# Improved reconstruction attacks using range query leakage

Marie-Sarah Lacharité

Brice Minaud

Kenny Paterson

Information Security Group
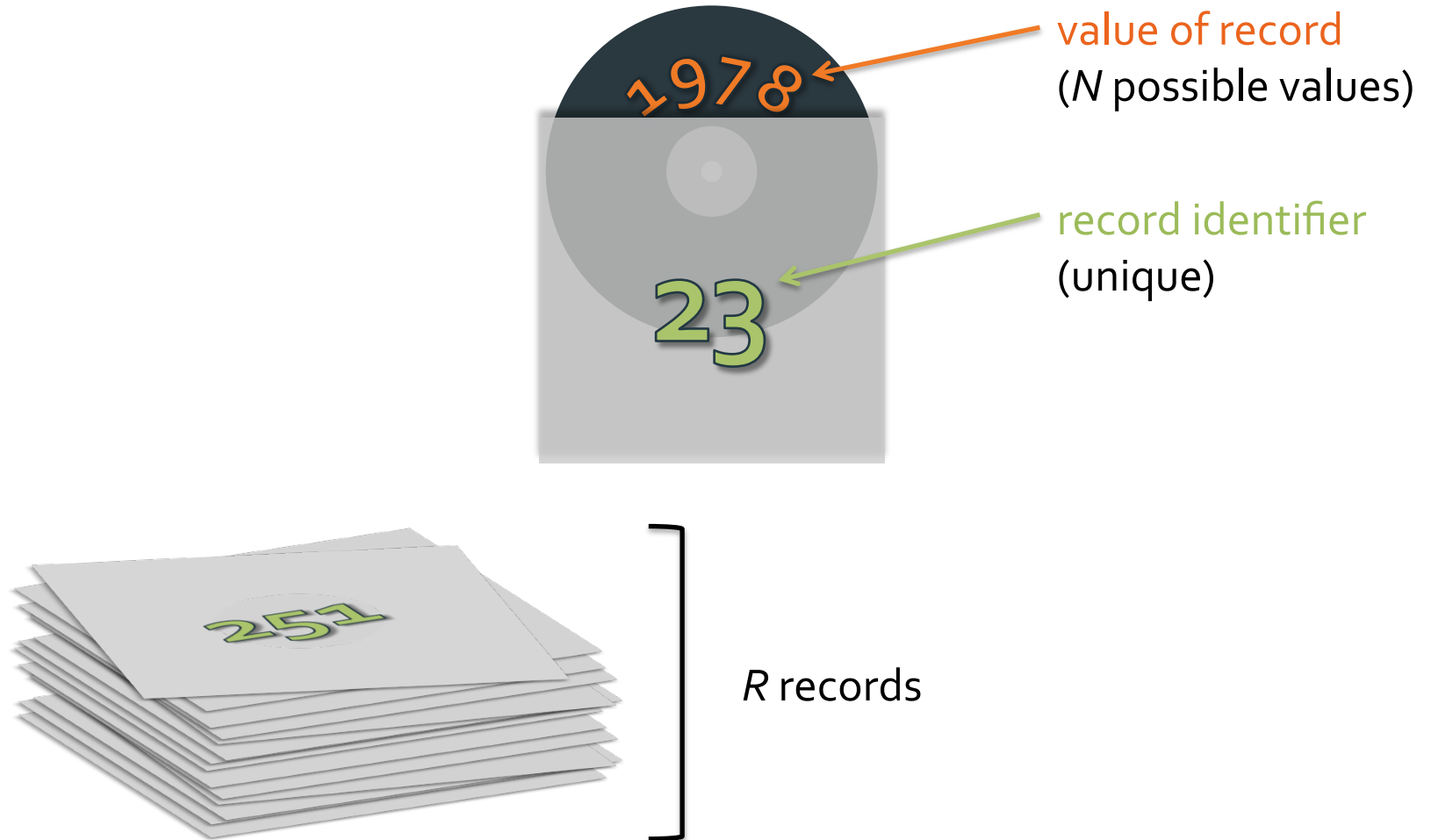
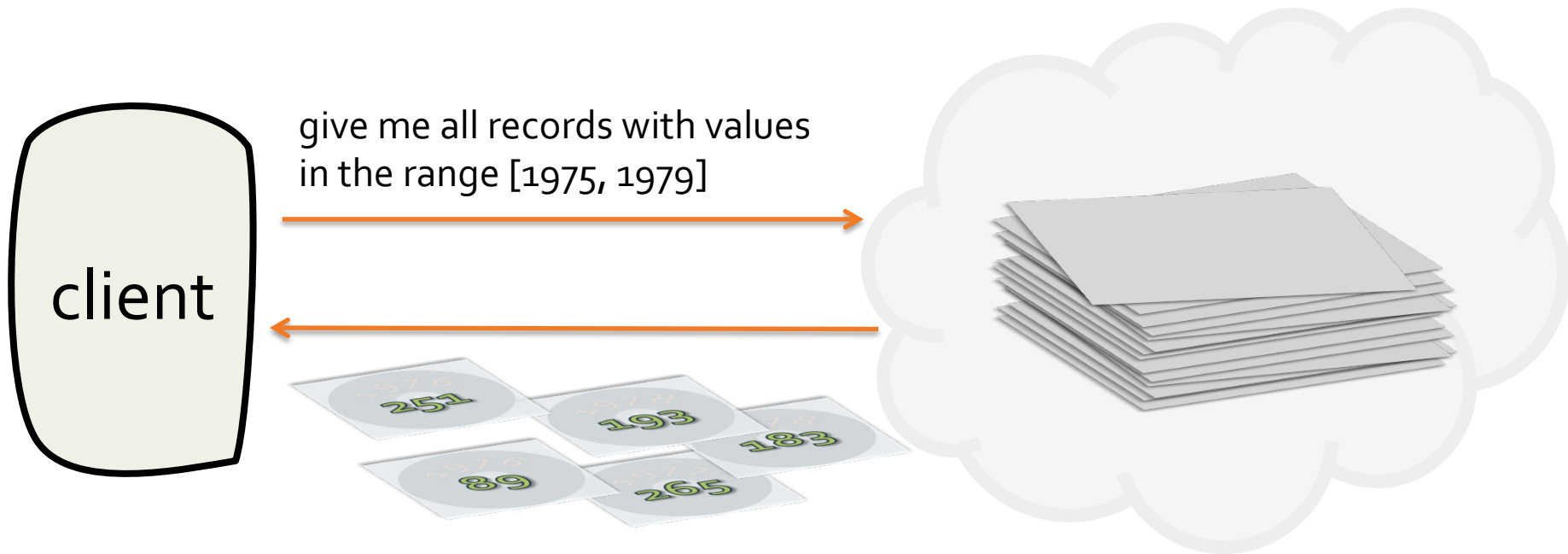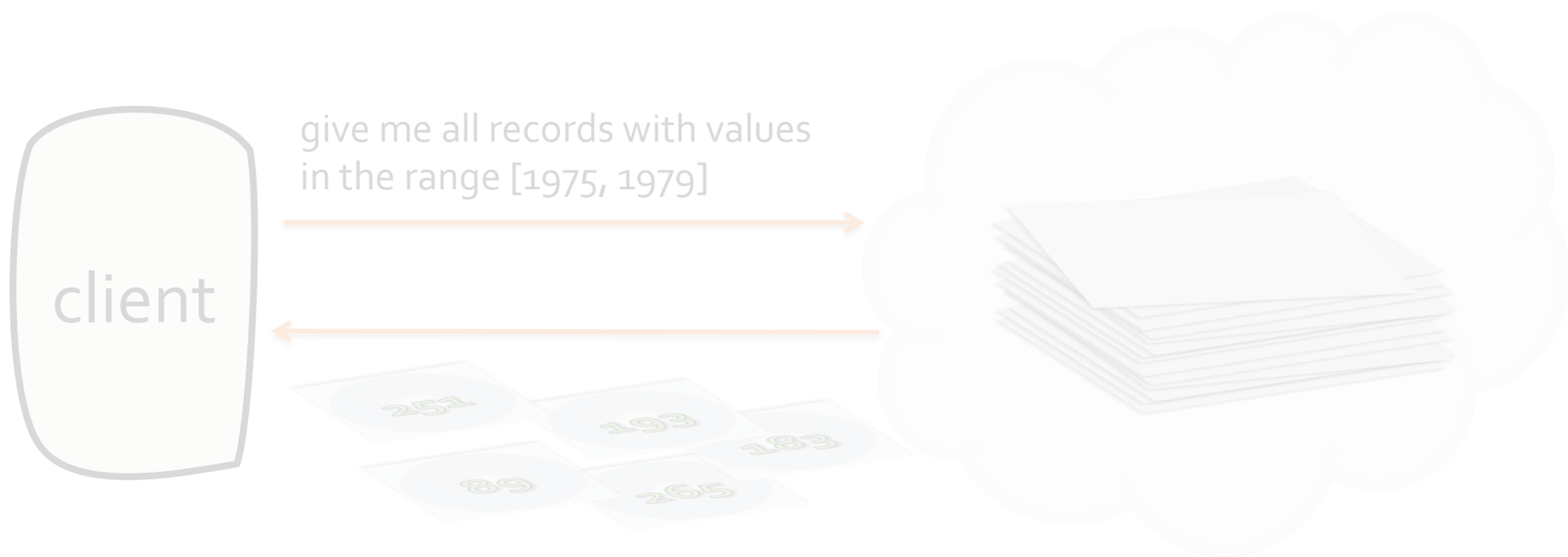ROYAL HOLLOWAY UNIVERSITY OF LONDON

