ($b = 0$) or $m'_1, \ldots, m'_k$ ($b = 1$)
- request evaluation of $f$
- request decryption

If there is an error, it's more likely that $m_1, \ldots, m_k$ were encrypted

**Distinguishing attack**

# CORRECTNESS IN PRACTICE

In practice (all / most libraries):
- Failure probability from $2^{-15}$ to $2^{-50}$
- It is derived from heuristic error analysis   (probabilities without randomness)

Why?
1) Leads to more efficient schemes
2) For the primary use-case of FHE, IND-CPA (passive) security suffices

**Next: how to exploit decryption errors to mount IND-CPA-D attacks on exact schemes!**

# ROADMAP

# REMINDERS ON BFV

**Plaintext space**: elements of $R_p = \mathbb{Z}_p[x] / x^N + 1$
- add in //

A **ciphertext** is of the form $(a, b) \in R_q^2$ s.t. $a \cdot s + b = \left(\frac{q}{p}\right) \cdot m + e$

- $s \in R_q$ is the secret key
- $m$ is the plaintext
- $e$ is the error
- $R_q = \mathbb{Z}_q[x] / x^N + 1$

To **decrypt**: $(a, b) \mapsto \left\lfloor (a \cdot s + b \bmod q) / \left(\frac{q}{p}\right) \right\rceil$

# AN ATTACK ON BFV

**Theory**

To get correctness,
bound the contributions of all errors
for all possible inputs

**Practice   (sometimes)**

Use heuristic bounds

$$\text{Noise}(ct_1 + ct_2) \approx \sqrt{\text{Noise}(ct_1)^2 + \text{Noise}(ct_2)^2}$$

# AN ATTACK ON BFV

**Theory**

To get correctness,
bound the contributions of all errors
for all possible inputs

**Practice (sometimes)**

Use heuristic bounds

$$\text{Noise}(\text{ct}_1 + \text{ct}_2) \approx \sqrt{\text{Noise}(\text{ct}_1)^2 + \text{Noise}(\text{ct}_2)^2}$$

For $i = 1 \dots k$: $\quad x_{i+1} \leftarrow x_i + x_i$

Estimate noise $\approx 2^{k/2}$
        => The computation is deemed legitimate
Real noise $\approx 2^k$

Start with $\text{ct} = \text{Enc}(0)$

**Key recovery**

msb                                                    lsb

| $\text{ct}_0$ | $m = 0$ | 0 | $e$ |

| $\text{ct}_i$ | 0 | 0 | $2^i \cdot e$ | 0 |

| $\text{ct}_k$ | $2^k \cdot e$ | 0 |

# AN ATTACK ON BFV

Q. Guo, D. Nabokov, E. Suvanto, T. Johansson: *Key recovery attacks on approximate homomorphic encryption with non-worst-case noise flooding countermeasures.* USENIX'24

M. Checri, R. Sirdey, A. Boudguiga, J.-P. Bultel: On the practical CPAD security of "exact" and threshold FHE schemes and libraries. Eprint 2024/116

Adaptation of [GNSJ24] to BFV
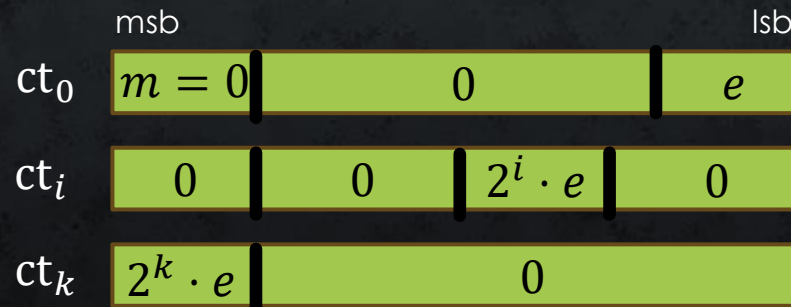
Concurrently obtained in [CSBB24]

