# Security in Outsourced Databases
# (Query Answer Assurance)
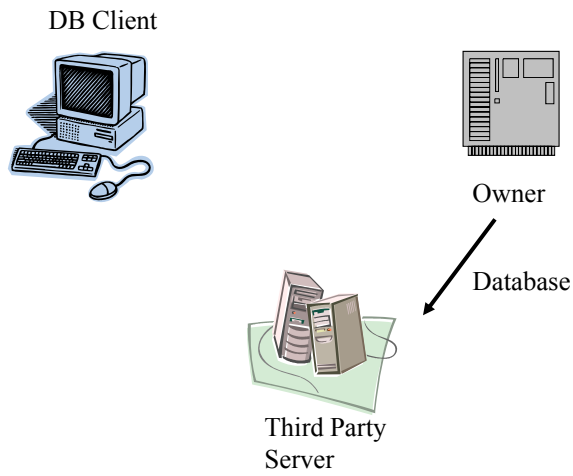
1

# Traditional Client-Server Arch.

DB Client

Query

Results

Owner

- Client queries are satisfied by a trusted server
- Secure the server
- Secure the communication channel, e.g. use SSL

2

# Data Publishing
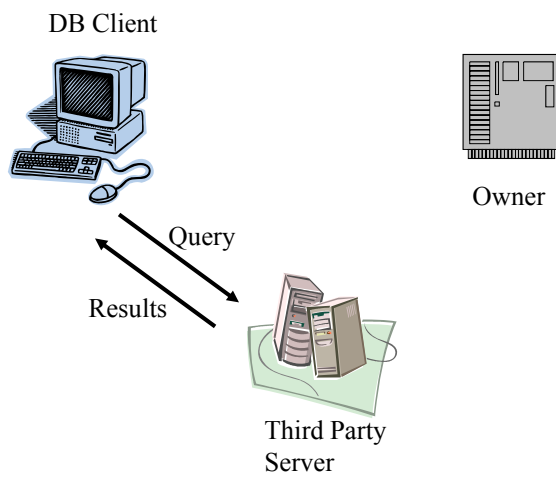## (Database-as-a-Service)

DB Client

Owner

Database

Third Party
Server

3

---

# Data Publishing

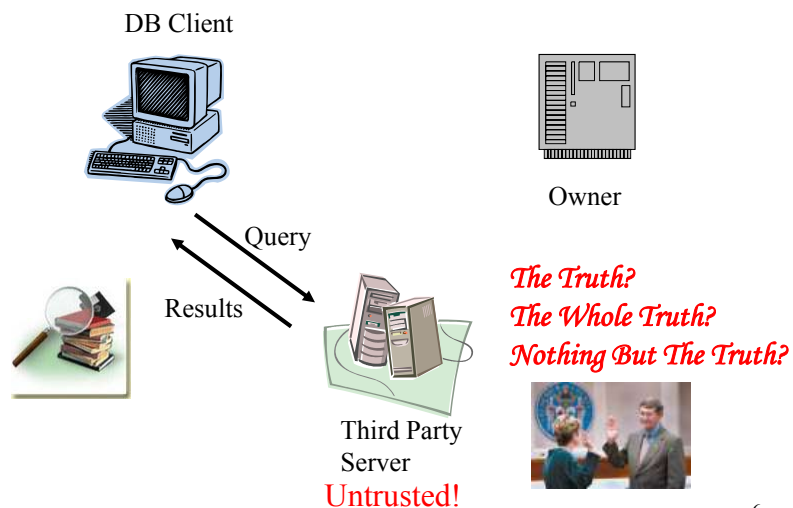DB Client

Owner

Query

Results

Third Party
Server

4

# Data Publishing

- Pushes business logic and data processing from corporate data centers to third party servers at the "edge" of the network
  - Distribution of (part of) the database to edge servers
  - Edge servers perform query processing
- **Why?**
  - Most organizations need DBMSs
  - DBMSs extremely complex to deploy, setup, maintain
  - Require skilled DBAs (at very high cost!)
- Advantages
  - Cuts down network latency and produces faster responses
  - Cheaper way to achieve scalability
  - Lowers dependency on corporate data center (removes single point of failure)
  - Reduced cost to client
    - Get what you need, pay for what you use and not for: hardware, software infrastructure or personnel to deploy, maintain, upgrade…
  - Reduced overall cost
    - cost amortization across users
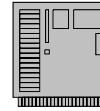  - Better service
    - leveraging experts

5

---

# The Challenge

DB Client



Owner

Query

Results

*The Truth?*
*The Whole Truth?*
*Nothing But The Truth?*

Third Party
Server
Untrusted!

6

# The Challenge

DB Client

Sel * FROM Emp
WHERE Sal < 5000

Owner

| ID | Name | Sal | Dept |
|----|------|------|------|
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 1 | D | 8010 | 1 |
| 4 | B | 2200 | 3 |
| 3 | E | 7000 | 2 |

Server

7

# The Challenge

DB Client

Sel * FROM Emp
WHERE Sal < 5000

Owner

Result =

| | | | |
|---|---|------|---|
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 4 | B | 2200 | 3 |

| ID | Name | Sal | Dept |
|----|------|------|------|
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 1 | D | 8010 | 1 |
| 4 | B | 2200 | 3 |
| 3 | E | 7000 | 2 |

Server

8

4

# Security Concerns

DB Client

$$\text{Result} = \begin{array}{cccc} 5 & A & 2000 & 1 \\ 2 & C & 3500 & 2 \\ 4 & B & 2200 & 3 \end{array}$$

Query

Result'

Server

| ID | Name | Sal | Dept |
|----|------|------|------|
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 1 | D | 8010 | 1 |
| 4 | B | 2200 | 3 |
| 3 | E | 7000 | 2 |

9

---

# Security Concerns

DB Client

$$\text{Result} = \begin{array}{cccc} 5 & A & 2000 & 1 \\ 2 & C & 3500 & 2 \\ 4 & B & 2200 & 3 \end{array}$$

Query

Result'

| 5 | A | 2000 | 1 |
|---|---|------|---|
| 2 | C | 3500 | 2 |
| 4 | B | 2200 | 3 |

Server is trustworthy!

