

CS 6431: Encrypted Databases

Instructor: Tom Ristenpart

<https://github.com/tomrist/cs6431-fall2021>



**CORNELL
TECH**

HOME OF THE
**JACOBS
INSTITUTE**



Minimizing data exposure risk

- Access controls
 - Least privilege principle applied to database
 - Monitoring & logging accesses
- Encrypted databases
 - Encryption at rest
 - Property-revealing encryption
- Data privacy protections
 - De-identification
 - K-anonymity
 - Differential privacy

Threat models for databases



Client
(e.g., laptop,
workstation,
web server)



Person	Age	Income
Alice	22	100,000
Bob	29	25,000
Charlie	31	40,000

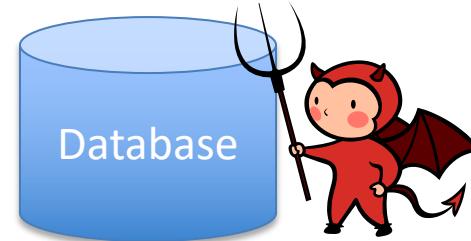
Malicious client accesses:

- Missing / misconfigured access control
- Break-in to authenticated client
- SQL injection attacks

Threat models for databases



Client
(e.g., laptop,
workstation,
web server)

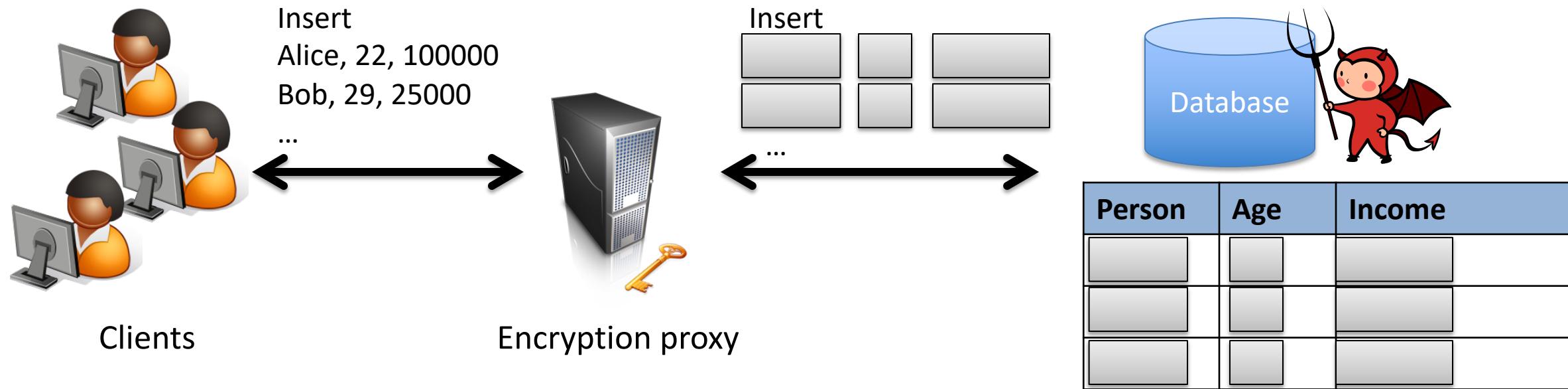


Person	Age	Income
Alice	22	100,000
Bob	29	25,000
Charlie	31	40,000

Compromise of database system:

- Loss or theft of persistent storage (hard disks)
- Malicious insider attack
- Subpoena / government access
- Complete compromise of system hosting database

One approach: database encryption



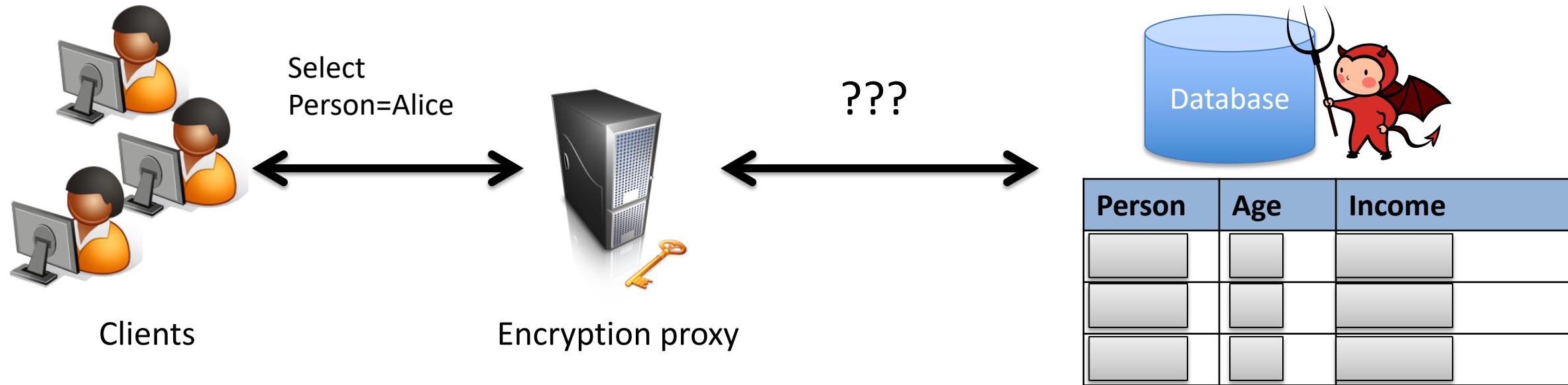
Encrypt data using secret key before insertion into database

Decrypt records fetched via queries

Encryption vs. DB functionality tension!
How can we fetch encrypted records?

Proxy can handle many clients (e.g., entire organization)

One approach: database encryption



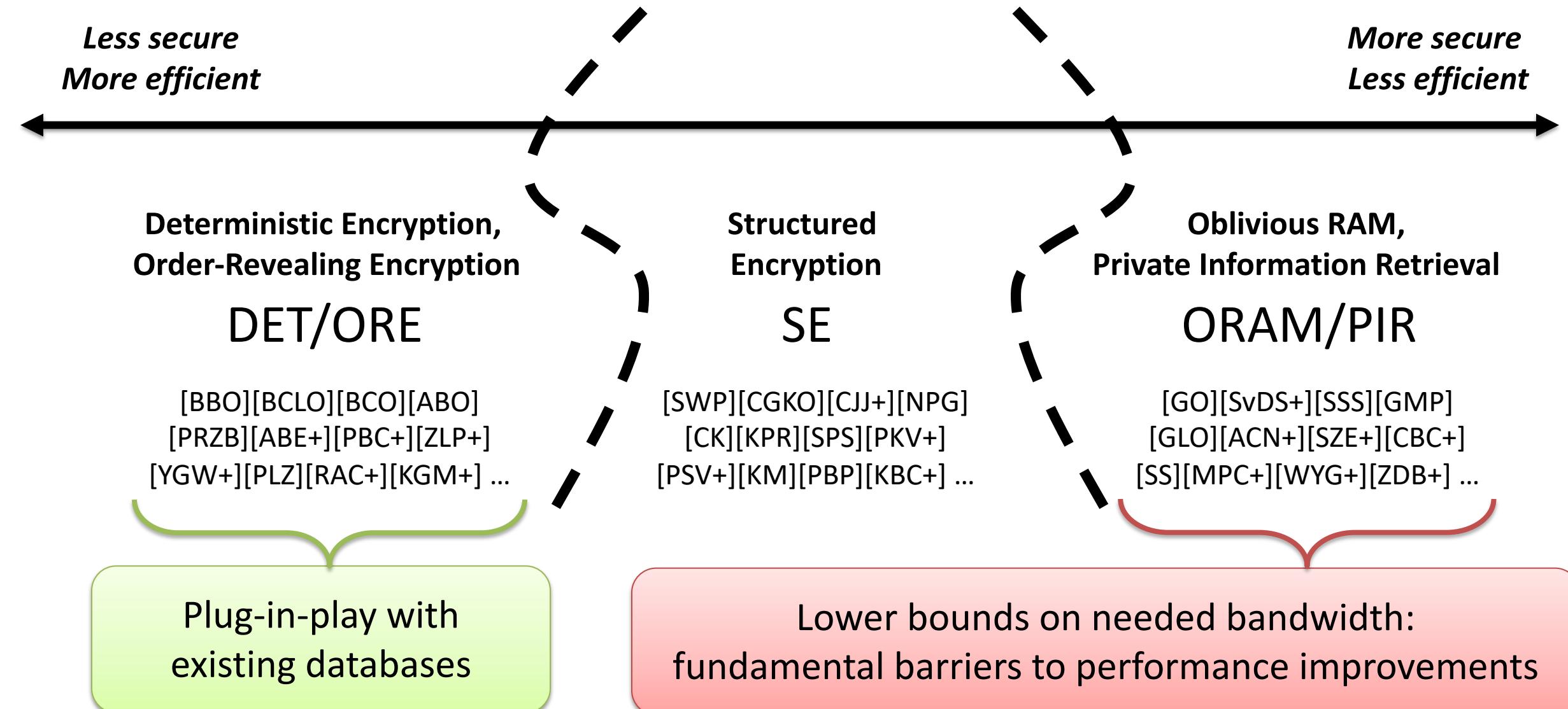
Encrypt data using secret key before insertion into database

Decrypt records fetched via queries

Encryption vs. DB functionality tension!
How can we fetch encrypted records?

