

# Secure Data Management at UW

Gerome Miklau and Dan Suciu

Trustworthy computing seminar  
Autumn 2004

1

## Today's talk

### 1 Controlled data publishing

secrecy: crypto + XML [VLDB 2003]

### 2 Analyzing information disclosure

secrecy: theory + relations [SIGMOD 2004]  
[ICDT 2005]

### 3 Tamper-resistant databases

integrity: crypto + relations [Current work]

2

2

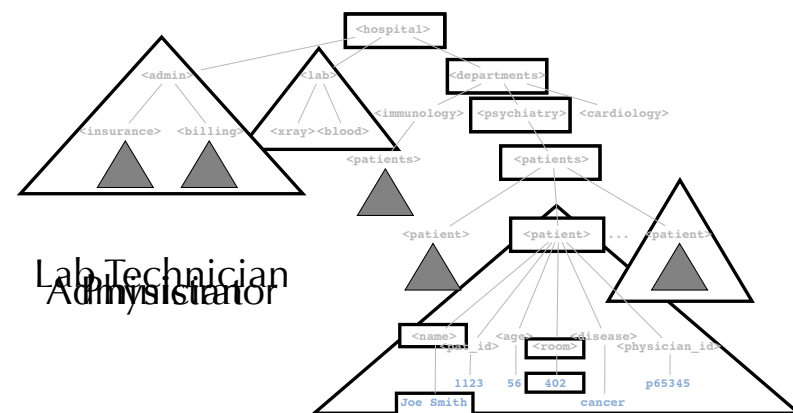
## Ingredients for data sharing

- ✓ File formats and tools  
XML, XML query languages
- ✓ Distributed processing  
Networked data sources  
Mediator systems, distributed systems
- 📁 Access control  
controlled distribution of data

