# Relational Database Encryption

Zheguang Zhao

Computer Science Department, Brown University, Rhode Island
zheguang.zhao@gmail.com

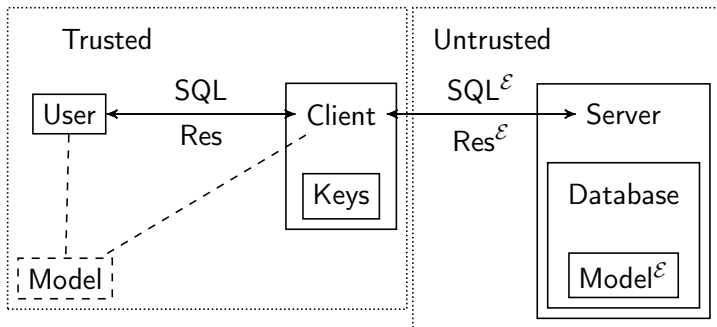Thesis Proposal, 12th December, 2019

# Part I

## Introduction

# Introduction

# 3 Problems

- Signal Problem: Secure application on untrusted OS
- Equifax Problem: Data within the private cloud
- Amazon Problem: Applications running on the public cloud

# Fundamental Problem: Relational Database Encryption



- Security
- Functionality
- Efficiency
- Legacy compliance

# Related Works

| Scheme | Leak. | Eff. | Legacy | Func. |
|---|---|---|---|---|
| FHE [14] | Less | Low | No | IntComp. |
| ORAM [11] | Less | Low | No | EqSel,Prefix,Range |
| PPE [6] | More | High | Y/N | $\subseteq$SQL(?) |
| STE [7, 10, 12] | Less | High | No | KeyVal,EqSel,EqJn,Prefix,Range |

Table: General observation in schemes based on various symmetric searchable encryption primitives.

## Thesis Problem

Use STE to construct secure, legacy-compliant, functional and efficient relational database encryption scheme.
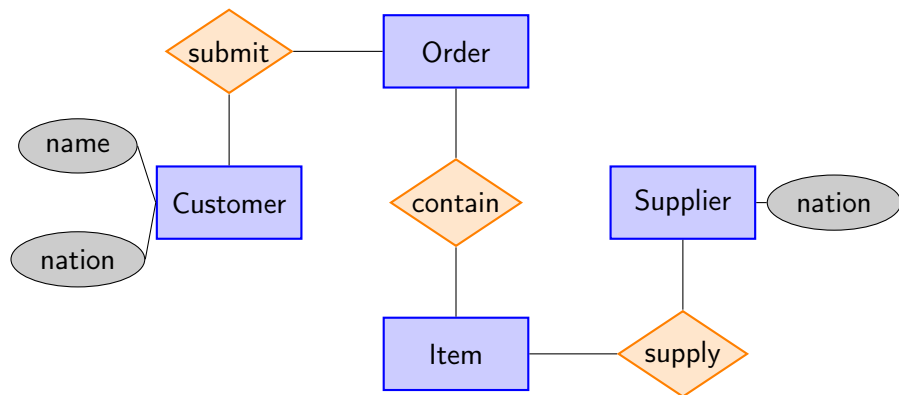
# Relational Model



- Contains entities and relationships
- First-order logic

# Cryptographic Primitives

- Cryptographic hash function $\mathcal{H}$
- Pseudorandom function $\mathcal{F}$
- Pseudorandom permutation $\mathcal{P}$
- Psuedorandom ciphertext under chosen-plaintext attack (RCPA-secure symmetric encryption $\mathcal{E}$)
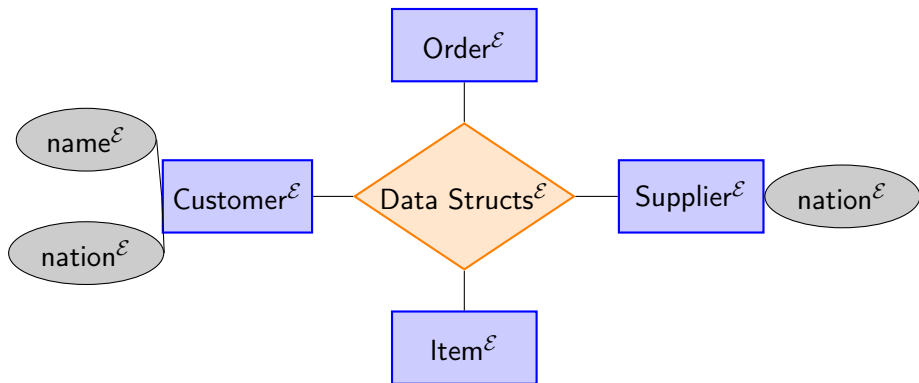- Switching Lemma ($\mathcal{F} \approx \mathcal{P}$)

# Introduction

# Encrypted Data Structures

## Security against Leakage

The scheme $\mathbf{dex} = (\mathbf{Setup}, \mathbf{ProcessQuery})$ is $\mathcal{L}$-secure if $\forall$ efficient A $\exists$ efficient S such that

$$\left| \Pr\left( \mathsf{Real}^A_{\mathbf{dex}}(\lambda) = 1 \right) - \Pr\left( \mathsf{Ideal}^A_{\mathbf{dex}, \mathcal{L}, S}(\lambda) = 1 \right) \right| \leq \mathsf{negl}(\lambda)$$