Proxy can handle many clients (e.g., entire organization)

Decades of academic work on this tension

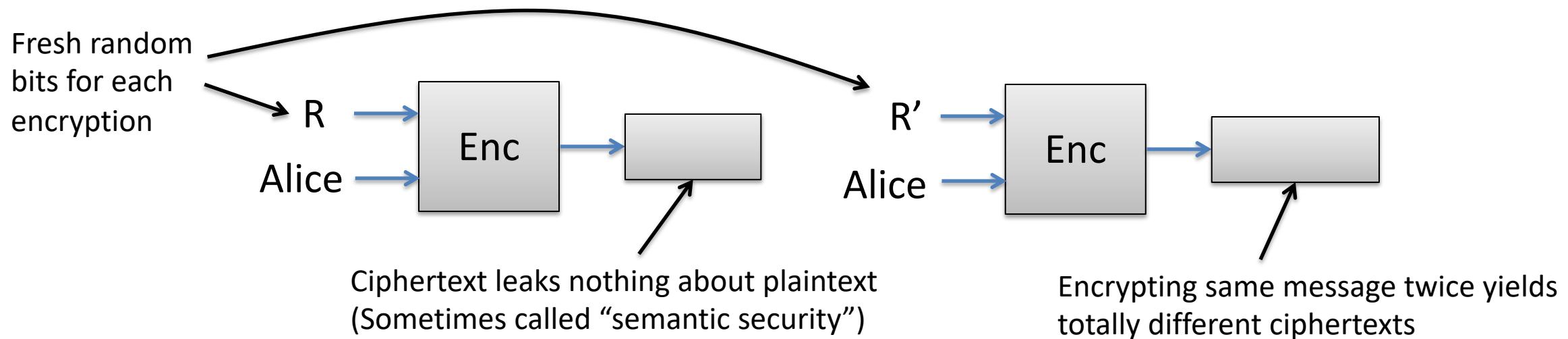


Deterministic encryption (DET)

The goal:

Encryption **reveals plaintext equality** to allow efficient record search

Conventional encryption *does not* reveal plaintext equality, by design:



Deterministic encryption (DET)

The goal:

Encryption **reveals plaintext equality** to allow efficient record search

Deterministic encryption: same plaintext maps to same ciphertext

$$\text{DET}(M) = \text{DET}(M') \text{ iff } M = M'$$



Examples:

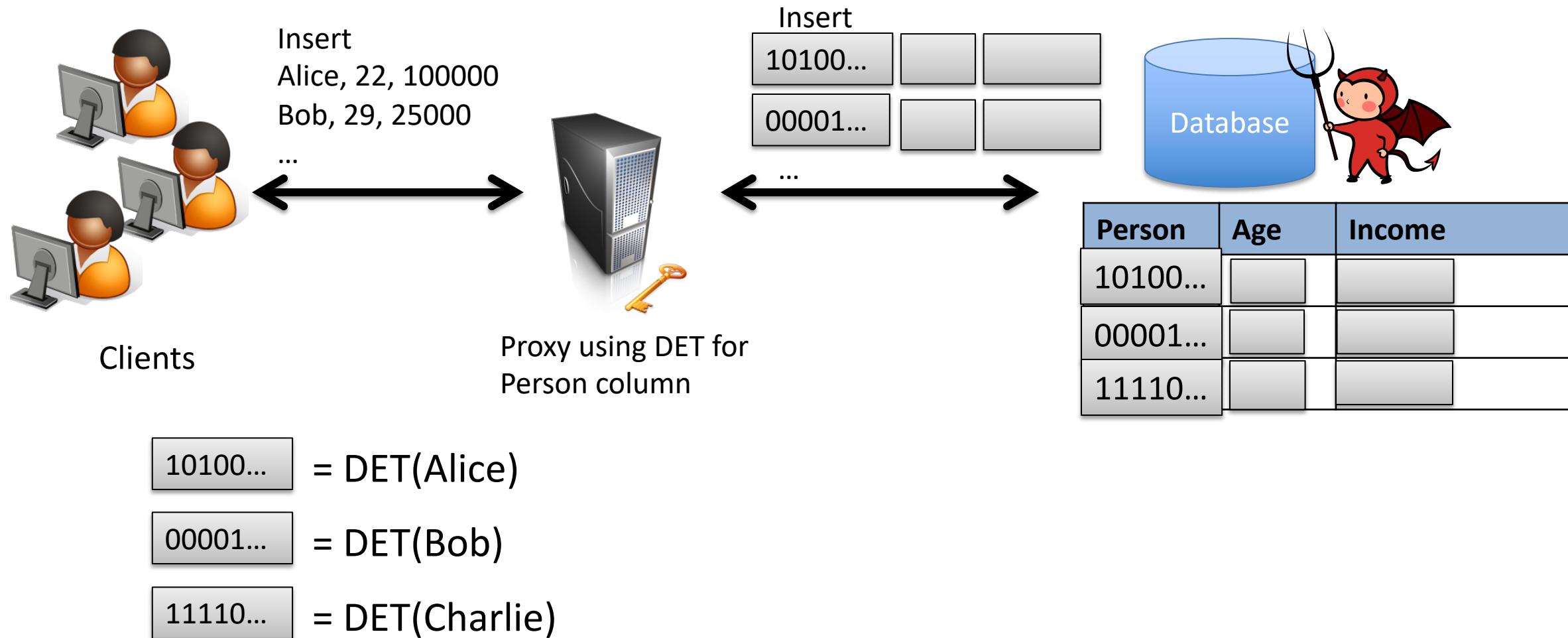
SIV mode, FFX

Encrypting same message twice yields
same ciphertext

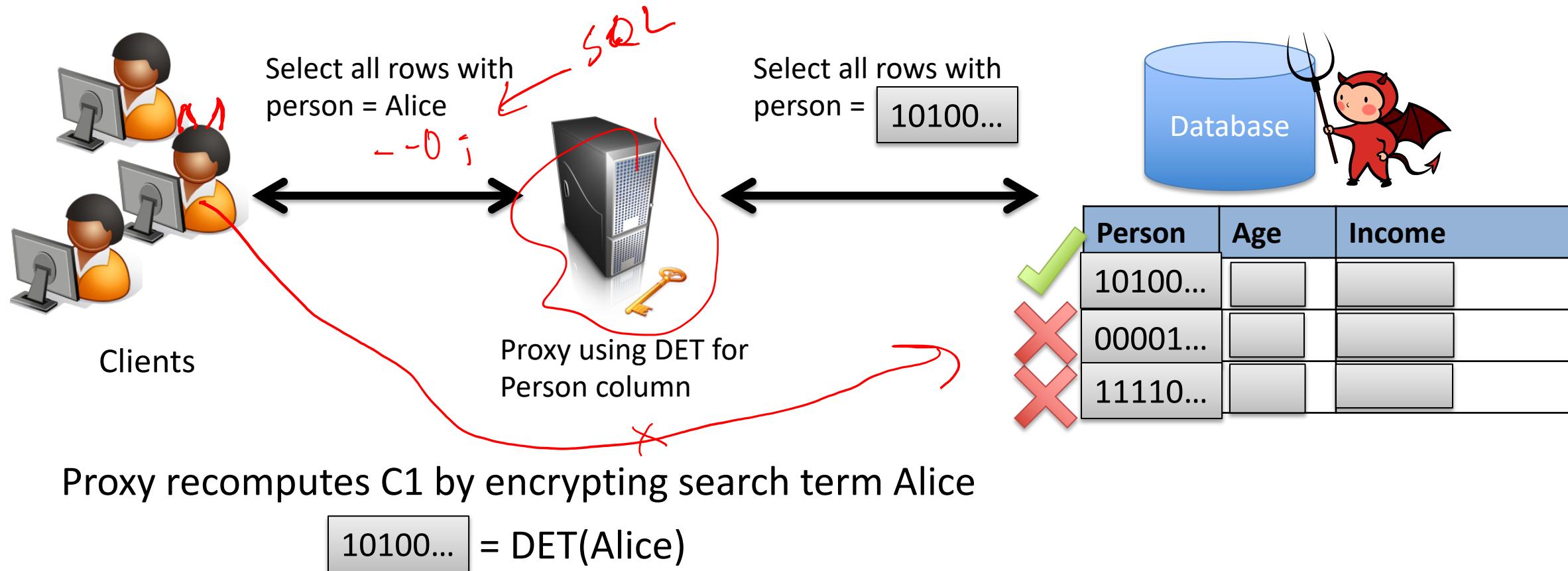
Leaks plaintext equality

Allows equality search (find all rows with name “Alice”)

