# SPIRAL: Fast High-Rate Single-Server Private Information Retrieval
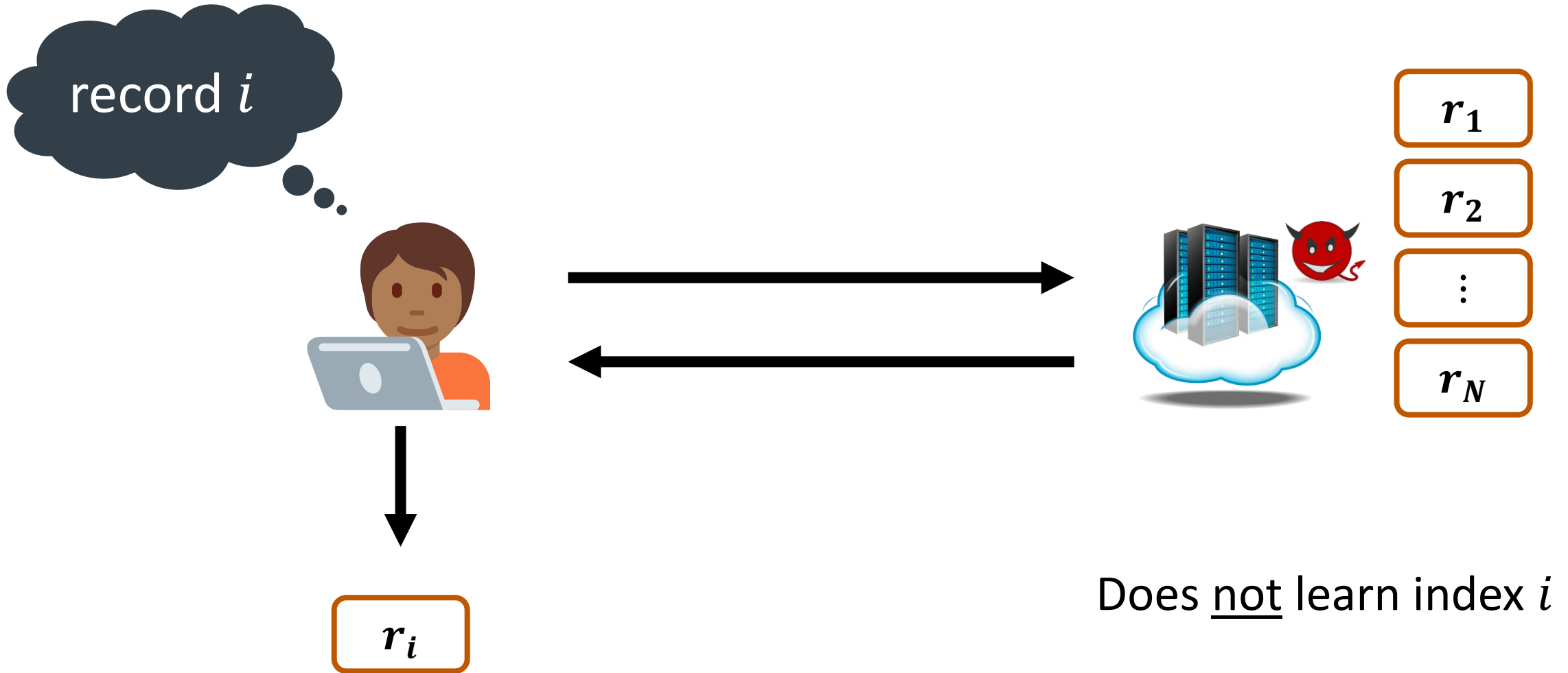
Samir Menon and <u>David Wu</u>
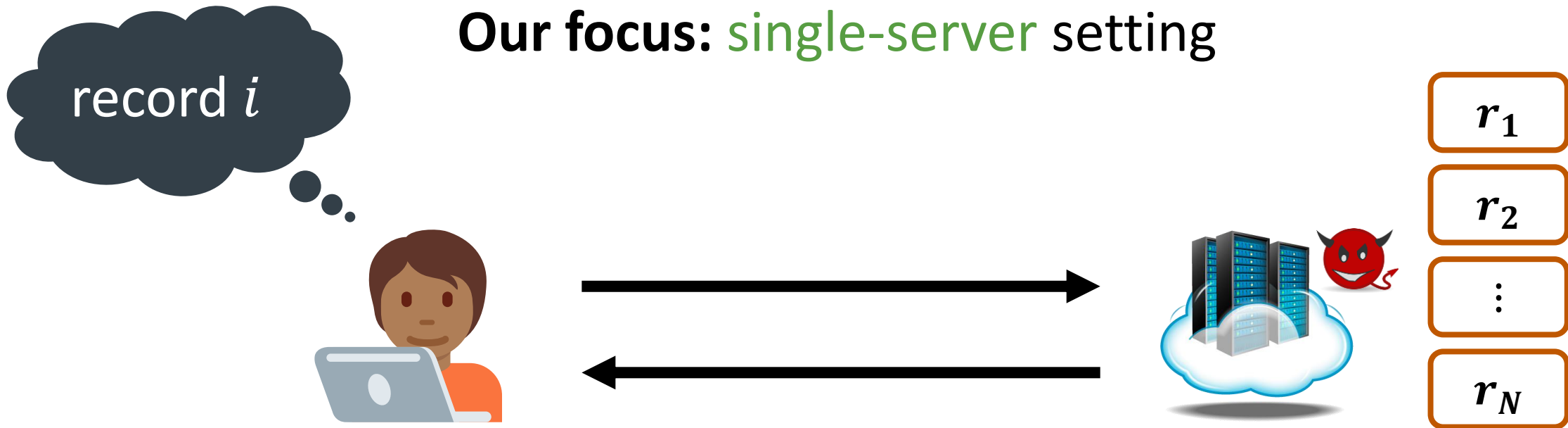
# Private Information Retrieval (PIR)

[CGKS95]



record $i$

$r_1$

$r_2$

$\vdots$

$r_N$

$r_i$

Does <u>not</u> learn index $i$

# Private Information Retrieval (PIR)

**Our focus:** single-server setting

record $i$

$r_1$

$r_2$

$\vdots$

$r_N$
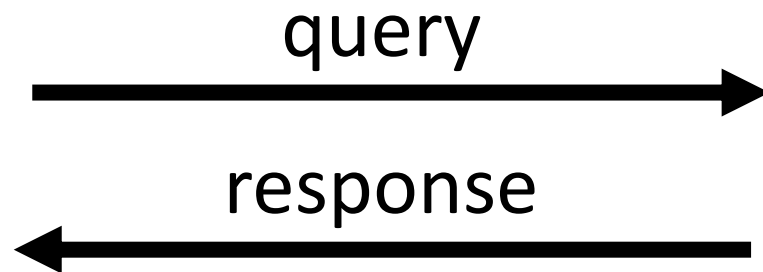
Basic building block in many privacy-preserving protocols

💬 Metadata-private messaging

🌐 Private DNS

👥 Contact discovery

😷 Private contact tracing

🔍 Safe browsing

🧭 Private navigation

# Efficiency Metrics

**1** **Query size**

query

response

**2** **Server Throughput**

$$\frac{\text{database size}}{\text{server computation time}}$$

*"measures how fast the server can respond as a function of database size"*

# Efficiency Metrics

**1** **Query size**

query →

response ←

Without preprocessing, server must perform a linear scan over the database

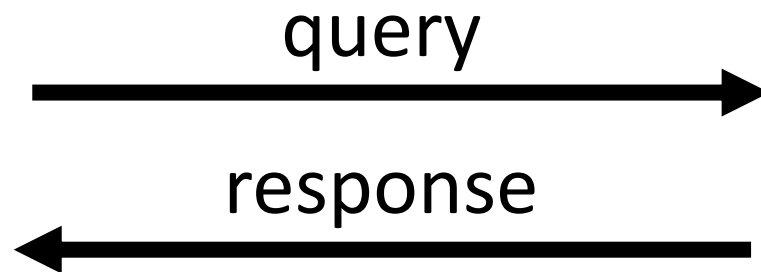**2** **Server Throughput**

$$\frac{\text{database size}}{\text{server computation time}}$$

*"measures how fast the server can respond as a function of database size"*

# Efficiency Metrics

Client generates a *reusable* set of public parameters

public parameters →

**4** **Public parameter size**

**1** **Query size**

query →

← response

**3** **Rate**

$$\frac{\text{record size}}{\text{response size}}$$

*"measures communication overhead in responses"*

**2** **Server Throughput**

$$\frac{\text{database size}}{\text{server computation time}}$$

*"measures how fast the server can respond as a function of database size"*

# The SPIRAL Family of PIR Protocols

Techniques to translate between FHE schemes enables new trade-offs in single-server PIR

Automatic parameter selection based on database configuration

## Base version of SPIRAL

**Query size:** 14 KB     4.5× smaller
**Rate:** 0.41     2.1× higher
**Throughput:** 333 MB/s     2.9× higher

(Database with $2^{14}$ records of size 100 KB)

**Cost:** 3.4× larger public parameters (17 MB)

## Streaming versions of SPIRAL

**Rate:** 0.81     3.4× smaller responses
**Throughput:** 1.9 GB/s     12.3× higher

## Best previous protocol:

**Rate:** 0.24
**Throughput:** 158 MB/s

# The SPIRAL Family of PIR Protocols

Techniques to translate between FHE sch[...]

Automatic parameter selection based on[...]

Higher throughput than running software AES over database
(Primary operation: 64-bit integer arithmetic)

## Base version of SPIRAL

| | | |
|---|---|---|
| **Query size:** | 14 KB | 4.5× smaller |
| **Rate:** | 0.41 | 2.1× higher |
| **Throughput:** | 333 MB/s | 2.9× higher |

(Database with $2^{14}$ records of size 100 KB)

**Cost:** 3.4× larger public parameters (17 MB)

## Streaming versions of SPIRAL

| | | |
|---|---|---|
| **Rate:** | 0.81 | 3.4× smaller responses |
| **Throughput:** | 1.9 GB/s | 12.3× higher |

## Best previous protocol:

| | |
|---|---|
| **Rate:** | 0.24 |
| **Throughput:** | 158 MB/s |

# The SPIRAL Family of PIR Protocols

Techniques to translate between FHE schemes enable

Automatic parameter selection based on database co

Cost of privately streaming a 2 GB movie from database of $2^{14}$ movies estimated to be $1.9\times$ more expensive than <u>no-privacy</u> baseline (based on AWS compute costs)

## Base version of SPIRAL

**Query size:** 14 KB     $4.5\times$ smaller
**Rate:** 0.41     $2.1\times$ higher
**Throughput:** 333 MB/s     $2.9\times$ higher

(Database with $2^{14}$ records of size 100 KB)

**Cost:** $3.4\times$ larger public parameters (17 MB)

## Streaming versions of SPIRAL

**Rate:** 0.81     $3.4\times$ smaller responses
**Throughput:** 1.9 GB/s     $12.3\times$ higher

## Best previous protocol:

**Rate:** 0.24
**Throughput:** 158 MB/s

**Starting point:** a $\sqrt{N}$ construction ($N$ = number of records)

| | | | |
|---|---|---|---|
| $r_{11}$ | $r_{12}$ | $r_{13}$ | $r_{14}$ |
| $r_{21}$ | $r_{22}$ | $r_{23}$ | $r_{24}$ |
| $r_{31}$ | $r_{32}$ | $r_{33}$ | $r_{34}$ |
| $r_{41}$ | $r_{42}$ | $r_{43}$ | $r_{44}$ |

Arrange the database as a
$\sqrt{N}$-by-$\sqrt{N}$ matrix

**Starting point:** a $\sqrt{N}$ construction ($N$ = number of records)

| | | | | | |
|---|---|---|---|---|---|
| $r_{11}$ | $r_{12}$ | $r_{13}$ | $r_{14}$ | $\times$ | 0 |
| $r_{21}$ | $r_{22}$ | $r_{23}$ | $r_{24}$ | $\times$ | 0 |
| $r_{31}$ | $r_{32}$ | $r_{33}$ | $r_{34}$ | $\times$ | 1 |
| $r_{41}$ | $r_{42}$ | $r_{43}$ | $r_{44}$ | $\times$ | 0 |

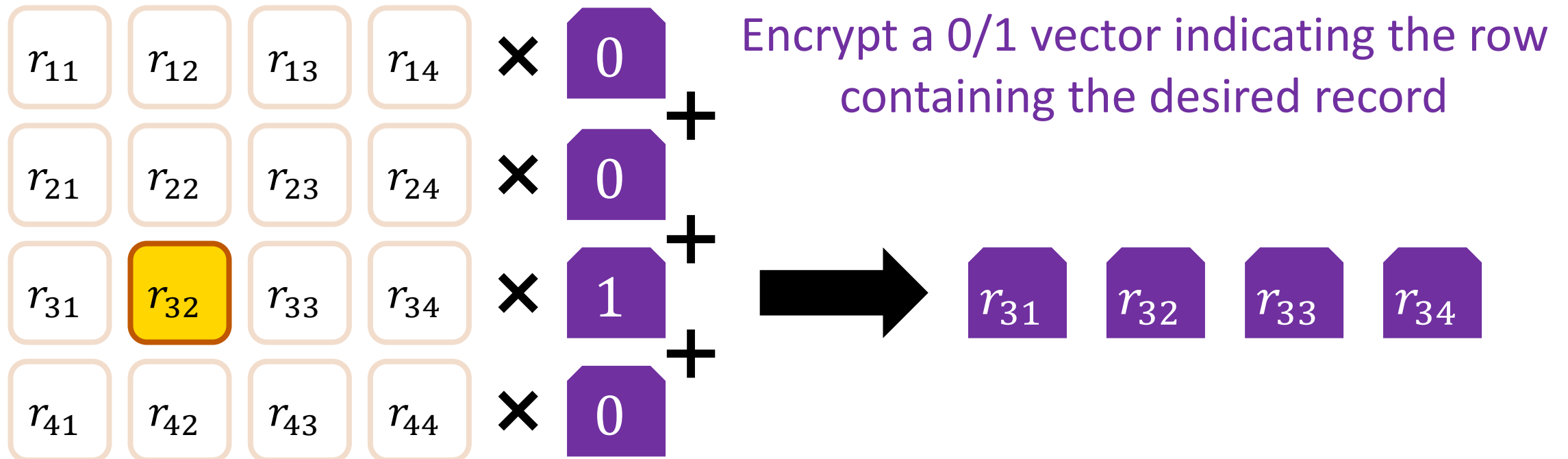Encrypt a 0/1 vector indicating the row containing the desired record