3

3

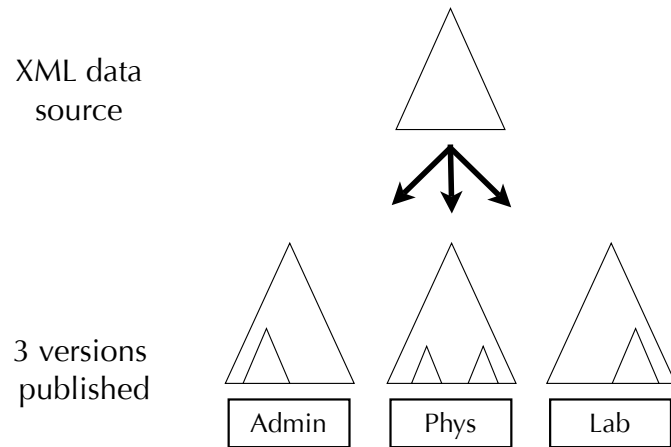
## XML data and access rights



4

4

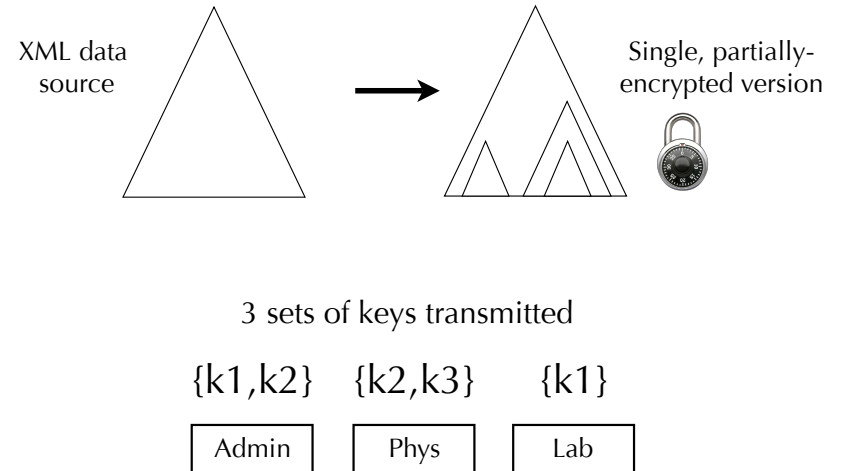
# Publishing Many Views



5

5

# Publishing one view + keys



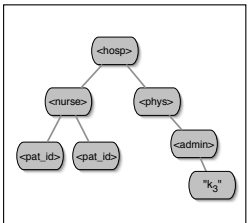
6

6

# Publishing protected data

① Policy Query	② Tree Protection	③ XML Output
----------------	-------------------	--------------

```
FOR $x in /doc/subj/subj
  $y in /doc/psy/psych
WHERE $x/exam-psych/id = $y/id
KEY   getKey($y)
TARGET $x
```



```
<EncryptedData>
  <EncryptionMethod Alg="AES"
    KeySize="128"/>
  <KeyInfo>
    <name>k1</name>
  </KeyInfo>
  <CipherData>
    <CipherValues>
      qZk+NkcGgWq6PiVxeFDCbjz
      DCbjzQ2jkcGgWq6PIVxeFFD
    </CipherValues>
  </CipherData>
</EncryptedData>
```

Evaluate policy on data

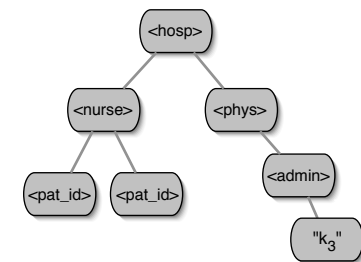
Encrypt, write as XML

7

7

# Tree Protection

- XML data tree guarded with formulas over keys.
- Node is accessible if "satisfying" keys are supplied.



Grammar for key formulas  
 $\sigma := T \mid F \mid k \mid \sigma \wedge \sigma' \mid \sigma \vee \sigma'$

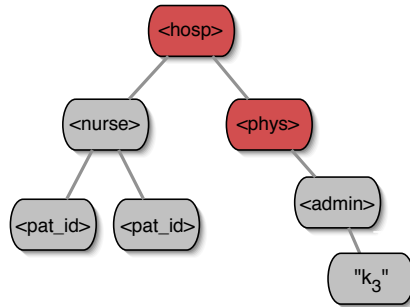
8

8

# Access Function

Given a set of keys  $K$ ,  $\mathbf{access}(K)$  computes the accessible nodes in a tree protection.

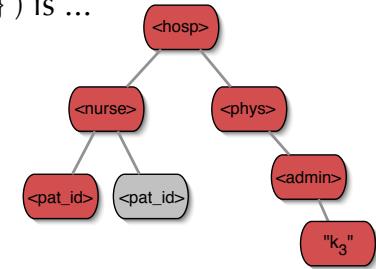
$\mathbf{access}(\{k_1\})$  is:



# Access Function

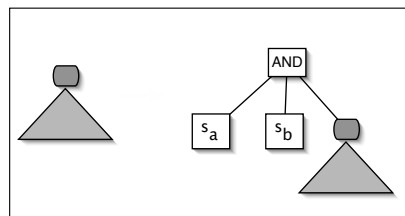
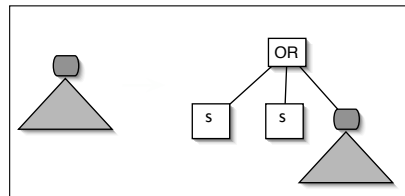
Access function can discover data values and use them as keys.

Example:  $\mathbf{access}(\{k_1, k_2\})$  is ...