Using DET in encrypted DBs



Using DET in encrypted DBs



Deterministic encryption enables fast retrieval of matching rows, but...

ciphertexts leak plaintext equality, allowing *leakage abuse attacks*

CryptDB paper [Popa et al. 2011]

- One of early academic works to suggest using property-revealing encryption for SQL-style DBs
 - “CryptDB’s approach is to execute queries over encrypted data”
- Popularized idea, particularly in academic circles & tech media
 - Described as efficient alternative to fully homomorphic encryption
- Combine existing DET, OPE schemes (plus one new one for database joins) with onioning (wrapping with randomized encryption)

Companies have deployed DE and/or ORE



Sometimes referred to as enabling
“computing on encrypted data”

Encryption *reveals just enough* about plaintexts
to perform specific DB operations

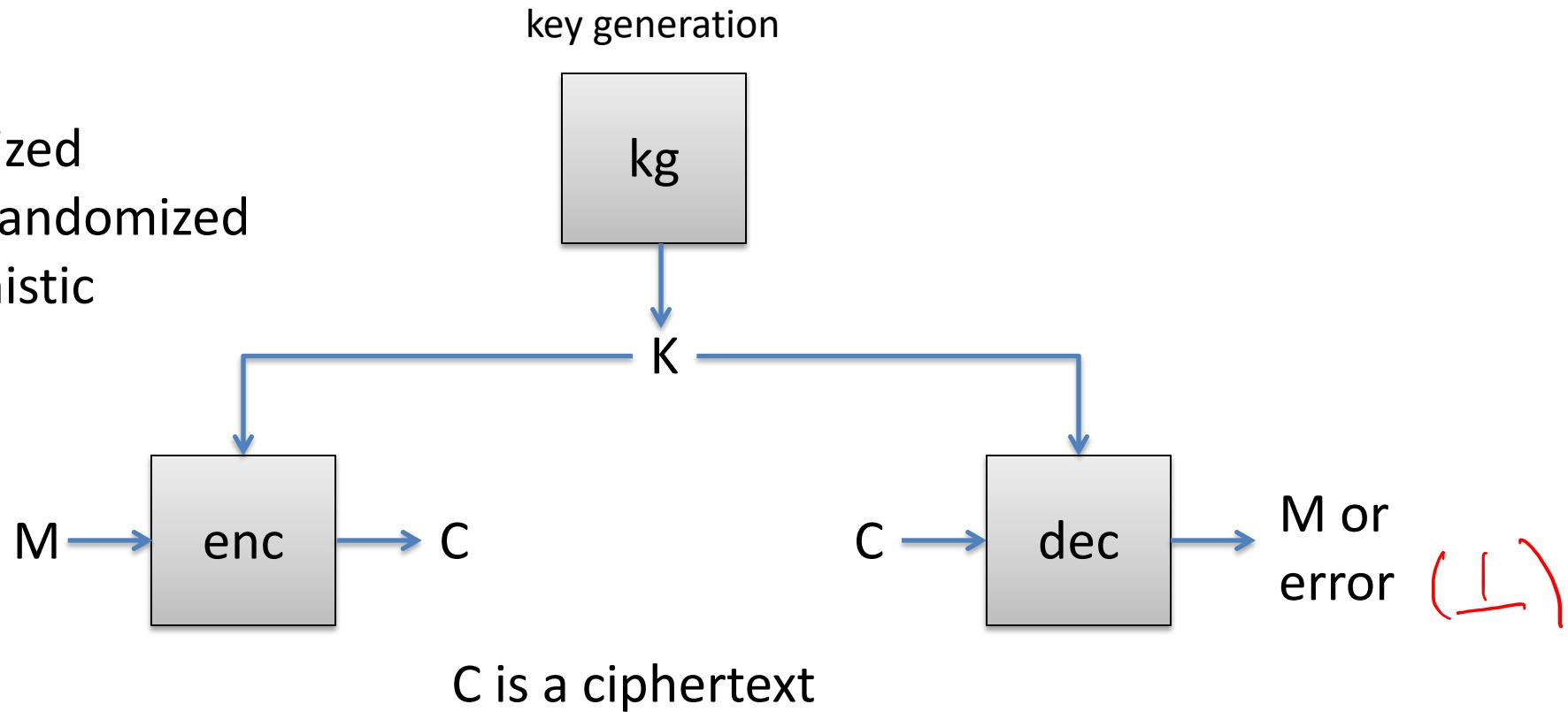
“*Property-revealing encryption*”

Formal security analyses *proving security**

* for definitions of security that allow
leaking information about plaintexts

Symmetric encryption (SE) definition

kg is randomized
enc **may** be randomized
dec deterministic



Correctness: $\Pr[\text{dec}(K, \text{enc}(K, M)) = M] = 1$, probability over enc's randomness

Write $C \leftarrow_{\$} \text{enc}_K(M)$ to signify running enc with fresh randomness

Assume $|C| = \text{clen}(|M|)$ for any M

Real-or-random indistinguishability (IND\$)

Standard security goal for symmetric encryption

Can't distinguish encryption of chosen messages from random string of bits

$\text{REAL}_{\text{SE}}^{\mathcal{A}}$

$K \leftarrow \$ \text{kg}$

$b' \leftarrow \$ \mathcal{A}^{\text{Enc}}$

Ret b'

$\text{Enc}(M)$ *K ⊕ m*

$C \leftarrow \$ \text{enc}_K(M)$

Ret C

$\text{RAND}_{\text{SE}}^{\mathcal{A}}$

$b' \leftarrow \$ \mathcal{A}^{\text{Enc}}$

Ret b'

$\text{Enc}(M)$

$C \leftarrow \$ \{0, 1\}^{\text{clen}(|M|)}$

Ret C

Adversary can repeatedly query Enc on same message

OTP encryption

$$\text{enc}_K(m) = K \oplus m$$

$$\text{dec}_K(c) = c \oplus K$$

Adversary is a randomized algorithm with access to Enc oracle

Advantage function measures efficacy of an adversary:

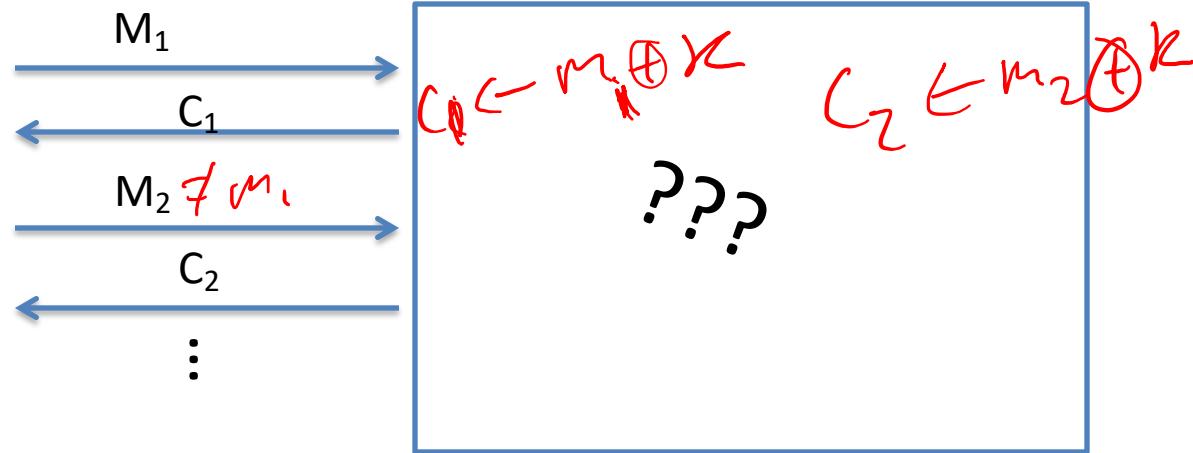
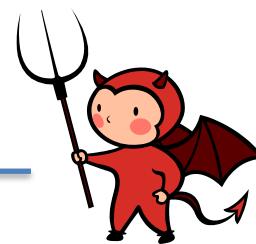
$$\text{Adv}_{\text{SE}}^{\text{ind}\$}(\mathcal{A}) = |\Pr[\text{REAL}_{\text{SE}}^{\mathcal{A}} \Rightarrow 1] - \Pr[\text{RAND}_{\text{SE}}^{\mathcal{A}} \Rightarrow 1]|$$