# Application Setting

# Storing Records in the Cloud

1978

value of record
(*N* possible values)

23

record identifier
(unique)

251

*R* records

give me all records with values
in the range [1975, 1979]

client

251  193  183  89  265

give me all records with values
in the range [1975, 1979]

client

251

193

183

89

265

**record identifiers**

{251, 89, 193, 265, 183}

OPE, ORE schemes, POPE, [HK16], Blind seer, [Lu12], [FJKNRS15],…

give me all records with values
in the range [1975, 1979]

*rank*

client

a+1
b

**record identifiers**

{251, 89, 193, 265, 183}

FH-OPE, Lewi-Wu, Arx, Cipherbase, EncKV,...

# Assumptions

1. Data is **dense:** all values appear in at least one record.

2. Queries are **uniformly distributed**.

**Target**: full reconstruction: find the value associated with each record.

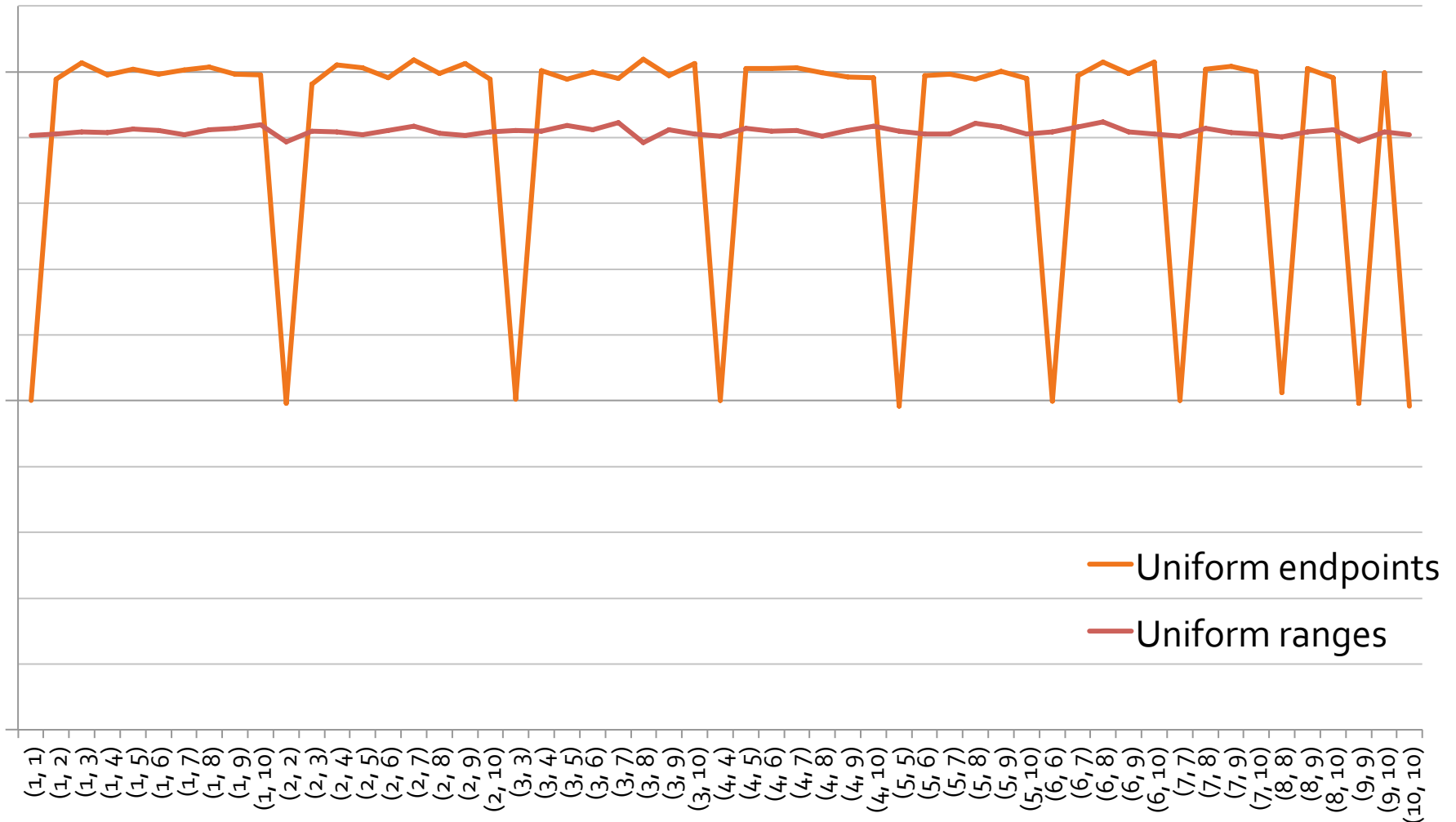**Best previous result (Kellaris et al., CCS 2016):**

