

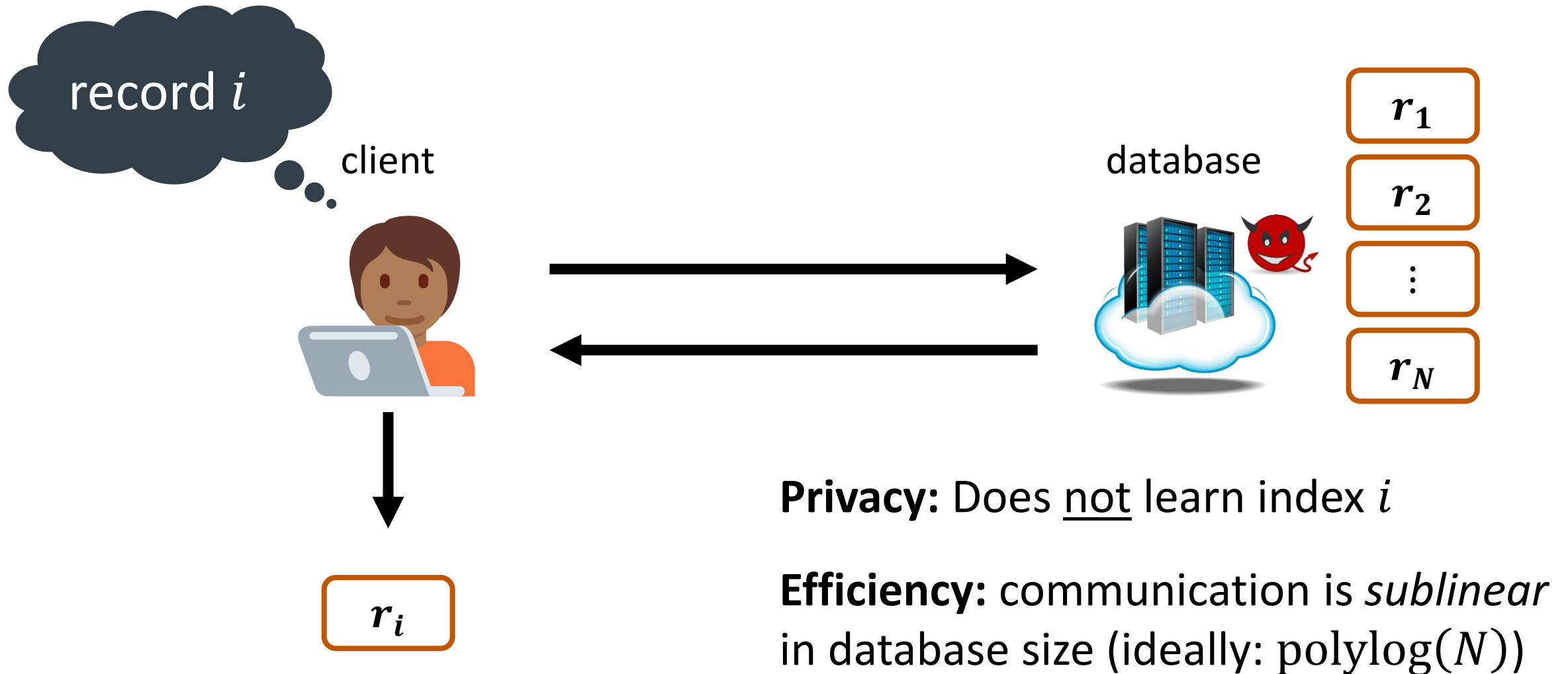
SPIRAL: Fast High-Rate Single-Server Private Information Retrieval

Samir Menon and David Wu

July 2023

Private Information Retrieval (PIR)

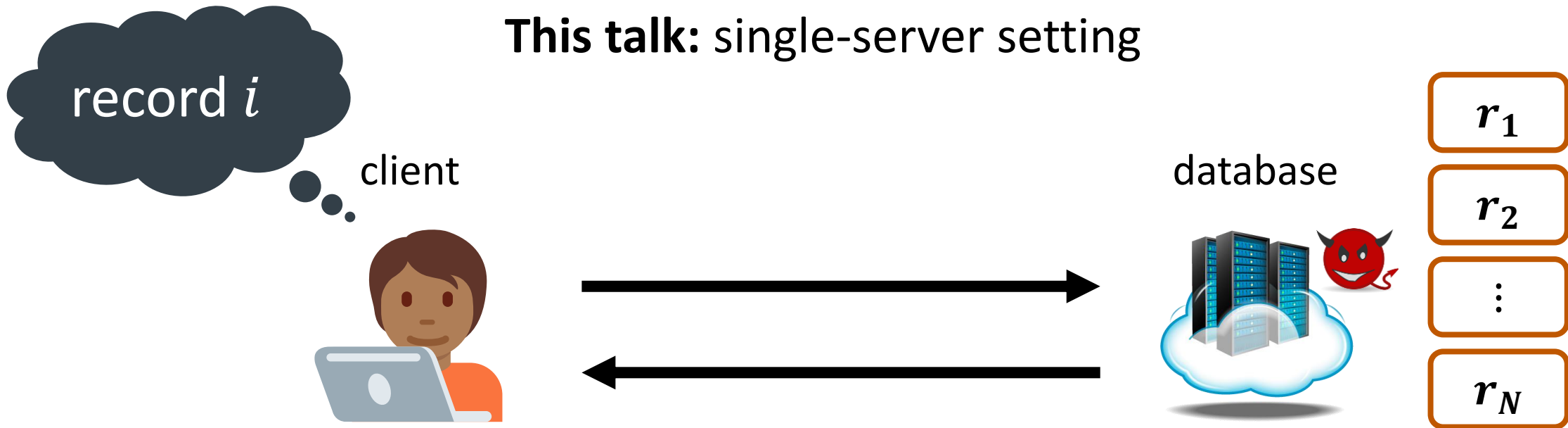
[CGKS95]



Private Information Retrieval (PIR)

[CGKS95]

This talk: single-server setting



Basic building block in many privacy-preserving protocols

- 💬 Metadata-private messaging
- 🌐 Private DNS
- 👤 Contact discovery
- 🙊 Private contact tracing
- 🎬 Private content delivery
- 🧭 Private navigation

PIR Metrics of Interest

Rate: *“measures communication overhead in responses”*

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

Throughput: *“measures how fast the server can respond as a function of database size”*

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

Constructions do not rely on server preprocessing so server always performs a linear scan

Goal: make linear scan as fast as possible

The SPIRAL Family of PIR Protocols

Leverages techniques to translate between homomorphic encryption schemes

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)

Independent of database and query and can be reused across multiple queries

Comparisons against schemes that do not require server preprocessing (i.e., server hints)

In particular, these exclude subsequent schemes such as FrodoPIR [DPC23], SimplePIR [HHCMV23], and Piano [ZPSZ23]

The SPIRAL Family of PIR Protocols

Leverages techniques to translate between homomorphic encryption schemes

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)

Independent of database and query and can be reused across multiple queries

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

PIR from Homomorphic Encryption

[K097]