Server

| ID | Name | Sal | Dept |
|----|------|------|------|
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 1 | D | 8010 | 1 |
| 4 | B | 2200 | 3 |
| 3 | E | 7000 | 2 |

10

# Security Concerns

Result =
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 4 | B | 2200 | 3 |

DB Client

Query

Result'
| 5 | A | **3500** | 1 |
| 2 | **D** | 3500 | 2 |
| 4 | B | 2200 | **1** |

Server is malicious!
Records are tampered

Server

| ID | Name | Sal | Dept |
|----|------|-----|------|
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 1 | D | 8010 | 1 |
| 4 | B | 2200 | 3 |
| 3 | E | 7000 | 2 |

11

---

# Security Concerns

Result =
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 4 | B | 2200 | 3 |

DB Client

Query

Result'
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |

Server is malicious!
Answers are dropped
(Incompleteness)

Server

| ID | Name | Sal | Dept |
|----|------|-----|------|
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 1 | D | 8010 | 1 |
| 4 | B | 2200 | 3 |
| 3 | E | 7000 | 2 |

12

# Security Concerns

DB Client

$$\text{Result} = \begin{array}{cccc} 5 & A & 2000 & 1 \\ 2 & C & 3500 & 2 \\ 4 & B & 2200 & 3 \end{array}$$

Query

Result'

| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 4 | B | 2200 | 3 |
| **1** | **D** | **1500** | **2** |
| **6** | **E** | **3400** | **1** |

Server

| ID | Name | Sal | Dept |
|----|------|------|------|
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 1 | D | 8010 | 1 |
| 4 | B | 2200 | 3 |
| 3 | E | 7000 | 2 |

Server is malicious!
Spurious answers are added

13

---

# Data Security Challenge:

Design Objectives:
- *Authenticity*: Every entry originated from the owner
- *Completeness*: No result entry is omitted from the answer
- *Precision*: Minimum information leakage
- *Security*: Computationally infeasible to cheat
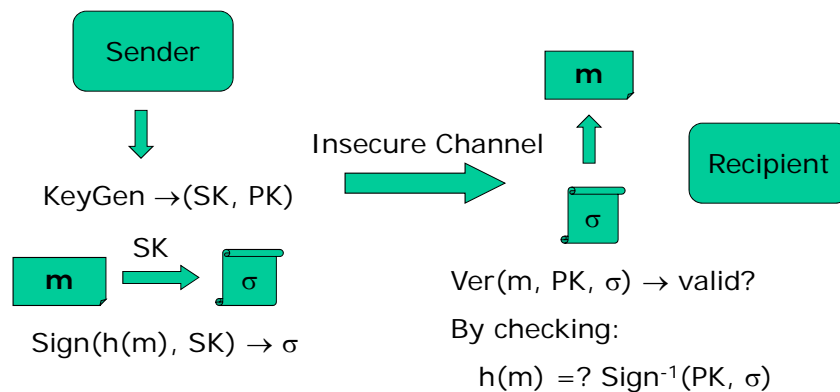- *Efficiency*: Polynomial proof

14

# Collision-resistant (one-way) hash functions

- Given x, easy to compute h(x); given h(x), difficult to determine x
- i.e., it is computationally hard to find $x_1$ and $x_2$ s.t. $h(x_1)=h(x_2)$
- Computational hard? Based on well established assumptions such as discrete logarithms
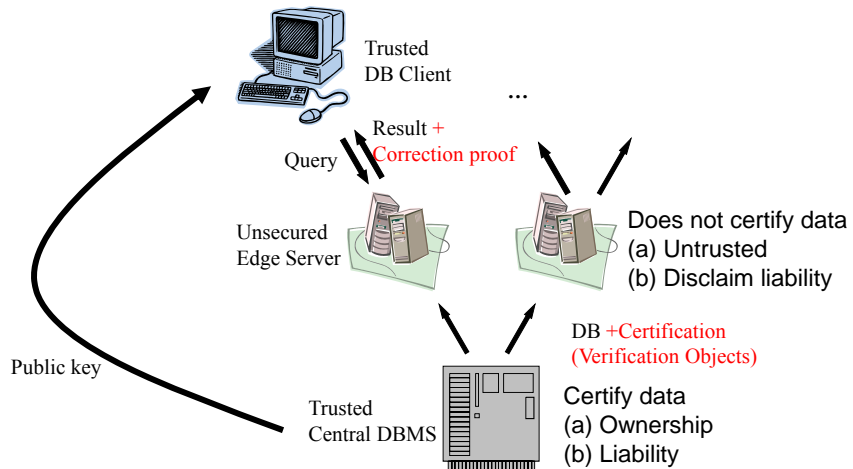- E.g., SHA, MD5

15

# Public key digital signature schemes

Cryptographic tool for authenticating the signed message as well as its origin, e.g., RSA, DSA

Sender

m

Insecure Channel

Recipient

KeyGen →(SK, PK)

σ

SK

m → σ

Ver(m, PK, σ) → valid?

Sign(h(m), SK) → σ

By checking:

$h(m) =? \text{Sign}^{-1}(PK, \sigma)$

16

# Authentic Publication Scheme

Trusted
DB Client

...