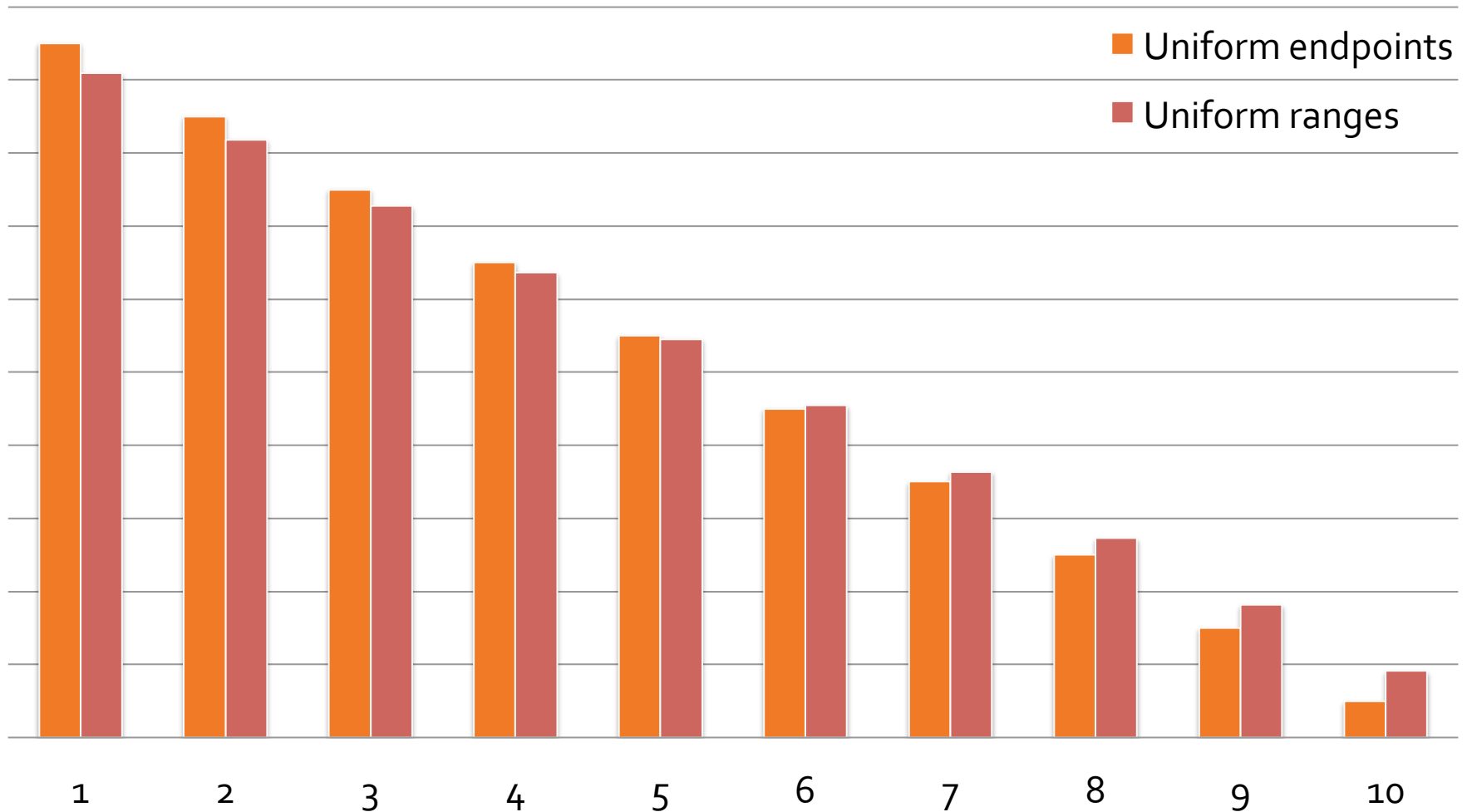Full reconstruction by analysing access pattern leakage from $O(N^2 \log N)$ queries.