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

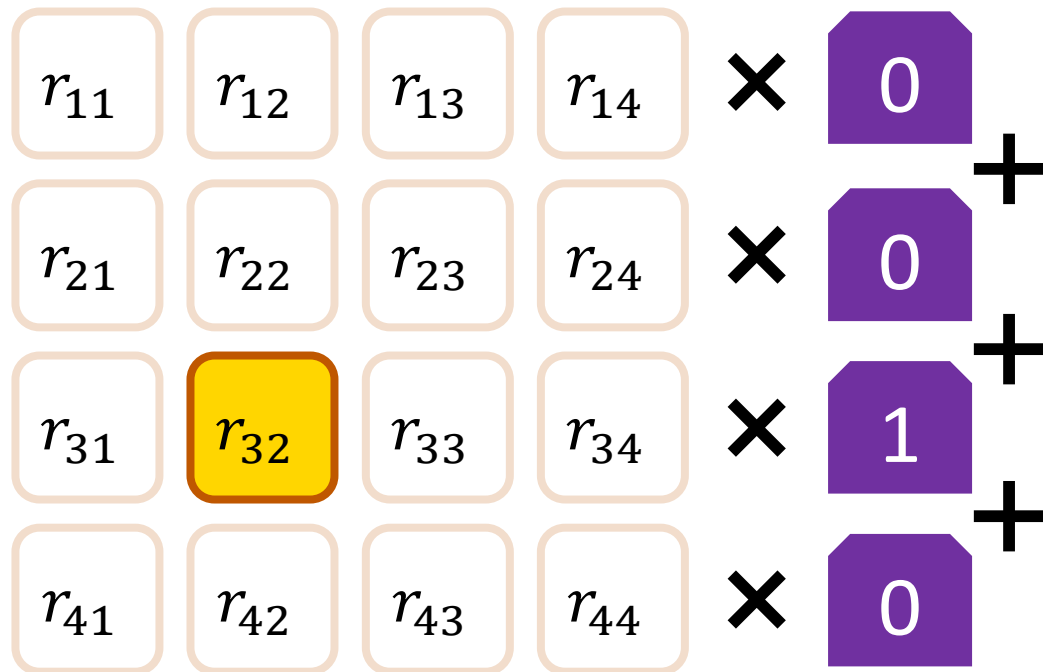
r_{11}	r_{12}	r_{13}	r_{14}
r_{21}	r_{22}	r_{23}	r_{24}
r_{31}	r_{32}	r_{33}	r_{34}
r_{41}	r_{42}	r_{43}	r_{44}

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

PIR from Homomorphic Encryption

[K097]

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



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

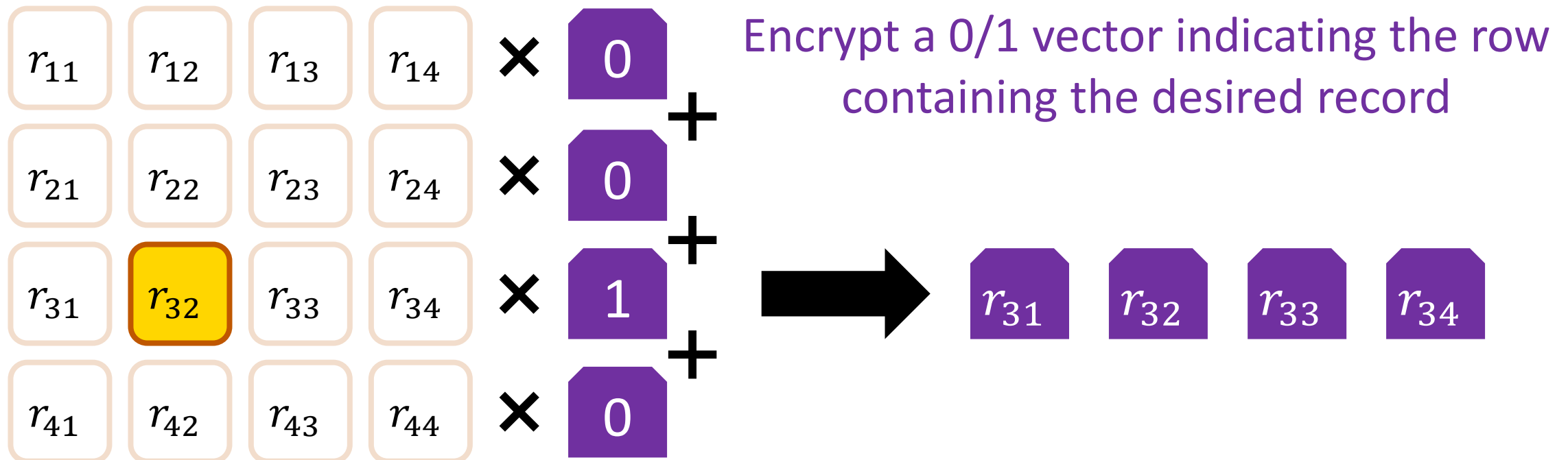
Arrange the database as a \sqrt{N} -by- \sqrt{N} matrix

Homomorphically compute product between query vector and database matrix

PIR from Homomorphic Encryption

[K097]

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



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

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

PIR from Homomorphic Encryption

[K097]

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



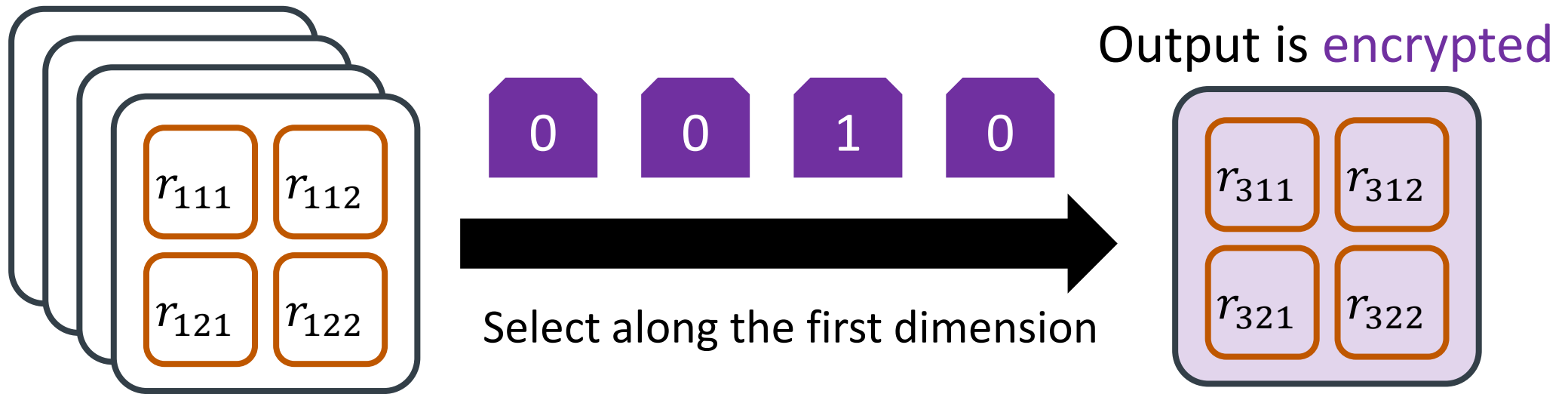
Response size: $O_\lambda(\sqrt{N})$

*Homomorphically compute product
between query vector and database matrix*

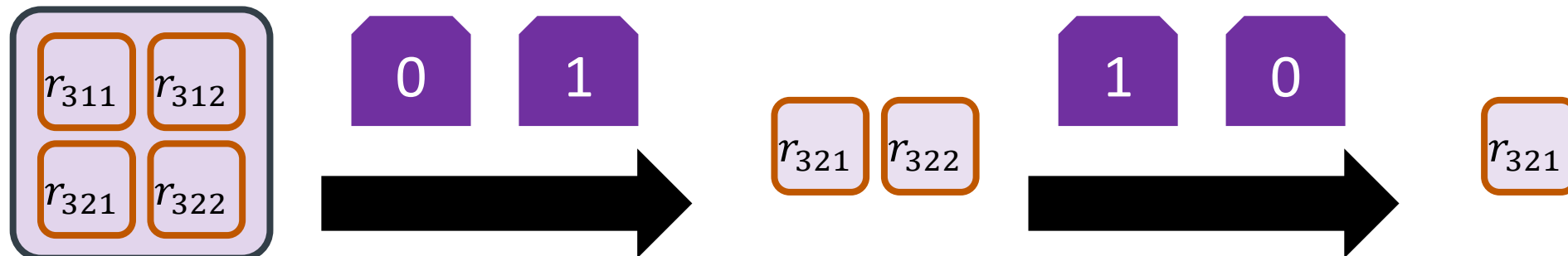
PIR from Homomorphic Encryption

[K097]

Sub- \sqrt{N} communication: view the database as **hypercube**



Approach: Use **homomorphic multiplication** [GH19, PT20, ALPRSSY21, MCR21]



SPIRAL: Composing FHE Schemes

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

Ciphertexts in lattice-based schemes are noisy encodings

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

Scheme 1: Regev's encryption scheme [Reg04]

Small ciphertexts (amortized); only supports additive homomorphism

18 KB plaintext \Rightarrow 43 KB ciphertext (2.4 \times expansion)

1 MB plaintext \Rightarrow 1.3 MB ciphertext (1.3 \times expansion)

allows the use of
smaller lattice
dimension and modulus

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

Large ciphertexts; supports homomorphic multiplication (with additive noise growth)

1 bit plaintext \Rightarrow 2.5 **MB** ciphertext

Can we get the best of both worlds?