| $\mathbf{Real_{dex}}(\lambda)$ | $\mathbf{Ideal_{dex, \mathcal{L}, S}}(\lambda)$ |
|---|---|
| $\mathbf{Init}(\mathsf{RDB})$ | $\mathbf{Init}(\mathsf{RDB})$ |
|   ❶ $(\mathsf{ERDB}, k) \sim \mathbf{Setup}(\mathsf{RDB}, \lambda)$ |   ❶ $L \sim \mathcal{L}(\mathsf{RDB})$ |
|   ❷ return ERDB |   ❷ $\mathsf{ERDB} \sim S_{\mathbf{Init}}(L)$ |
| $\mathbf{Query}(Q_{\mathsf{RDB}})$ |   ❸ return ERDB |
|   ❶ $(\mathsf{Res}, Q_{\mathsf{ERDB}}) \sim \mathbf{ProessQuery}(k, Q_{\mathsf{RDB}})$ on ERDB | $\mathbf{Query}(Q_{\mathsf{RDB}})$ |
|   ❷ return $Q_{\mathsf{ERDB}}$ |   ❶ $L \sim \mathcal{L}(L, Q_{\mathsf{RDB}})$ |
| $\mathbf{Update}(Q_{\mathsf{RDB}})$ |   ❷ $(\mathsf{Res}, Q_{\mathsf{ERDB}}) \sim S_{\mathbf{Query}}(L)$ |
|   ❶ $(\mathsf{Res}, Q_{\mathsf{ERDB}}) \sim \mathbf{ProcessQuery}(k, Q_{\mathsf{RDB}})$ on ERDB |   ❸ return $Q_{\mathsf{ERDB}}$ |
|   ❷ return ERDB | $\mathbf{Update}(Q_{\mathsf{RDB}})$ |
| $\mathbf{Final}(b)$ |   ❶ $L \sim \mathcal{L}(L, Q_{\mathsf{RDB}})$ |
|   ❶ output $b$ |   ❷ $(\mathsf{Res}, Q_{\mathsf{ERDB}}) \sim S_{\mathbf{Update}}(L)$ |
| |   ❸ return $Q_{\mathsf{ERDB}}$ |
| | $\mathbf{Final}(b)$ |
| |   ❶ output $b$ |

# Part II

## Construction

# Construction

- Express data structure & computation in domain-specific language

---

[1]Or with transitive closure and windowing function

# Emulation

- Express data structure & computation in domain-specific language
- Varying expressibility
  1. SPC algebra: conjunctive queries
  2. Relational algebra: $\subseteq$ first-order logic
  3. Datalog: least fixed-point logic
  4. SQL: Turing-complete with procedral extension[1]

---

[1] Or with transitive closure and windowing function

# Emulation

- Express data structure & computation in domain-specific language
- Varying expressibility
  1. SPC algebra: conjunctive queries
  2. Relational algebra: $\subseteq$first-order logic
  3. Datalog: least fixed-point logic
  4. SQL: Turing-complete with procedral extension[1]
- Examples
  - Datalog, SQLite: no procedures.
  - Azure SQL Data Warehouse, SparkSQL: no transitive closure.

---

[1]Or with transitive closure and windowing function

# Emulation

- Express data structure & computation in domain-specific language
- Varying expressibility
    1. SPC algebra: conjunctive queries
    2. Relational algebra: $\subseteq$ first-order logic
    3. Datalog: least fixed-point logic
    4. SQL: Turing-complete with procedral extension[1]
- Examples
    - Datalog, SQLite: no procedures.
    - Azure SQL Data Warehouse, SparkSQL: no transitive closure.

## Problem

Minimum language for expressing an (encryted) data structure

---

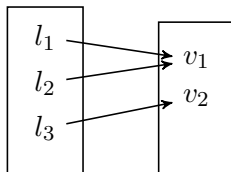[1]Or with transitive closure and windowing function

# Construction

## Map

A map is a binary relation of labels and values
$M : \{(l \to v) \mid l \in L, v \in V\}$ which associates each label to a unique value.
The operation $M[l] = v$ if $(l \to v) \in M$.

Example:



$M[l_2] = v_1$

Emulation:

| label | value |
|:-----:|:-----:|
| $l_1$ | $v_1$ |
| $l_2$ | $v_1$ |
| $l_3$ | $v_2$ |

$\pi_{\texttt{value}}\sigma_{\texttt{label}=l_2}T_M \implies \begin{bmatrix} v_1 \end{bmatrix}$

# Construction

# Multi-Map

## Multi-Map

A multi-map is a map[1] between labels and sets of values
$MM : \{(l \rightarrow V_l \mid l \in L, V_l \subseteq V)\}$. The operation $MM[l] = V_l$ if
$(l \rightarrow V_l) \in MM$.

Example:



$MM[l_1] = \{v_1, v_2\}$

---

[1]Elsewhere the term "map" will be reserved to a map that is *not* a multi-map.

### Multi-Map

A multi-map is a map[1] between labels and sets of values
$MM : \{(l \to V_l \mid l \in L, V_l \subseteq V)\}$. The operation $MM[l] = V_l$ if
$(l \to V_l) \in MM$.

Example:



$MM[l_1] = \{v_1, v_2\}$

Emulation Idea 1: Padding
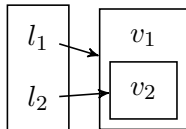
| label | value |
|-------|-------|
| $l_1$ | $v_1, v_2$ |
| $l_2$ | $v_2$ |

Bad: Set of sets is unrelational.

---

[1]Elsewhere the term "map" will be reserved to a map that is *not* a multi-map.

Emulation of Idea 2: Flattening

| label | value |
|:-----:|:-----:|
| $l_1$ | $v_1$ |
| $l_1$ | $v_2$ |
| $l_2$ | $v_2$ |

$$\pi_{\texttt{value}}\sigma_{\texttt{label}=l_1}T_{\mathsf{MM}} = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

Idea 2: Flattening



Bad: No longer a map[1].

Emulation of Idea 2: Flattening

| label | value |
|:-----:|:-----:|
| $l_1$ | $v_1$ |
| $l_1$ | $v_2$ |
| $l_2$ | $v_2$ |

$$\pi_{\texttt{value}}\sigma_{\texttt{label}=l_1}T_{\textsf{MM}} = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

---

[1]However it meets a non-map definition of multi-map.

Idea 3a: two maps

## Label transformation

$f : \mathbb{Z}^+ \times L \to L'$ is a one-to-one function.

Idea 3a: two maps

# Multi-Map

## Label transformation

$f : \mathbb{Z}^+ \times L \to L'$ is a one-to-one function.

Idea 3a: two maps



Idea 3b: three maps



The $i$ in $f(i, l)$ means the $i$th map of $L \to L'$.