Probabilities are over random coins used in games

Real-or-random indistinguishability (IND\$)

Adversary gets to submit messages to oracle.
Doesn't know which oracle it has access to

If $C_1 \oplus C_2 = m_1 \oplus m_2$



Intuition:

Ciphertexts are indistinguishable from random bits. Contain no information about plaintext

- Such information-theoretic security only possible if $|K| = |M|$ (e.g., one-time pad)

Instead, restrict to **computationally efficient** adversaries

- Breakthrough paradigm shift from 1980s (e.g., Goldwasser, Micali paper)

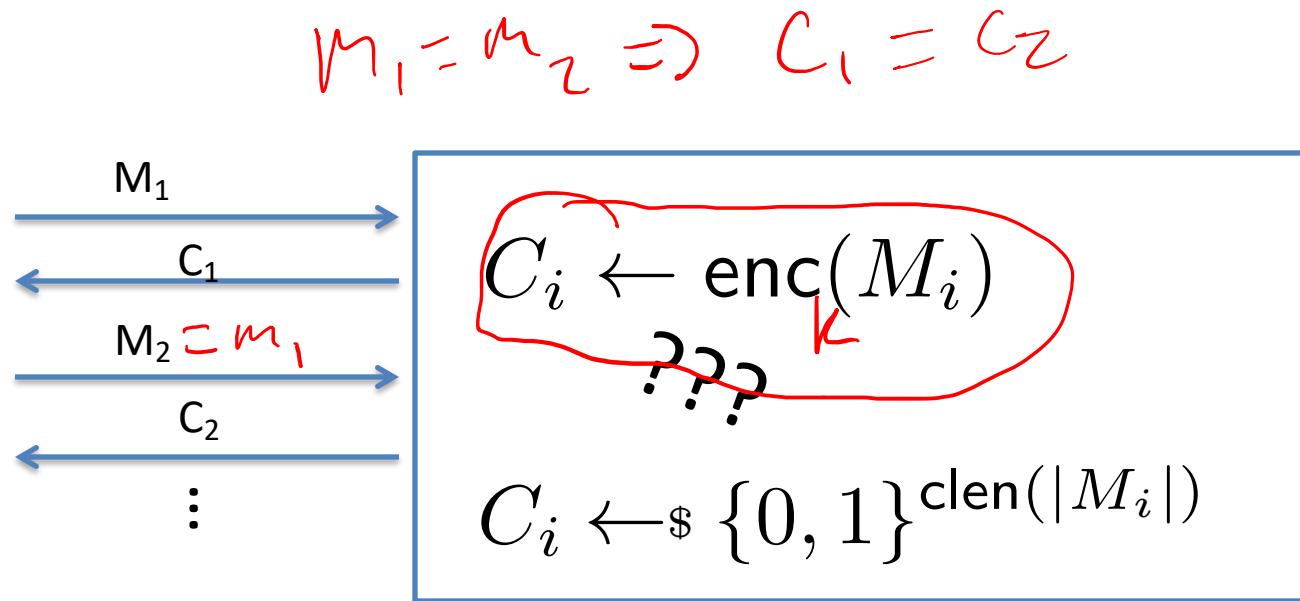
Deterministic encryption is IND\$-insecure

Adversary gets to submit messages to oracle.

Doesn't know which oracle it has access to

If $C_1 = C_2$ then Ret 1
Ret 0

$(m_1 \neq m_2)$



How can adversary choose queries in order to distinguish between REAL and RAND?

REAL: getting back deterministic encryptions of plaintext

RAND: getting back randomly chosen bit string no matter what is input

$$\Pr[\text{REAL}_{\text{DET}} \Rightarrow q] = 1$$

$$\Pr[\text{RAND} \dots]$$

$$= \frac{1}{2^{\text{clen}(|M_1|)}}$$

Real-or-random indistinguishability for deterministic encryption

$$\text{REAL}_{\text{SE}}^{\mathcal{A}}$$
$$K \xleftarrow{\$} \mathbf{kg}$$
$$b' \xleftarrow{\$} \mathcal{A}^{\text{Enc}}$$

Ret b'

$$\text{Enc}(M)$$
$$C \xleftarrow{\$} \mathsf{enc}_K(M)$$

Ret C

$$\text{RAND}_{\text{SE}}^{\mathcal{A}}$$
$$b' \xleftarrow{\$} \mathcal{A}^{\text{Enc}}$$

Ret b'

$$\text{Enc}(M)$$
$$C \xleftarrow{\$} \{0, 1\}^{\mathsf{clen}(|M|)}$$

Ret C

Adversary can repeatedly
query Enc on same message

Real-or-random indistinguishability for deterministic encryption

$\text{REAL}_{\text{SE}}^{\mathcal{A}}$

$K \leftarrow \$ \text{kg}$

$b' \leftarrow \$ \mathcal{A}^{\text{Enc}}$

Ret b'

$\text{Enc}(M)$

$C \leftarrow \$ \text{enc}_K(M)$

Ret C

$\text{RAND}_{\text{SE}}^{\mathcal{A}}$

$b' \leftarrow \$ \mathcal{A}^{\text{Enc}}$

Ret b'

$\text{Enc}(M)$

$C \leftarrow \$ \{0, 1\}^{\text{clen}(|M|)}$

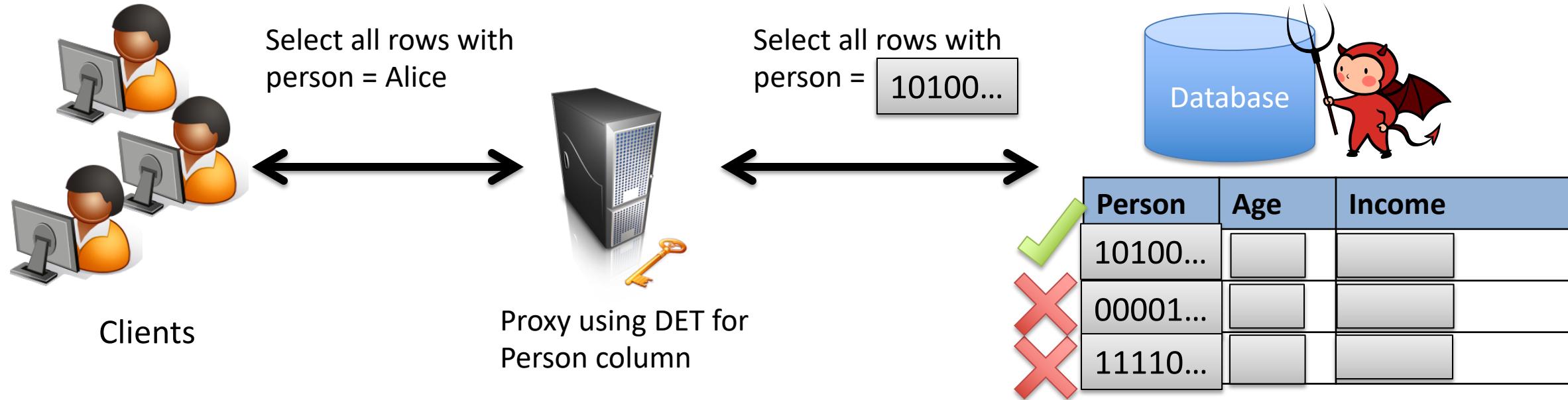
Ret C

Adversary **cannot** repeatedly query Enc on same message

Can formalize via pseudocode changes to game or as restriction on class of adversaries considered

Can show schemes that achieve this: Synthetic IV (SIV) from [Rogaway, Shrimpton 2006]
What does this notion of security allow in terms of plaintext leakage?