SPIRAL: Composing FHE Schemes

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

Ciphertexts in lattice-based schemes are noisy encodings

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

Scheme 1: Regev's encryption scheme [Reg04]

Small ciphertexts (amortized); only supports additive homomorphism

18 KB plaintext \Rightarrow 43 KB ciphertext (2.4 \times expansion)

1 MB plaintext \Rightarrow 1.3 MB ciphertext (1.3 \times expansion)

allows the use of
smaller lattice
dimension and modulus

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

Large ciphertexts; supports homomorphic multiplication (with additive noise growth)

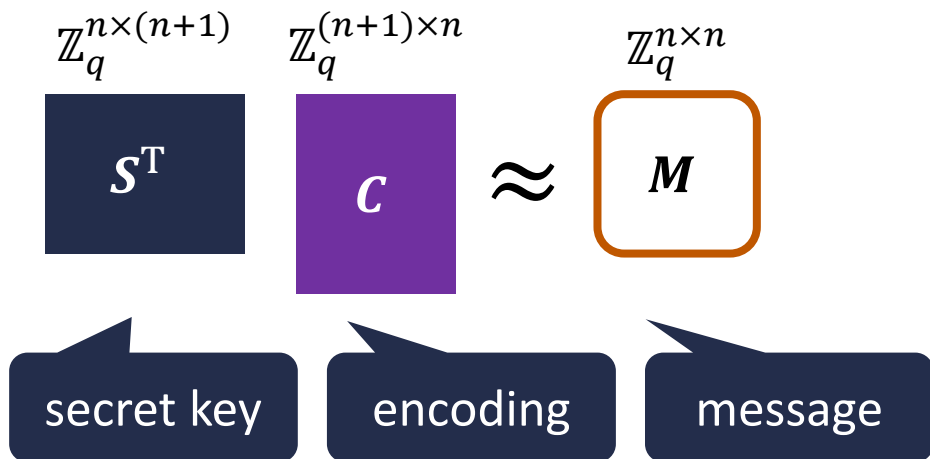
1 bit plaintext \Rightarrow 2.5 **MB** ciphertext

SPIRAL: Use GSW for homomorphic multiplication, Regev for communication

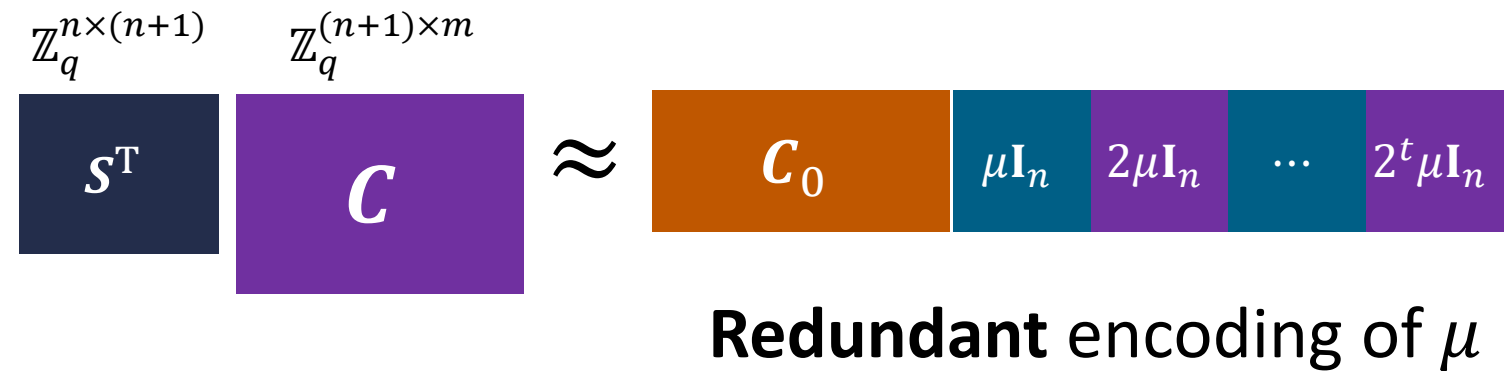
Regev-GSW Homomorphisms

[CGGI18]

Regev encoding of $M \in \mathbb{Z}_q^{n \times n}$:



GSW encoding of $\mu \in \mathbb{Z}_q$:

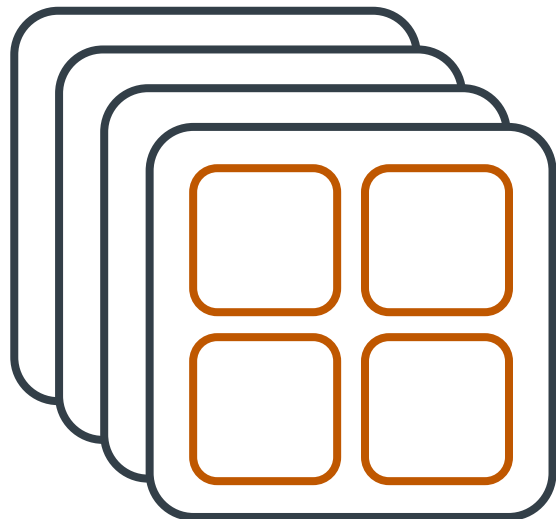


Key property: given Regev encoding of message M and GSW encoding of scalar μ , can efficiently derive a Regev encoding of $\mu \cdot M$

Simplification: actual scheme operates over polynomial rings

The Gentry-Halevi Blueprint

[GH19]