- Full reconstruction with O($N \log N$) queries

  - in fact, expected $N \cdot (3 + \log N)$.

- Approximate reconstruction with relative accuracy ε from O($N \cdot (\log 1/ε)$) queries

  - in fact, expected $5/4 \cdot N \cdot (\log 1/ε) + O(N)$.

- Approximate reconstruction using an *auxiliary distribution* and rank leakage.

  - more efficient in practice, evaluation via simulation.

  - applies in the non-dense case too, giving a new attack on OPE/ORE schemes.
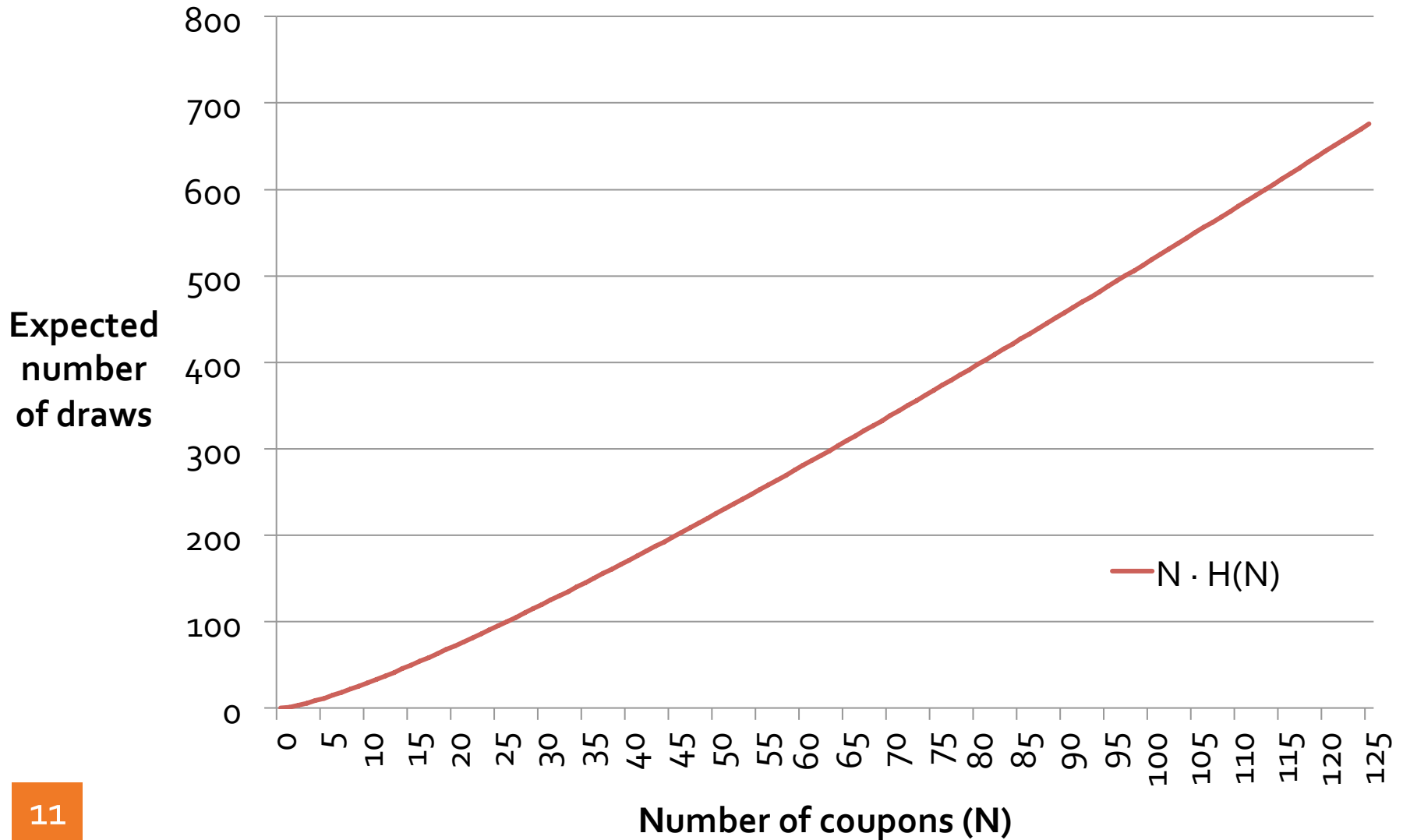
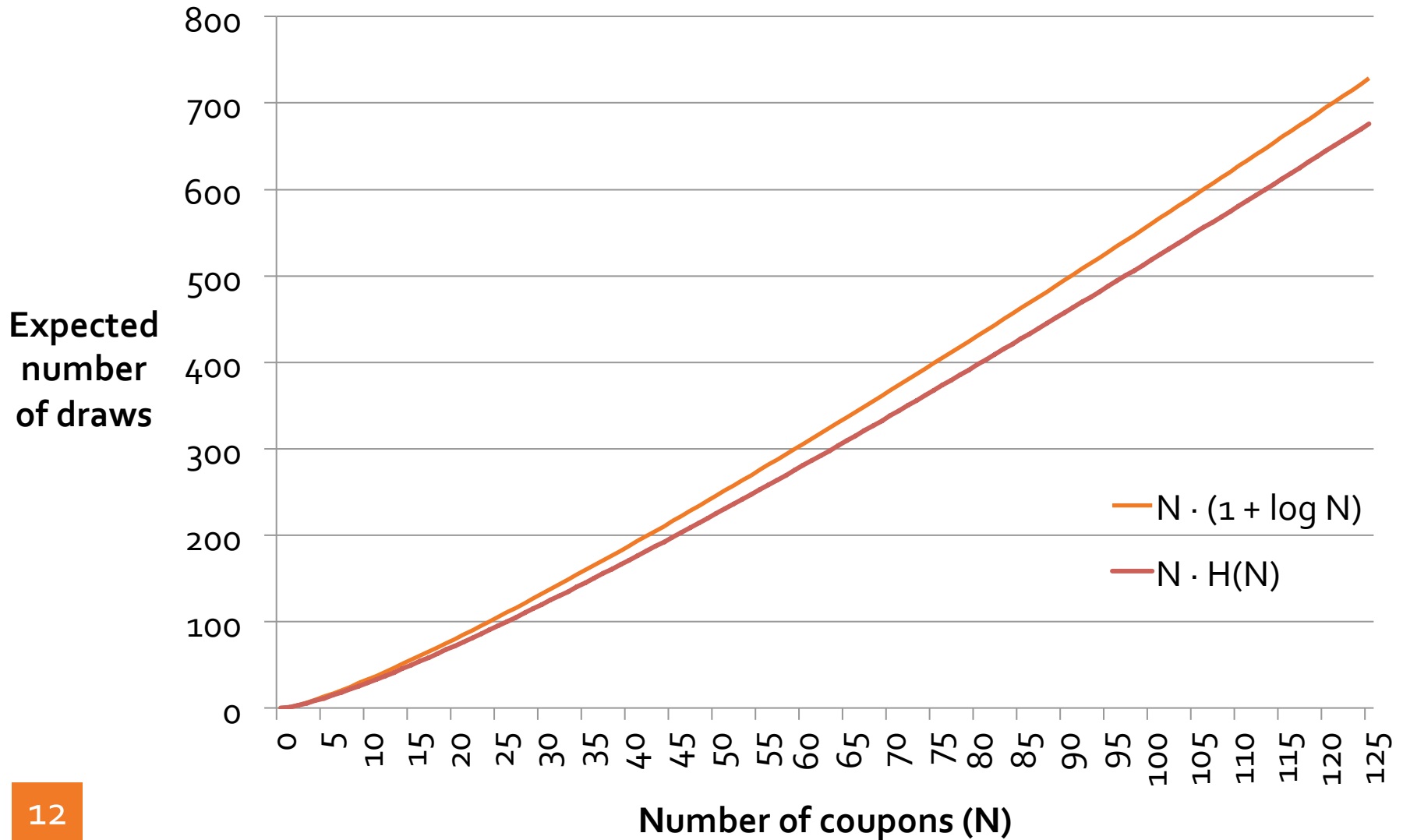# Uniform Queries: Uniform Endpoints vs. Uniform Ranges (*N*=10)

Distribution of Left Endpoints:
Uniform Endpoints vs. Uniform Ranges (*N*=10)

Legend: Uniform endpoints, Uniform ranges

X-axis: 1 2 3 4 5 6 7 8 9 10

# Coupon Collector's Problem



Expected number of draws vs. Number of coupons (N), showing the curve $N \cdot H(N)$.