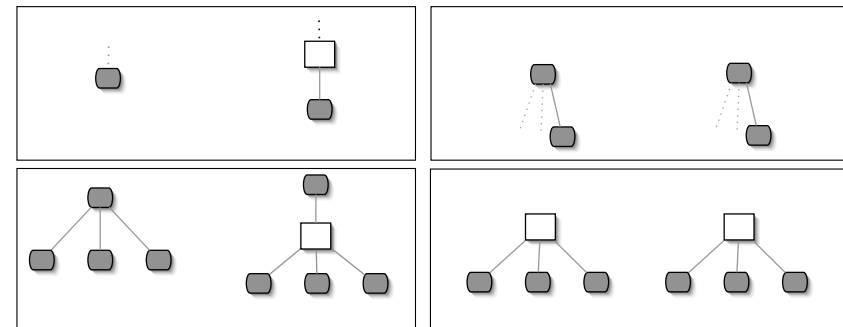
# Normalization Rules

- a tree protection is normalized if every formula is atomic
- soundness: access function invariant under rewritings



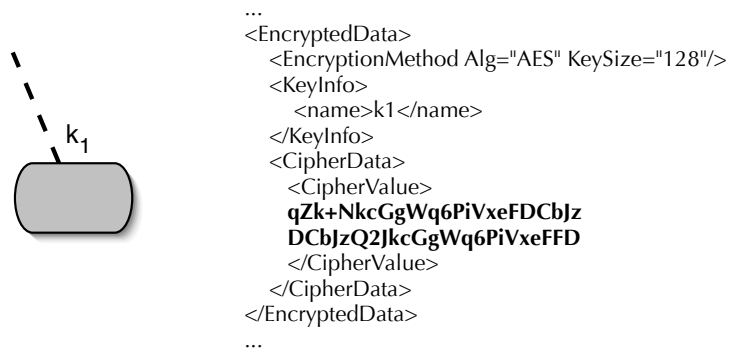
Benaloh and Leichter. CRYPTO 1988

# Rewrite rules



Rewrite rules can be used to optimize formula expressions.

# Implementing a Protection



Leaf node

Encrypted leaf node, following  
XML Encryption Recommendation

13

13

# Security

## Security Property

- given document D, set of keys K:
- ✓ if node  $x \in \text{access}(K)$  then  $x$  is efficiently computable from D.
  - ☐ if node  $x \notin \text{access}(K)$  then the best algorithm for finding  $x$  requires guessing keys.

Additional disclosure:

- number of children of a node
- size of ciphertext, duplicate subtrees,
- policy information.

14

14

# Today's talk

## 1 Controlled data publishing

secrecy: crypto + XML

## 2 Analyzing information disclosure

secrecy: theory + relations

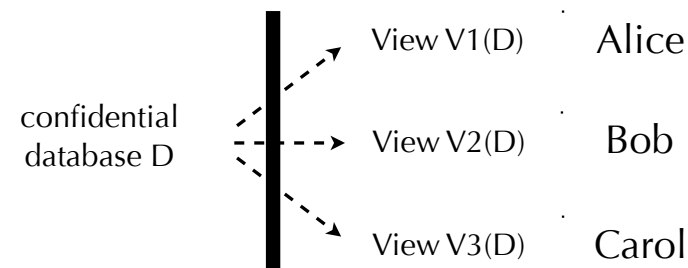
## 3 Tamper-resistant databases

integrity: crypto + relations

15

15

# Protecting data using views



16

16

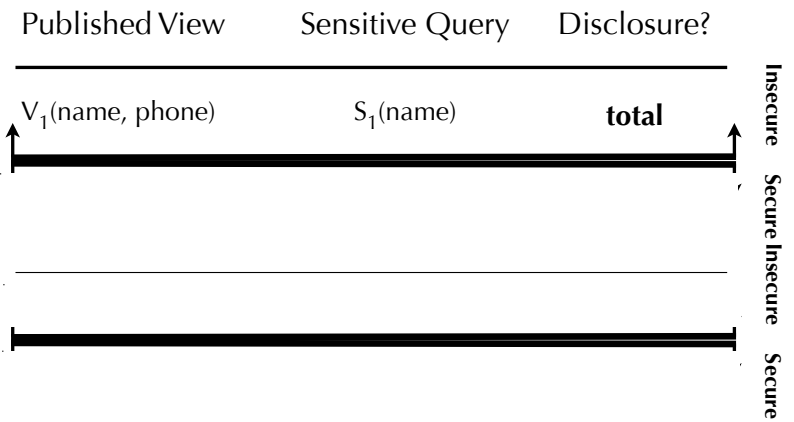
# Problem statement

## Query - View security

- a published view  $V$
- a sensitive query  $S$
- Does  $V$  reveal anything about  $S$  ?