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

Query contains 2^{ν_1} Regev encodings



Technically
uses matrices

Indicator for index along first dimension

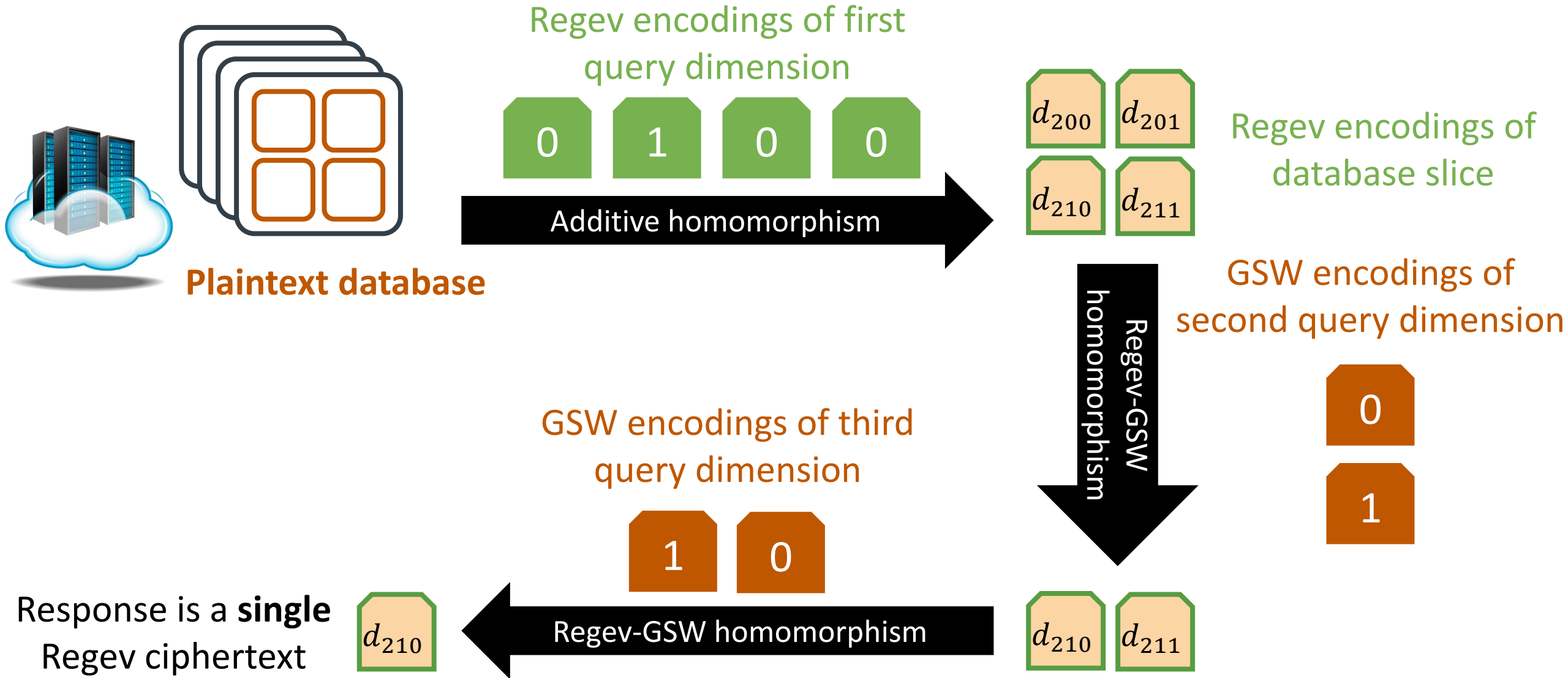
Query contains ν_2 GSW encodings



Indicator for index along subsequent dimensions

The Gentry-Halevi Blueprint

[GH19]



The Gentry-Halevi Blueprint

[GH19]

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

Drawback: large queries

Query contains 2^{ν_1} Regev encodings



Indicator for index along first dimension

Query contains ν_2 GSW encodings



Indicator for index along subsequent dimensions

Estimated size:
4 MB/ciphertext

Estimated query size:
30 MB

The Gentry-Halevi Blueprint

[GH19]

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

Drawback: large queries

Query contains 2^{ν_1} Regev encodings



Indicator for index along first dimension

Query contains ν_2 GSW encodings



Indicator for index along subsequent dimensions

SealPIR query size:
66 KB

Estimated query size:
30 MB

The SPIRAL Protocol

Key idea: Expand Regev encodings into GSW encodings

OnionPIR [MCR21]: use Regev-GSW homomorphism for the **scalar** case

This work: General toolkit to translate between scalar/matrix Regev and GSW

Transformations useful for **query compression** and **response packing**

Assembling GSW Encodings

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

$$\mu \mathbf{S}^T \mathbf{G} = \begin{array}{|c|c|c|c|c|c|} \hline \mathbf{C}_0 & \mu \mathbf{I}_n & 2\mu \mathbf{I}_n & 2^2 \mu \mathbf{I}_n & \dots & 2^t \mu \mathbf{I}_n \\ \hline \end{array}$$

$$\mathbf{C} = \begin{array}{|c|c|c|c|c|c|} \hline \mathbf{A} & \mathbf{B}_0 & \mathbf{B}_1 & \mathbf{B}_2 & \dots & \mathbf{B}_t \\ \hline \end{array}$$

Break \mathbf{C} into *blocks*

Assembling GSW Encodings

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

$$\mu \mathbf{S}^T \mathbf{G} = \left[\mathbf{C}_0 \mid \mu \mathbf{I}_n \mid 2\mu \mathbf{I}_n \mid 2^2 \mu \mathbf{I}_n \mid \dots \mid 2^t \mu \mathbf{I}_n \right]$$

\approx \approx

$$\mathbf{S}^T \mathbf{C} = \left[\mathbf{S}^T \mathbf{A} \mid \mathbf{S}^T \mathbf{B}_0 \mid \mathbf{S}^T \mathbf{B}_1 \mid \mathbf{S}^T \mathbf{B}_2 \mid \dots \mid \mathbf{S}^T \mathbf{B}_t \right]$$

Break \mathbf{C} into *blocks*

Leverage “key-switching”

Standard Regev encodings of $\mu, 2\mu, \dots, 2^t \mu$

Query Compression in SPIRAL

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

Query contains 2^{ν_1} matrix Regev encodings



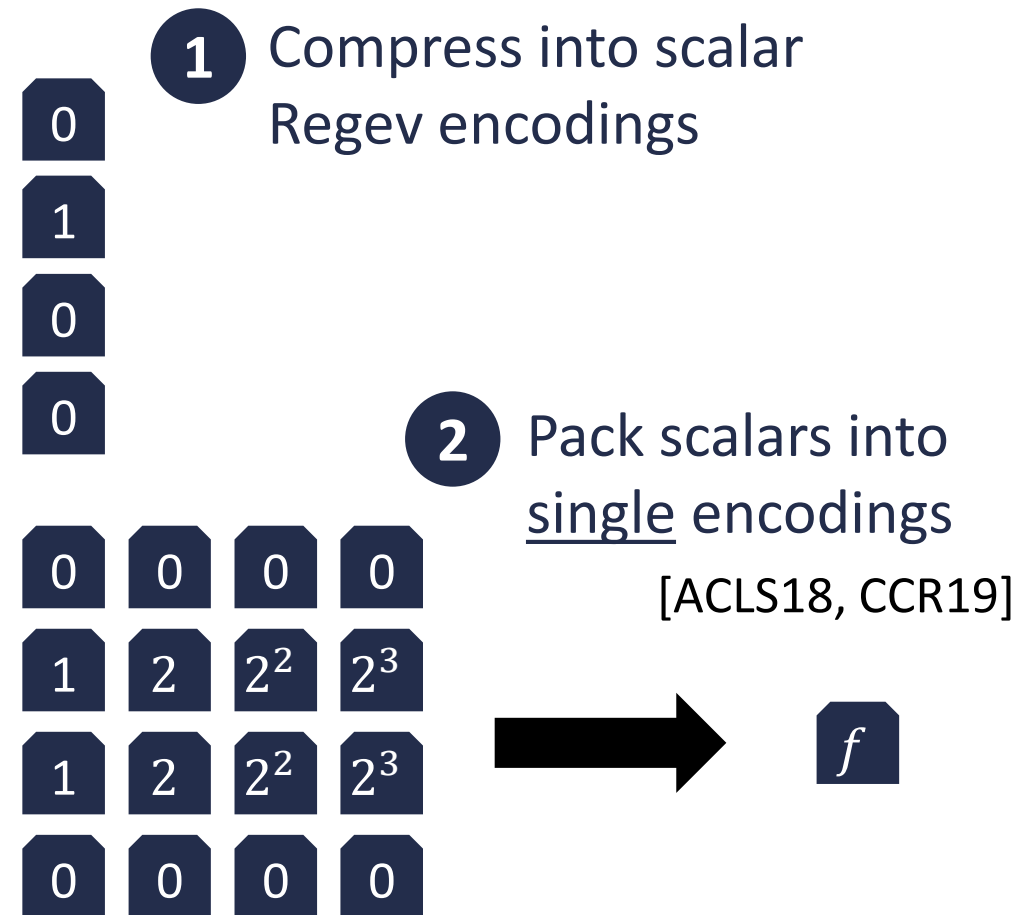
Indicator for index along first dimension

Query contains ν_2 GSW encodings



Indicator for index along subsequent dimensions

Similar techniques possible for *response* compression [see paper]