Result +
Correction proof

Query

Unsecured
Edge Server

Does not certify data
(a) Untrusted
(b) Disclaim liability

DB +Certification
(Verification Objects)

Public key

Trusted
Central DBMS

Certify data
(a) Ownership
(b) Liability

---

# Naïve Scheme

Each attribute has a signed digest
Each tuple has a signed digest

Relation R

| $D_T$ | $(A_1, D_1)$ | ... | $(A_i, D_i)$ | | ... |
|---|---|---|---|---|---|
| | | | | | |
| | | | | | |

$D_T$ – Signed tuple digest
$D_{Ai}$ – attribute digest

# Naïve Scheme

Query: SELECT $A_3$, $A_4$, … FROM R

Filtered attributes

Result tuples

| $D_T$ | $A_3$ | $A_4$ | **...** | $D_1$ | $D_2$ | $D_5$ | **...** |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| | | | | | | | |

$D_T$ – Signed tuple digest
$D_i$ – attribute digest of $A_i$

19

---

# Naïve Scheme (Example)

| A1 | B1 | C1 | a1 | b1 | c1 | T1 |
|---|---|---|---|---|---|---|
| A2 | B2 | C2 | a2 | b2 | c2 | T2 |
| A3 | B3 | C3 | a3 | b3 | c3 | T3 |

$T = sign(g(h(A)|h(B)|h(C))$
    g and h are collision-resistant hash functions
$ai = h(Ai)$

Retrieve whole of first tuple:
    Server returns A1, B1, C1, T1; Client can compute h(A1), h(B1) and
    h(C1), and verify T1 from A1, B1 and C1

Retrieve only attributes A1 and B1 of first tuple:
    Server returns A1, B1, c1 and T1; Client has no access to C1, so
    c1 has to be provided

20

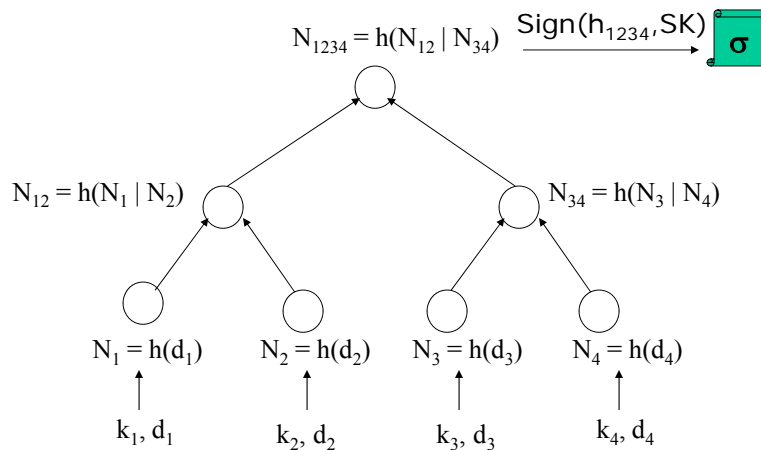Issues??

# Using Merke Hash Tree (MHT)

- For each tuple t, a tuple hash h(t) is computed
$$h(t) = h(h(t.A1) \mid h(t.A2) \mid \ldots \mid h(t.An))$$
- Assume a total order on attribute A of a relation R with $|R|$ tuples (e.g., based on the primary key)
  - MHT(R,A) is a binary tree with $|R|$ leaf nodes and hash values h(i) associated with node i
  - If i is a leaf node, then h(i) = h(ti), ti is the ith tuple in the order
  - If i is an internal node, then h(i) = h(h(l), h(r)) where l and r are the left and right children of node i.
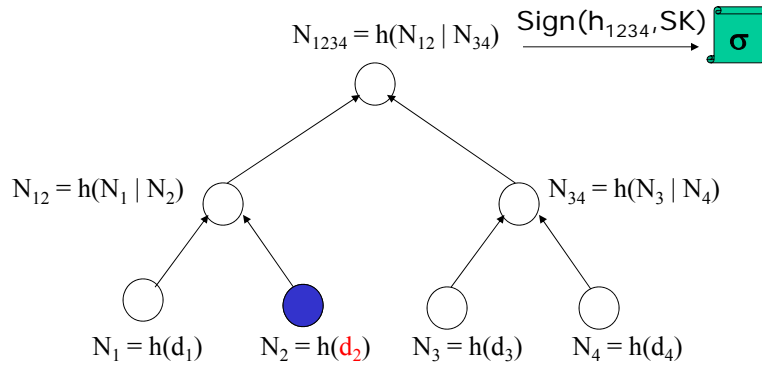  - The root hash is the digest of all values in the Merkle-hash tree MHT(R,A).

21

# Merkle Hash Tree



$N_{1234} = h(N_{12} \mid N_{34})$    $\text{Sign}(h_{1234}, \text{SK})$   $\sigma$

$N_{12} = h(N_1 \mid N_2)$       $N_{34} = h(N_3 \mid N_4)$

$N_1 = h(d_1)$   $N_2 = h(d_2)$   $N_3 = h(d_3)$   $N_4 = h(d_4)$

$k_1, d_1$    $k_2, d_2$    $k_3, d_3$    $k_4, d_4$

Ordering attribute: $k_1 < k_2 < k_3 < k_4$; $d_i$ are tuples
Owner needs to sign root node ($N_{1234}$)