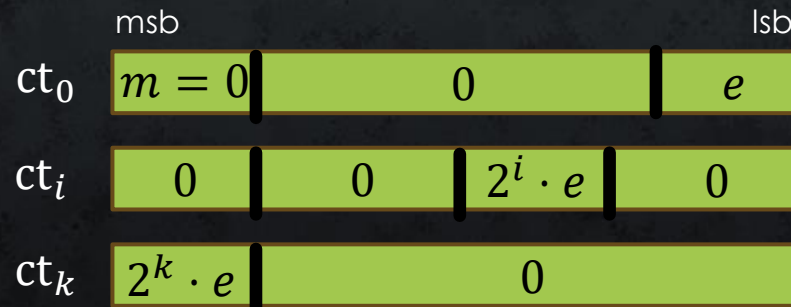For $i = 1 \dots k$:     $x_{i+1} \leftarrow x_i + x_i$

Estimate noise $\approx 2^{k/2}$
          => The computation is deemed legitimate
Real noise $\approx 2^k$

Start with $\text{ct} = \text{Enc}(0)$

msb                                              lsb

$\text{ct}_0$  | $m = 0$ | $0$ | $e$ |

$\text{ct}_i$  | $0$ | $0$ | $2^i \cdot e$ | $0$ |

$\text{ct}_k$  | $2^k \cdot e$ | $0$ |

**Key recovery**

# DOES IT WORK ON OPENFHE?

**OpenFHE:**

- claims to get IND-CPA-D security for CKKS,
- Has measures in place for correctness of exact schemes.

# DOES IT WORK ON OPENFHE?

**OpenFHE:**
- claims to get IND-CPA-D security for CKKS,
- Has measures in place for correctness of exact schemes.

We tested the attack on **OpenFHE**'s BFV,

With: $\quad N = 2^{12}, \quad p = 2^{16} + 1, \quad q = 2^{60}, \quad \sigma \approx 2^{7.41}$

Start with an encryption of 0, and iterate $k = 44$ times

Estimated error probability $\approx 2^{-2^{27.5}}$

But decryption gives the initial noise, and we recover $s$

Only additions => attack is instantaneous

# WHY DOES IT WORK ON OPENFHE?

| Practice (sometimes) | OpenFHE |
|---|---|
| Heuristic bounds | Triangular inequality |
| $\text{Noise}(ct_1 + ct_2) \approx \sqrt{\text{Noise}(ct_1)^2 + \text{Noise}(ct_2)^2}$ | $\text{Noise}(ct_1 + ct_2) \leq \text{Noise}(ct_1) + \text{Noise}(ct_2)$ |

But the attack **does** succeed!

# WHY DOES IT WORK ON OPENFHE?

<div style="border:2px solid; background:#f5b942">

**Practice   (sometimes)**

Heuristic bounds

$$\text{Noise}(\text{ct}_1 + \text{ct}_2) \approx \sqrt{\text{Noise}(\text{ct}_1)^2 + \text{Noise}(\text{ct}_2)^2}$$

</div>

<div style="border:2px solid; background:#f5b942">

**OpenFHE**

Triangular inequality

$$\text{Noise}(\text{ct}_1 + \text{ct}_2) \leq \text{Noise}(\text{ct}_1) + \text{Noise}(\text{ct}_2)$$

</div>

But the attack **does** succeed!

There is an error in the handling of addition error bounds in OpenFHE.

For $k$ additions, OpenFHE multiplies the error by $k$.

For $i = 1 \ldots k$:    $x_{i+1} \leftarrow x_i + x_i$                    $k$ additions but error grows as $2^k$
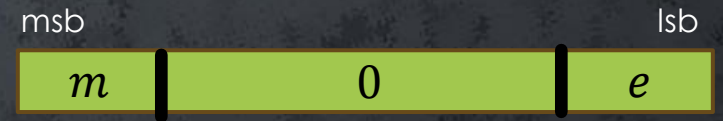
# ROADMAP

1- Motivation

2- Attacks against CKKS

3- IND-CPA-D versus IND-CPA for exact schemes

4- An attack against BFV/BGV addition

**5- Attacks against bootstrapping algorithms**

6- Concluding remarks

# REMINDERS ON DM/CGGI

**Plaintext space**:  elements of   $\{0,1\}$

- Binary gates



A **ciphertext** is of the form $(a, b) \in \mathbb{Z}_q^n \times \mathbb{Z}_q$   s.t.   $\langle a, s \rangle + b = \left(\frac{q}{8}\right) \cdot m + e$