The SPIRAL Protocol

record i



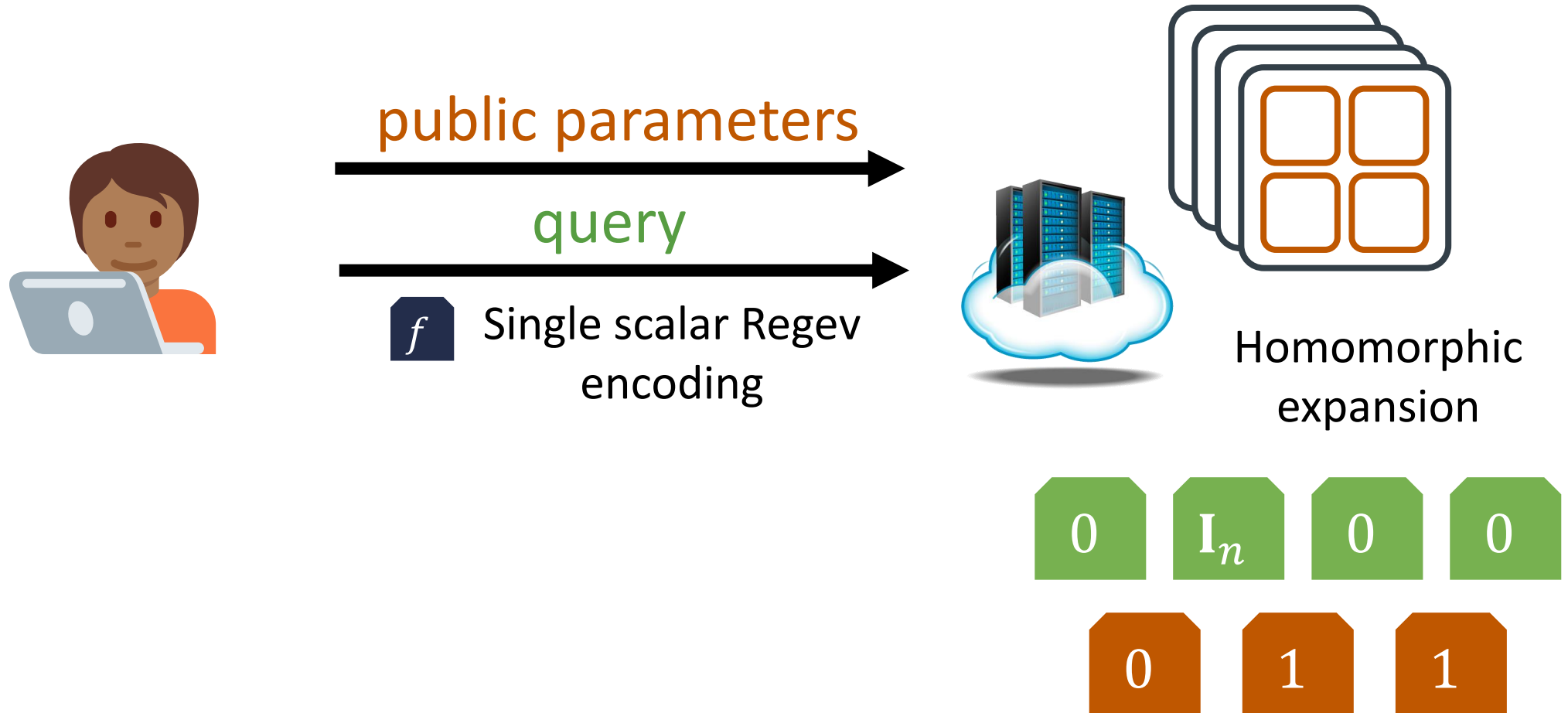
public parameters



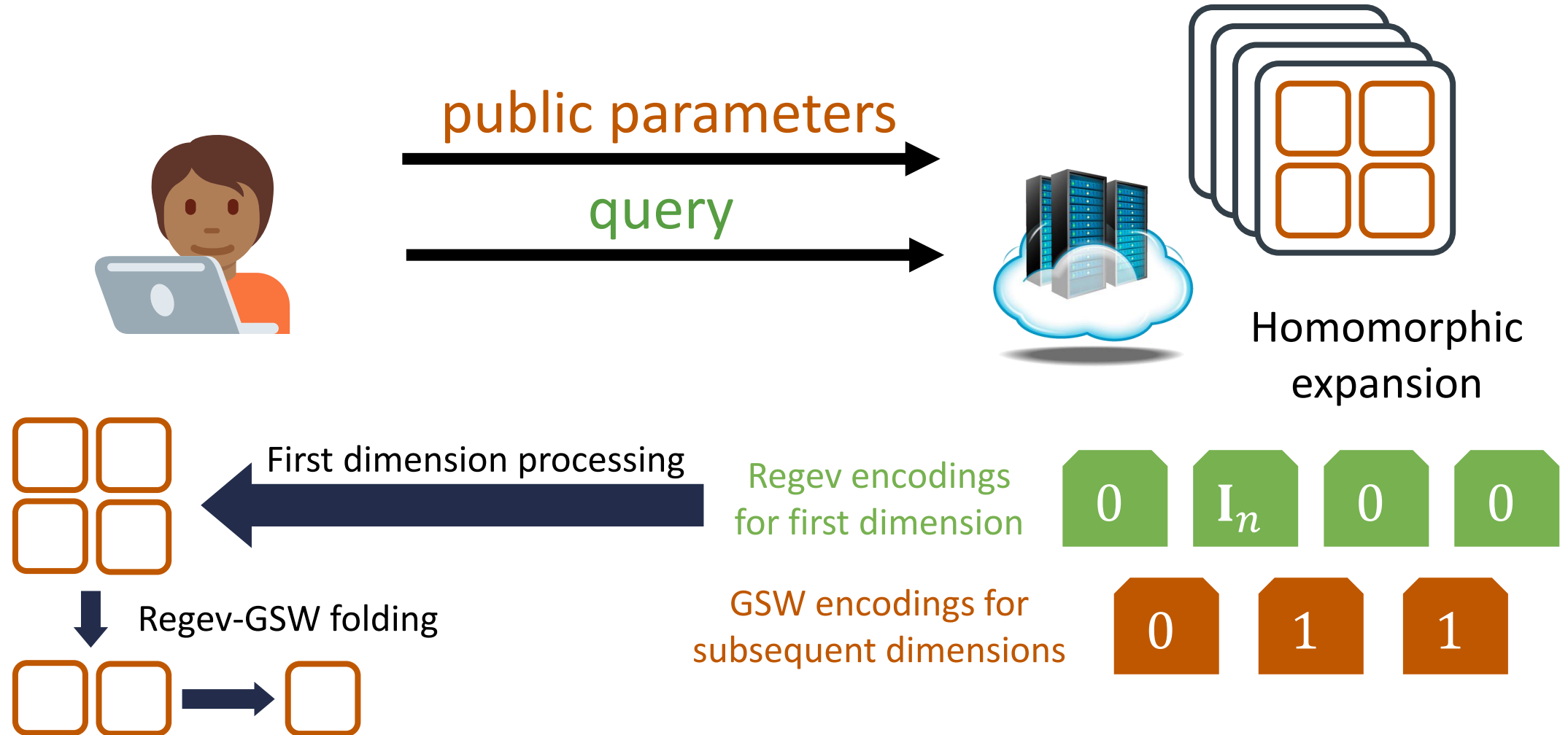
Key-switching matrices for
ciphertext expansion and
translation