# Multi-Map

- Knowing $f$, reduce to one map

- Idea 3c: One map

# Multi-Map

- Knowing $f$, reduce to one map
- Emulate as table $T_{\mathsf{MM}}$

- Idea 3c: One map



- $T_{\mathsf{MM}}$ for $\mathsf{M} : f \to V$

| label | value |
|---|---|
| $f(1, l_1)$ | $v_1$ |
| $f(1, l_2)$ | $v_2$ |
| $f(2, l_1)$ | $v_2$ |

# Multi-Map

- Knowing $f$, reduce to one map
- Emulate as table $T_{\mathsf{MM}}$
- Express $\mathsf{MM}[l]$ as union

  $\mathsf{MM}[l_1] \rightarrow \{\mathsf{M}[f(1,l_1)]\} \cup \{\mathsf{M}[f(2,l_1)]\}$

  $\mathsf{MM}[l_2] \rightarrow \{\mathsf{M}[f(1,l_2)]\}$

- Idea 3c: One map



- $T_{\mathsf{MM}}$ for $\mathsf{M} : f \rightarrow V$

| label | value |
|-------|-------|
| $f(1,l_1)$ | $v_1$ |
| $f(1,l_2)$ | $v_2$ |
| $f(2,l_1)$ | $v_2$ |

# Multi-Map

- Knowing $f$, reduce to one map
- Emulate as table $T_{\mathsf{MM}}$
- Express $\mathsf{MM}[l]$ as union

  $\mathsf{MM}[l_1] \rightarrow \{\mathsf{M}[f(1, l_1)]\} \cup \{\mathsf{M}[f(2, l_1)]\}$

  $\mathsf{MM}[l_2] \rightarrow \{\mathsf{M}[f(1, l_2)]\}$

- Is iteration a viable idea?
- In relational algebra?

- Idea 3c: One map



- $T_{\mathsf{MM}}$ for $\mathsf{M} : f \rightarrow V$

| label | value |
|-------|-------|
| $f(1, l_1)$ | $v_1$ |
| $f(1, l_2)$ | $v_2$ |
| $f(2, l_1)$ | $v_2$ |

# Multi-Map

## Multi-Map as Recursive Map

A multi-map $MM : L \to 2^V$ is a map $M : f(\mathbb{Z}^+, L) \to V$ where $f$ is a one-to-one function. The operation $MM[l] = \text{rec}(1, l)$ where

$$\text{rec}(i, l) = \begin{cases} \{M[f(i, l)]\} \cup \text{rec}(i + 1, l) & \text{if} \quad \exists f(i, l) \in \text{Dom}(M) \\ \emptyset & \text{else} \end{cases}$$

# Multi-Map

## Multi-Map as Recursive Map

A multi-map $\text{MM} : L \to 2^V$ is a map $\text{M} : f(\mathbb{Z}^+, L) \to V$ where $f$ is a one-to-one function. The operation $\text{MM}[l] = \text{rec}(1, l)$ where

$$\text{rec}(i, l) = \begin{cases} \{\text{M}[f(i, l)]\} \cup \text{rec}(i + 1, l) & \text{if} \quad \exists f(i, l) \in \text{Dom}(\text{M}) \\ \emptyset & \text{else} \end{cases}$$

## Negative Result

Relational algebra cannot express the multi-map as the recursive map.

# Relational Algebra with Transitive Closure

- Transitive closure is needed
- Undecidablility impacts query optimization (Later).

## Syntax of RCTE

Transitive closure can be expressed throuh the recursive common table
expression (RCTE) as in SQL-99 [3]

```
With Recursive  View  As
     BaseSubquery
   Union All
     RecursiveSubquery
```

# Relational Algebra with Transitive Closure

- But the semantics of RCTE has only ad-hoc definition in the literature. Here we provide an algebraic one.

### Semantics of RCTE

The semantics of the RCTE defines that the recursive view is equivalent to

$$\text{View} = \text{rec}_\text{rcte}(1)$$

where the recursive function is defined as

$$\text{rec}_\text{rcte}(i) = \begin{cases} \Delta\text{View}_i \cup \text{rec}_\text{rcte}(i+1) & \text{if} \quad \Delta\text{View}_i \neq \emptyset \\ \emptyset & \text{else} \end{cases}$$

with the view increment

$$\Delta\text{View}_i = \begin{cases} \text{BaseSubquery} & \text{if} \quad i = 1 \\ \text{RecursiveSubquery} \mid \Delta\text{View}_{i-1} & \text{else} \end{cases}$$

# Relational Algebra with Transitive Closure

- The RCTE semantics can express the recursive map. Why?
- Simply compare

$$\mathsf{rec}(i, l) = \begin{cases} \{\mathsf{M}[f(i,l)]\} \cup \mathsf{rec}(i+1, l) & \text{if} \quad \exists f(i,l) \in \mathrm{Dom}(\mathsf{M}) \\ \emptyset & \text{else} \end{cases}$$

and

$$\mathsf{rec}_{\mathsf{rcte}}(i) = \begin{cases} \Delta\mathsf{View}_i \cup \mathsf{rec}_{\mathsf{rcte}}(i+1) & \text{if} \quad \Delta\mathsf{View}_i \neq \emptyset \\ \emptyset & \text{else} \end{cases}$$

# Multi-Map Emulation

$T_{\mathsf{MM}}$ for rec. map

| label | value |
|-------|-------|
| $f(1, l_1)$ | $v_1$ |
| $f(1, l_2)$ | $v_2$ |
| $f(2, l_1)$ | $v_2$ |

# Multi-Map Emulation

- Operation MM[$l$] is emulated by $\mathbf{RCTE}(l)$ as

  With Recursive  View  As
  $$T_{\mathsf{MM}} \bowtie_{\mathtt{label}=f(i,l)} \{\langle i : 1 \rangle\}$$
  Union All
  $$T_{\mathsf{MM}} \bowtie_{\mathtt{label}=f(i,l)} \pi_{i+1 \to i}\mathsf{View}$$
  $\pi_{\mathtt{value}}\mathsf{View}$