Using DET in encrypted DBs



Proxy recomputes C1 by encrypting search term Alice

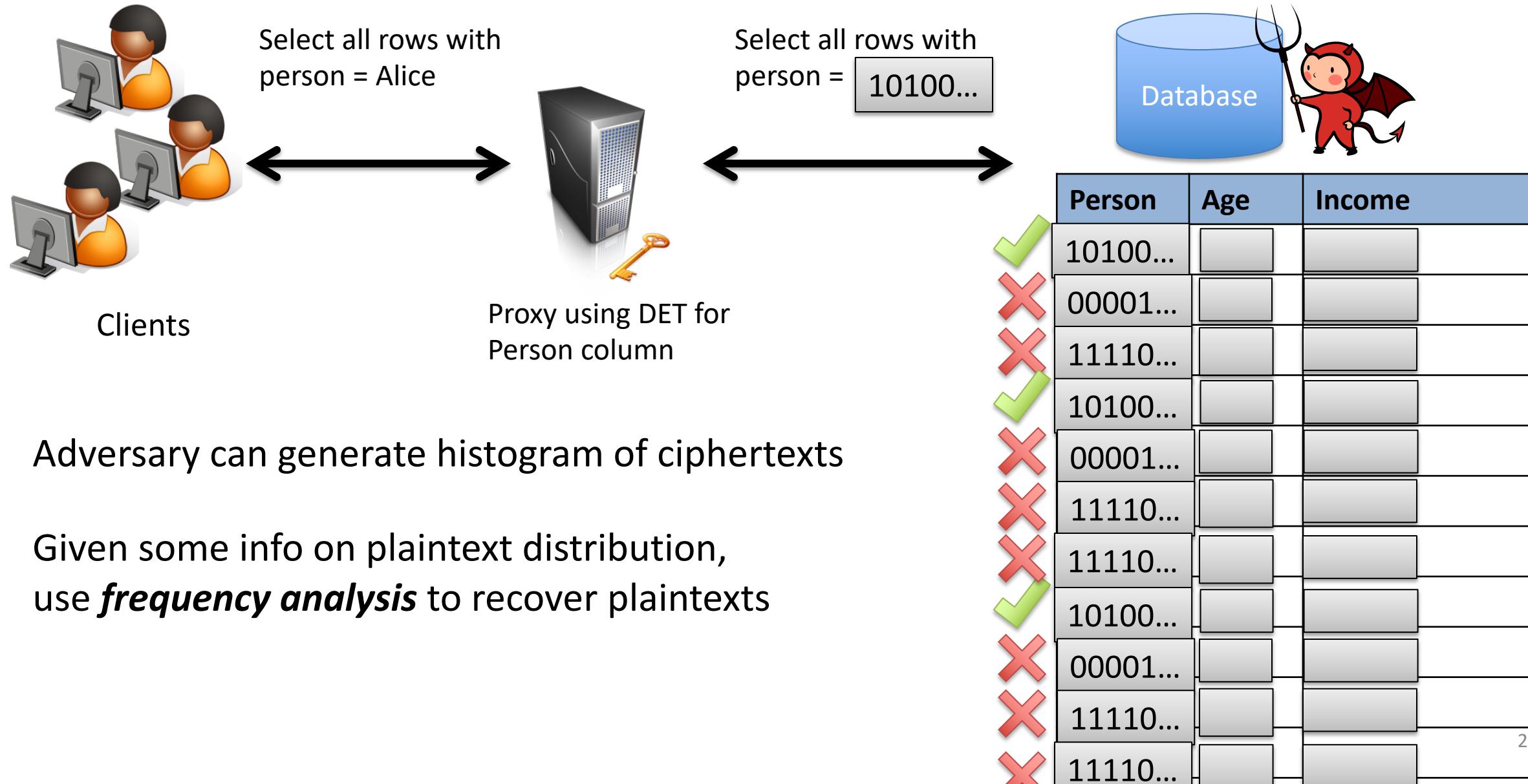
$$10100\ldots = \text{DET}(\text{Alice})$$

Deterministic encryption enables fast retrieval of matching rows, but...

ciphertexts leak plaintext equality, allowing *leakage abuse attacks*

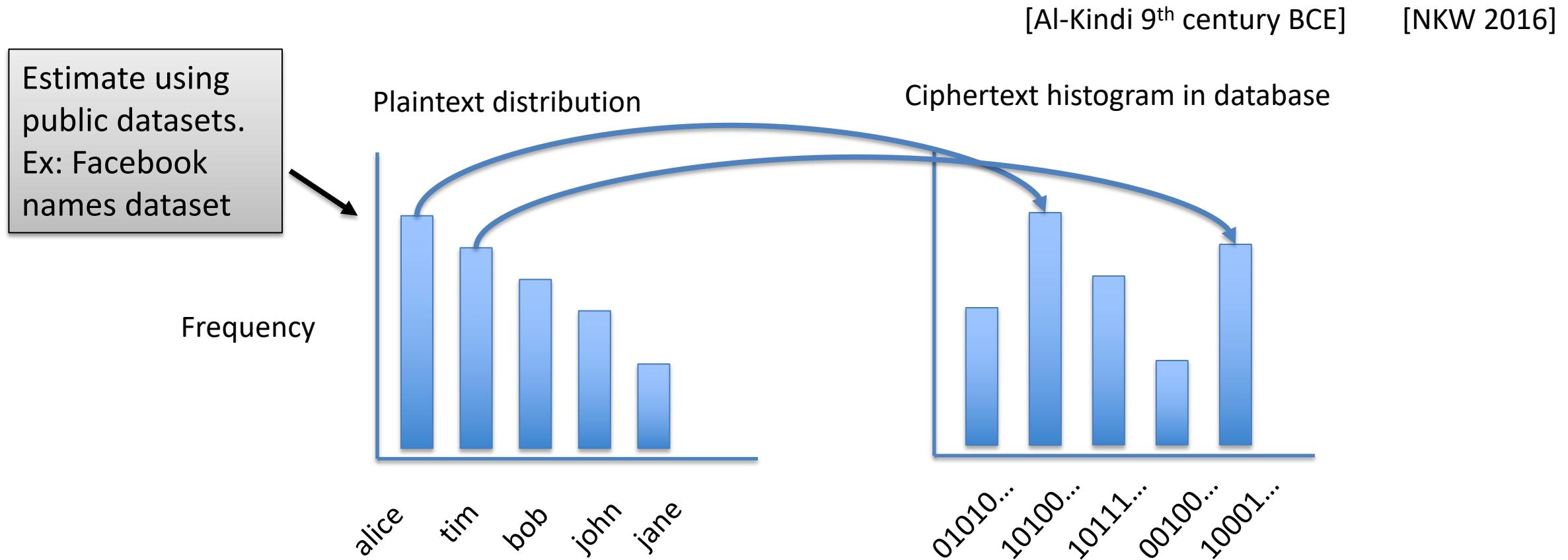
Leakage-Abuse Attacks against DET

[Naveed,
Kamara, Wright
2015]



Frequency attacks against DET

Idea: if attacker knows something about distribution of plaintexts, we can exploit it



For deterministic encryption this is all information adversary gets from database
Frequency analysis optimal attack [Lacharite, Paterson 2015]

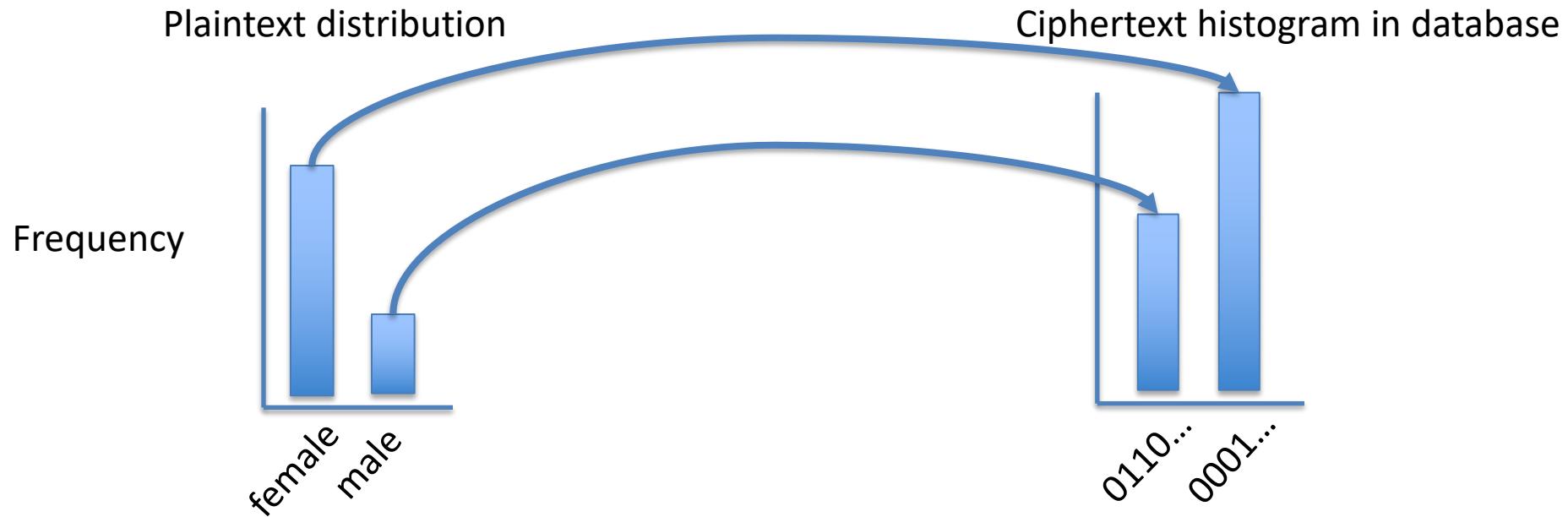
DET frequency analysis can be very effective