Arrange the database as a $\sqrt{N}$-by-$\sqrt{N}$ matrix

*Homomorphically* compute product between query vector and database matrix

# PIR from Homomorphic Encryption

**Starting point:** a $\sqrt{N}$ construction ($N$ = number of records)

| $r_{11}$ | $r_{12}$ | $r_{13}$ | $r_{14}$ | $\times$ | 0 |
| $r_{21}$ | $r_{22}$ | $r_{23}$ | $r_{24}$ | $\times$ | 0 |
| $r_{31}$ | $r_{32}$ | $r_{33}$ | $r_{34}$ | $\times$ | 1 |
| $r_{41}$ | $r_{42}$ | $r_{43}$ | $r_{44}$ | $\times$ | 0 |

$+$ $+$ $+$

Encrypt a 0/1 vector indicating the row containing the desired record

$\longrightarrow$ $r_{31}$ $r_{32}$ $r_{33}$ $r_{34}$

Arrange the database as a $\sqrt{N}$-by-$\sqrt{N}$ matrix

Database is in the clear, so *additive* homomorphism suffices

# PIR from Homomorphic Encryption

**Starting point:** a $\sqrt{N}$ construction ($N$ = number of records)

Client decrypts to learn records

Encrypt a 0/1 vector indicating the row containing the desired record

$r_{31}$  $r_{32}$  $r_{33}$  $r_{34}$

**Response size:** $\sqrt{N} \cdot \mathrm{poly}(\lambda)$

*Homomorphically* compute product between query vector and database matrix

# PIR from Homomorphic Encryption

**Starting point:** a $\sqrt{N}$ construction ($N$ = number of records)

Client decrypts to learn records

Encrypt a 0/1 vector indicating the row containing the desired record

$r_{31}$  $r_{32}$  $r_{33}$  $r_{34}$

**Response size:** $\sqrt{N} \cdot \text{poly}(\lambda)$

ciphertext size
($\lambda$ is security parameter)

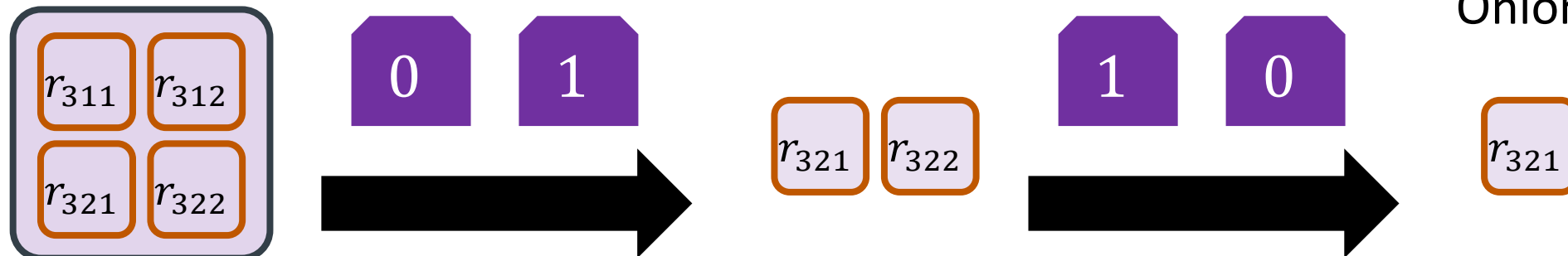*Homomorphically* compute product between query vector and database matrix

# PIR from Homomorphic Encryption

Beyond $\sqrt{N}$ communication: view the database as hypercube

Output is encrypted



Select along the first dimension

**Approach:** Use homomorphic multiplication

Gentry-Halevi [GH19]

OnionPIR [MCR21]

# SPIRAL: Composing FHE Schemes

Follows Gentry-Halevi blueprint of composing **two** lattice-based FHE schemes:

FHE ciphertexts are noisy encodings

Homomorphic operations increase noise; more noise = larger parameters = less efficiency

**Scheme 1:** Regev's encryption scheme [Reg04]

High-rate; only supports additive homomorphism

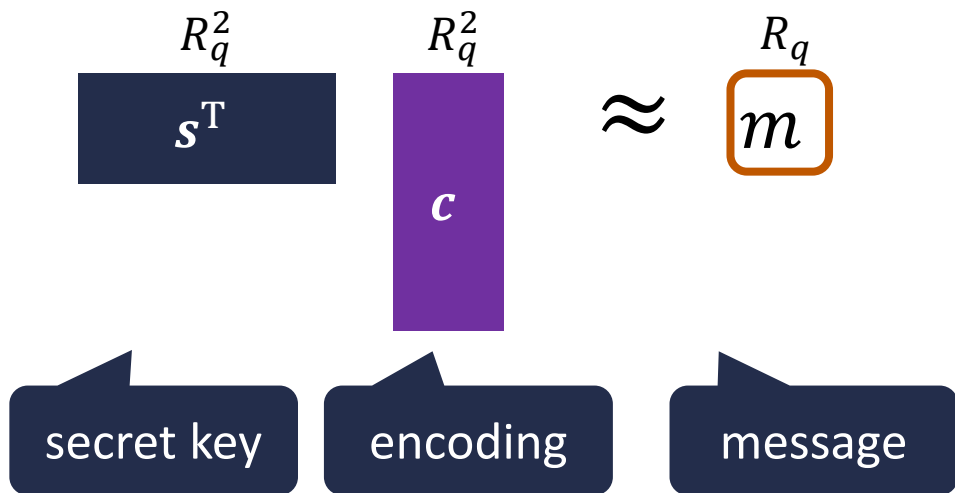**Scheme 2:** Gentry-Sahai-Waters encryption scheme [GSW13]

Low rate; supports homomorphic multiplication (with <u>additive</u> noise growth)

**Goal:** get the best of *both* worlds

# Regev Encodings (over Rings)

Regev <u>encoding</u> of a scalar $m \in R$:

$$\underset{\substack{R_q^2 \\ s^{\mathrm{T}}}}{} \quad \underset{\substack{R_q^2 \\ c}}{} \approx \underset{R_q}{\boxed{m}}$$

secret key    encoding    message

- Secret key allows recovery of noisy version of original message
- To support decryption of "small" values $t \in R_p$, we encode $t$ as $(q/p)t$
- Decryption recovers noisy version of $(q/p)t$ and rounding yields $t$

$$\text{rate} = \frac{\log p}{2 \log q} < \frac{1}{2}$$

**OnionPIR:** rate = 0.24

All elements are polynomials in the ring $R = \mathbb{Z}[x]/(x^d + 1)$ where $d = 2^k$

# Matrix Regev Encodings (over Rings)

Regev <u>encoding</u> of a matrix $M \in R_q^{n \times n}$:

$$\overset{R_q^{n \times (n+1)}}{S^{\mathrm{T}}} \quad \overset{R_q^{(n+1) \times n}}{C} \approx \overset{R_q^{n \times n}}{M}$$

**Idea:** "Reuse" encryption randomness

$$\text{rate} = \frac{n^2 \log p}{n(n+1) \log q} = \frac{n^2}{n^2 + n} \frac{\log p}{\log q}$$

**Additively homomorphic:**

$$S^{\mathrm{T}} C_1 \approx M_1$$

$$S^{\mathrm{T}} C_2 \approx M_2$$

$$S^{\mathrm{T}} (C_1 + C_2) \approx M_1 + M_2$$

All elements are polynomials in the ring $R = \mathbb{Z}[x]/(x^d + 1)$ where $d = 2^k$

# Gentry-Sahai-Waters Encodings

GSW <u>encoding</u> of a bit $\mu \in \{0,1\}$:

**Gadget matrix [MP12]:**

$$R_q^{n \times (n+1)} \quad R_q^{(n+1) \times n} \qquad R_q^{n \times (n+1)} \quad R_q^{(n+1) \times m}$$

$$\boxed{S^{\mathrm{T}}} \; \boxed{C} \;\approx\; \boxed{\mu}\; \boxed{S^{\mathrm{T}}} \; \boxed{G}$$

$$G = \begin{bmatrix} g^{\mathrm{T}} & & \\ & \ddots & \\ & & g^{\mathrm{T}} \end{bmatrix}$$

$$g^{\mathrm{T}} = [\,1 \quad 2 \quad 2^2 \quad \cdots \quad 2^{\lfloor \log_z q \rfloor}\,]$$
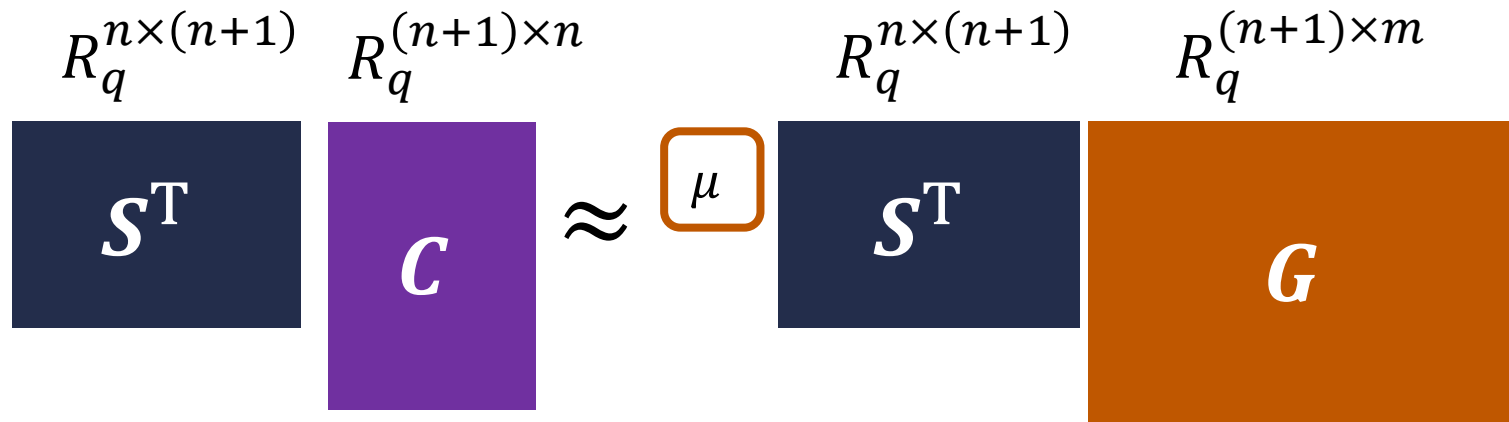
$$m = (n+1)\log q$$

"Powers-of-2" matrix

Construction will use other decomposition bases

All elements are polynomials in the ring $R = \mathbb{Z}[x]/(x^d + 1)$ where $d = 2^k$

# Gentry-Sahai-Waters Encodings

GSW <u>encoding</u> of a bit $\mu \in \{0,1\}$:

**Gadget matrix [MP12]:**

$$R_q^{n \times (n+1)} \quad R_q^{(n+1) \times n}$$

$$\boxed{S^{\mathrm{T}}} \boxed{C} \approx \boxed{\mu} \boxed{S^{\mathrm{T}}} \boxed{G}$$

$$R_q^{n \times (n+1)} \quad R_q^{(n+1) \times m}$$

$$G = \begin{bmatrix} g^{\mathrm{T}} & & \\ & \ddots & \\ & & g^{\mathrm{T}} \end{bmatrix}$$

$$g^{\mathrm{T}} = \begin{bmatrix} 1 & 2 & 2^2 & \dots & 2^{\lfloor \log_z q \rfloor} \end{bmatrix}$$

$$m = (n+1) \log q$$

"Powers-of-2" matrix

Construction will use other decomposition bases

**Main property:** for every vector $v \in \mathbb{Z}_q^{n+1}$, can define $G^{-1}(v) \in \{0,1\}^m$ where $GG^{-1}(v) = v$
*"binary decomposition"*

All elements are polynomials in the ring $R = \mathbb{Z}[x]/(x^d + 1)$ where $d = 2^k$

GSW <u>encoding</u> of a bit $\mu \in \{0,1\}$:

**Gadget matrix** [MP12]:

$$R_q^{n \times (n+1)} \quad R_q^{(n+1) \times n} \qquad\qquad R_q^{n \times (n+1)} \quad R_q^{(n+1) \times m}$$

$$\boxed{S^{\mathrm{T}}} \; \boxed{C} \approx \boxed{\mu} \; \boxed{S^{\mathrm{T}}} \; \boxed{G}$$

$$G = \begin{bmatrix} g^{\mathrm{T}} & & \\ & \ddots & \\ & & g^{\mathrm{T}} \end{bmatrix}$$

$$g^{\mathrm{T}} = [1 \quad 2 \quad 2^2 \quad \dots \quad 2^{\lfloor \log_z q \rfloor}]$$

$$\text{rate} = \frac{1}{d(n+1)^2 \log q}$$

$$m = (n+1)\log q$$

"Powers-of-2" matrix

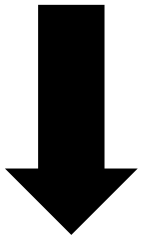**Concretely:** $d = 2048, n \geq 1, q = 2^{56}$

Construction will use other decomposition bases

All elements are polynomials in the ring $R = \mathbb{Z}[x]/(x^d + 1)$ where $d = 2^k$

# Regev-GSW Homomorphism

$$S^{\mathrm{T}} C_{\mathrm{Reg}} \approx M$$

$$S^{\mathrm{T}} C_{\mathrm{GSW}} \approx \mu S^{\mathrm{T}} G$$

With noise terms:
$$S^{\mathrm{T}} C_{\mathrm{GSW}} G^{-1}(C_{\mathrm{Reg}}) = \mu M + E_{\mathrm{GSW}} G^{-1}(C_{\mathrm{Reg}}) + \mu E_{\mathrm{Reg}}$$

**Asymmetric noise growth:** if all GSW ciphertexts are "fresh," then noise accumulation is additive in the number of multiplications

$$S^{\mathrm{T}} C_{\mathrm{GSW}} G^{-1}(C_{\mathrm{Reg}}) \approx \mu S^{\mathrm{T}} C_{\mathrm{Reg}} \approx \mu M$$

$C_{\mathrm{GSW}} G^{-1}(C_{\mathrm{Reg}})$ is a Regev encoding of $\mu M$

# The Gentry-Halevi Blueprint

Database is represented as $2^{\nu_1} \times \underbrace{2 \times 2 \times \cdots \times 2}_{2^{\nu_2}}$ hypercube

Query contains $2^{\nu_1}$ <u>matrix</u> Regev ciphertexts

| 0 | $\mathbf{I}_n$ | 0 | 0 | 0 | 0 |

Indicator for index along first dimension

Query contains $\nu_2$ GSW ciphertexts

| 0 | 1 | 1 | 0 |

Indicator for index along subsequent dimensions

Each GSW ciphertext participates in only <u>one</u> multiplication with a Regev ciphertext!

Response is a <u>single</u> matrix Regev ciphertext

# The Gentry-Halevi Blueprint

Database is represented as $2^{\nu_1} \times \underbrace{2 \times 2 \times \cdots \times 2}_{2^{\nu_2}}$ hypercube

**Drawback:** large queries

Can compress using polynomial encoding method of Angel et al.
[ACLS18]

Query contains $2^{\nu_1}$ matrix Regev ciphertexts

| 0 | $\mathbf{I}_n$ | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|

Indicator for index along first dimension

**Estimated size:**
4 MB/ciphertext

Query contains $\nu_2$ GSW ciphertexts

| 0 | 1 | 1 | 0 |
|---|---|---|---|

**Estimated query size:**
30 MB

Indicator for index along subsequent dimensions

# The Gentry-Halevi Blueprint

Database is represented as $2^{\nu_1} \times \underbrace{2 \times 2 \times \cdots \times 2}_{2^{\nu_2}}$ hypercube

**Drawback:** large queries

Can compress using polynomial encoding method of Angel et al. [ACLS18]

Query contains $2^{\nu_1}$ matrix Regev ciphertexts

| 0 | $\mathbf{I}_n$ | 0 | 0 | 0 | 0 |

Indicator for index along first dimension

SealPIR query size: 66 KB

**Estimated query size:** 30 MB

Query contains $\nu_2$ GSW ciphertexts

| 0 | 1 | 1 | 0 |

Indicator for index along subsequent dimensions

# OnionPIR

**High-level:** Gentry-Halevi approach with *scalar* Regev ciphertexts ($n = 1$)

Leverages Chen et al. approach [CCR19] to "assemble" GSW ciphertext using Regev-GSW multiplication

Regev ciphertexts can be packed using polynomial encoding method [ACLS18, CCR19]

Use of scalar Regev ciphertexts reduces the rate to $\approx 0.24$
(over $4\times$ response overhead)

# This Work: Translating Between Regev and GSW

**"Best of both worlds":** Small queries (as in OnionPIR) with the high rate/throughput of the Gentry-Halevi scheme

| | | |
|---|---|---|
| **Query size:** | 14 KB | 2000× smaller than Gentry-Halevi (4.5× smaller than OnionPIR) |
| **Rate:** | 0.41 | 2.1× higher than OnionPIR |
| **Throughput:** | 333 MB/s | 2.9× higher than OnionPIR |

(Database with $2^{14}$ records of size 100 KB)

Comparable improvements for other database configurations; more speed-ups in streaming setting

**Cost:** 3.4× larger public parameters for extra translation keys

Leverage simple key-switching techniques for query and response compression

Scalar Regev → Matrix Regev
Matrix Regev → GSW

Scalar Regev → Matrix Regev

Query compression

Response compression

(for large records)

# Scalar Regev → Matrix Regev

**Input:** encoding $\boldsymbol{c}$ where $\boldsymbol{s}_1^{\mathrm{T}} \boldsymbol{c} \approx m$

$$\boldsymbol{s}_1^{\mathrm{T}} = [-\tilde{s}_0 \mid 1] \in R_q^2$$

$$\boldsymbol{c}^{\mathrm{T}} = [c_0 \mid c_1] \in R_q^2$$

**Output:** encoding $\boldsymbol{C}$ where $\boldsymbol{S}_2^{\mathrm{T}} \boldsymbol{C} \approx m \mathbf{I}_n$



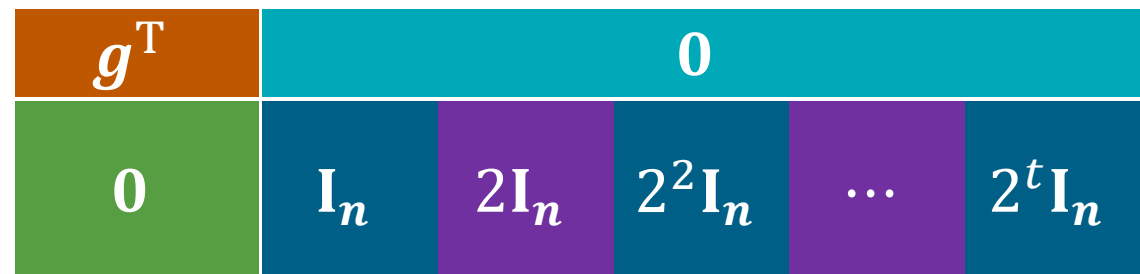$$\boldsymbol{S}_2^{\mathrm{T}} \boldsymbol{C} = m \boldsymbol{I}_n$$

Can replace with $\boldsymbol{S}_2$ with arbitrary secret key using standard key-switching techniques

# Matrix Regev → GSW

**Goal:** use Regev encodings to construct $C$ such that $S^{\mathrm{T}}C \approx \mu S^{\mathrm{T}}G$

$$S^{\mathrm{T}} = [-s \mid \mathbf{I}_n] \in R_q^{n \times (n+1)}$$

$$G = \begin{bmatrix} g^{\mathrm{T}} & & \\ & \ddots & \\ & & g^{\mathrm{T}} \end{bmatrix} \xrightarrow{\text{rearrange}}$$



$t = \log q$

$$\mu S^{\mathrm{T}}G = \begin{array}{|c|c|c|c|c|c|} \hline -\mu s g^{\mathrm{T}} & \mu \mathbf{I}_n & 2\mu \mathbf{I}_n & 2^2 \mu \mathbf{I}_n & \cdots & 2^t \mu \mathbf{I}_n \\ \hline \end{array}$$

$$C = \begin{array}{|c|c|c|c|c|c|} \hline A & B_0 & B_1 & B_2 & \cdots & B_t \\ \hline \end{array}$$

Break $C$ into *blocks*

# Matrix Regev → GSW

**Goal:** use Regev encodings to construct $C$ such that $S^{\mathrm{T}}C \approx \mu S^{\mathrm{T}}G$



$$S^{\mathrm{T}}C = \boxed{S^{\mathrm{T}}A} \quad \boxed{S^{\mathrm{T}}B_0 \mid S^{\mathrm{T}}B_1 \mid S^{\mathrm{T}}B_2 \mid \cdots \mid S^{\mathrm{T}}B_t}$$

$$\mu S^{\mathrm{T}}G = \boxed{-\mu sg^{\mathrm{T}}} \quad \boxed{\mu\mathbf{I}_n \mid 2\mu\mathbf{I}_n \mid 2^2\mu\mathbf{I}_n \mid \cdots \mid 2^t\mu\mathbf{I}_n}$$

$B_0, \ldots, B_t$ are matrix Regev ciphertexts encrypting $\mu\mathbf{I}_n, 2\mu\mathbf{I}_n, \ldots, 2^t\mu\mathbf{I}_n$

Can derive from scalar Regev encodings of $\mu, 2\mu, \ldots, 2^t\mu$

# Matrix Regev → GSW

**Goal:** use Regev encodings to construct $C$ such that $S^{\mathrm{T}} C \approx \mu S^{\mathrm{T}} G$

$$S^{\mathrm{T}} C = \boxed{S^{\mathrm{T}} A} \quad \begin{array}{|c|c|c|c|c|} \hline S^{\mathrm{T}} B_0 & S^{\mathrm{T}} B_1 & S^{\mathrm{T}} B_2 & \cdots & S^{\mathrm{T}} B_t \\ \hline \end{array}$$

Write $S^{\mathrm{T}} = [-s \mid \mathbf{I}_n]$

Let $s_{\mathrm{Reg}}$ be the key for a Regev encoding scheme

$$\mu S^{\mathrm{T}} G = \boxed{-\mu s g^{\mathrm{T}}} \quad \begin{array}{|c|c|c|c|c|} \hline \mu \mathbf{I}_n & 2\mu \mathbf{I}_n & 2^2 \mu \mathbf{I}_n & \cdots & 2^t \mu \mathbf{I}_n \\ \hline \end{array}$$

Construct key-switching matrix $W$:
$$S^{\mathrm{T}} W \approx -s \left( s_{\mathrm{Reg}}^{\mathrm{T}} \otimes g^{\mathrm{T}} \right)$$

Let $c_0, \ldots c_t$ be encodings of $\mu, \ldots, 2^t \mu$ under $s_{\mathrm{Reg}}$: $s_{\mathrm{Reg}}^{\mathrm{T}} c_i \approx 2^i \mu$

Let $C = [c_0 \mid \cdots \mid c_t]$

Then, $S^{\mathrm{T}} W g^{-1}(C) \approx -s \left( s_{\mathrm{Reg}}^{\mathrm{T}} \otimes g^{\mathrm{T}} \right) g^{-1}(C) \approx -s [\mu \mid \cdots \mid 2^t \mu] = -\mu s g^{\mathrm{T}}$

# Matrix Regev → GSW

**Goal:** use Regev encodings to construct $\boldsymbol{C}$ such that $\boldsymbol{S}^{\mathrm{T}}\boldsymbol{C} \approx \mu \boldsymbol{S}^{\mathrm{T}}\boldsymbol{G}$

$$\boldsymbol{S}^{\mathrm{T}}\boldsymbol{C} = \left[ \begin{array}{c|ccccc} \boldsymbol{S}^{\mathrm{T}}\boldsymbol{A} & \boldsymbol{S}^{\mathrm{T}}\boldsymbol{B}_0 & \boldsymbol{S}^{\mathrm{T}}\boldsymbol{B}_1 & \boldsymbol{S}^{\mathrm{T}}\boldsymbol{B}_2 & \cdots & \boldsymbol{S}^{\mathrm{T}}\boldsymbol{B}_t \end{array} \right]$$

Write $\boldsymbol{S}^{\mathrm{T}} = [\,-\boldsymbol{s} \mid \mathbf{I}_n\,]$

Let $\boldsymbol{s}_{\mathrm{Reg}}$ be the key for a Regev encoding scheme

$$\mu \boldsymbol{S}^{\mathrm{T}}\boldsymbol{G} = \left[ \begin{array}{c|ccccc} -\mu \boldsymbol{s}\boldsymbol{g}^{\mathrm{T}} & \mu \mathbf{I}_n & 2\mu \mathbf{I}_n & 2^2\mu \mathbf{I}_n & \cdots & 2^t\mu \mathbf{I}_n \end{array} \right]$$

Construct key-switching matrix $\boldsymbol{W}$:
$$\boldsymbol{S}^{\mathrm{T}}\boldsymbol{W} \approx -\boldsymbol{s}\left( \boldsymbol{s}_{\mathrm{Reg}}^{\mathrm{T}} \otimes \boldsymbol{g}^{\mathrm{T}} \right)$$

Let $\boldsymbol{c}_0, \ldots$ be encodings of $\ldots 2^t \ldots$ under $\boldsymbol{s}_{\mathrm{Reg}}$: $\boldsymbol{s}_{\mathrm{Reg}}^{\mathrm{T}} \boldsymbol{c}_i \approx 2^i \mu$