$T_{\mathsf{MM}}$ for rec. map

| label | value |
|-------|-------|
| $f(1,l_1)$ | $v_1$ |
| $f(1,l_2)$ | $v_2$ |
| $f(2,l_1)$ | $v_2$ |

# Multi-Map Emulation

$T_{MM}$ for rec. map

| label | value |
|-------|-------|
| $f(1, l_1)$ | $v_1$ |
| $f(1, l_2)$ | $v_2$ |
| $f(2, l_1)$ | $v_2$ |

- Operation $MM[l]$ is emulated by $\mathbf{RCTE}(l)$ as

  ```
  With Recursive  View  As
  ```
  $$T_{MM} \bowtie_{\mathtt{label}=f(i,l)} \{\langle i : 1 \rangle\}$$
  ```
      Union All
  ```
  $$T_{MM} \bowtie_{\mathtt{label}=f(i,l)} \pi_{i+1 \to i}\mathsf{View}$$
  $$\pi_{\mathtt{value}}\mathsf{View}$$

- $\mathbf{RCTE}[l_1]$:

$$\{i = 1, f(1, l_1), v_1\} \cup \{i = 2, f(2, l_1), v_2\} \cup \{i = 3, f(3, l_1), \mathsf{nil}\} \nearrow \emptyset$$

- Projects as $\begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$

# Construction

- Multi-Map



- Recursive map

## Encrypted Multi-Map

- Correct:
  $\Pr(\mathsf{MM}[l] \equiv \mathsf{EMM}[l]) = 1 - \mathsf{negl}(\lambda)$

- Multi-Map



- Recursive map

# Encrypted Multi-Map

- Correct:
  $\Pr(\mathsf{MM}[l] \equiv \mathsf{EMM}[l]) = 1 - \mathsf{negl}(\lambda)$
- Leakage: # edges, queries, results

- Multi-Map



- Recursive map

# Encrypted Multi-Map

- Correct:
  $\Pr(\mathsf{MM}[l] \equiv \mathsf{EMM}[l]) = 1 - \mathsf{negl}(\lambda)$
- Leakage: # edges, queries, results
- Structure: recursive map homomorphism
  - perspective on correctness and security
  - for emulation

- Multi-Map



- Recursive map

# Encrypted Multi-Map

## $\Pi_{\mathrm{bas}}$ Encrypted Multi-Map

The scheme $\Pi_{\mathrm{bas}}$ [8] can be seen as a transformation of a recursive map representation of MM

- Client: $l' = \mathcal{F}(k, l)$, $\mathsf{trpd}_j(l') = \mathcal{F}(k, l'||j)^1$ , $\mathcal{E}(\cdot, \cdot)$
- Server: $\mathcal{F}_s(\mathsf{trpd}_1, i)$, $\mathcal{D}(\mathsf{trpd}_2, \cdot)^2$

where

- $\mathcal{F}, \mathcal{F}_s$ are PRFs, and $\mathcal{E}, \mathcal{D}$ are part of RCPA-secure symmetric key encryption.

---

[1] Output splitting is used in [8] instead.
[2] Response revealing.

# Encrypted Multi-Map

Figure: Multi-Map



Figure: Recursive Map



Figure: Encrypted Multi-Map

$$l_1' = \mathcal{F}(k, l_1) \dashrightarrow \mathcal{F}_s(\mathsf{trpd}_1(l_1'), 1) \rightarrow \mathcal{E}(\mathsf{trpd}_2(l_1'), v_1)$$

$$l_2' = \mathcal{F}(k, l_2) \dashrightarrow \mathcal{F}_s(\mathsf{trpd}_1(l_2'), 1) \rightarrow \mathcal{E}(\mathsf{trpd}_2(l_2'), v_2)$$

$$l_1' = \mathcal{F}(k, l_1) \dashrightarrow \mathcal{F}_s(\mathsf{trpd}_1(l_1'), 2) \rightarrow \mathcal{E}(\mathsf{trpd}_2(l_1'), v_2)$$

# Encrypted Multi-Map

## $\Pi_{bas}$ Security and Correctness

Provided by Cash *et al.* [7].

## $\Pi_{bas}$ Structure

$\Pi_{bas}$ preserves the *equivalent* recursive map for the multi-map.

- Homomorphism: by proving comp. ind. 1-1 on the label transformation. What about value transformation?
- Elementary equivalence: by Ehrenfreucht-Fraïssé Game [1, 2], with extension to pseudorandom objects.

# Emulation of Encrypted Multi-Map

| label | value |
|---|---|
| $\mathcal{F}_s(\mathsf{trpd}_1(l_1'), 1)$ | $\mathcal{E}(\mathsf{trpd}_2(l_1'), v_1)$ |
| $\mathcal{F}_s(\mathsf{trpd}_1(l_2'), 1)$ | $\mathcal{E}(\mathsf{trpd}_2(l_2'), v_2)$ |
| $\mathcal{F}_s(\mathsf{trpd}_1(l_1'), 2)$ | $\mathcal{E}(\mathsf{trpd}_2(l_1'), v_2)$ |

Table: $T_{\mathsf{EMM}}$ for $\Pi_{\mathrm{bas}}$. $l' = \mathcal{F}(k,l), \mathsf{trpd}_j(l') = \mathcal{F}(k, l'||j)$

- Operation EMM[$l$] is emulated by $\mathbf{RCTE}(l)$ as

$$\text{With Recursive } \text{View } \text{As}$$
$$T_{\mathsf{EMM}} \bowtie_{\texttt{label}=\mathcal{F}_s(\mathsf{trpd}_1(l'))} \{\langle i : 1 \rangle\}$$
$$\text{Union All}$$
$$T_{\mathsf{EMM}} \bowtie_{\texttt{label}=\mathcal{F}_s(\mathsf{trpd}_2(l'))} \pi_{i+1 \to i}\text{View}$$
$$\pi_{\texttt{value}}\text{View}$$

# Construction

# Independence

Table: Customer$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $c_1$ | Alice$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_2$ | Bob$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_3$ | Alice$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |

Table: Supplier$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $s_1$ | Li-Ning$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |
| $s_2$ | Nike$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |

# Independence

Table: Customer$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $c_1$ | Alice$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_2$ | Bob$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_3$ | Alice$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |

Table: Supplier$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $s_1$ | Li-Ning$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |
| $s_2$ | Nike$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |

Table: EMM$_\sigma$

| label | value |
|-------|-------|
| $\mathcal{F}(\mathsf{trpd}_1(\text{C.Name=Alice}'), 1)$ | $\mathcal{E}((\mathsf{trpd}_2(\text{C.Name=Alice}')), c_1)$ |
| $f(\text{C.Name=Bob}, 1)$ | $e(\text{C.Name=Bob}, c_2)$ |
| $f(\text{C.Name=Alice}, 2)$ | $e(\text{C.Name=AliceS}, c_3)$ |

Figure: MM$_\sigma$

# Independence

### Table: Customer$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $c_1$ | Alice$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_2$ | Bob$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_3$ | Alice$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |

### Table: Supplier$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $s_1$ | Li-Ning$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |
| $s_2$ | Nike$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |

### Table: EMM$_\sigma$

| label | value |
|-------|-------|
| $\mathcal{F}(\text{trpd}_1(\text{C.Name=Alice}'), 1)$ | $\mathcal{E}((\text{trpd}_2(\text{C.Name=Alice}')), c_1)$ |
| $f(\text{C.Name=Bob}, 1)$ | $e(\text{C.Name=Bob}, c_2)$ |
| $f(\text{C.Name=Alice}, 2)$ | $e(\text{C.Name=AliceS}, c_3)$ |

### Figure: MM$_\sigma$



- $\pi_{\text{Country}}\sigma_{\text{Name=Alice}}\text{Customer} \equiv$
  $\pi_{\text{Country}^{\mathcal{F}}} \underbrace{\text{Customer}^{\mathcal{F}} \ltimes_{\text{rid=value}} \mathbf{RCTE}_\sigma(\text{C.Name=Alice})}_{\sigma^{\mathcal{E}}_{\text{C.Name=Alice}}\text{Customer}^{\mathcal{F}}}$

# Independence

Table: Customer$^{\mathcal{F}}$

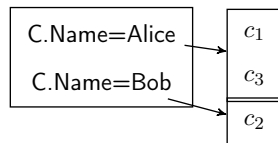| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $c_1$ | Alice$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_2$ | Bob$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_3$ | Alice$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |

Table: Supplier$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $s_1$ | Li-Ning$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |
| $s_2$ | Nike$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |

Figure: MM$_{\bowtie}$

C.Country = S.Country

| $(c_1, s_2)$ |
|---|
| $(c_2, s_2)$ |
| $(c_3, s_1)$ |

# Independence

Table: Customer$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $c_1$ | Alice$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_2$ | Bob$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_3$ | Alice$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |

Table: Supplier$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $s_1$ | Li-Ning$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |
| $s_2$ | Nike$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |

Figure: MM$_{\bowtie}$

C.Country = S.Country

| |
|---|
| $(c_1, s_2)$ |
| $(c_2, s_2)$ |
| $(c_3, s_1)$ |

Table: EMM$_{\bowtie}$

| label | value$_1$ | value$_2$ |
|-------|-----------|-----------|
| $f(\text{C.Country} = \text{S.Country}, 1)$ | $e(\text{C.Country} = \text{S.Country}, c_1)$ | $e(\text{C.Country} = \text{S.Country}, s_2)$ |
| $f(\text{C.Country} = \text{S.Country}, 2)$ | $e(\text{C.Country} = \text{S.Country}, c_2)$ | $e(\text{C.Country} = \text{S.Country}, s_2)$ |
| $f(\text{C.Country} = \text{S.Country}, 3)$ | $e(\text{C.Country} = \text{S.Country}, c_3)$ | $e(\text{C.Country} = \text{S.Country}, s_1)$ |

# Independence

Table: Customer$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $c_1$ | Alice$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_2$ | Bob$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_3$ | Alice$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |

Table: Supplier$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $s_1$ | Li-Ning$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |
| $s_2$ | Nike$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |

Figure: $\text{MM}_{\bowtie}$



C.Country = S.Country $\longrightarrow$

$(c_1, s_2)$

$(c_2, s_2)$

$(c_3, s_1)$

Table: $\text{EMM}_{\bowtie}$

| label | value$_1$ | value$_2$ |
|-------|-----------|-----------|
| $f(\text{C.Country} = \text{S.Country}, 1)$ | $e(\text{C.Country} = \text{S.Country}, c_1)$ | $e(\text{C.Country} = \text{S.Country}, s_2)$ |
| $f(\text{C.Country} = \text{S.Country}, 2)$ | $e(\text{C.Country} = \text{S.Country}, c_2)$ | $e(\text{C.Country} = \text{S.Country}, s_2)$ |
| $f(\text{C.Country} = \text{S.Country}, 3)$ | $e(\text{C.Country} = \text{S.Country}, c_3)$ | $e(\text{C.Country} = \text{S.Country}, s_1)$ |

- Customer $\bowtie_{\text{C.Country}=\text{S.Country}}$ Supplier $\equiv$
  $\underbrace{\mathbf{RCTE}_{\bowtie}(\text{C.Country}=\text{S.Country}) \bowtie_{\texttt{value}_1=\texttt{rid}} \text{C} \bowtie_{\texttt{value}_2=\texttt{rid}} \text{S}}_{\text{C}^{\mathcal{E}} \bowtie_{\text{C.Country}=\text{S.Country}}^{\mathcal{E}} \text{S}^{\mathcal{E}}}$

Table: Customer$^\mathcal{F}$

| rid | Name$^\mathcal{F}$ | Country$^\mathcal{F}$ |
|-----|--------------------|-----------------------|
| $c_1$ | Alice$^\mathcal{E}$ | US$^\mathcal{E}$ |
| $c_2$ | Bob$^\mathcal{E}$ | US$^\mathcal{E}$ |
| $c_3$ | Alice$^\mathcal{E}$ | China$^\mathcal{E}$ |