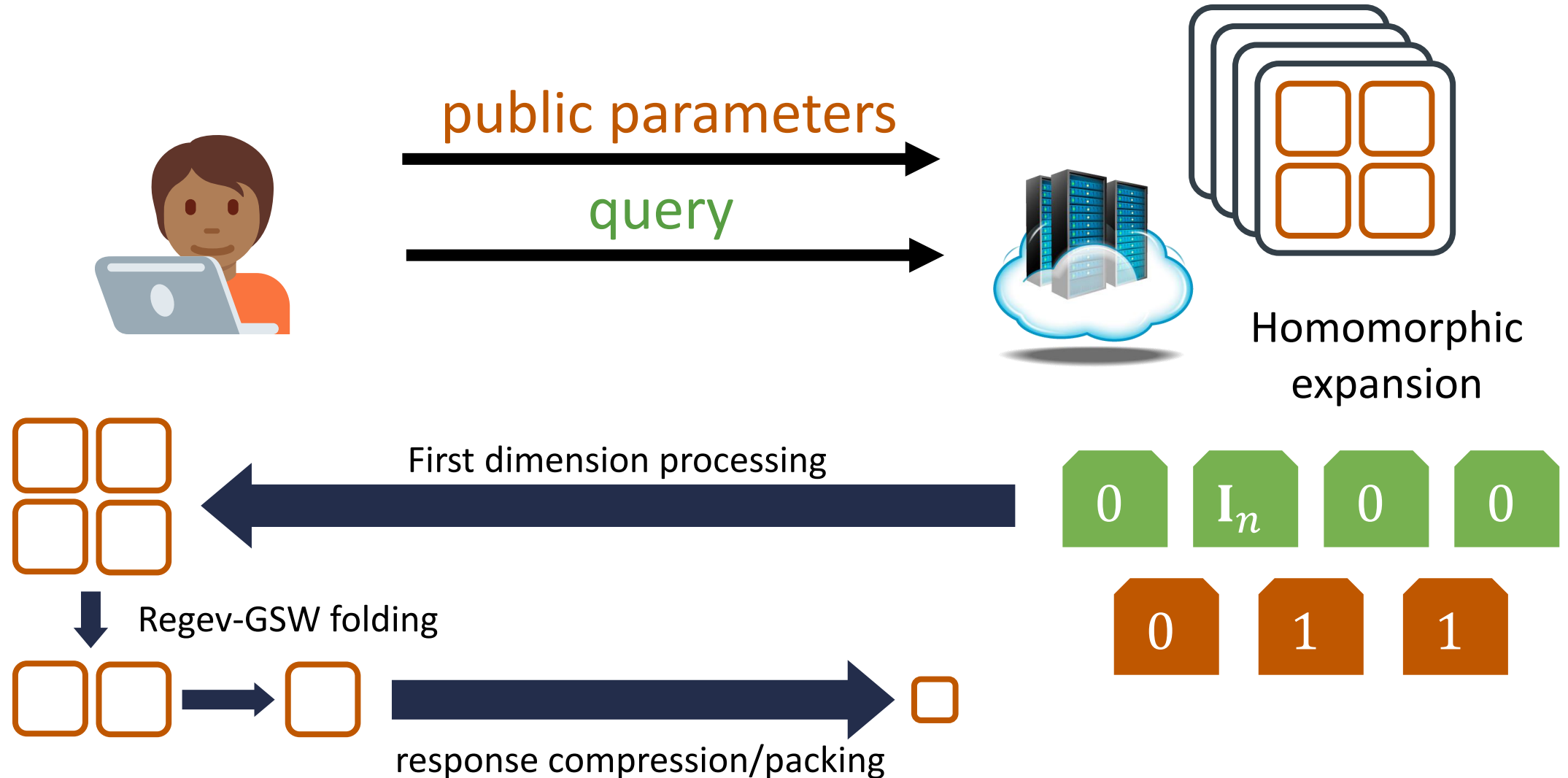
The SPIRAL Protocol



The SPIRAL Protocol



The SPIRAL Protocol



Basic Comparisons

Database	Metric	SealPIR	FastPIR	OnionPIR	SPIRAL
2¹⁸ 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×

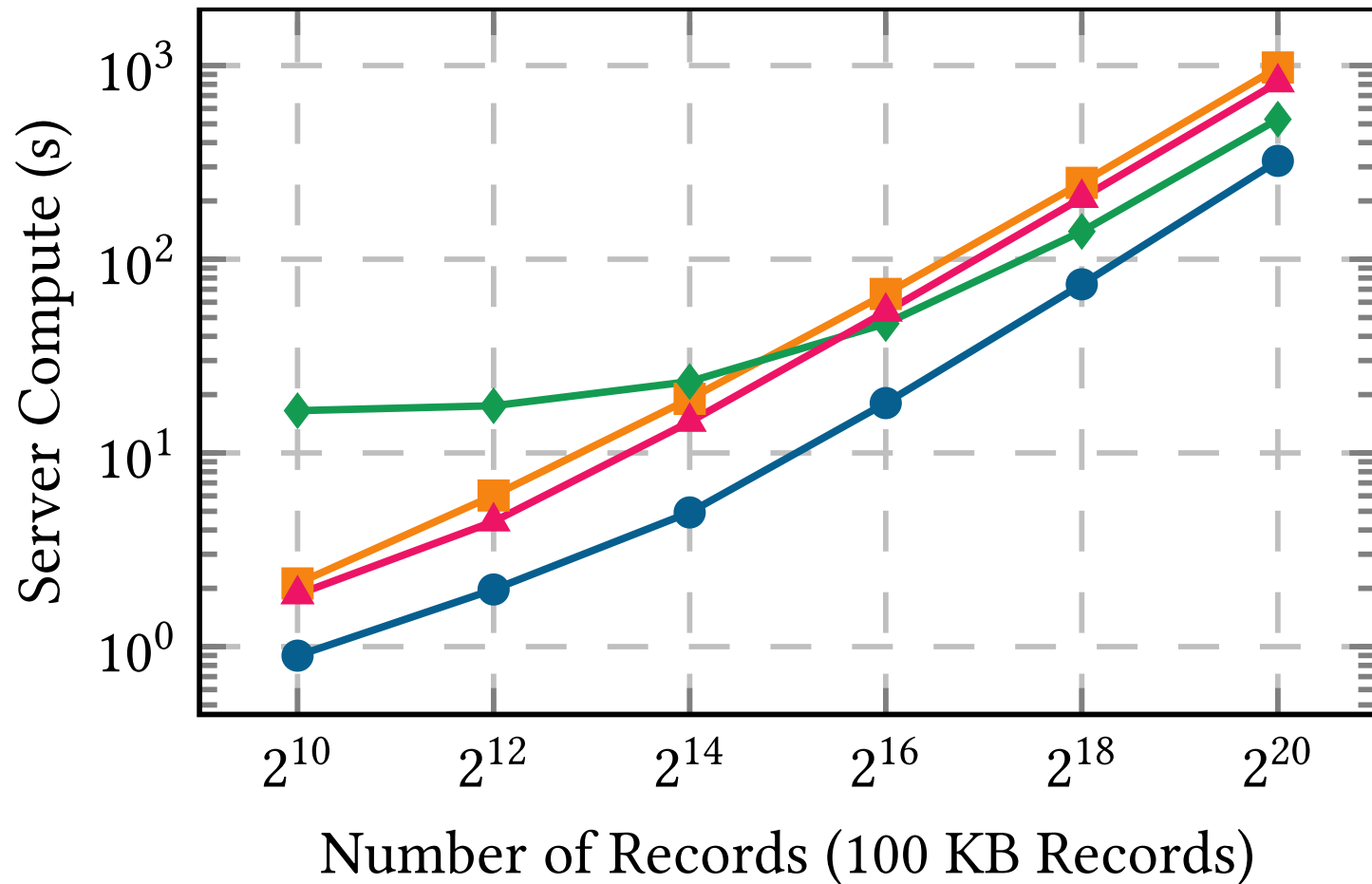
reduce response size by 2×

reduce compute time by 2×

Comparisons against schemes that do not require server preprocessing (i.e., server hints)

In particular, these exclude subsequent schemes such as FrodoPIR, SimplePIR, and Piano

Basic Comparisons (with Large 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

—●— SPIRAL —■— SealPIR —◆— FastPIR —▲— OnionPIR

The Streaming Setting

Streaming setting: same query reused over multiple databases

Private video stream (database D_i contains i^{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)