Let $\boldsymbol{C} = \ldots$

> Define $\boldsymbol{A} = \boldsymbol{W}\boldsymbol{g}^{-1}(\boldsymbol{C})$

Then, $\boldsymbol{S}^{\mathrm{T}}\boldsymbol{W}\boldsymbol{g}^{-1}(\boldsymbol{C}) \approx -\boldsymbol{s}\left( \boldsymbol{s}_{\mathrm{Reg}}^{\mathrm{T}} \otimes \boldsymbol{g}^{\mathrm{T}} \right) \boldsymbol{g}^{-1}(\boldsymbol{C}) \approx -\boldsymbol{s}[\,\mu \mid \cdots \mid 2^t\mu\,] = -\mu \boldsymbol{s}\boldsymbol{g}^{\mathrm{T}}$

# Matrix Regev → GSW

$\boldsymbol{W}\boldsymbol{g}^{-1}([\boldsymbol{c_0}|\cdots|\boldsymbol{c_t}])$

Concatenate blocks to obtain GSW encoding of $\mu$

| $\boldsymbol{c_0}$ | $\mu$ |
| $\boldsymbol{c_1}$ | $2\mu$ |
| $\vdots$ | |
| $\boldsymbol{c_t}$ | $2^t\mu$ |

$-\mu \boldsymbol{s}\boldsymbol{g}^{\mathrm{T}}$

Scalar Regev to Matrix Regev

$\mu \mathbf{I}_n$

$\vdots$

$2^t\mu \mathbf{I}_n$

$-\mu \boldsymbol{s}\boldsymbol{g}^{\mathrm{T}}$  $\mu \mathbf{I}_n$  $\cdots$  $2^t\mu \mathbf{I}_n$
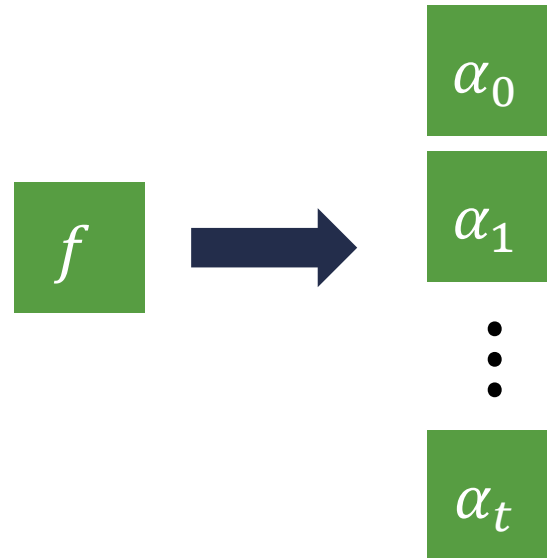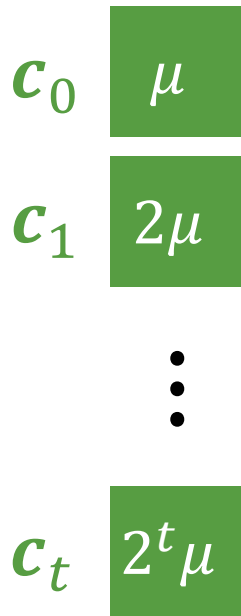
Ciphertext contains
$(n+1)^2(t+1)$ elements of $R_q$

scalar Regev encodings: elements of $R_q^2$

matrix Regev encodings: elements of $R_q^{(n+1)\times n}$

**Takeaway:** instead of sending $(n+1)^2(t+1)$ ring elements per GSW ciphertext, only need to send $2(t+1)$

# Further Compression via Polynomial Encodings

[ACLS18, CCR19]: let $f(x) = \alpha_0 + \alpha_1 x + \cdots + \alpha_t \cdot x^t$ with $t < d$

$c_0$ $\mu$

$c_1$ $2\mu$

$\vdots$

$c_t$ $2^t\mu$

$f$ $\rightarrow$ $\alpha_0$ $\alpha_1$ $\vdots$ $\alpha_t$

Expands a Regev encoding of a polynomial into Regev encodings of its coefficients

**Cost:** additional (reusable) public parameters needed for Regev-to-GSW translation

**Takeaway:** We can pack $(\mu, 2\mu, \ldots 2^t\mu)$ into a <u>single</u> polynomial

As long as $t + 1 < d$, client and communicate a GSW ciphertext with a <u>single</u> Regev encoding (2 ring elements)

$(n+1)^2(t+1)$ ring elements

$\downarrow$

2 ring elements

# Query Expansion in Spiral

Database is represented as $2^{\nu_1} \times \underbrace{2 \times 2 \times \cdots \times 2}_{2^{\nu_2}}$ hypercube
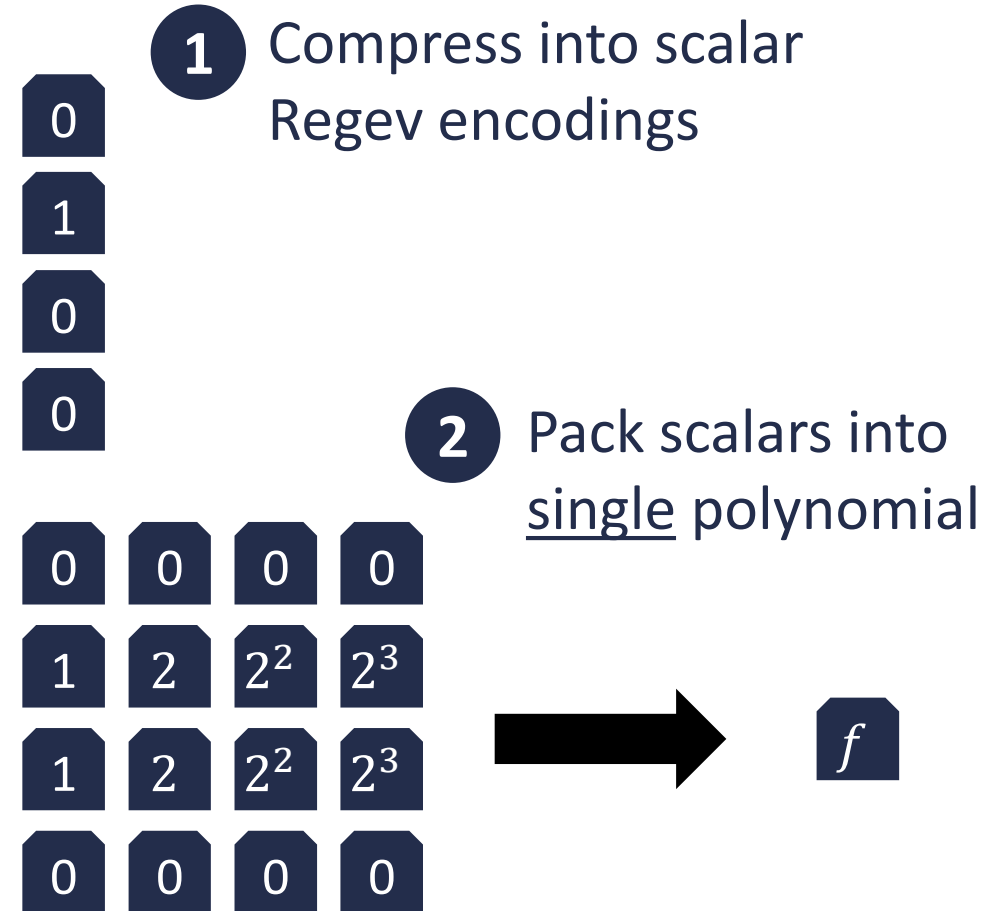
Query contains $2^{\nu_1}$ matrix Regev ciphertexts

| 0 | $\mathbf{I}_n$ | 0 | 0 |

Indicator for index along first dimension

Query contains $\nu_2$ GSW ciphertexts

| 0 | 1 | 1 | 0 |

Indicator for index along subsequent dimensions

**1** Compress into scalar Regev encodings

| 0 |
| 1 |
| 0 |
| 0 |

**2** Pack scalars into <u>single</u> polynomial

| 0 | 0 | 0 | 0 |
| 1 | 2 | $2^2$ | $2^3$ |
| 1 | 2 | $2^2$ | $2^3$ |
| 0 | 0 | 0 | 0 |

$f$

# Query Expansion in Spiral

public parameters

Moving costs from online to offline phase

query

response

offline and one-time cost

online cost

SPIRAL also achieves higher rate and throughput

**Trade-off:** larger public parameters, smaller queries

SealPIR:     3 MB
OnionPIR:  5 MB
SPIRAL:      18 MB

SealPIR:      66 KB          Gentry-Halevi: ≈30 MB
OnionPIR:  63 KB          SPIRAL: 14 KB

# Response Compression via Modulus Switching

PIR response consists of a single <u>matrix</u> Regev encoding

$$[-\boldsymbol{s} \mid \mathbf{I}_n] \quad \boldsymbol{C} \quad \approx \quad \frac{q}{p} \cdot \boldsymbol{M}$$

Modulus $q$ must be large enough to support target number of homomorphic operations

$$\text{rate} \propto \frac{\log p}{\log q}$$

Standard technique in FHE: *modulus reduction*

Rescale ciphertext by $\frac{q'}{q}$ where $q' < q$

$$\text{rate} \propto \frac{\log p}{\log q'}$$

Rescaling introduces small amount of noise (from rounding)

**This work:** Observe that rounding error $\boldsymbol{E}$ is scaled by $[-\boldsymbol{s} \mid \mathbf{I}_n]$

$$[-\boldsymbol{s} \mid \mathbf{I}_n] \quad \boldsymbol{E}$$

# Response Compression via Modulus Switching

PIR response consists of a single <u>matrix</u> Regev encoding

$$[-\boldsymbol{s} \mid \mathbf{I}_n]\ \boldsymbol{C} \approx \boxed{\frac{q}{p} \cdot \boldsymbol{M}}$$

Modulus $q$ must be large enough to support target number of homomorphic operations

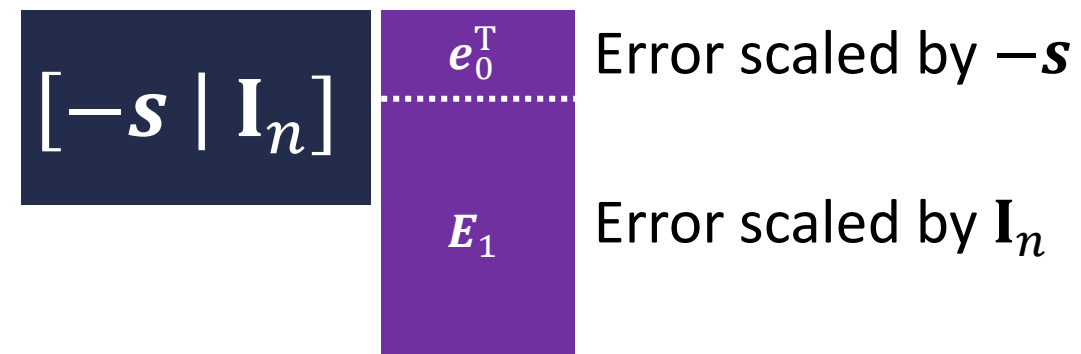$$\text{rate} \propto \frac{\log p}{\log q}$$

Standard technique in FHE: *modulus reduction*

Rescale ciphertext by $\frac{q'}{q}$ where $q' < q$

$$\text{rate} \propto \frac{\log p}{\log q'}$$

Rescaling introduces small amount of noise (from rounding)

**This work:** Observe that rounding error $\boldsymbol{E}$ is scaled by $[-\boldsymbol{s} \mid \mathbf{I}_n]$

$$[-\boldsymbol{s} \mid \mathbf{I}_n]$$

$\boldsymbol{e}_0^{\mathrm{T}}$  Error scaled by $-\boldsymbol{s}$

$\boldsymbol{E}_1$  Error scaled by $\mathbf{I}_n$

# Response Compression via Modulus Switching

PIR response consists of a single <u>matrix</u> Regev encoding

$$[-\boldsymbol{s} \mid \mathbf{I}_n] \quad \boxed{C} \approx \boxed{\frac{q}{p} \cdot \boldsymbol{M}}$$

**Observation:** At least half of the error components are scaled by identity matrix!

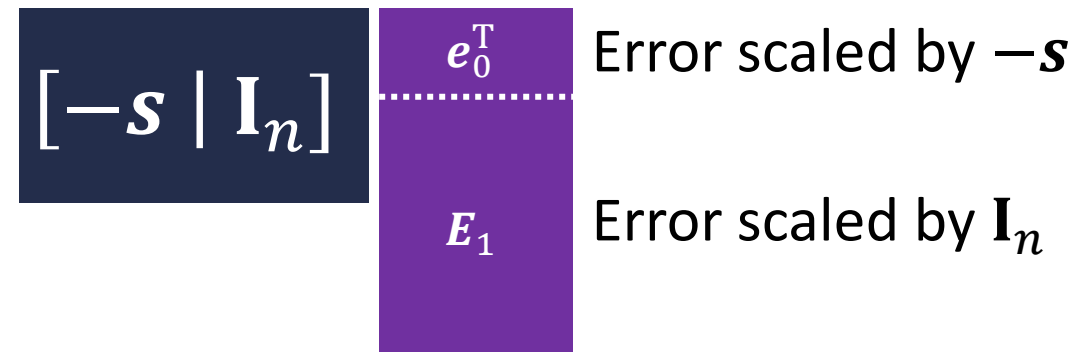**Approach:** Use two different moduli to rescale the ciphertext