22

# MHT: Point Search

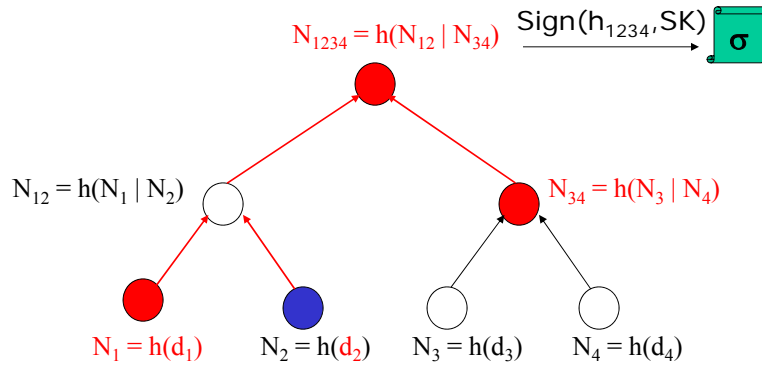$N_{1234} = h(N_{12} \mid N_{34})$ $\xrightarrow{\text{Sign}(h_{1234},\text{SK})}$ $\sigma$

$N_{12} = h(N_1 \mid N_2)$

$N_{34} = h(N_3 \mid N_4)$

$N_1 = h(d_1)$ $N_2 = h(d_2)$ $N_3 = h(d_3)$ $N_4 = h(d_4)$

Query: Retrieve tuple $d_2$

23

# MHT: Point Search

$N_{1234} = h(N_{12} \mid N_{34})$ $\xrightarrow{\text{Sign}(h_{1234},\text{SK})}$ $\sigma$

$N_{12} = h(N_1 \mid N_2)$

$N_{34} = h(N_3 \mid N_4)$

$N_1 = h(d_1)$ $N_2 = h(d_2)$ $N_3 = h(d_3)$ $N_4 = h(d_4)$

Edge server returns $d_2$, $N_1$, $N_{34}$ and signed $N_{1234}$
Client computes $N_{1234} = h(h(h(d_2)|N_1), N_{34})$ and verify
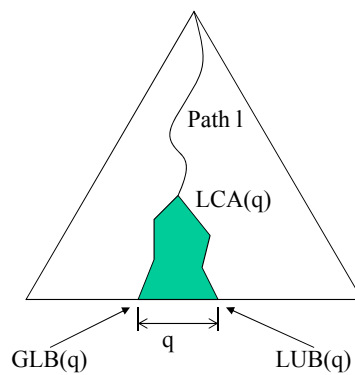that the signed value is correct

24

# MHT: Point Search

$$N_{1234} = h(N_{12} \mid N_{34}) \qquad \xrightarrow{\text{Sign}(h_{1234}, \text{SK})} \quad \sigma$$

$N_{12} = h(N_1 \mid N_2)$

$N_{34} = h(N_3 \mid N_4)$

$N_1 = h(d_1)$     $N_2 = h(d_2)$     $N_3 = h(d_3)$     $N_4 = h(d_4)$

Edge server returns $d_2$, $N_1$, $N_{34}$ and signed $N_{1234}$ (and the structure)
Client computes $N_{1234} = h(h(h(d_2)|N_1), N_{34})$ and verify that the
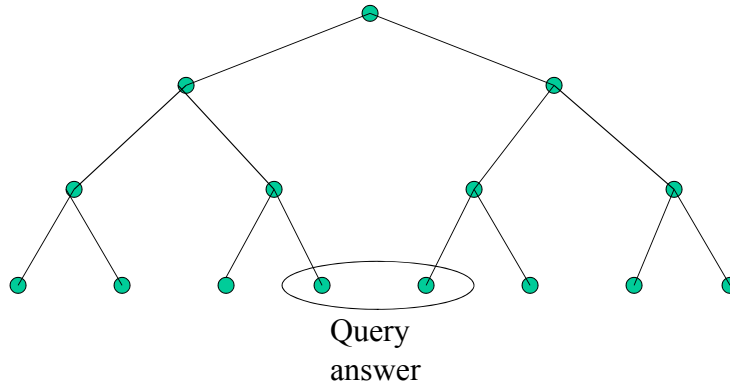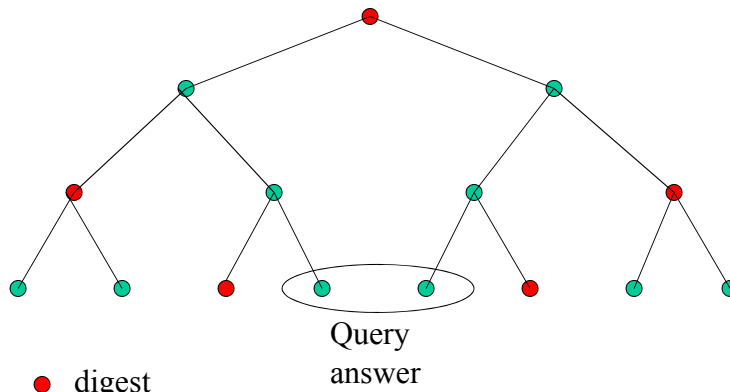signed value is correct     25

---

# Range Queries

Path l

LCA(q)

GLB(q)     q     LUB(q)

26

# Example: Range queries

Query
answer

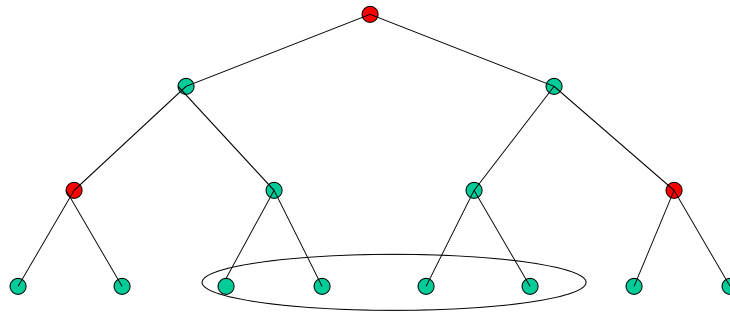What are returned?