Non-uniform, small-domain values easy for adversary to recover

[Bindschaedler et al. 2018] simulated attacks:

Recover 100% of DET-encrypted values in hospital datasets

Most columns binary or small number of possibilities



DET frequency analysis can be very effective

Non-uniform, small-domain values easy for adversary to recover

[Bindschaedler et al. 2018] simulated attacks:

Recovers 100% of DET-encrypted values in hospital datasets

Most columns binary or small number of possibilities

More uniform, larger-domain values can be harder for adversary to recover

[Naveed et al. 2016] simulated attacks:

Recovers 10% of DET-encrypted age values in hospital datasets

Data closer to uniform
Good plaintext distribution estimates harder to get

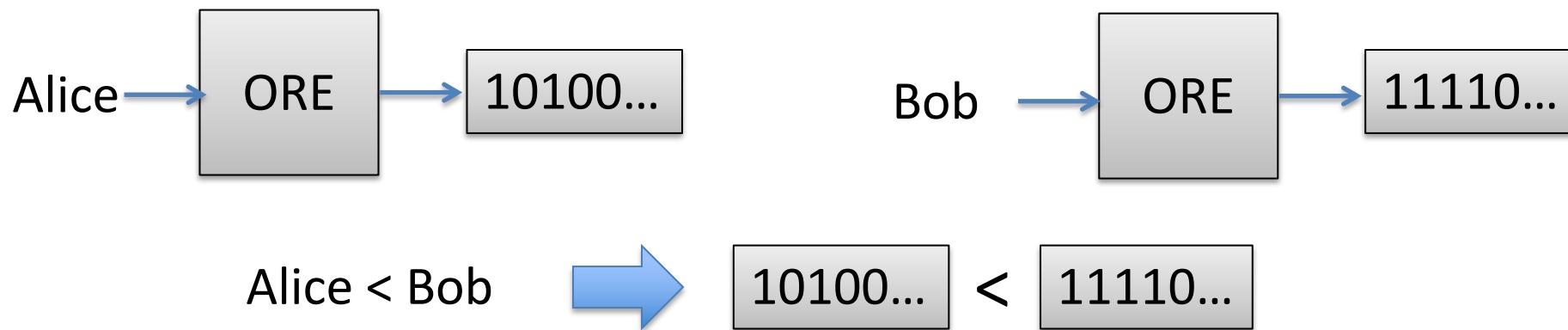
Order-revealing encryption (ORE)

The goal:

Encryption *reveals plaintext equality and ordering* to allow efficient range searches

Deterministic encryption such that ciphertexts reveal plaintext order

$$\text{ORE}(M) \leq \text{ORE}(M') \text{ iff } M \leq M'$$



Practical schemes by: [Boldyreva, Chennette, Lee, O'Neill 2009] (BCLO)
[Chenette, Lewi, Weis, Wu 2015] (CLWW)

CLWW construction

Uses a pseudorandom function (PRF), e.g.: $F_K(X) = \text{HMAC}(K, X)$ for secret key K

Let n -bit input message be $M = M_1 M_2 \dots M_n$

$$u_i = F_K(i \| M_1 M_2 \dots M_{i-1} \| 0^{n-i}) + M_i \bmod 3$$

Ciphertext is u_1, u_2, \dots, u_n

Given two ciphertexts $\begin{matrix} u_1, u_2, \dots, u_n \\ u'_1, u'_2, \dots, u'_n \end{matrix}$ how to infer ordering of M, M' ?

Find index j s.t. $u_j \neq u'_j$

If $u_j > u'_j$ then output $M > M'$

Otherwise output $M' < M$

Leaks more than order! What does it leak?

Trivial attacks against ORE: BCLO leakage

Some ORE schemes leak a lot more than just order of plaintexts

[BCO 11] showed that BCLO leaks about half of plaintext bits trivially

Plaintext	Ciphertext	m_c
michael	cyrzjipnouushzh	michaekypfbkfr
david	aenpse cevvpkmr	david jwbvhec
robert	emlqrnycvblqqnd	robert lwyeorrr
john	ccnncczzpruvjhd	johmzzzsfbunn
james	bzkxrq gzortby	james zyovtq
daniel	aelfspocabjdvjc	daniel jgaginu
richard	ekrzjmjhjxykbba	richardkmfnwwx
jose	ccqrlzzziozokby	josdzzzxvfruqq
mark	cwmlfzzzjxhlkh	marjzzzxqyduv
christopher	zokwwbrbibyouo	christotnqfolw

Apply a few arithmetic operations to ciphertext

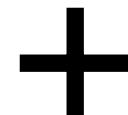
[Grubbs et al. 2017] attacks

New attack framework that combines:

frequency
analysis



ordering
leakage



scheme-specific
leakage

Scheme(s)	First names	Last names
Kerschbaum [27]	26%	6%
Popa et al. [36], Kerschbaum [28]	84%	38%
BCLO [12, 13]	99%	97%
CLWW [18]	98%	75%
BCLO + CLWW [18]	85%	44%
Baseline Guessing	4%	1%

Percentage of row values recovered
from encrypted database

Academic literature on leakage abuse attacks

Improvements on attacks

Varying leakage types

- Multiple-column leakage
- Nearest neighbor access patterns
- “Volume” leakage (# of values returned for range query)

Important for guiding practice:



[Islam, Kuzu, Kantarcioglu 2013]

[Cash, Grubbs, Perry, Ristenpart 2015]

[Naveed, Kamara, Wright 2015]

[Durak, DuBuisson, Cash 2016]

[Poillot, Wright 2016]

[Grubbs, McPherson, Naveed, Ristenpart, Shmatikov 2016]

[Zhang, Katz, Papamanthou 2016]

[Kellaris, Kolios, Nissim, O'Neill 2016]

[Grubbs, Sekniqi, Bindschaedler, Naveed, Ristenpart 2017]

[Grubbs, Ristenpart, Shmatikov 2017]

[Bindschaedler, Grubbs, Cash, Ristenpart, Shmatikov 2018]

[Grubbs, Lacharité, Minaud, Paterson 2018]

[Grubbs, Lacharité, Minaud, Paterson 2019]