Standard technique in FHE: *modulus reduction*

Rescale ciphertext by $\frac{q'}{q}$ where $q' < q$

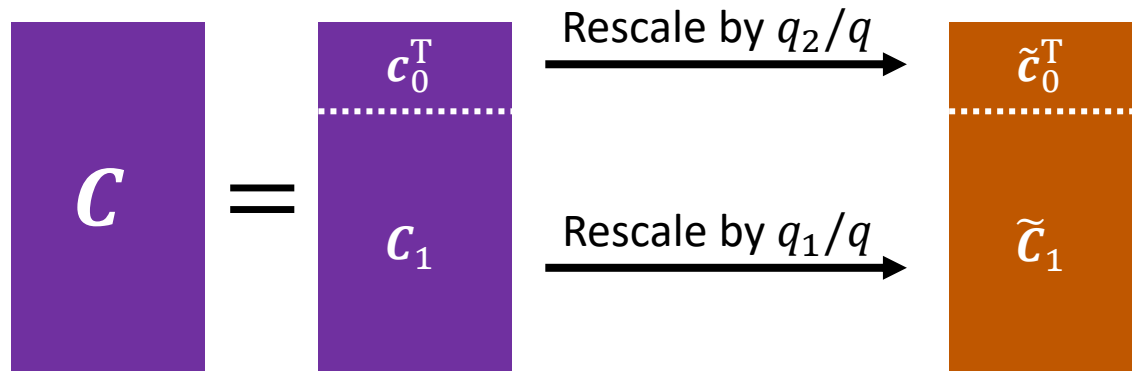$$\text{rate} \propto \frac{\log p}{\log q'}$$

Rescaling introduces small amount of noise (from rounding)

**This work:** Observe that rounding error $\boldsymbol{E}$ is scaled by $[-\boldsymbol{s} \mid \mathbf{I}_n]$

$$[-\boldsymbol{s} \mid \mathbf{I}_n]$$

$e_0^{\mathrm{T}}$ — Error scaled by $-\boldsymbol{s}$

$E_1$ — Error scaled by $\mathbf{I}_n$

# Response Compression via Modulus Switching

PIR response consists of a single <u>matrix</u> Regev encoding



$$C = \begin{array}{c} c_0^{\mathrm{T}} \\ \hline c_1 \end{array}$$

Rescale by $q_2/q$ → $\tilde{c}_0^{\mathrm{T}}$

Rescale by $q_1/q$ → $\tilde{C}_1$

**Observation:** At least half of the error components are scaled by identity matrix!

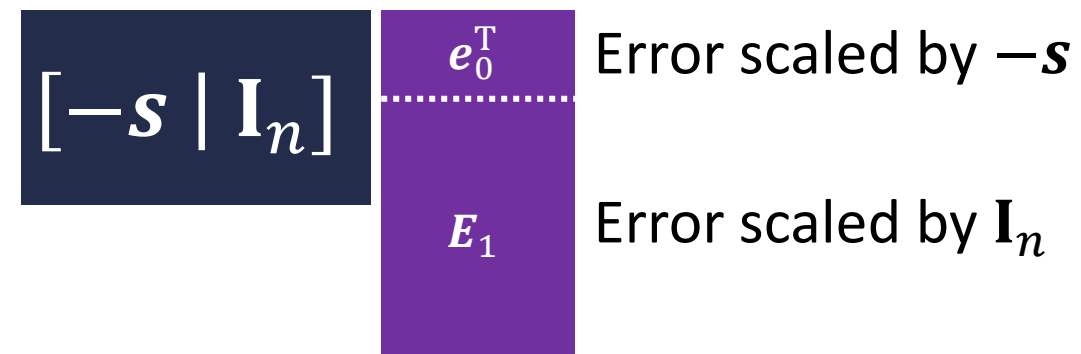**Approach:** Use two different moduli to rescale the ciphertext

Standard technique in FHE: *modulus reduction*

Rescale ciphertext by $\dfrac{q'}{q}$ where $q' < q$

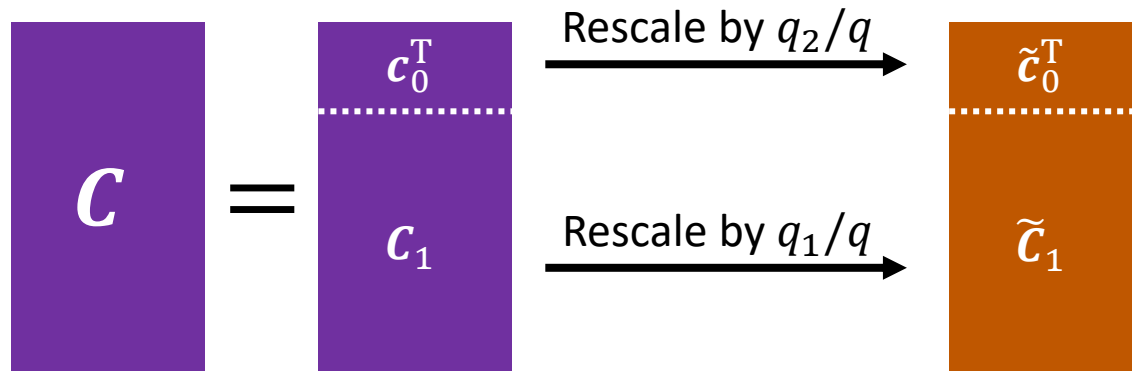$$\text{rate} \propto \frac{\log p}{\log q'}$$

Rescaling introduces small amount of noise (from rounding)

**This work:** Observe that rounding error $E$ is scaled by $[-s \mid \mathbf{I}_n]$

$$[-s \mid \mathbf{I}_n] \quad \begin{array}{c} e_0^{\mathrm{T}} \\ \hline E_1 \end{array}$$

Error scaled by $-s$

Error scaled by $\mathbf{I}_n$

# Response Compression via Modulus Switching

PIR response consists of a single <u>matrix</u> Regev encoding



Rescale by $q_2/q$

Rescale by $q_1/q$

$C = $ with components $c_0^{\mathrm{T}}$, $c_1$ rescaled to $\tilde{c}_0^{\mathrm{T}}$, $\tilde{c}_1$

**Observation:** At least half of the error components are scaled by identity matrix!
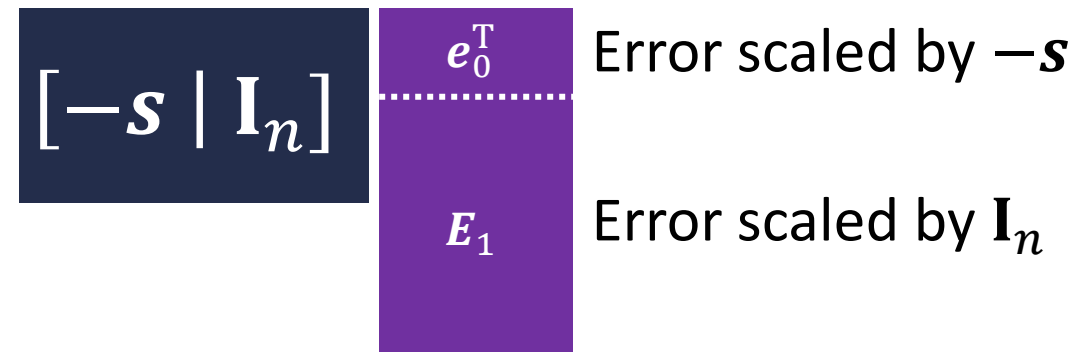
**Approach:** Use two different moduli to rescale the ciphertext

$$\text{rate} = \frac{n^2 \log p}{n^2 \log q_1 + n \log q_2}$$

- SealPIR                              0.01
- Gentry-Halevi (estimated)   0.44
- OnionPIR                          0.24

**Overall rate:**   0.34 (with vanilla modulus switching)
0.81 (with split modulus switching)

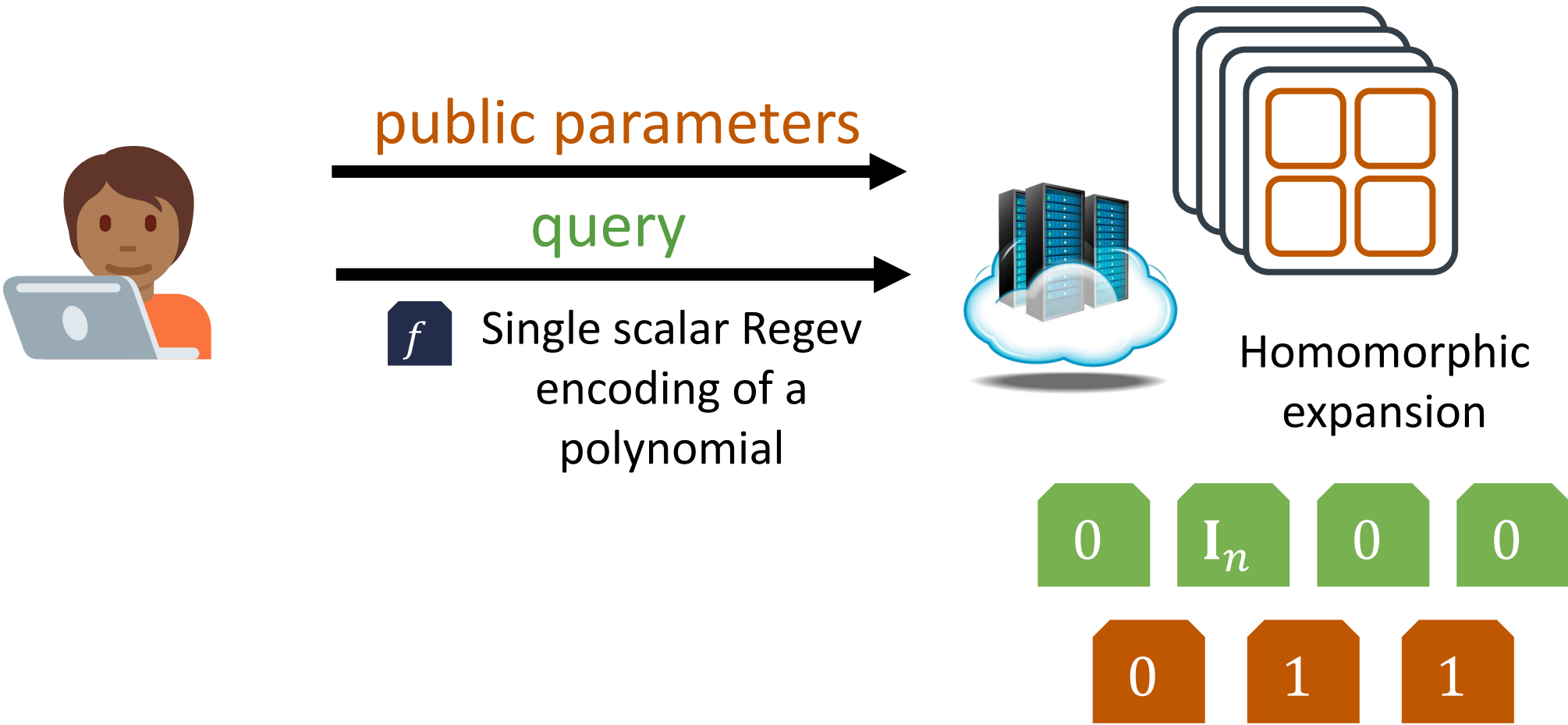**This work:** Observe that rounding error $E$ is scaled by $[-s \mid \mathbf{I}_n]$

$[-s \mid \mathbf{I}_n]$ with $e_0^{\mathrm{T}}$ → Error scaled by $-s$

$E_1$ → Error scaled by $\mathbf{I}_n$

# Vanilla SPIRAL



record $i$

public parameters

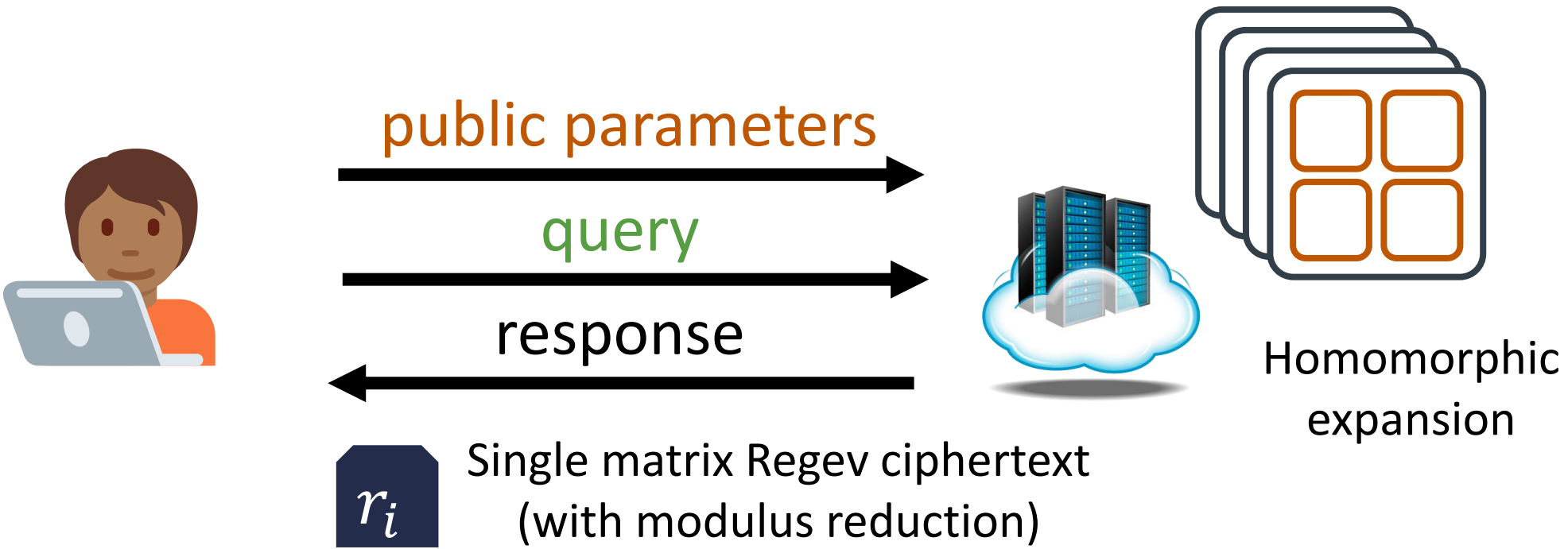Key-switching matrices for ciphertext expansion and translation