# Example: Range queries

Query
answer

● digest

What are returned?

# Example: Range queries



● digest

What are returned?
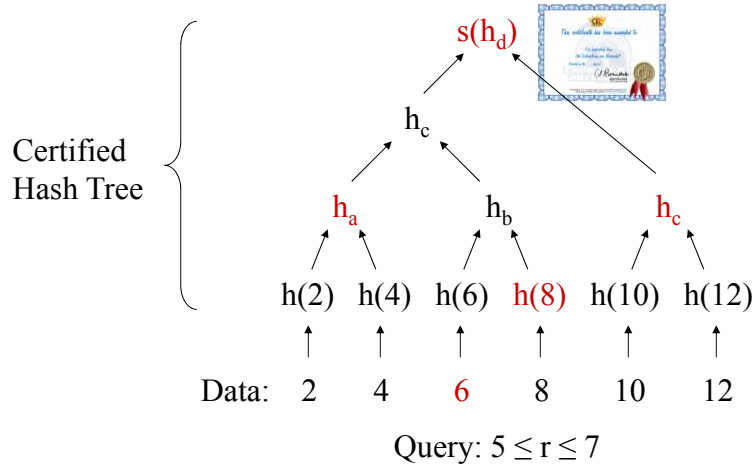
# Proving Authenticity is Easy



$s(h_d)$

$h_c$

Certified Hash Tree

$h_a = h(h(2)|h(4))$        $h_b$        $h_c$

$h(2)$  $h(4)$  $h(6)$  $h(8)$  $h(10)$  $h(12)$

Data:   2      4      6      8      10      12

Query: $5 \leq r \leq 7$

# Proving Authenticity is Easy

$s(h_d)$

$h_c$

$h_a$     $h_b$     $h_c$

Certified
Hash Tree

h(2)  h(4)  h(6)  h(8)  h(10)  h(12)

Data:   2     4     6     8     10     12

Query: $5 \leq r \leq 7$

31

# Proving Completeness is Easy But …

$s(h_d)$

$h_c$

$h_a$     $h_b$     $h_c$

Certified
Hash Tree

h(2)  h(4)  h(6)  h(8)  h(10)  h(12)

Data:   2     4     6     8     10     12

Query: $5 \leq r \leq 7$

32

# Precision may be compromised!



Certified Hash Tree

$s(h_d)$

$h_c$

$h_a$     $h_b$     $h_c$

$h(2)$  $h(4)$  $h(6)$  $h(8)$  $h(10)$  $h(12)$

Data:   2    4    6    8    10    12

Query: $5 \leq r \leq 7$

- Compromise precision: Disclose left and right neighbors
- May violate access control policy

33

---

# Example

- Access control: U can only see records with salary < 8000
- Results are records 2, 3, and 5.
- If system does not return record 1, U will not know that the answer is complete since it is possible that there is a record with Sal > 7000 but < 8000 that is not returned.
- If system returns record 1, then it violates the access control policy!
- Need an authentication mechanism that verifies completeness without compromising access control rules

| ID | Name | Sal | Dept |
|----|------|------|------|
| 5 | A | 2000 | 1 |
| 2 | C | 3500 | 2 |
| 1 | D | 8010 | 1 |
| 4 | B | 2200 | 3 |
| 3 | E | 7000 | 2 |

34

17

# What's the problem?

- A Merkle hash tree is needed for every sort-order on a table
- VO (Verification Object – the data used for verification) needs to contain links all the way to the root,
  - VO grows linearly to query result and logarithmic to base table size
- Projections may have to be performed by clients
- No provision for dynamic updates on the database
- Weak in terms of access control
  - Attributes that are supposed to be filtered out must also be returned for verification

35

# A signature-chain-based scheme: Let's start simple …

- Consider a sorted list of distinct integers, $R = \{r_1 \ldots, r_{i-1}, r_i, r_{i+1}, \ldots r_n\}$
- Retrieve record whose value is greater than or equal to $\alpha$
  - $\alpha \leq r$ (i.e., $\sigma_{\alpha \leq r}(R)$ )
- Result $Q = \{r_a, r_{a+1}, \ldots r_b \}$, i.e., $r_{a-1} < \alpha \leq r_a < r_{a+1} < \ldots r_{b = r_n}$
- Result is complete iff:
  - Contiguity: Each pair of successive entries $r_i$, $r_{i+1}$ in Q also appears in R (based on Signature Chain)
  - Terminal: Last element of Q is also last element of R, i.e., $r_b = r_n$ (based on Signature Chain)
  - Origin: $r_a$ is the first element in R that satisfies the query condition, i.e., $r_{a-1} < \alpha \leq r_a$ (based on Private Boundary Proof)

36

# Signature Chain

- For each data value, there is an associated signature
  - Computed from its own value, and that of its left and right neighbors
  - $sig(r_i) = s(h(g(r_{i-1}) \mid g(r_i) \mid g(r_{i+1})))$

$$\cdots \ r_{i-1} \ \text{[chain]} \ r_i \ \text{[chain]} \ r_{i+1} \ \text{[chain]} \ r_{i+2} \ \cdots$$

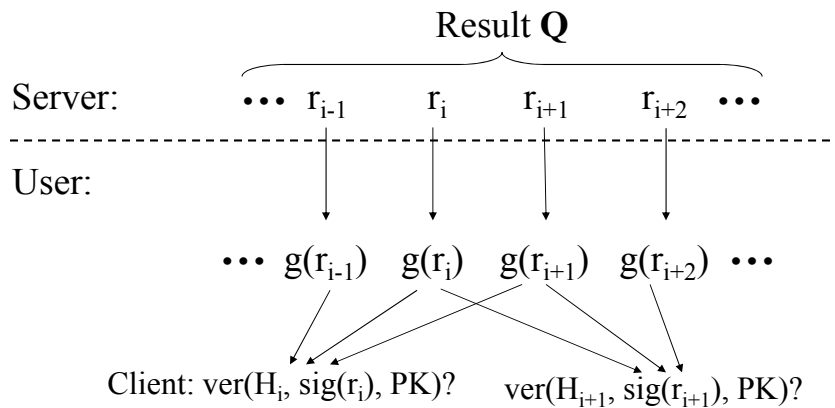- Owner stores the $(r_i, sig(r_i))$ pair in the server
- During querying, server returns (answer, signature) pairs and more …(verification objects) …