- $s \in \mathbb{Z}_q^n$  is the secret key        • $e$ is the error
- $m$ is the plaintext bit

To **decrypt**:   $(a, b) \mapsto \left\lfloor (\langle a, s \rangle + b \mod q) / \left(\frac{q}{8}\right) \right\rceil$

# DM/CGGI BOOTSTRAPPING

LWE ctxt with key $s$
Modulo $q$

**ModSwitch** →

LWE ctxt with key $s$
Modulo $2N$

↑ **KeySwitch**

↓ **BlindRotate**

LWE ctxt with key $s'$
Modulo $q$

← **SampleExtract**

$RLWE_N$ ctxt with key $s'$
Modulo $q$

25

# DM/CGGI BOOTSTRAPPING

LWE ctxt with key $s$
Modulo $q$
<u>Noise variance</u>:  $\sigma_{br}^2 + \sigma_{ks}^2$

ModSwitch →

LWE ctxt with key $s$
Modulo $2N$
<u>Noise variance</u>:  $\sigma_{br}^2 + \sigma_{ks}^2 + \sigma_{ms}^2$

KeySwitch ↑

BlindRotate ↓

LWE ctxt with key $s'$
Modulo $q$
<u>Noise variance</u>:  $\sigma_{br}^2$

← SampleExtract

$RLWE_N$ ctxt with key $s'$
Modulo $q$
<u>Noise variance</u>:  $\sigma_{br}^2$

25

# DM/CGGI GATE BOOTSTRAPPING

Two LWE ctxts with key $s$
Modulo $q$
<u>Noise variance</u>:  $\sigma_{br}^2 + \sigma_{ks}^2$

Add and

ModSwitch

LWE ctxt with key $s$
Modulo $2N$
<u>Noise variance</u>:  $4\sigma_{br}^2 + 4\sigma_{ks}^2 + \sigma_{ms}^2$

KeySwitch

BlindRotate

# EXPLOITING DECRYPTION ERROR

Add and

ModSwitch

LWE ctxt with key $s$
Modulo $2N$
Noise variance: $4\sigma_{br}^2 + 4\sigma_{ks}^2 + \sigma_{ms}^2$

BlindRotate

- Gate bootstrapping fails
if the noise spills over the ptxt

- After ModSwitch is where noise is largest

- If gate bootstrapping fails,
then the ModSwitch error must be large

# EXPLOITING MODSWITCH ERROR

**ModSwitch**:     $\text{ct} \bmod q \;\mapsto\; \text{ct}' = \left\lfloor \left(\frac{2N}{q}\right) \cdot \text{ct} \right\rceil \bmod 2N$

$\langle \text{ct}, \text{sk} \rangle = e \;\Rightarrow\; \langle \text{ct}', \text{sk} \rangle = \langle e_{\text{rnd}}, \text{sk} \rangle + e \,,$     where   $e_{\text{rnd}}$   is known

A failure tells that   $\langle e_{\text{rnd}}, \text{sk} \rangle + e \geq \frac{2N}{16} \,,$  for a known $e_{\text{rnd}}$

Attack can be completed with statistical analysis

# IN PRACTICE

M. Dahl, D. Demmler, S. E. Kazdadi, A. Meyre, J.-B. Orfila, D. Rotaru, N. P. Smart, S. Tap, M. Walter: *Noah's ark: efficient threshold-FHE using noise flooding*. WAHC'23

We considered Zama's TFHE-rs

- For the default parameters, decryption error probability is $\approx 2^{-40}$

- We simulated that 256 decryption errors suffices

- Mounting the attack would take $\approx 2^{16}$ CPU years

- There are parameter sets with much poorer correctness

- The attack extends the [DDK+23] threshold-FHE scheme

# AN ATTACK ON CKKS BOOTSTRAPPING

CKKS BTS has 4 steps:
1. S2C
2. ModRaise
3. C2S
4. EvalMod

# AN ATTACK ON CKKS BOOTSTRAPPING

CKKS BTS has 4 steps:
1. S2C
2. ModRaise
3. C2S
4. EvalMod

Polynomial approximation to the mod-1 function, over a given number $2K + 1$ of periods.
- Higher $K$ => more costly
- Smaller $K$ => higher probability of error