# Vanilla SPIRAL

public parameters

query

*f* Single scalar Regev encoding of a polynomial

Homomorphic expansion

| 0 | $\mathbf{I}_n$ | 0 | 0 |

| 0 | 1 | 1 |

# Vanilla SPIRAL

# Vanilla SPIRAL

public parameters

query

response

Homomorphic expansion

$r_i$ Single matrix Regev ciphertext
(with modulus reduction)

Many parameter choices in SPIRAL:
    Plaintext matrix dimension
    Plaintext modulus
    Decomposition bases for key-switching
    Database arrangement

Trade-offs in public parameter size, query size, server throughput, and rate

Use estimated running time + compute cost to choose parameters for an input database configuration

Automatic parameter selection tool

# Basic Comparisons

| Database | Metric | SealPIR | FastPIR | OnionPIR | SPIRAL |
|---|---|---|---|---|---|
| $2^{18}$ records 30 KB records (7.9 GB database) | Public Param. Size | 3 MB | 1 MB | 5 MB | 18 MB |
| | Query Size | 66 KB | 8 MB | 63 KB | 14 KB |
| | Response Size | 3 MB | 262 KB | 127 KB | 84 KB |
| | Server Compute | 74.91 s | 50.5 s | 52.7 s | 24.5 s |
| | Rate: | | | 0.24 | 0.36 |
| | Throughput: | | | 149 MB/s | 322 MB/s |

Database configuration preferred by OnionPIR

**Compared to OnionPIR:**
reduce query size by 4.5×        increase public parameter size by 3.6×
reduce response size by 2×
reduce compute time by 2×

# Basic Comparisons (with Larger Records)



**Throughput for 100 GB database ($2^{20}$ records):**
- Spiral:        310 MB/s (322 s)
- SealPIR:       102 MB/s (977 s)
- FastPIR:       189 MB/s (528 s)
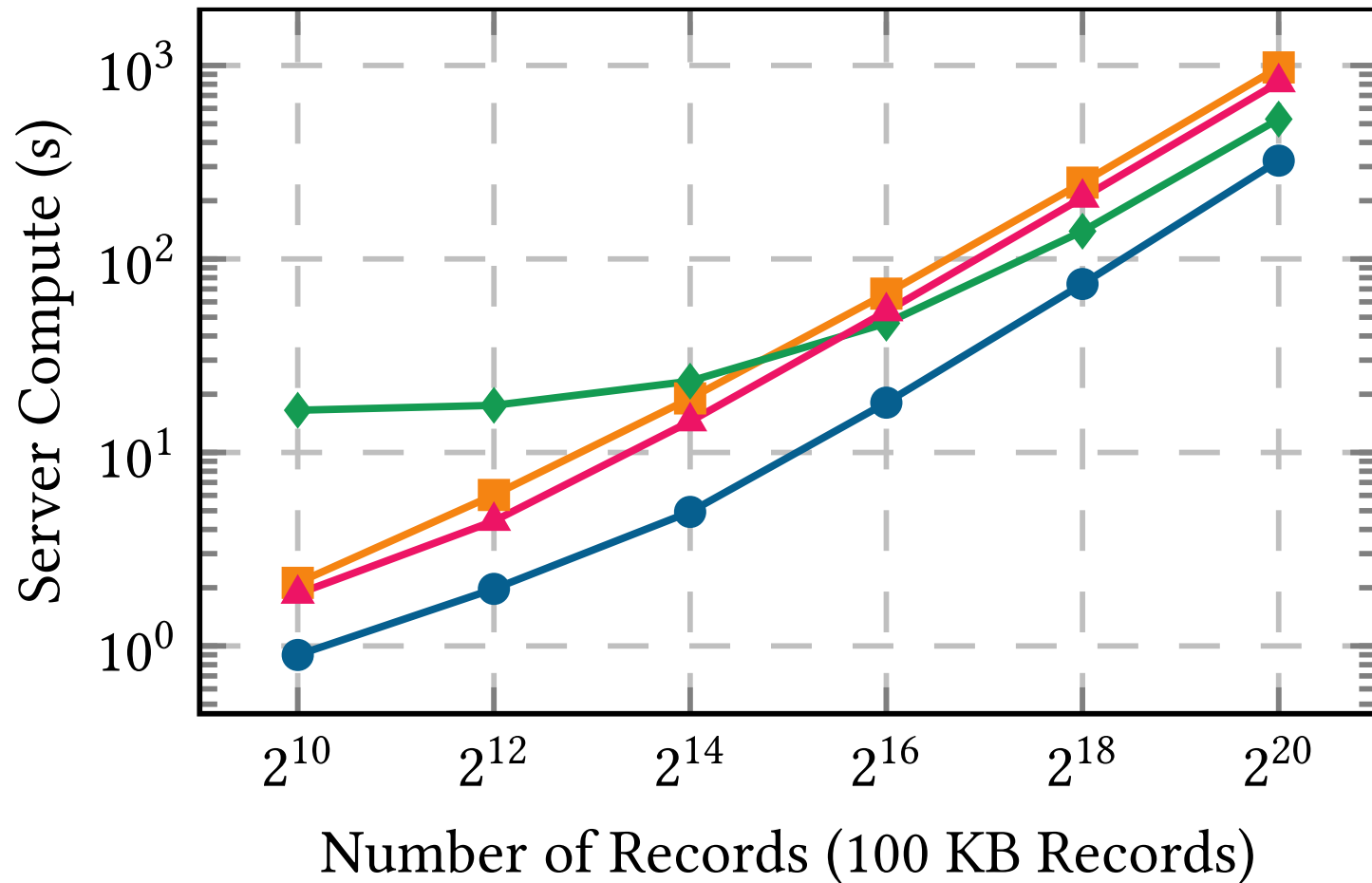- OnionPIR:      122 MB/s (817 s)

Spiral also has smaller query size and response size, but larger public parameters

All measurements based on single-thread/single-core processing

Server cost is <u>linear</u> in database size

# Basic Comparisons (with Larger Records)



**Client costs:**
- Generating <u>reusable</u> public parameters is the most expensive operation, but still < 1 s
- Query generation and response decoding are fast (30 ms and < 1 ms)

**Server costs:**
- Query expansion typically takes ≈ 1 second (less than 1.5% of overall compute when number of records is large)
- Parameter selection favors configurations that evenly distributes the work between first layer processing and ciphertext folding

(*see paper for detailed microbenchmarks*)

# The Streaming Setting: SPIRALSTREAM
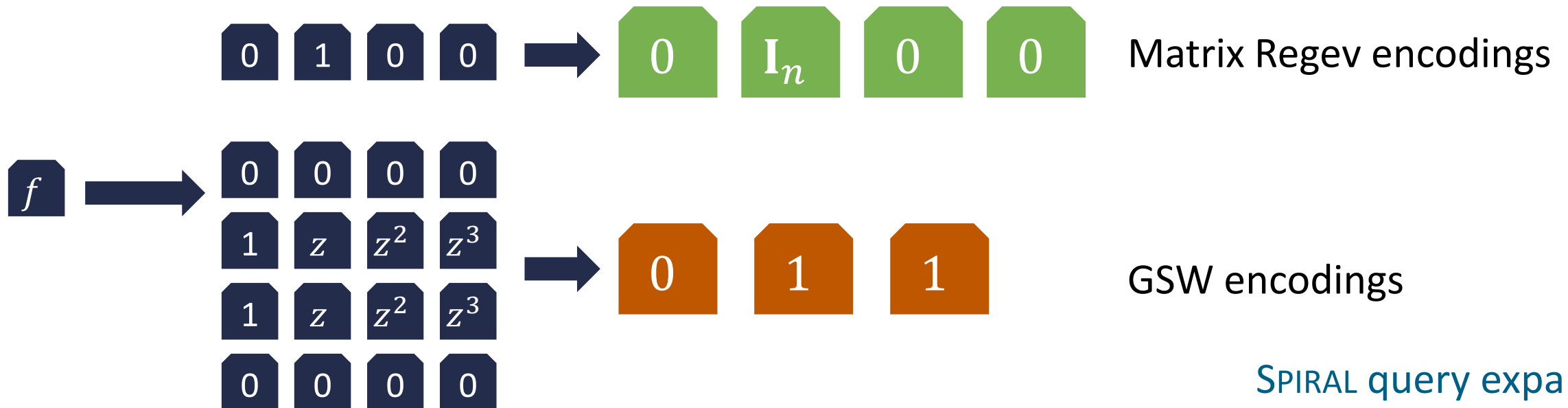
**Streaming setting:** <u>same</u> query reused over multiple databases

Private video stream (database $D_i$ contains $i^{\text{th}}$ block of media)      [GCMSAW16]

Private voice calls (repeated polling of the same "mailbox")      [AS16, AYAAG21]

**Goal:** minimize online costs (i.e., server compute, response size)

Consider larger public parameters or query size (amortized over lifetime of stream)



Matrix Regev encodings

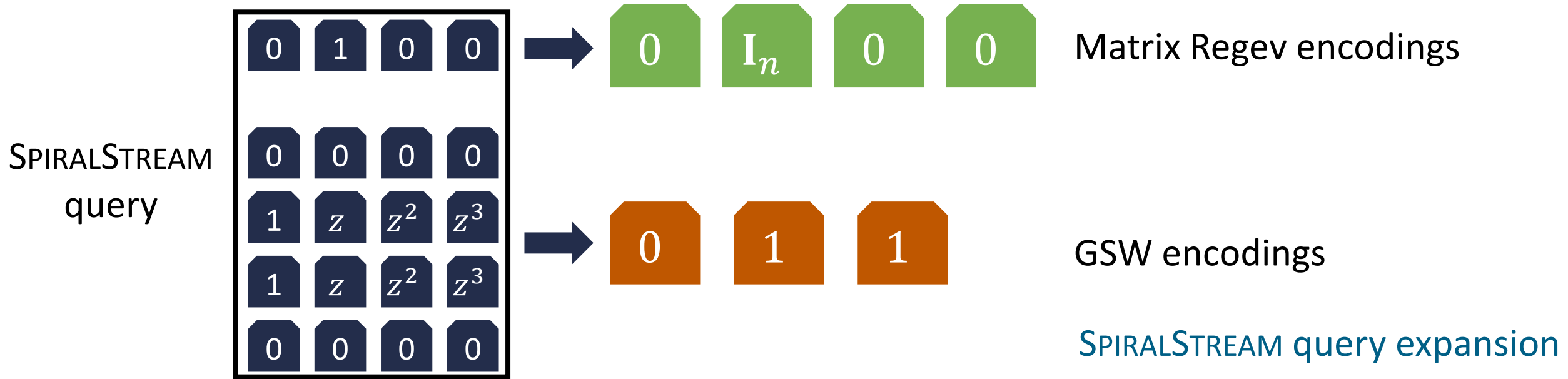GSW encodings

SPIRAL query expansion

# The Streaming Setting: SPIRALSTREAM

Removing the initial expansion <u>significantly</u> reduces the noise growth from query expansion

Decreases size of public parameters (no more automorphism keys)

Better control of noise growth $\Rightarrow$ higher server throughput and higher rate