$h^i (r) = h^{i-1} ( h (r) ) \qquad h^0 (r) = h (r) \qquad g(r) = h^{U-r-1} (r)$
$U$ = max value outside of domain (known to all users)
$s$ is a signature function using owner's private key

37

---

# Signature Chain Ensures Contiguity

Server returns $(r_i , sig(r_i))$-pairs

Result **Q**

Server: $\cdots \ r_{i-1} \quad r_i \quad r_{i+1} \quad r_{i+2} \ \cdots$

User:

$\cdots \ g(r_{i-1}) \quad g(r_i) \quad g(r_{i+1}) \quad g(r_{i+2}) \ \cdots$

Client: $ver(H_i, sig(r_i), PK)?$ $\quad ver(H_{i+1}, sig(r_{i+1}), PK)?$

$H_i = h(g(r_{i-1}) \mid g(r_i) \mid g(r_{i+1}))$
Signature chain: $sig(r_i) = s(h(g(r_{i-1}) \mid g(r_i) \mid g(r_{i+1})))$
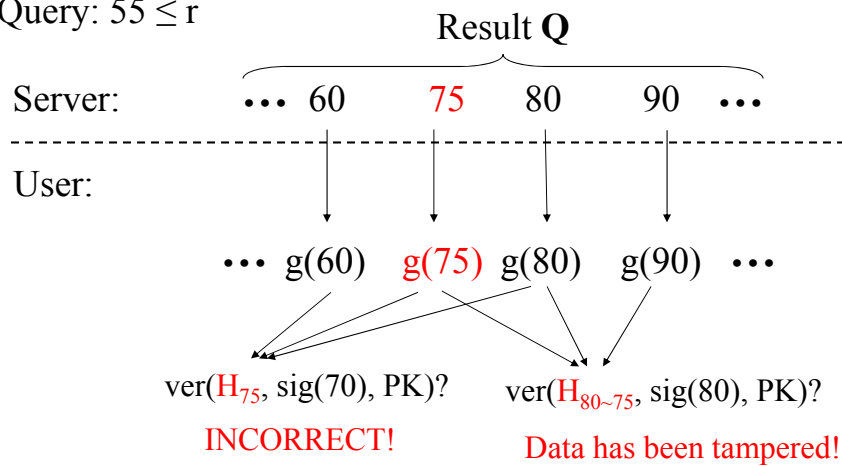
38

## Signature Chain Ensures Contiguity

Query: $55 \leq r$

Result **Q**

Server:  $\cdots$  60      70      80      90  $\cdots$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

User:

$\cdots$  g(60)   g(70)  g(80)   g(90)  $\cdots$

ver($H_{70}$, sig(70), PK)?          ver($H_{80}$, sig(80), PK)?

39

---

## Signature Chain Ensures Contiguity

Query: $55 \leq r$

Result **Q**

Server:  $\cdots$  60      75      80      90  $\cdots$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

User:

$\cdots$  g(60)   g(75)  g(80)   g(90)  $\cdots$

ver($H_{75}$, sig(70), PK)?          ver($H_{80\sim75}$, sig(80), PK)?

INCORRECT!                Data has been tampered!

40

# Signature Chain Ensures Contiguity

Query: $55 \leq r$

Result **Q**

Server: $\cdots$ 60  70  80  90 $\cdots$

User:

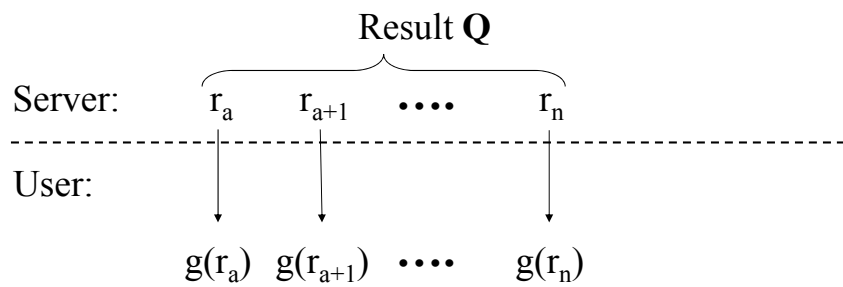$\cdots$ g(60)  g(80)  g(90) $\cdots$

ver($H_{80}$, sig(80), PK)?
$H_{60\sim80}$ will be computed (without 70) -
will not match sig(80). INCORRECT!!!!