# Spectrum of disclosure

Employee (name, dept, phone)



# Existing techniques

- Logical inference
  - Answering queries using views
  - Statistical databases
- exact disclosure  
 aggregation of numerical attributes

# Our basic strategy

Compare Mallory's knowledge about sensitive query  $S$ :

$$\text{Knowledge about } S, \text{ a priori} \stackrel{?}{=} \text{Knowledge about } S, \text{ given } V$$

When these are equal,  $V$  provides **no information** about  $S$

# The adversary's knowledge

Knowledge about S, *a priori*  $\stackrel{?}{=}$  Knowledge about S, given V

- |  |                     |
|--|---------------------|
| <ul style="list-style-type: none"> <li>• Schema</li> <li>• Domain</li> <li>• Probability of each database</li> </ul> | + V<br>+ V(I) = ans |
|--|---------------------|

Probability of each answer to S  $\stackrel{?}{=}$  Probability of each answer to S

# Probabilities of databases

- Fix schema R, and domain.
- For each tuple t:  $0 \leq \text{Pr}[t] \leq 1$
- The probability of a database instance I is:

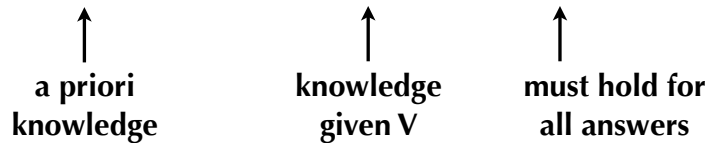
$$\text{Pr}[I] = \prod_{t \in I} \text{Pr}[t] \times \prod_{t \notin I} (1 - \text{Pr}[t])$$

**assumption: tuple-independence**

# Definition: query-view security

S and V are secure (denoted  $S|V$ ) if:

$$\text{Pr}[S(I)=x] = \text{Pr}[S(I)=x | V(I)=y] \quad \forall x \forall y$$

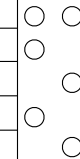


- independence of probabilistic events
- inspired by Shannon's perfect secrecy [1949]

# Concrete example

- relation Edge(X,Y)
- nodes={a,b}
- tuple probability = 1/2
- possible graphs:

1			
2	(a,a)		
3	(a,b)		
4	(b,a)		
5	(a,a) (a,b)		
6	(a,a) (b,a)		
	⋮		
15	(a,a) (a,b) (b,a)		
16	(a,a) (a,b) (b,a) (b,b)		



**Sensitive query:**

$S(x) :- \text{Edge}(x,y)$

for  $S(I) = \{ (a) \}$ ,

$$\text{Pr}[S(I) = \{ (a) \}] = 3/16$$

**Published view:**

$V(y) :- \text{Edge}(x,y)$

for  $V(D) = \{ (a) \}$ ,

$$\text{Pr}[S(I) = \{ (a) \} | V(I) = \{ (a) \}] = 1/3$$

## Goal: logical criterion

Brute force → Logical criterion

For each  
answer to S  
answer to V  
Compare probabilities

depends on domain &  
probability distribution

Analyze the view  
expressions S and V

25

25

## Deciding query-view security

Theorem

Given query S and view V, deciding whether  
 $S \mid V$  is  $\Pi_2^P$ -complete

when S and V secure, then they are secure:

- for any (sufficiently large) domain
- for any probability distribution

26

26

## Other consequences

Supplemental knowledge

- security standard easily extended
- compare

Encrypted view

- no sensitive query is secure

Supplemental info: cardinality

- no sensitive query is secure

27

27

## Measuring disclosure

When query-view security fails:

for some x,y:

$$\Pr[ S(I)=x ] \neq \Pr[ S(I)=x \mid V(I)=y ]$$

↖ **how do we evaluate  
the difference?**

- magnitude, or closeness to 1?
- aggregate over many answers
- see paper for limited case.