Table: Supplier$^\mathcal{F}$

| rid | Name$^\mathcal{F}$ | Country$^\mathcal{F}$ |
|-----|--------------------|-----------------------|
| $s_1$ | Li-Ning$^\mathcal{E}$ | China$^\mathcal{E}$ |
| $s_2$ | Nike$^\mathcal{E}$ | US$^\mathcal{E}$ |



$\bowtie$C.Country=S.Country

$O(\sigma T)$  $O(T)$

$\sigma_{\text{Name=Bob}}$    S

C

Hints from thought experiment:

- On efficiency: $O(\sigma T^2)$ where $\sigma$ is the selectivity
- On optimal leakage: only the encrypted result set

# Independence



Figure: $MM_\sigma$

| C.Name=Alice | $c_1$ |
| C.Name=Bob | $c_3$ |
| | $c_2$ |

Figure: $MM_{\bowtie}$

| C.Country = S.Country | $(c_1, s_2)$ |
| | $(c_2, s_2)$ |
| | $(c_3, s_1)$ |

$\bowtie_{\text{C.Country=S.Country}}$
$O(\sigma T)$ $O(T)$
$\sigma_{\text{Name=Bob}}$ $S$
$C$

$\bowtie_{\text{value}_2=\text{rid}}$
$\bowtie_{\text{rid=value}_1}$ $S^{\mathcal{F}}$
$O(\sigma T)$ $O(T^2)$
$\ltimes_{\text{rid=value}}$ $\mathbf{RCTE}_{\bowtie}(\text{C.Country=S.Country})$
$\mathbf{RCTE}_\sigma(\text{C.Name=Bob})$

# Independence

Table: Customer$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $c_1$ | Alice$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_2$ | Bob$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |
| $c_3$ | Alice$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |

Table: Supplier$^{\mathcal{F}}$

| rid | Name$^{\mathcal{F}}$ | Country$^{\mathcal{F}}$ |
|-----|------|---------|
| $s_1$ | Li-Ning$^{\mathcal{E}}$ | China$^{\mathcal{E}}$ |
| $s_2$ | Nike$^{\mathcal{E}}$ | US$^{\mathcal{E}}$ |

- Also similar issue with compound formula, such as
  $C.Name = Alice \wedge C.Country = US$ for filters, or
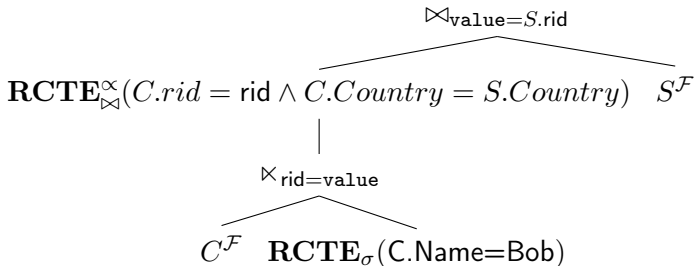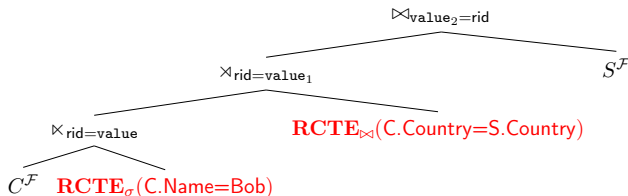  $C.Country = S.Country \wedge C.Name = S.Name$.

## Suboptimality

Independent operators are suboptimal in both efficiency and leakage for correlated filter and join predicates.

# Construction

# Dependence

- Introduce dependence between encrypted operators



$\bowtie_{\text{C.Country=S.Country}}$

$\sigma_{\text{Name=Bob}}$    S;

C

$\bowtie_{\texttt{value}_2=\texttt{rid}}$

$\ltimes_{\texttt{rid}=\texttt{value}_1}$     $S^{\mathcal{F}}$

$\ltimes_{\texttt{rid}=\texttt{value}}$     $\mathbf{RCTE}_{\bowtie}(\text{C.Country=S.Country})$

$C^{\mathcal{F}}$   $\mathbf{RCTE}_{\sigma}(\text{C.Name=Bob})$

$\bowtie_{\texttt{value}=S.\text{rid}}$

$\mathbf{RCTE}_{\bowtie}^{\propto}(C.rid = \text{rid} \wedge C.Country = S.Country)$    $S^{\mathcal{F}}$

$\ltimes_{\text{rid}=\texttt{value}}$

$C^{\mathcal{F}}$   $\mathbf{RCTE}_{\sigma}(\text{C.Name=Bob})$

# Dependence

Figure: $\mathsf{MM}_\sigma$

Figure: $\mathsf{MM}_{\bowtie}^\propto$

# Dependence

Figure: $\mathsf{MM}_\sigma$

| C.Name=Alice | | $c_1$ |
|---|---|---|
| | | $c_3$ |
| C.Name=Bob | | |
| | | $c_2$ |

Figure: $\mathsf{MM}_\bowtie^\propto$

| C.Country=S.Country $\wedge$ $C.\mathsf{rid} = c_1$ | $s_1$ |
|---|---|
| C.Country=S.Country $\wedge$ $C.\mathsf{rid} = c_2$ | $s_2$ |
| C.Country=S.Country $\wedge$ $C.\mathsf{rid} = c_3$ | |

- Problem: Client has to keep $O(T)$ trapdoors for $\mathsf{MM}_\bowtie^\propto$, and interaction
- Solution:
  - Client: master trapdoor $\mathsf{trpd}_\bowtie = \mathcal{F}(k, \mathsf{C.Country=S.Country})$
  - Server: derive $O(T)$ trapdoors for each $c_i$ as $\mathsf{trpd}_j = \mathcal{F}(\mathsf{trpd}_\bowtie, c_i || j)$ for $j = 1, 2$

# Construction

# Normal Form

Figure: General Multi-Map



Figure: $MM_\sigma$

# Normal Form

Figure: General Multi-Map



Figure: $MM_\sigma$



## Multi-Map Range Partition under 1NF

The range of $MM_\sigma$ *always* forms a partition in the space of $C^{\mathcal{F}}.\mathsf{rid} = \{c_i\}_{i=1}^{T}$, because of 1st normal form

- $C.Name$ is an elemantary set
- $C^{\mathcal{F}}.\mathsf{rid} = \{c_i\}_{i=1}^{T}$ is a candidate key $\rightarrow$ uniquely identifies a row

By contrast

- The general multi-map may have overlapping value sets in its range
- EMM suitable for document keyword model, but overkill relational model?