# Coupon Collector's Problem
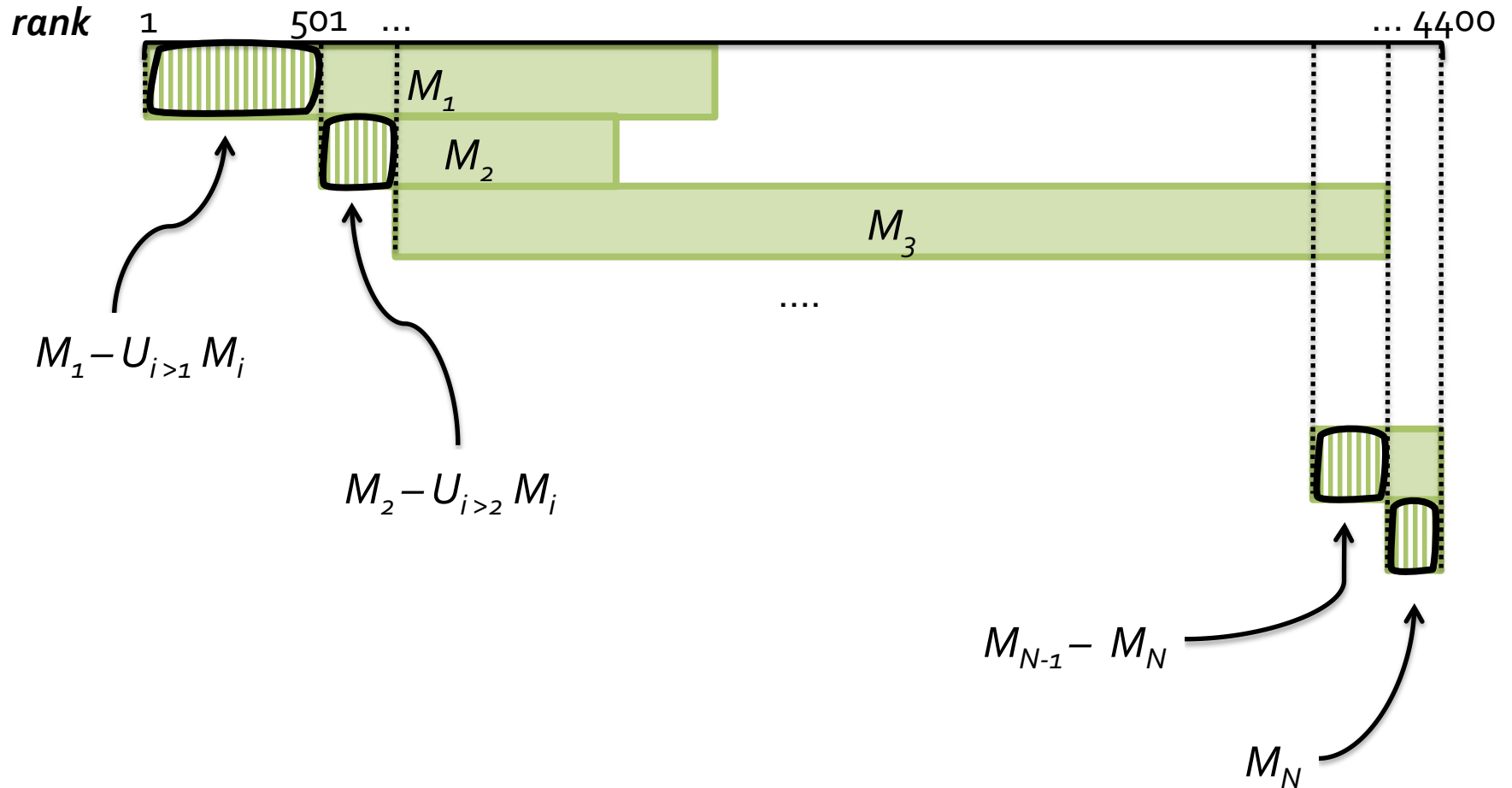
# Attack 1: Full Reconstruction

# Motivating Example (with Rank Leakage)

- Suppose **left endpoints** of query intervals are chosen **uniformly at random.**

- Wish to observe at least 1 query with each of the $N$ possible left endpoints.