41

---

# How To Ensure $r_n$ Is The Last Record?

Result **Q**

Server: $r_a$  $r_{a+1}$  $\cdots$  $r_n$

User:

g($r_a$)  g($r_{a+1}$)  $\cdots$  g($r_n$)

Create a fictitious record $r_{n+1}$ that is
larger than the largest value but smaller
than U
• $sig(r_{n+1}) = s(h(g(r_n)|g(r_{n+1})|h(U)))$

42

# Signature Chain Ensures Terminal

Result **Q**

Server: $r_a$ $r_{a+1}$ •••• $r_n$ $g(r_{n+1})$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
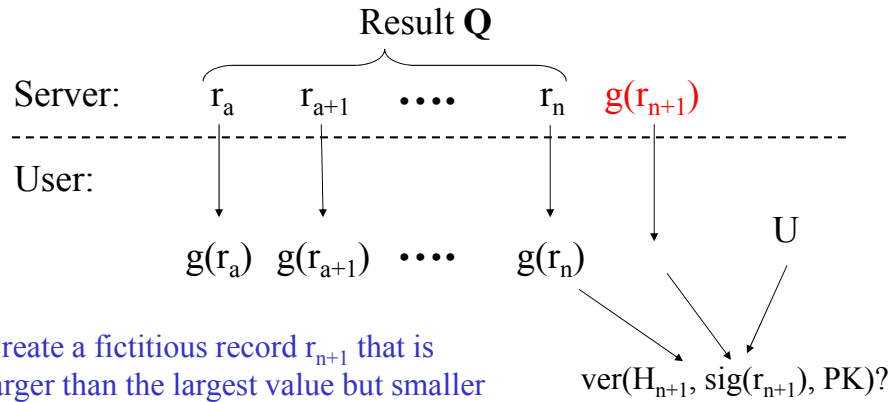
User:

$g(r_a)$ $g(r_{a+1})$ •••• $g(r_n)$

U

Create a fictitious record $r_{n+1}$ that is larger than the largest value but smaller than U

• $sig(r_{n+1}) = s(h(g(r_n)|g(r_{n+1})|h(U)))$

• server returns $g(r_{n+1})$ instead of $r_{n+1}$

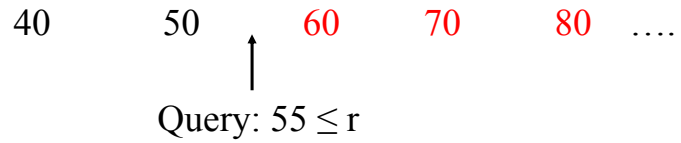$ver(H_{n+1}, sig(r_{n+1}), PK)?$

43

---

# How to prove Origin (without revealing the boundary point)??

40        50        60        70        80    ….

44

---

# How to prove Origin??

40        50        60        70        80    ….

↑

Query: $55 \leq r$
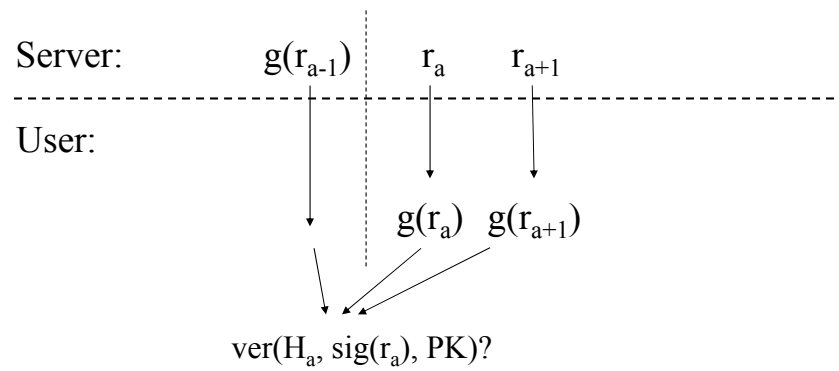
A naïve solution is to return 50. By proving that 50 is chained to 60, we know that no answer has been dropped. But, this reveals the value of 50.

45

# How about this …

Server:        $g(r_{a-1})$    $r_a$    $r_{a+1}$

User:

$g(r_a)$  $g(r_{a+1})$

$ver(H_a, sig(r_a), PK)$?

Query: $\alpha \leq r$

46

# The basic idea fails …

**Cheat!**

Server:   50   60   g(70)   |   80   90

User:

$g(r_a)$   $g(r_{a+1})$

$ver(H_a, sig(r_a), PK)?$   ⬅ User can't detect!

Query: $55 \le r$

47

---

# Private Boundary Proof Ensures Origin

Server:   $h^{\alpha - r_{a-1} - 1}(r_{a-1})$   $r_a$   $r_{a+1}$

User:   **??**

$g(r_{a-1})$   $g(r_a)$   $g(r_{a+1})$

$ver(H_a, sig(r_a), PK)?$

$h^i(r) = h^{i-1}(h(r))$
$g(r) = h^{U-r-1}(r)$

Query: $\alpha \le r$

48

## Private Boundary Proof Ensures Origin

Server: $\quad h^{\alpha - r_{a-1} - 1}(r_{a-1}) \quad r_a \quad r_{a+1}$

A **collaborative scheme** to compute the hash value

User: **hash U - α times**

$g(r_{a-1}) \quad g(r_a) \quad g(r_{a+1})$

$\text{ver}(H_a, \text{sig}(r_a), PK)?$

$$h^i(r) = h^{i-1}(h(r))$$
$$g(r) = h^{U-r-1}(r)$$
$$= h^{U-\alpha}(h^{\alpha-r-1}(r))$$

Query: $\alpha \leq r$

49

---

## Back to our example

**Cheat!**

Require that the inverse of $h^i$ for $i < 0$ be undefined

Server: $\quad h^{\alpha-70-1}(70) \quad 80 \quad 90$

User: hash U - 55 times

**Wrong!** Undefined! $\quad g(r_a) \quad g(r_{a+1})$

$\text{ver}(H_a, \text{sig}(r_a), PK)?$ ← User detects cheating

Query: $55 \leq r$

50

# Back to our example

Server:     $h^{55-50-1}$(50)    60    70   ....

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

User:     hash U - 55 times

g(50)     g(60)   g(70)

ver($H_{60}$, sig(60), PK)?