28

28

## Partial disclosure

- Intuition
  - domains are large
  - databases are small (and of known size)
- New probabilistic model
  - each tuple  $t$  equally-likely...
  - $\text{prob}[t]$  s.t. database size constant
- Practical security
  - $\lim \Pr( S | V ) = 0$  as domain  $\rightarrow \infty$

[ Dalvi, Miklau, Suciu. ICDT 2005 (to appear) ]

29

29

## Today's talk

- 1 Controlled data publishing  
secrecy: crypto + XML
- 2 Analyzing information disclosure  
secrecy: theory + relations
- 3 Tamper-resistant databases  
integrity: crypto + relations

30

30

## Integrity

### Data integrity

an assurance that data has not been modified by an unauthorized party.

### Consistency

an assurance that the items in a collection are "fresh".

### Query integrity

an assurance that a query answer is correct and complete.

31

31

## Tampering threats to DB systems

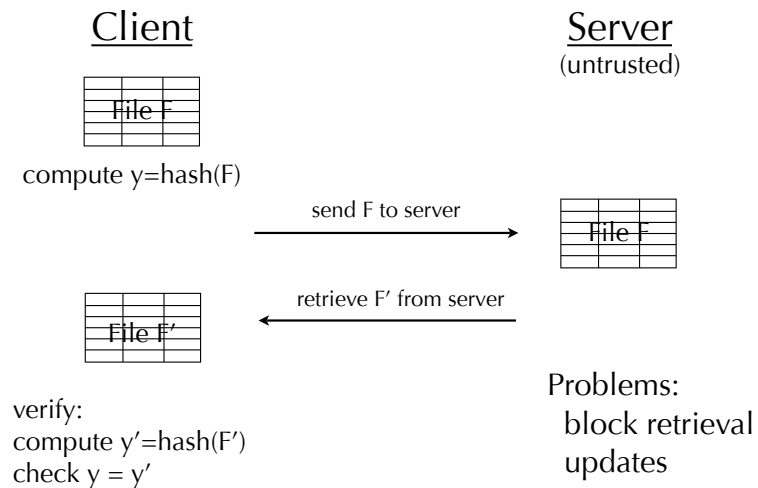
- DB access control vulnerabilities
  - specification failure
  - enforcement failure
  - subversion (e.g. sql injection, weak authentication)
- DB extensions - user defined functions
- general OS and network threats
- privileged parties: OS admin, DB admin
- service provider model

32

32



# Data integrity with hashing

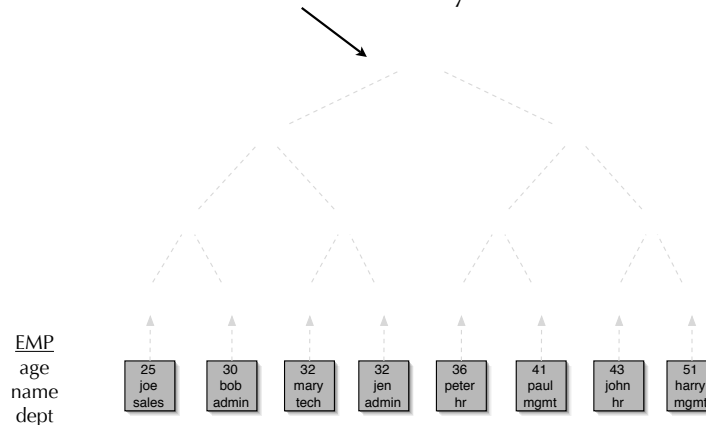


33

33

# Merkle hash tree on relation

Client stores root hash locally



[Merkle 1989] [Devanbu, et al. 2001]