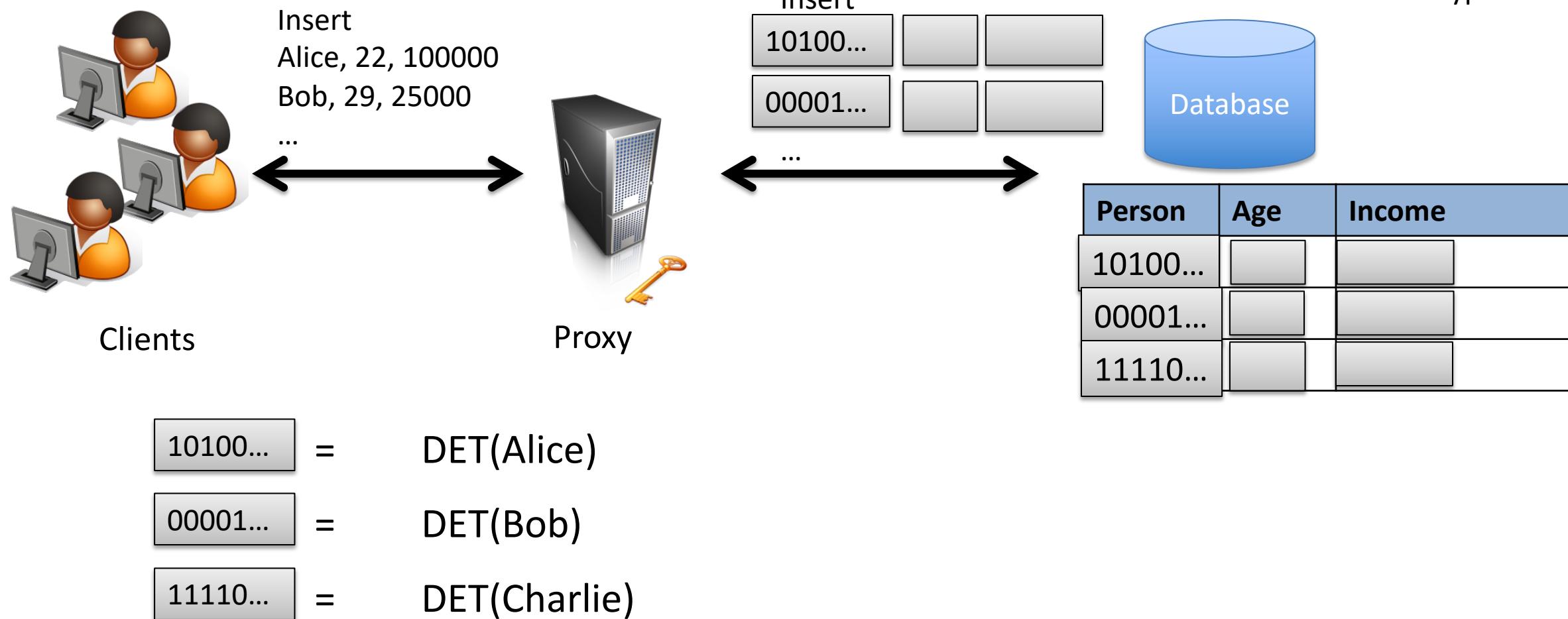
[Kornaropoulos, Papamanthou, Tamassia 2019]

[Kornaropoulos, Papamanthou, Tamassia 2020]

...

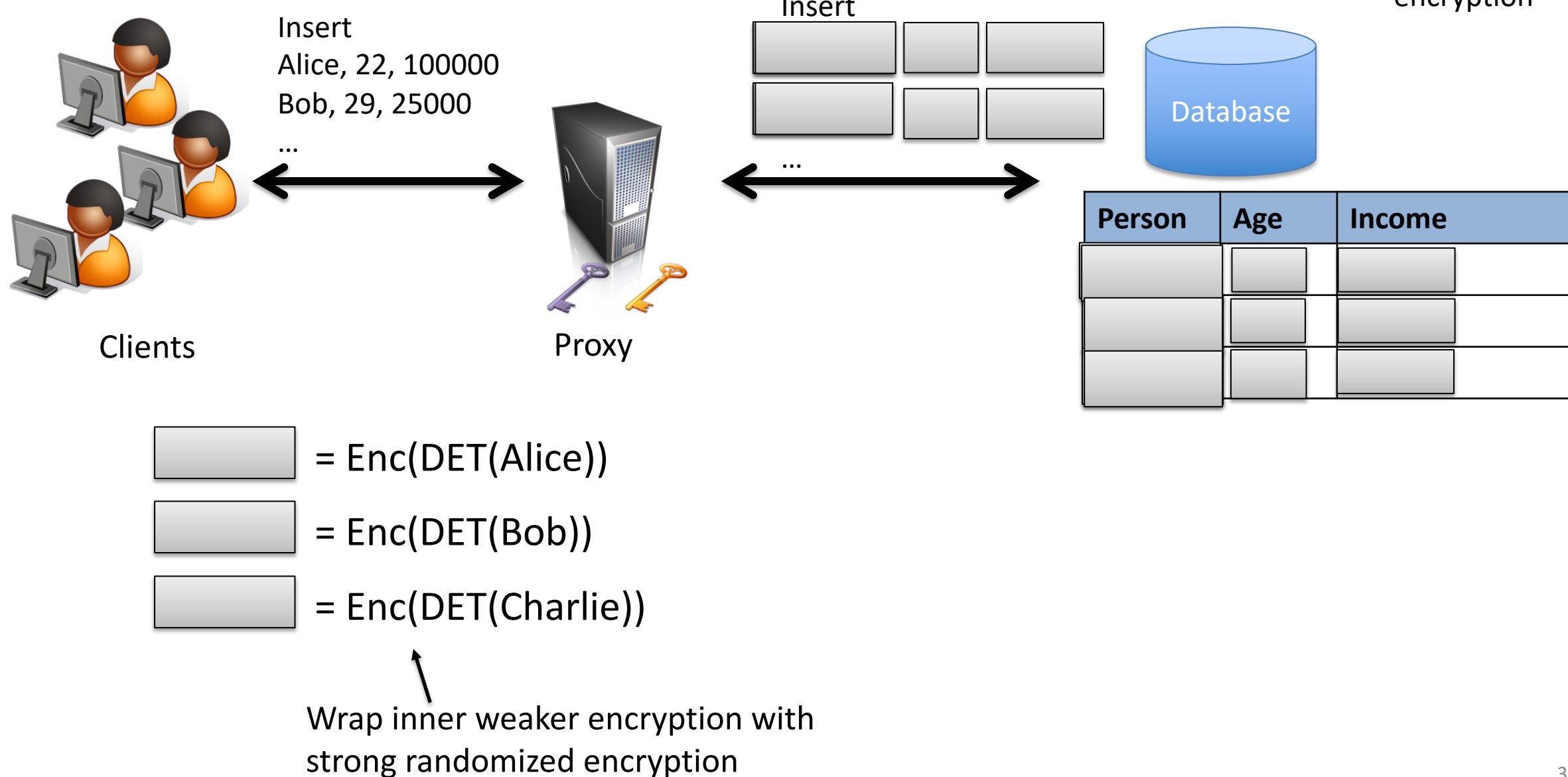
Encryption wrapping (aka onioning)

Example of
structured
encryption



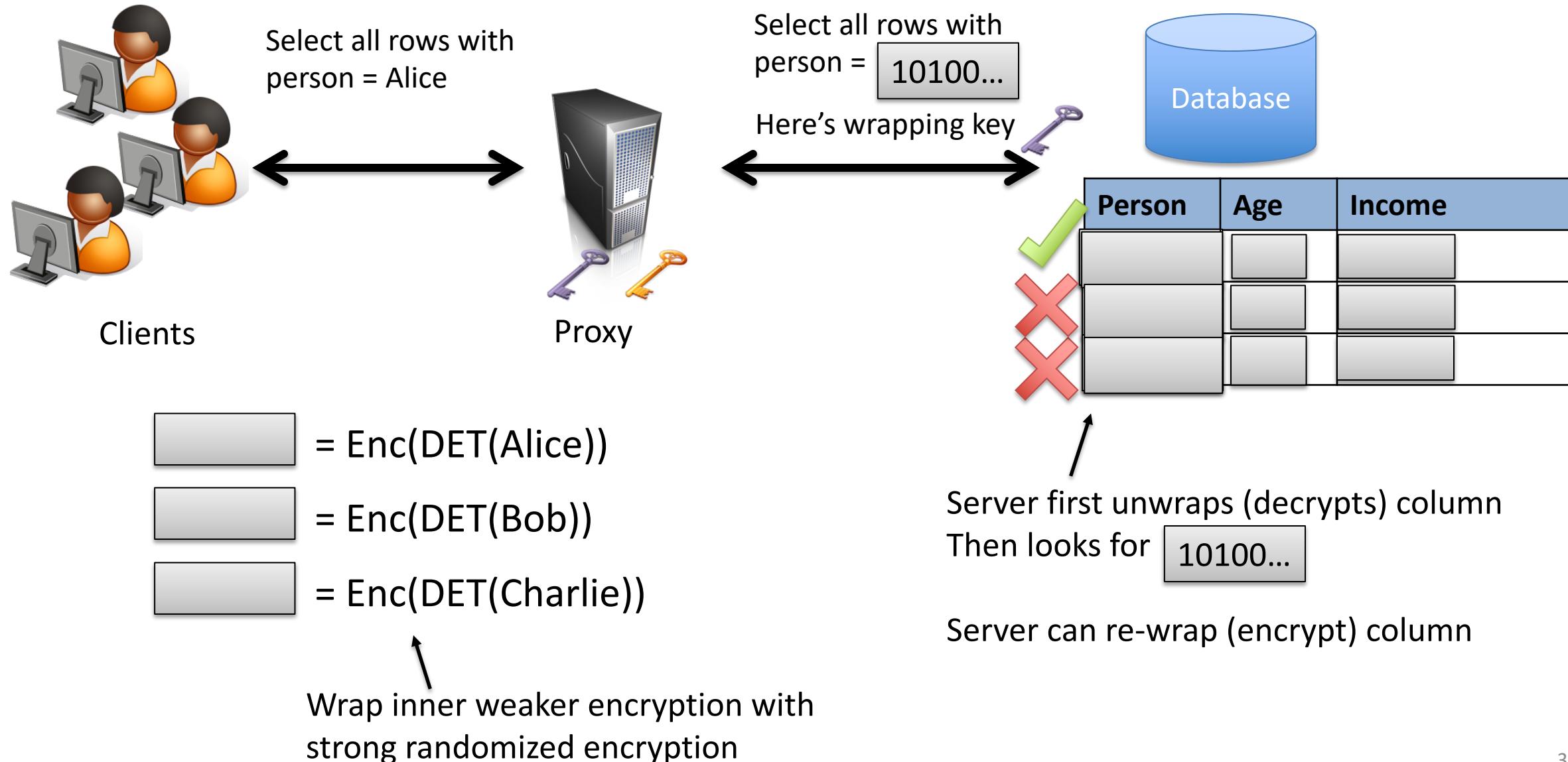
Encryption wrapping (aka onioning)