# AN ATTACK ON CKKS BOOTSTRAPPING

CKKS BTS has 4 steps:
1. S2C
2. ModRaise
3. C2S
4. EvalMod

Polynomial approximation to the mod-1 function, over a given number $2K + 1$ of periods.
- Higher $K$ => more costly
- Smaller $K$ => higher probability of error

**Input of EvalMod is not in the approximation range => Output is nonsense**

**When that happens, we have an equation**

$$\langle x, \mathbf{sk} \rangle + \mathbf{e} \geq \textbf{\textit{bound}}, \quad \text{where } x \text{ is known.}$$

(like the DM/CGGI attack)

**Example: OpenFHE**
(claims INDCPA-D security for CKKS)

Probability of error ranges from $2^{-22}$ to $2^{-57}$

# ROADMAP

1- Motivation

2- Attacks against CKKS

3- IND-CPA-D versus IND-CPA for exact schemes

4- An attack against BFV/BGV

5- An attack against DM/CGGI

**6- Concluding remarks**

# TAKE-AWAY

**IND-CPA security:**

**one cannot distinguish between encryptions of two different plaintexts**

**IND-CPA-D security:**

**Same, but the attacker may ask for decryption of ciphertexts for which it is supposed to know the underlying plaintext**

**IND-CPA-D attacks on exact schemes**

BGV / BFV
DM / CGGI
(Exact) CKKS

**"when applied to standard (exact) encryption schemes, IND-CPA-D is perfectly equivalent to IND-CPA"**

# All competitive FHE schemes can suffer from IND-CPA-D attacks

# ATTACKS OF DIFFERENT NATURES

| Attack | Scheme | Decryption oracle or validity oracle? | Key recovery or distinguishing? |
|---|---|---|---|
| [LM21] | CKKS | Decryption | Key recovery |
| [LMSS22] | CKKS with limited decryption noise | Decryption | Distinguishing |
| [GNST24] | CKKS with heuristic error analysis | Decryption | Key recovery |
| Our work | FHE with imperfect correctness | Validity oracle | Distinguishing |
| Our work & [CSBB24] | BFV/BGV with heuristic error analysis | Can be adapted to validity oracle | Key recovery |
| Our work | DM/CGGI with large decryption error | Validity oracle | Key recovery |
| Our work | Exact CKKS | Validity oracle | Key recovery |

# ATTACKS OF DIFFERENT NATURES

| Attack | Scheme | Decryption oracle or validity oracle? | Key recovery or distinguishing? |
|---|---|---|---|
| [LM21] | CKKS | Decryption | Key recovery |
| [LMSS22] | CKKS with limited decryption noise | Decryption | Distinguishing |
| [GNST24] | CKKS with heuristic error analysis | Decryption | Key recovery |
| Our work | FHE with imperfect correctness | Validity oracle | Distinguishing |
| Our work & [CSBB24] | BFV/BGV with heuristic error analysis | Can be adapted to validity oracle | Key recovery |
| Our work | DM/CGGI with large decryption error | Validity oracle | Key recovery |
| Our work | Exact CKKS | Validity oracle | Key recovery |

**The situation is arguably worse for exact schemes!**

# COUNTERMEASURES

For all schemes:
- **tiny failure probability**
- **no heuristic** noise analysis

For (approximate) CKKS:
- **high-precision** computation
- followed by **noise flooding**

efficiency

# COUNTERMEASURES

For all schemes:
- **tiny failure probability**
- **no heuristic** noise analysis

For (approximate) CKKS:
- **high-precision** computation
- followed by **noise flooding**

efficiency

And be very diligent with the implementation:

- IND-CPA:     be cautious about **KeyGen & Enc**
- IND-CPA-D:   be cautious about **KeyGen, Enc, Eval & Dec**

# QUESTIONS?