Query: $55 \leq r$

51

---

# Putting the Pieces Together

Result **Q**

Distributor:    $h^{\alpha-r_{a-1}-1}(r_{a-1})$   $r_a$    $r_{a+1}$   •••   $r_n$   $g(r_{n+1})$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

User:     hash U - α times

$g(r_{a-1})$    $g(r_a)$   $g(r_{a+1})$   •••   $g(r_n)$

ver($H_a$, sig($r_a$), PK)?      ver($H_{a+1}$, sig($r_{a+1}$), PK)?

Query: $\alpha \leq r$

52

# Other cases

- $\alpha \leq r$
- $\beta \geq r$ ( Result = $\{r_a, r_{a+1}, \ldots r_b\}$, i.e., $r_a, \ldots r_b \leq \beta < r_{b+1}$
  - Need to verify that $r_{b+1} > \beta$
  - Define $g(r) = \mathbf{h}^{\,r-L-1}\,(\mathbf{r}) = \mathbf{h}^{\,\beta - L}\,(\,\mathbf{h}^{\,r-\beta-1}\,(\mathbf{r})\,)$ where L is a value outside of the minimum value of the domain
- So, we have $\alpha \leq r \leq \beta$
- $r = \alpha \equiv \alpha \leq r \leq \alpha$
- $\alpha < r < \beta \equiv \alpha+1 \leq r \leq \beta-1$
- $\alpha \neq r \equiv (L < r < \alpha) \cup (\alpha < r < R)$

53

---

# NULL Answers??

- Consider Q: $\alpha \leq r$.
- $Q = \emptyset$ because $r_n < \alpha$.
  - Server returns $\mathbf{h}^{\,\alpha - r_n - 1}\,(\mathbf{r}),\ \mathbf{g(r_{n+1})},\ \mathbf{sig(r_{n+1})}$
  - User computes $\mathbf{h}^{\,U - \alpha}\,(\,\mathbf{h}^{\,\alpha - r_n - 1}\,(\mathbf{r})\,)$ and verifies $ver(H_{n+1}, sig(r_{n+1}), PK)$?

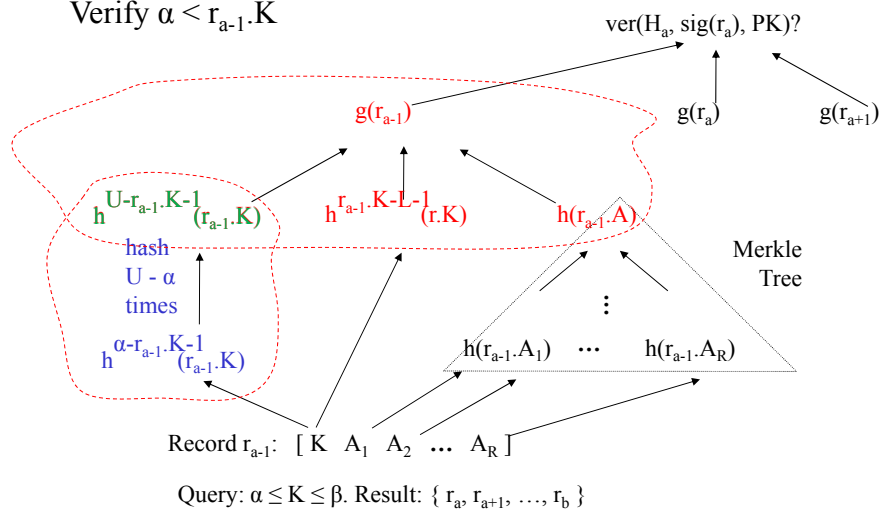- How about $r_i < \alpha \leq \beta < r_{i+1}$ ?

54

27

# One More Vulnerability

- User can discover $r_{a-1}$ through brute force enumeration of numbers below $r_a$
- Solution:
  - Record $[K, A_1, .., A_m]$, $K$ = ordering attribute
  - $g(r_i.K \mid r_i.A_1 \mid \ldots \mid r_i.A_m)$
  - Brute-force attack is no longer feasible

55

# Completeness Verification for Range Queries



Verify $\alpha < r_{a-1}.K$

$ver(H_a, sig(r_a), PK)?$

$g(r_{a-1})$

$g(r_a)$  $g(r_{a+1})$

$h^{U-r_{a-1}.K-1}(r_{a-1}.K)$   $h^{r_{a-1}.K-L-1}(r.K)$   $h(r_{a-1}.A)$

Merkle Tree

hash
$U - \alpha$
times

$h^{\alpha-r_{a-1}.K-1}(r_{a-1}.K)$

$h(r_{a-1}.A_1)$ $\cdots$ $h(r_{a-1}.A_R)$

Record $r_{a-1}$:  $[\, K \quad A_1 \quad A_2 \; \ldots \; A_R \,]$

Query: $\alpha \leq K \leq \beta$. Result: $\{\, r_a, r_{a+1}, \ldots, r_b \,\}$

56

28

# Other queries

- SP Query
  - Based on MHT(r.A)
  - Ordering attribute has to be returned (even if it is not part of the target attributes). Why?
  - For attributes that are filtered out, digests may need to be returned
- SPJ Query
  - R.Ai = S.Aj  (Ai is foreign-key in R, Aj is primary key in S)
    - Referential integrity constraint mandates that every instance of R.Ai must have a matching entry in S.Aj
    - So, only need to deal with selection conditions on R.Ai or S.Aj
    - Create a signature chain for R.Ai

57

# What else?

- What about data freshness?
- More efficient scheme
- Ad-hoc joins
- Aggregates
- Multi-dimensional data
- Computation
- Complete (complex) queries

58

# Summary

- Malicious service provider may cheat
- Users need assurance on their query answers
- Merkle hash tree offers a good solution but …
- Signature chain guarantee completeness without violating access control policy

59