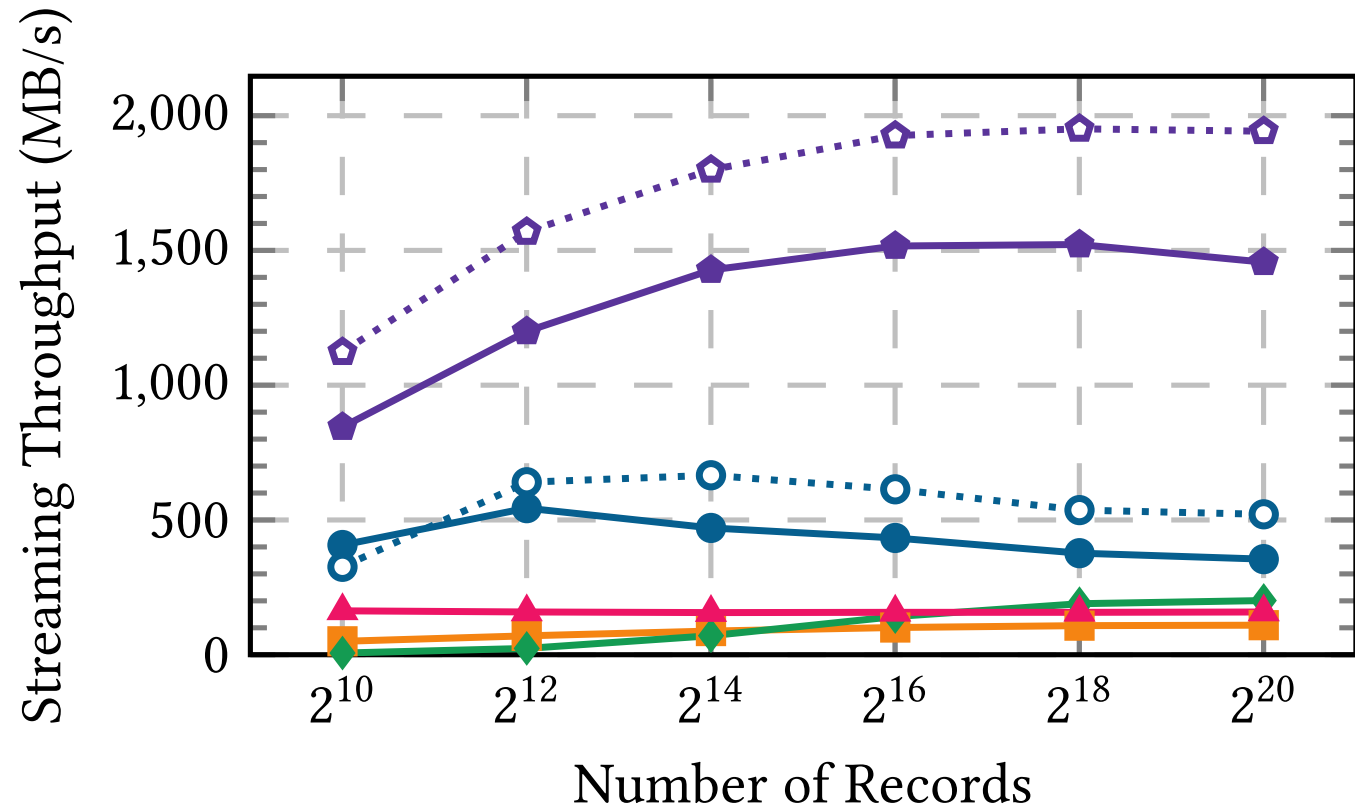
Consider larger public parameters or query size (amortized over lifetime of stream)

Approach: send all of the Regev encodings (and only use Regev-GSW translation)

The Streaming Setting

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

Packing outperforms non-packed protocol for streaming settings

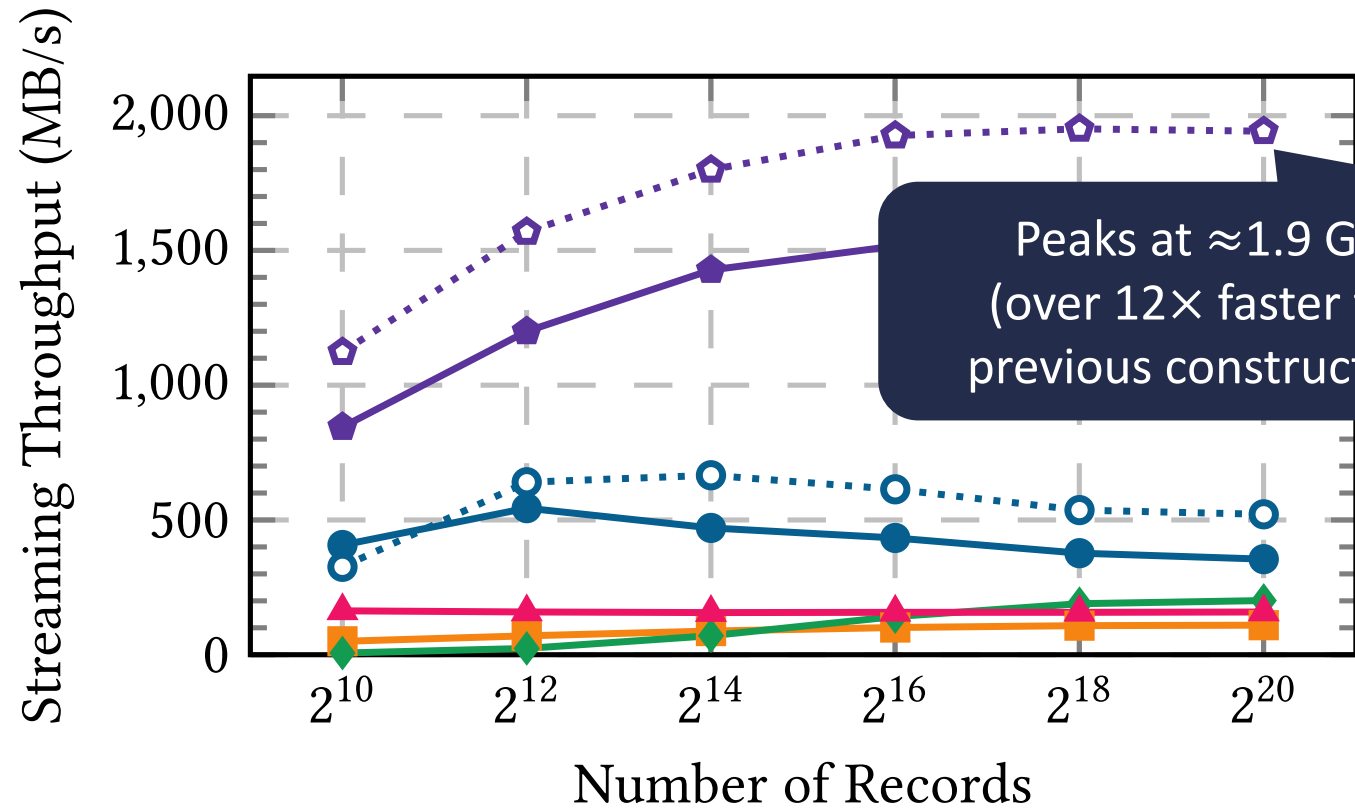


Packed versions rely on response compression (larger public parameters, higher throughput)