# Normal Form

Figure: Many-to-Many Join



$c_1$ → $s_1$
$c_2$ → $s_2$
$c_3$ → $s_3$

Figure: 1-to-Many Join



$c_1$ → $s_1$
$c_2$ → $s_2$
→ $s_3$

Figure: Many-to-1 Join



$c_1$ → $s_1$
$c_2$ → $s_2$
$c_3$ →

## Multi-Map Range Partition for Joins

The range partition generalize to $MM_{\bowtie}^{\propto}$ too.

# Normal Form

Figure: Many-to-Many Join



Figure: 1-to-Many Join



Figure: Many-to-1 Join



## Multi-Map Range Partition for Joins

The range partition generalize to $MM_\bowtie^\propto$ too.

3NF has only three key joins:

- Foreign-to-foreign key join $\in$ many-to-many join
- Primary-to-foreign key join $\in$ 1-to-many join
- Foreign-to-primary key join $\in$ many-to-one join

## Worst-case Optimal Space for Joins

The worst-case optimal space for joins in a 3NF data model is $O(T)$.

Adapt EMM for 1/3NF for efficiency/security?

# Normal Form

Adapt EMM for 1/3NF for efficiency/security?

## Semi-Encrypted Multi-Map

The semi-encrypted multi-map is $\Pi_{\mathrm{bas}}$ with its values in clear. Formally it is a transformation of MM

- Client: $\mathcal{F}(k, \cdot)$, $\mathsf{trpd}_j(l') = \mathcal{F}(k, l'||j)$
- Server: $\mathcal{F}_s(\cdot, i)$

# Normal Form



Figure: Encrypted Multi-Map

$$l'_1 = \mathcal{F}(k, l_1) \rightarrow \mathcal{F}_s(\mathsf{trpd}_1(l'_1), 1) \rightarrow \mathcal{E}(\mathsf{trpd}_2(l'_1), v_1)$$

$$l'_2 = \mathcal{F}(k, l_2) \rightarrow \mathcal{F}_s(\mathsf{trpd}_1(l'_2), 1) \rightarrow \mathcal{E}(\mathsf{trpd}_2(l'_2), v_2)$$

$$l'_1 = \mathcal{F}(k, l_1) \rightarrow \mathcal{F}_s(\mathsf{trpd}_1(l'_1), 2) \rightarrow \mathcal{E}(\mathsf{trpd}_2(l'_1), v_2)$$

Figure: Semi-Encrypted Multi-Map

$$l'_1 = \mathcal{F}(k, l_1) \rightarrow \mathcal{F}_s(\mathsf{trpd}_1(l'_1), 1) \rightarrow v_1$$

$$l'_2 = \mathcal{F}(k, l_2) \rightarrow \mathcal{F}_s(\mathsf{trpd}_1(l'_2), 1) \qquad v_2$$

$$l'_1 = \mathcal{F}(k, l_1) \rightarrow \mathcal{F}_s(\mathsf{trpd}_1(l'_1), 2)$$

Leaks "co-occurence" pattern. Insecure for document keyword model.

# Normal Form

## Security of SEMM under 1NF

The SEMM under 1NF leaks only dimensions of table (trivial co-ocurrence).
SEMM achieves the same security as $\Pi_{\mathrm{bas}}$ for range as row ids.

Proof sketch:

- Co-occurrence is the same for every row. Each $v_i$ uniquely identifies a row. So the range size is # rows. Number of in-edges for each row id $v_i$ is always equal to # attributes. All rows have the same # attributes.

Figure: SEMM is Secure under 1NF



$$l'_1 = \mathcal{F}(k, l_1) \longrightarrow \mathcal{F}_s(\mathsf{trpd}_1(l'_1), 1) \longrightarrow v_1$$
$$l'_2 = \mathcal{F}(k, l_2) \longrightarrow \mathcal{F}_s(\mathsf{trpd}_1(l'_2), 1)$$
$$l'_3 = \mathcal{F}(k, l_3) \longrightarrow \mathcal{F}_s(\mathsf{trpd}_1(l'_3), 1)$$
$$l'_1 = \mathcal{F}(k, l_1) \longrightarrow \mathcal{F}_s(\mathsf{trpd}_1(l'_1), 2)$$

$v_2$

# Normal Form

Implications

1. Can divide the SEMM and collocate with tables such that $SEMM.\texttt{value} = \texttt{rid}$.
2. Share the cleartext $SEMM.\texttt{value}$ for all SEMMs for the same table to reduce SEMM size.
3. Pre-indexing on $SEMM.\texttt{value}$ to reduce computation time.

# Normal Form

Implications

1. Can divide the SEMM and collocate with tables such that $SEMM.\texttt{value} = $ rid.
2. Share the cleartext $SEMM.\texttt{value}$ for all SEMMs for the same table to reduce SEMM size.
3. Pre-indexing on $SEMM.\texttt{value}$ to reduce computation time.

Achieve worst-case optimal space for joins

1. SEMM only index many-to-1 or 1-to-many joins.
2. Worst-case optimal space for joins: avoid storing many-to-many joins in SEMM, but factor them into two many-to-1 or 1-to-many joins.

# Construction

1. Worst-case optimal joins (efficiency). Idea from [7]
2. Range queries [10]
3. Updates with minimum interactions [12]
4. Consolidate the framework/proofs in this thesis
5. Query optimization
6. Security against malicious server

# Construction

- Prototype system is open-source [15]
- Based on the algebraic core of Apache Spark SQL [9], interface with any database endpoints e.g. PostgreSQL [13]
- Parallel database encryption
- Fusion of plaintext and encrypted operators.

Figure: The **dex** encrypted relational database.