Example of structured encryption



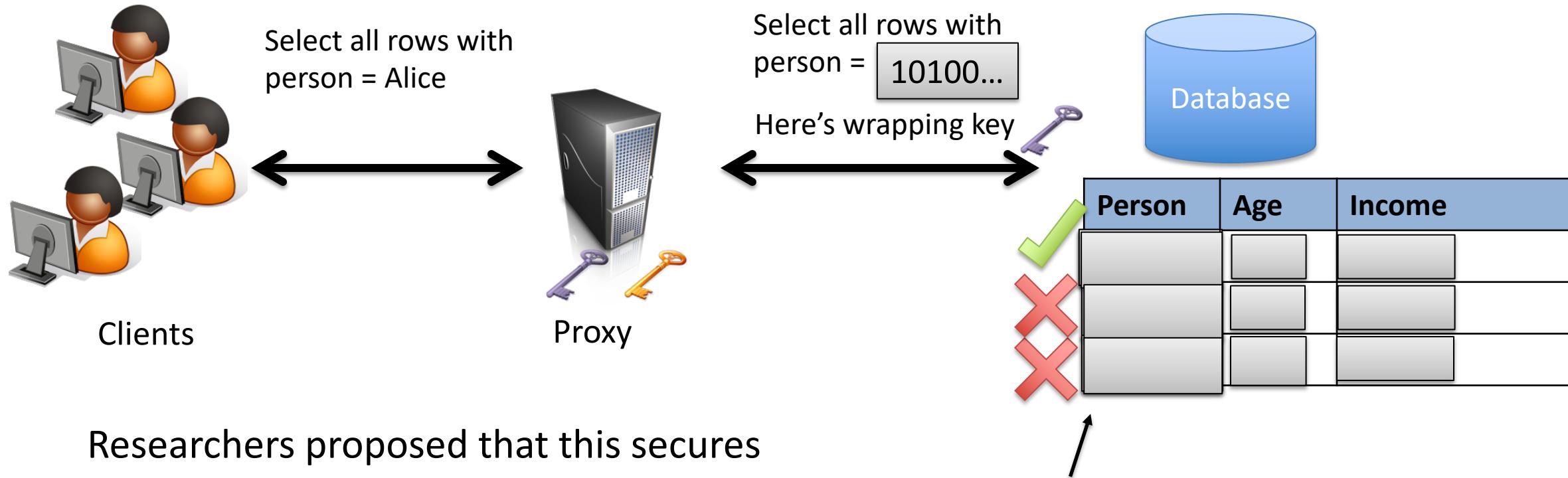
Encryption wrapping (aka onioning)

Example of structured encryption



Encryption wrapping (aka onioning)

Example of structured encryption



Researchers proposed that this secures against "*snapshot adversaries*:

- Attackers that only obtain dump of database entries learn nothing

Server first unwraps (decrypts) column
Then looks for 10100...

Server can re-wrap (encrypt) column

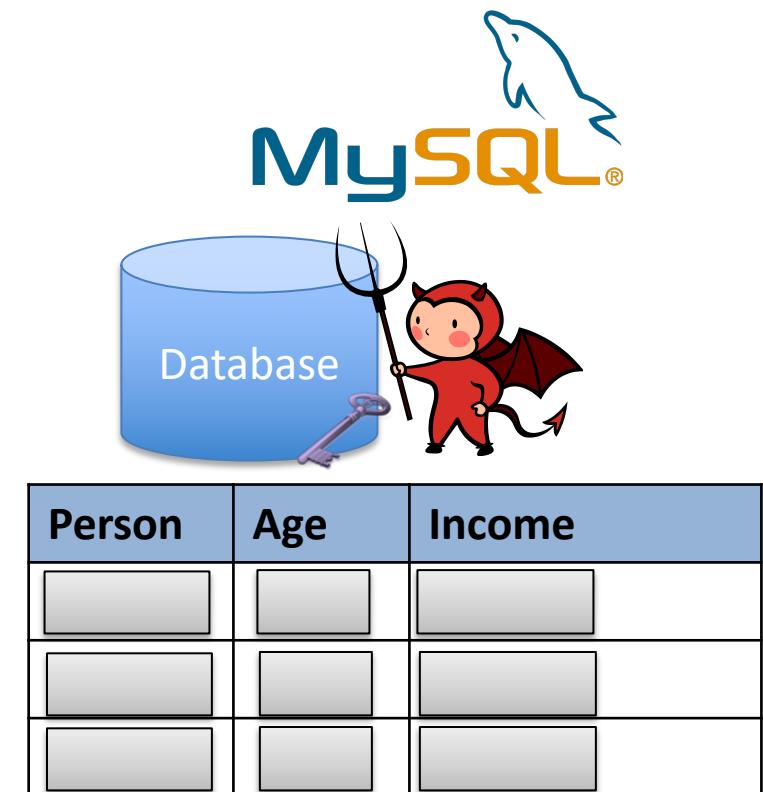
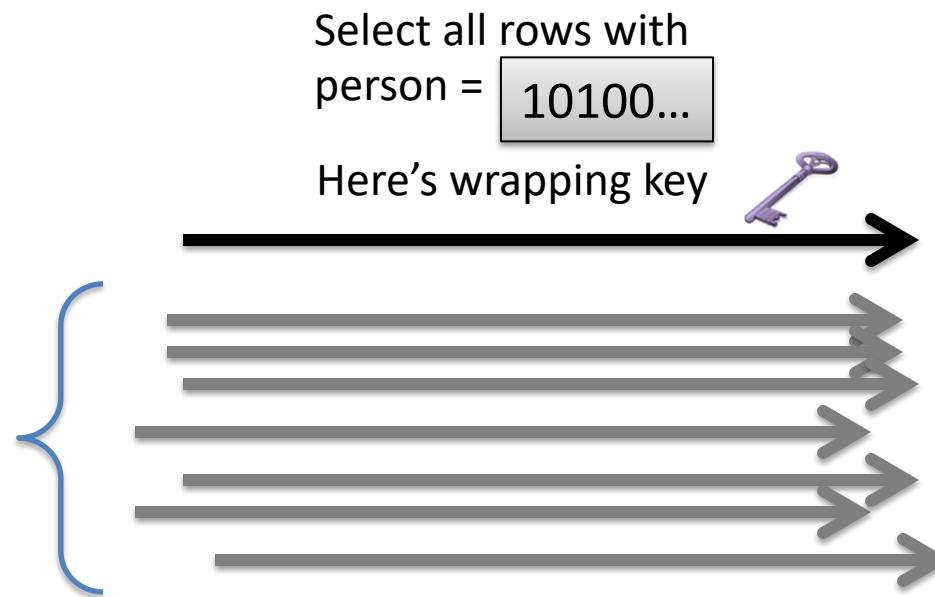
Currently: snapshot adversaries don't exist

Queries are left in memory, cached, and even written to logs in typical databases

Attackers in practice will get access to all of it

[Grubbs et al. 2017]
experiments

100,000 more operations,
not involving wrapping key



Wrapping key still in DB's memory, available to attacker!

Encrypted DBs

- Big research area:
 - Property-revealing encryption
 - Systems designs
 - Systems security
 - Leakage-abuse attacks
- Used widely in practice
- Illustrative pitfalls in terms of “provable” security
 - “CryptDB is a system that provides practical and provable confidentiality ...”