34

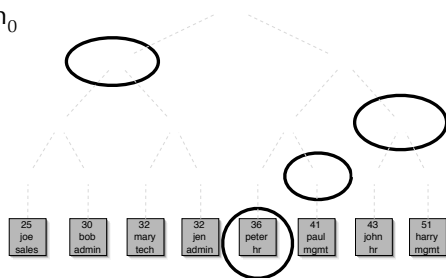
34

# Verifiable query (1)

CLIENT: Give me the record for 'peter'

SERVER: answer is  $t_5$   
hash path is  $h_{101}, h_{11}, h_0$

CLIENT: verify by computing  
new root hash  $h'_\epsilon$   
check  $h'_\epsilon = h_\epsilon$



This proves correctness and consistency

35

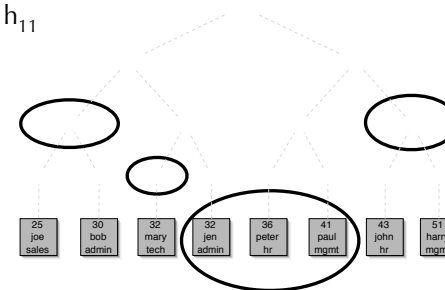
35

# Verifiable query (2)

CLIENT: Return all employees with age between 35 and 40

SERVER: answer is  $t_5$ . Also return  $t_4, t_6$ .  
hash path is  $h_{010}, h_{00}, h_{11}$

CLIENT: verify by computing  
new root hash  $h'_\epsilon$   
check  $h'_\epsilon = h_\epsilon$



This proves correctness, completeness, and consistency

36

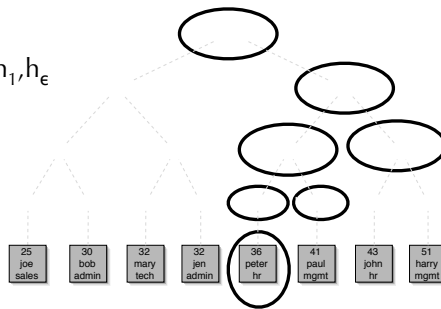
36

# Updating the database

CLIENT: Move Peter from department HR to MGMT

-- verify peter's record --

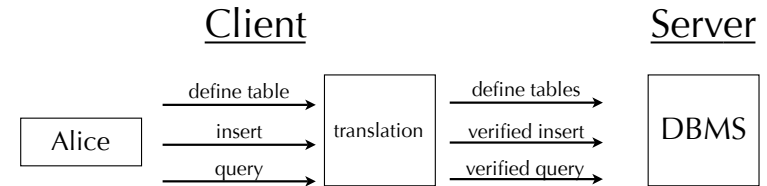
CLIENT: compute new  $h_{100}, h_{10}, h_1, h_\epsilon$   
send to server



37

37

# Implementing a tamper-resistant database



- Smart client, oblivious server
- Relational representation of hash tree
- Query definition
- Index selection

38

38

# Cost of integrity

- Reasonable communication overhead
- Reasonable client computation
- Modest storage overhead
- Good scalability
- Throughput:

preliminary results

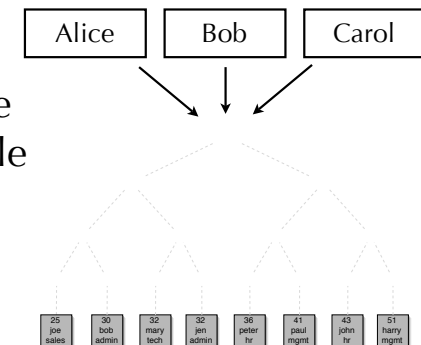
	integrity	no integrity	multiple
Query	2.0 ms	.4 ms	5.0
Range query	6.1 ms	1.3 ms	4.7
Insert	8.3 ms	.8 ms	10

39

39

# Multiple clients

How do we manage integrity with multiple users?



40

40

Questions ?