# System & Evaluation

TPC-H Benchmark [5] on the dependence scheme (without normal form).

| dex-cor | mean(ms)(m44xlarge) | rel.err. | slowdown vs. Postgres | slowdown vs. Postgres(t22xlarge) |
|---------|---------------------|----------|-----------------------|----------------------------------|
| q1 | 2149.5 | 1.38% | 1.4 | 1.0 |
| q10 | 217837.6 | 0.16% | 115.0 | 32.7 |
| q11 | 993.3 | 2.15% | 13.8 | 5.4 |
| q12 | 38976.8 | 0.25% | 32.7 | 26.9 |
| q13 | 61460.2 | 0.27% | 84.0 | 64.9 |
| q14 | 110169.4 | 0.15% | 40.5 | 29.7 |
| q15 | 100299.6 | 0.29% | 41.0 | 33.3 |
| q16 | 478.2 | 1.99% | 3.0 | 2.6 |
| q17 | 437.2 | 2.29% | 6.0 | 4.0 |
| q18 | 280982.8 | 0.28% | 66.6 | 50.0 |
| q19 | 31976.2 | 0.36% | 373.1 | 324.3 |
| q2 | 2445.4 | 0.52% | 15.8 | 12.3 |
| q20 | 27546.1 | 0.51% | 284.6 | 35.7 |
| q21 | 447845.6 | 0.22% | 441.9 | 354.6 |
| q22 | 55829.3 | 0.36% | 134.1 | 108.7 |
| q3 | 40152.1 | 0.38% | 22.8 | 17.2 |
| q4 | 119831.3 | 0.20% | 30.6 | 23.1 |
| q5 | 4645337.0 | 0.62% | 2231.8 | 1875.5 |
| q7 | 47910.1 | 0.30% | 131.8 | 81.6 |
| q8 | 117817.6 | 0.44% | 57.4 | 38.8 |
| q9a | 464302.6 | 0.20% | 10.8 | 12.3 |
| q9b | 356243.1 | 0.32% | 8.3 | 9.4 |

| table name | row est. | attrs | page est. | total(bytes) | index(bytes) | table(bytes) |
|---|---|---|---|---|---|---|
| region | 5 | 4(3) | 1 | 32 k (64 k) | 16 k(48 k) | 8192 |
| orders | $1.5e+06$ | 10(9) | 71484(26405) | 604 m(603 m) | 45 m(396 m) | 559 m(206 m) |
| supplier | 10000 | 8(7) | 447(226) | 3936 k(4680 k) | 328 k(2840 k) | 3600 k(1832 k) |
| customer | 150000 | 9(8) | 7620(3758) | 64 m(75 m) | 4640 k(46 m) | 60 m (29m) |
| partsupp | 800000 | 7(5) | 35580(18242) | 302 m(337 m) | 24 m(194 m) | 278 m (143 m) |
| nation | 25 | 5(4) | 1 | 32 k(80 k) | 16 k(64 k) | 8192 |
| lineitem | $6.00139e+06$ | 19(16) | 500251(117594) | 4090 m(3338 m) | 181 m(2419 m) | 3909 m(919 m) |
| part | 200000 | 10(9) | 9656(3832) | 82 m(85 m) | 6184 k(55 m) | 75 m(30 m) |
| $T_\sigma^\perp$ | $8.74646e+07$ | 2 | 1079906 | 13 G | 4926 m | 8439 m |
| $T_\bowtie^\infty$ | $5.45298e+07$ | 2 | 673291 | 8332 m | 3071 m | 5261 m |

Table: TPCH benchmark **dex-cor** versus plaintext Postgres storage size. The plaintext Postgres storage size is shown in parenthesis. TPC-H scale factor is 1.

# References I

1.  Fraïssé, R. *Sur quelques classifications des systemes de relations.* PhD thesis (1955).

2.  Ehrenfeucht, A. An application of games to the completeness problem for formalized theories. *Fund. Math* **49,** 13 (1961).

3.  Eisenberg, A. & Melton, J. SQL: 1999, formerly known as SQL3. *ACM Sigmod record* **28,** 131–138 (1999).

4.  Bellare, M. & Rogaway, P. *The security of triple encryption and a framework for code-based game-playing proofs.* in *Annual International Conference on the Theory and Applications of Cryptographic Techniques* (2006), 409–426.

5.  Council, T. P. P. TPC-H benchmark specification. *Published at http://www. tcp. org/hspec. html* **21,** 592–603 (2008).

# References II

6. Popa, R. A., Redfield, C., Zeldovich, N. & Balakrishnan, H. *CryptDB: protecting confidentiality with encrypted query processing.* in *Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles* (2011), 85–100.

7. Cash, D. *et al. Highly-scalable searchable symmetric encryption with support for boolean queries.* in *Annual Cryptology Conference* (2013), 353–373.

8. Cash, D. *et al. Dynamic searchable encryption in very-large databases: data structures and implementation..* in *NDSS* **14** (2014), 23–26.

9. Armbrust, M. *et al. Spark sql: Relational data processing in spark.* in *Proceedings of the 2015 ACM SIGMOD international conference on management of data* (2015), 1383–1394.

10. Faber, S. *et al.* *Rich queries on encrypted data: Beyond exact matches.* in *European Symposium on Research in Computer Security* (2015), 123–145.

11. Degitz, A., Köhler, J. & Hartenstein, H. *Access Pattern Confidentiality-Preserving Relational Databases: Deployment Concept and Efficiency Evaluation..* in *EDBT/ICDT Workshops* (2016).

12. Kamara, S. & Moataz, T. *SQL on structurally-encrypted databases.* in *International Conference on the Theory and Application of Cryptology and Information Security* (2018), 149–180.

13. Group, T. P. G. D. *PostgreSQL.* https://www.postgresql.org/. 2019.

14. Tan, B. H. M., Lee, H. T., Wang, H., Ren, S. Q. & Aung, K. M. M. Efficient Private Comparison Queries over Encrypted Databases using Fully Homomorphic Encryption with Finite Fields. *IACR Cryptology ePrint Archive* **2019,** 332 (2019).

15. Zhao, Z. *encrypted-spark*. https://github.com/zheguang/encrypted-spark/tree/dex. 2019.