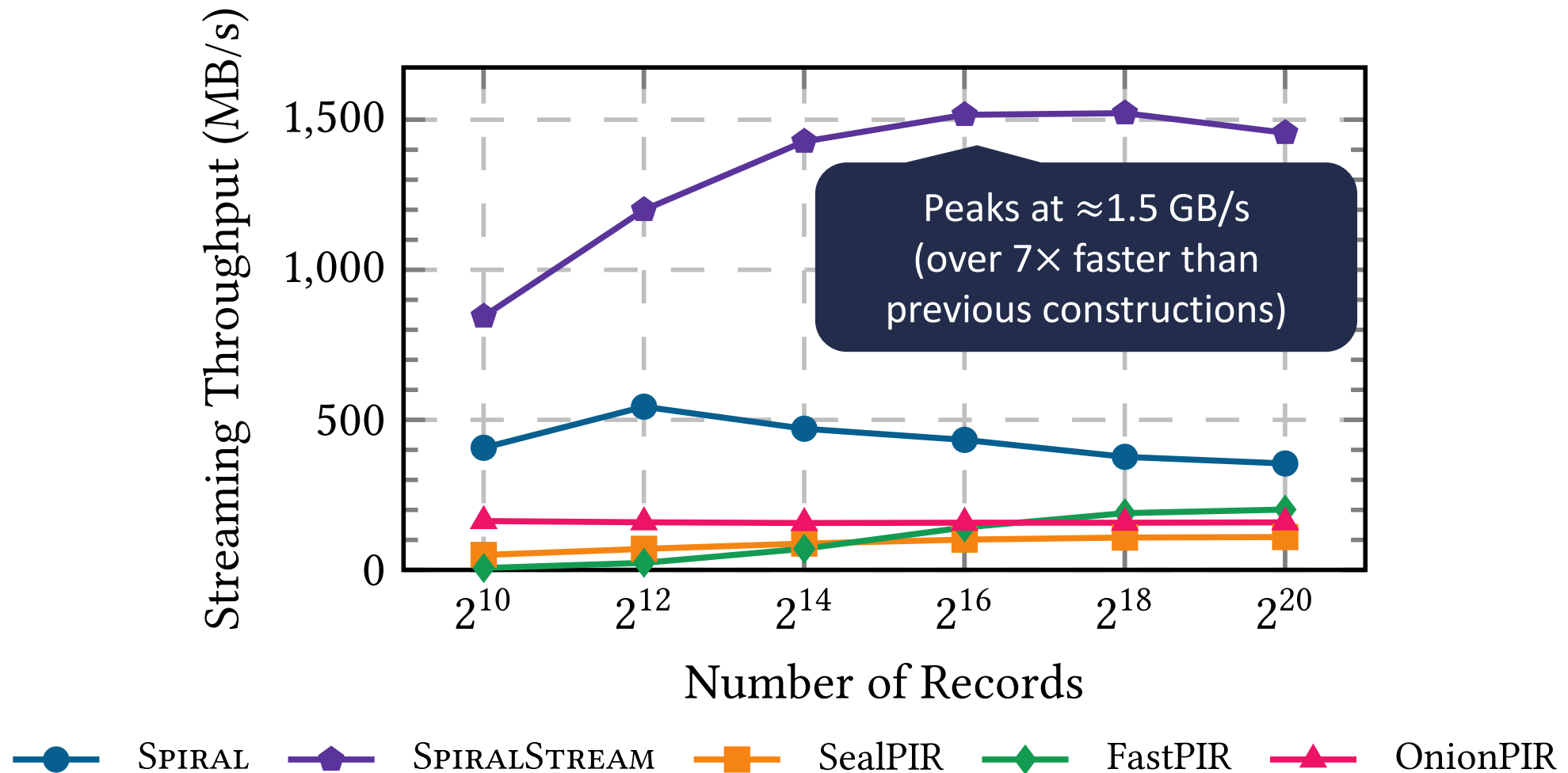Larger queries (more Regev encodings)

SPIRALSTREAM query

| | | | |
|---|---|---|---|
| 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 1 | $z$ | $z^2$ | $z^3$ |
| 1 | $z$ | $z^2$ | $z^3$ |
| 0 | 0 | 0 | 0 |

$\longrightarrow$ $\boxed{0}$ $\boxed{\mathbf{I}_n}$ $\boxed{0}$ $\boxed{0}$   Matrix Regev encodings

$\longrightarrow$ $\boxed{0}$ $\boxed{1}$ $\boxed{1}$   GSW encodings

SPIRALSTREAM query expansion

# The Streaming Setting: SPIRALSTREAM

| Database | Metric | OnionPIR | SPIRAL | SPIRALSTREAM |
|---|---|---|---|---|
| $2^{18}$ records 30 KB records (7.9 GB database) | Public Param. Size | 5 MB | 18 MB | 3 MB |
| | Query Size | 63 KB | 14 KB | 15 MB |
| | Response Size | 127 KB | 84 KB | 62 KB |
| | Server Compute | 52.7 s | 24.5 s | 9.0 s |
| | Rate: | 0.23 | 0.36 | 0.48 |
| | Throughput: | 149 MB/s | 322 MB/s | 874 MB/s |

25% reduction in response size
2.7× increase in throughput

# The Streaming Setting: SPIRALSTREAM

**Streaming throughput:** ignoring query expansion costs, assuming optimal record size for each system

# Higher Rate via Response Packing: SPIRALPACK

*Can we further reduce response size?*

$$\text{rate} = \frac{n^2 \log p}{n \log q_2 + n^2 \log q_1} \qquad q_1 = 4p$$

Increasing the plaintext dimension $n$ increases the rate

SPIRAL and SPIRALSTREAM use $n = 2$

   Higher values of $n$ increases <u>computational</u> cost

   Each Regev encoding is a $(n + 1) \times n$ matrix, so number of ring operations per
   homomorphic operation scale with $O(n^3)$    [Not using fast matrix multiplications here]

**SPIRALPACK:** Perform homomorphic operations with $n = 1$ and pack <u>responses</u>

# Higher Rate via Response Packing: SᴘɪʀᴀʟPᴀᴄᴋ

## Sᴘɪʀᴀʟ



Plaintext space: $R_p^{n \times n}$

Each record is $n \times n$ matrix

## SᴘɪʀᴀʟPᴀᴄᴋ



Split database into $n^2$ databases

$i^{\text{th}}$ database contains $i^{\text{th}}$ entry of record (elements of $R_p$)

Better throughput
Worse rate

Response consists of $n^2$ Regev encodings

# Higher Rate via Response Packing: SPIRALPACK

$n^2$ Regev ciphertexts with dimension 1

Consists of $2n^2$ ring elements

Variant of scalar Regev to matrix Regev transformation

Requires publishing $n$ key-switching matrices

Packing done only at the very end (cost does <u>not</u> scale with number of records)

1 Regev ciphertext with dimension $n$

Consists of $n(n + 1)$ ring elements

**SPIRALPACK:** higher throughput and rate (for sufficiently large records), larger public parameters

# Higher Rate via Response Packing: SPIRALPACK

| Database | Metric | OnionPIR | SPIRAL | | SPIRALSTREAM | |
|---|---|---|---|---|---|---|
| **$2^{18}$ records** **30 KB records** **(7.9 GB database)** | Public Param. Size | 5 MB | 18 MB $\rightarrow$ | 18 MB | 3 MB $\rightarrow$ | 16 MB |
| | Query Size | 63 KB | 14 KB $\rightarrow$ | 14 KB | 15 MB $\rightarrow$ | 30 MB |
| | Response Size | 127 KB | 84 KB $\rightarrow$ | 86 KB | 62 KB $\rightarrow$ | 96 KB |
| | Server Compute | 52.7 s | 24.5 s $\rightarrow$ | 17.7 s | 9.0 s $\rightarrow$ | 5.3 s |

- Small records $\Rightarrow$ can only take advantage of low packing dimension
- Higher throughputs since homomorphic operations cheaper
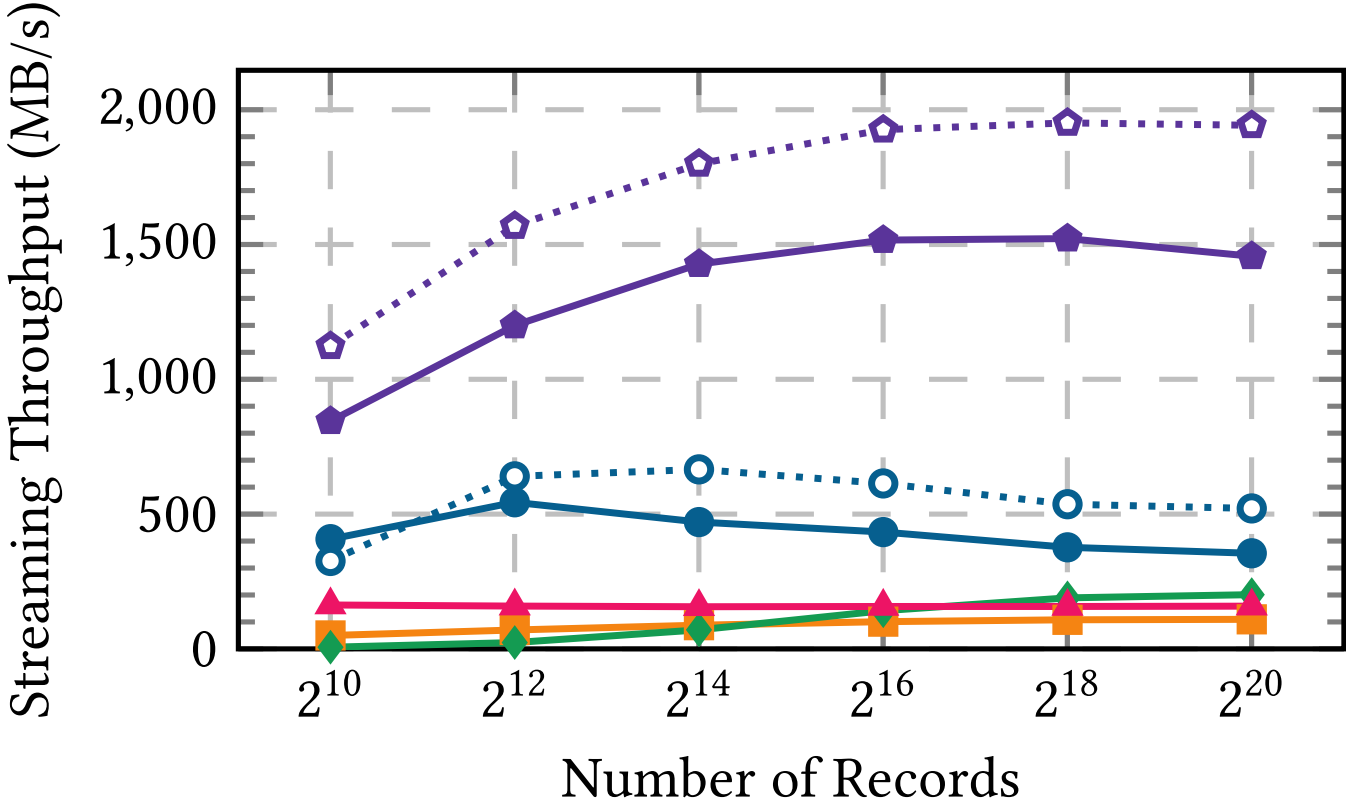- Responses larger due to extra noise from response packing

# Higher Rate via Response Packing: SPIRALPACK

| Database | Metric | OnionPIR | SPIRAL | | SPIRALSTREAM | |
|---|---|---|---|---|---|---|
| $2^{18}$ records 30 KB records (7.9 GB database) | Public Param. Size | 5 MB | 18 MB | $\rightarrow$ 18 MB | 3 MB | $\rightarrow$ 16 MB |
| | Query Size | 63 KB | 14 KB | $\rightarrow$ 14 KB | 15 MB | $\rightarrow$ 30 MB |
| | Response Size | 127 KB | 84 KB | $\rightarrow$ 86 KB | 62 KB | $\rightarrow$ 96 KB |
| | Server Compute | 52.7 s | 24.5 s | $\rightarrow$ 17.7 s | 9.0 s | $\rightarrow$ 5.3 s |
| $2^{14}$ records 100 KB records (1.6 GB database) | Public Param. Size | 5 MB | 17 MB | $\rightarrow$ 47 MB | 1 MB | $\rightarrow$ 24 MB |
| | Query Size | 63 KB | 14 KB | $\rightarrow$ 14 KB | 8 MB | $\rightarrow$ 30 MB |
| | Response Size | 508 KB | 242 KB | $\rightarrow$ 188 KB | 208 KB | $\rightarrow$ 150 KB |
| | Server Compute | 14.4 s | 4.92 s | $\rightarrow$ 4.58 s | 2.4 s | $\rightarrow$ 1.2 s |
| | Rate: | 0.20 | 0.41 | $\rightarrow$ 0.53 | 0.48 | $\rightarrow$ 0.67 |
| | Throughput: | 114 MB/s | 333 MB/s | $\rightarrow$ 358 MB/s | 683 MB/s | $\rightarrow$ 1.4 GB/s |

With 100 KB records, higher rate **and** throughput in exchange for larger public parameters