The Streaming Setting

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



Packing outperforms non-packed protocol for streaming settings

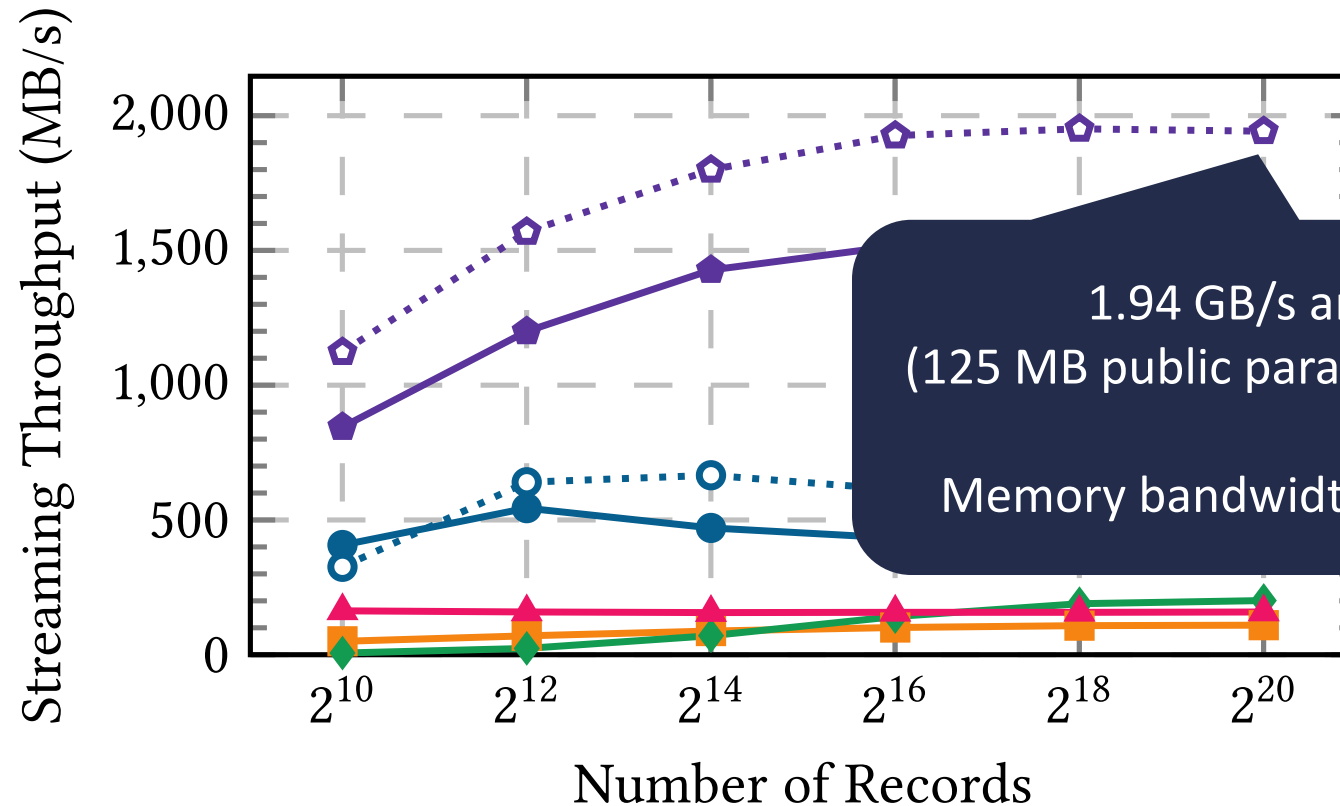
Peaks at ≈ 1.9 GB/s (over 12 \times faster than previous constructions)

Packed versions rely on response compression (larger public parameters, higher throughput)



The Streaming Setting

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



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

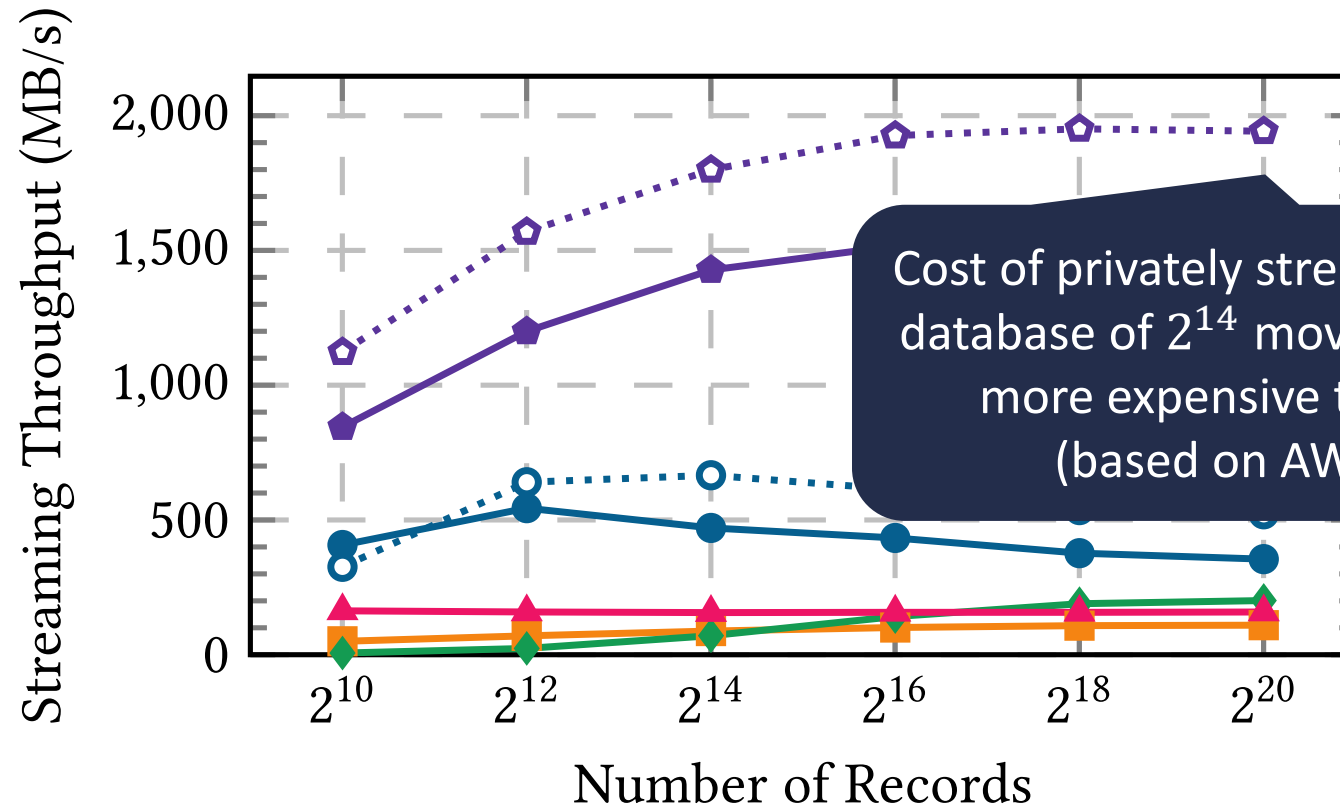
Memory bandwidth on system: ≈ 10 GB/s

Packing outperforms non-packed protocol for streaming settings



The Streaming Setting

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



Cost of privately streaming a 2 GB movie from database of 2^{14} movies estimated to be 1.9× more expensive than direct download (based on AWS compute costs)

Packing outperforms non-packed protocol for streaming settings

- SPIRAL
- SPIRALPACK
- ◆ SPIRALSTREAM
- ◇ SPIRALSTREAMPACK
- SealPIR
- ◆ FastPIR
- ▲ OnionPIR

The SPIRAL Family of PIR

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

Used for both query compression and response compression

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

Streaming versions of SPIRAL

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

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

Some Recent Developments in PIR

Server preprocessing (client downloads hint at beginning of protocol)

SimplePIR, DoublePIR [HHCMV23]

Very high throughput (nearly memory bandwidth!)

Well suited for databases with small records (a few bits)

Piano [ZPSZ23]

Sublinear server computational costs (can scale better to databases that are >100 GB)

Preprocessing phase requires *streaming* the entire database

Server preprocessing (*without* hint)

Doubly-efficient PIR [LMW23]

Server encodes the database to answer queries in sublinear time

Concrete efficiency not yet clear

Some Recent Developments in PIR

Server preprocessing (client downloads hint at time of query)

SimplePIR, DoublePIR [HHCMV23]

Very high throughput (nearly memory bandwidth)

Well suited for databases with small records

Piano [ZPSZ23]

Sublinear server computational costs (can scale better to databases that are >100 GB)

Preprocessing phase requires *streaming* the entire database

Server preprocessing (*without* hint)

Doubly-efficient PIR [LMW23]

Server encodes the database to answer queries

Concrete efficiency not yet clear

Many other directions!

- Protocols for batch queries [MR23]
- Supporting keyword search [PSY23]
- Authenticating the response [CNCWF23]

Takeaway: PIR is an exciting area to work in with many different trade-offs to explore!

The SPIRAL Family of PIR

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

Used for both query compression and response compression

Automatic parameter selection to choose lattice parameters based on database configuration

Paper: <https://eprint.iacr.org/2022/368>

Code: <https://github.com/menonsamir/spiral-rs>

Demo: <https://spiralwiki.com>

Thank you!