# Packing in the Streaming Setting

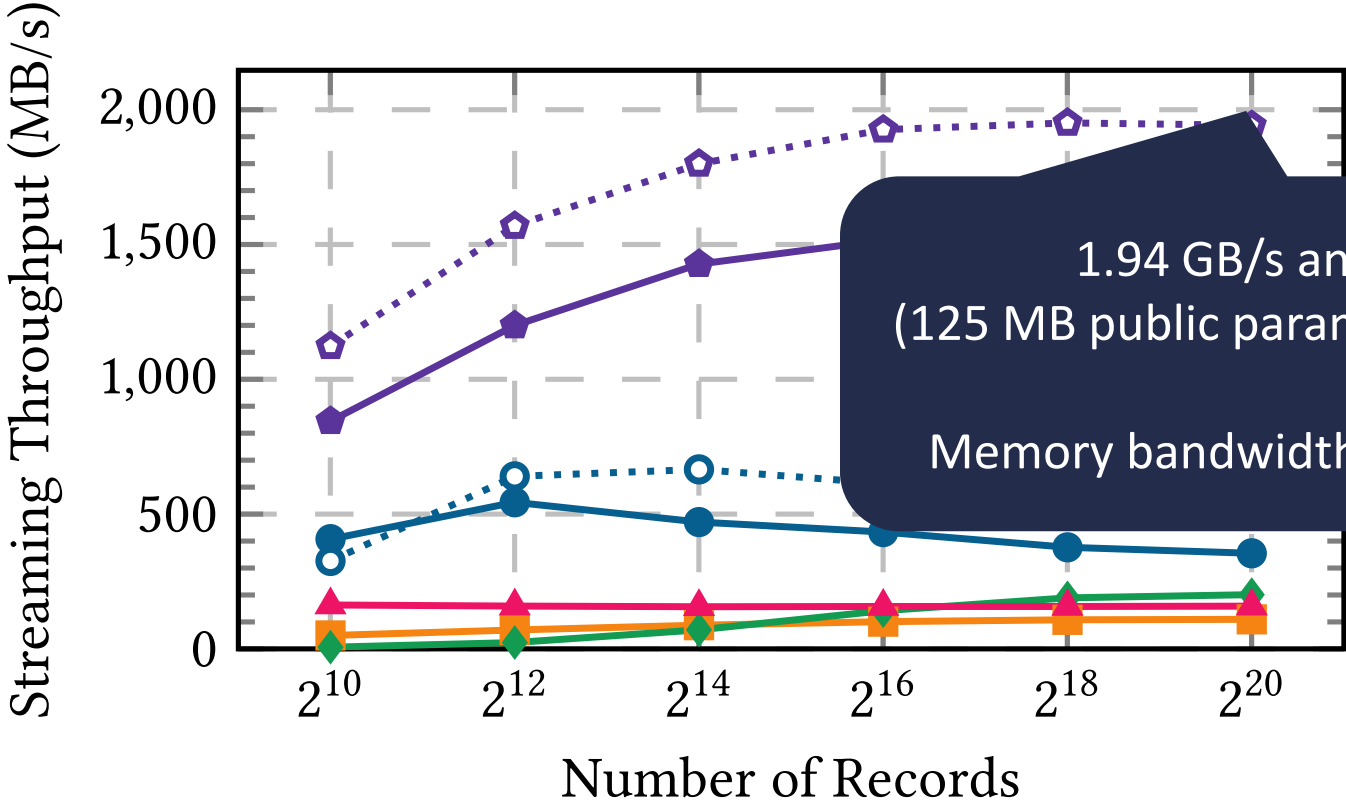**Streaming throughput:** ignoring query expansion costs, assuming optimal record size for each system



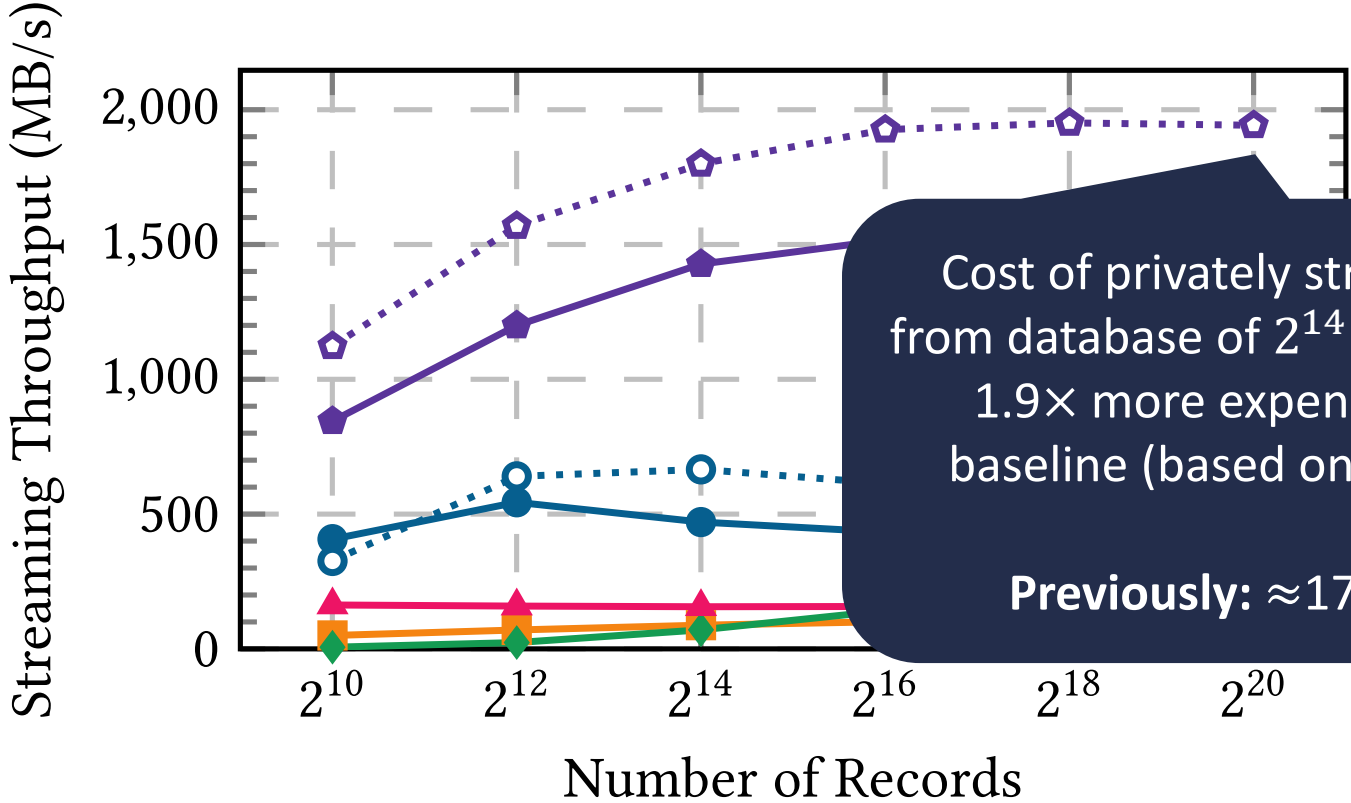Packing outperforms non-packed protocol for streaming settings

# Packing in the Streaming Setting

**Streaming throughput:** ignoring query expansion costs, assuming optimal record size for each system



Packing outperforms non-packed protocol for streaming settings

1.94 GB/s and a rate of 0.81
(125 MB public parameter and 30 MB query)

Memory bandwidth on system: ≈10 GB/s

# Packing in the Streaming Setting

**Streaming throughput:** ignoring query expansion costs, assuming optimal record size for each system



Packing outperforms non-packed protocol for streaming settings

Cost of privately streaming a 2 GB movie from database of $2^{14}$ movies estimated to be 1.9× more expensive than <u>no-privacy</u> baseline (based on AWS compute costs)
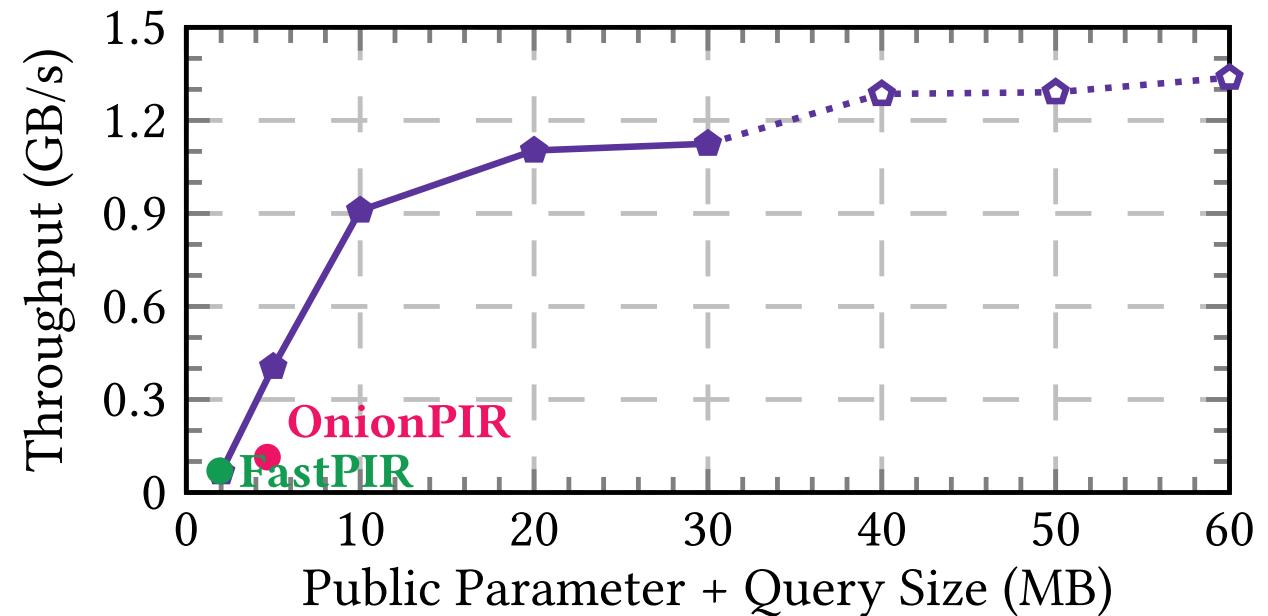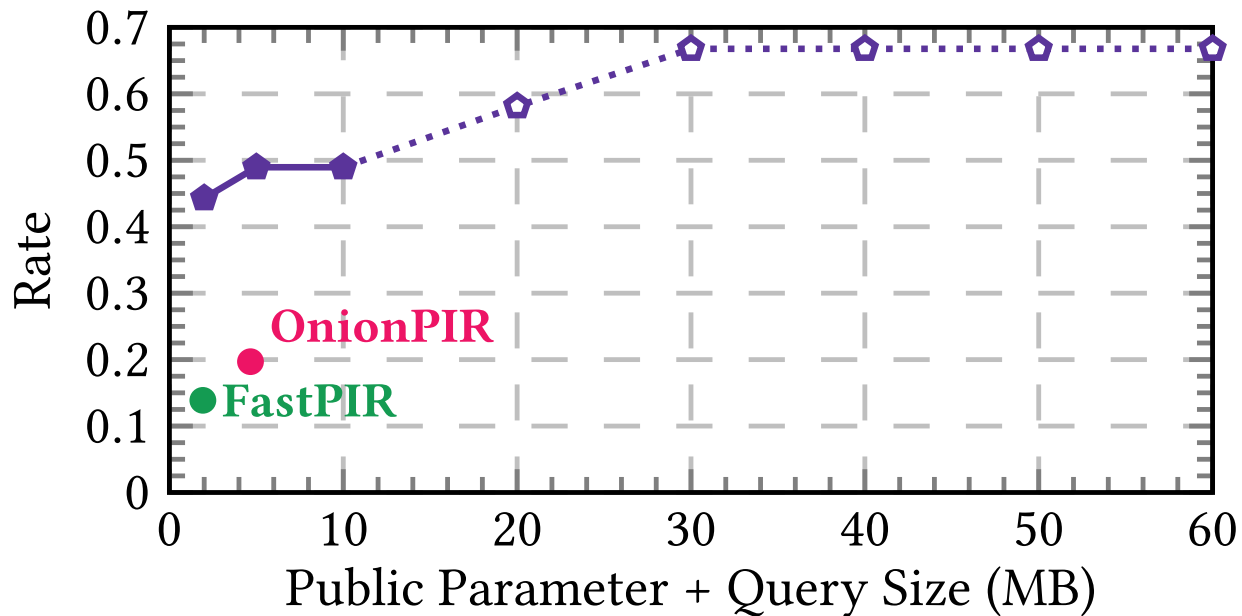
**Previously:** ≈17× more expensive

Legend:
- SPIRAL
- SPIRALPACK
- SPIRALSTREAM
- SPIRALSTREAMPACK
- SealPIR
- FastPIR
- OnionPIR

# A Systematic Way to Explore PIR Trade-Offs

Parameter selection tool can be used to minimize online cost
with constraints on public parameter and query sizes

(Database configuration: $2^{14} \times 100$ KB database)

# The SPIRAL Family of PIR

Techniques to translate between FHE schemes enables new trade-offs in single-server PIR

Scalar Regev → Matrix Regev
Regev → GSW

Query compression

Scalar Regev → Matrix Regev

Response compression
(for large records)

Automatic parameter selection to choose lattice parameters based on database configuration

## Base version of SPIRAL

| | | |
|---|---|---|
| **Query size:** | 14 KB | 4.5× smaller |
| **Rate:** | 0.41 | 2.1× higher |
| **Throughput:** | 333 MB/s | 2.9× higher |

(Database with $2^{14}$ records of size 100 KB)

## Streaming versions of SPIRAL

| | | |
|---|---|---|
| **Rate:** | 0.81 | 3.4× smaller |
| **Throughput:** | 1.9 GB/s | 12.3× higher |

Improvements primarily due to query and response compression

# The SPIRAL Family of PIR

Techniques to translate between FHE schemes enables new trade-offs in single-server PIR

Scalar Regev → Matrix Regev
Regev → GSW

Query compression

Scalar Regev → Matrix Regev

Response compression
(for large records)

Automatic parameter selection to choose lattice parameters based on database configuration

**Base version of SPIRAL**

| **Query size:** | 14 KB | 4.5× smaller |
| **Rate:** | 0.41 | 2.1× higher |
| **Throughput:** | 333 MB/s | 2.9× higher |

(Database with $2^{14}$ records of size 100 KB)

**Streaming versions of SPIRAL**

| **Rate:** | 0.81 | 3.4× smaller |
| **Throughput:** | 1.9 GB/s | 12.3× higher |

Improvements primarily due to fine-tuning scheme parameters for database configuration

# Future Directions

Leveraging FHE composition in other privacy-preserving systems

    Private set intersection (PSI)

    Oblivious RAM (ORAM)

Hardware acceleration for higher throughput

Leveraging preprocessing to achieve <u>sublinear</u> server computation

**Paper:** `https://eprint.iacr.org/2022/368`
**Code:** `https://github.com/menonsamir/spiral`

# Thank you!