- Expected number of queries needed is at most $N \cdot (1 + \log N)$.

| hidden | leaked | | |
|---|---|---|---|
| [x,y] | a = rank(x-1) | b = rank(y) | matching IDs |
| [20,25] | 1300 | 1500 | $M_{20}$ |
| [1,18] | 0 | 1200 | $M_1$ |
| [55,125] | 3100 | 4400 | $M_{55}$ |
| [2,10] | 500 | 800 | $M_2$ |
| [7,98] | 700 | 4200 | $M_7$ |

relabelled for convenience
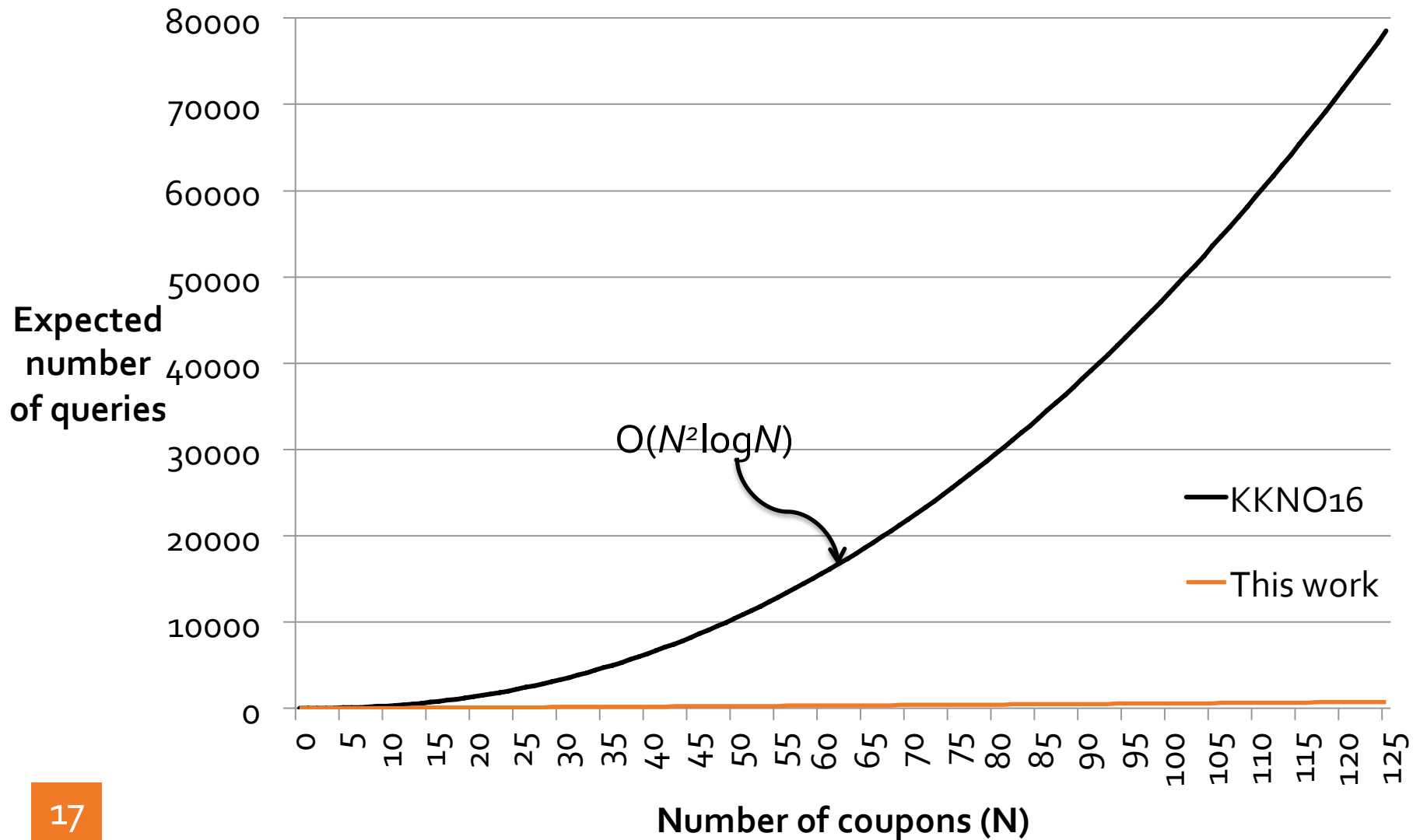
# Motivating Example (with Rank Leakage)

rank    1       501    ....                             ... 4400

$M_1$

$M_2$

$M_3$

....

$M_1 - U_{i>1}\, M_i$

$M_2 - U_{i>2}\, M_i$

$M_{N-1} - M_N$

$M_N$

# Full Reconstruction (with Rank Leakage)

- Now suppose queries have **ranges** chosen **uniformly at random.**

- We present a data-optimal algorithm (fails ⇨ full reconstruction is impossible).

- Expected number of sufficient queries is at most
$$N \cdot (2 + \log N) \text{ for } N \geq 27.$$

- Main idea: partition, then sort (easy with rank leakage, harder without).

- Expected number of necessary queries is at least
$$1/2 \cdot N \cdot \log N - O(N)$$

  **for any algorithm**.

**Expected number of queries**

$O(N^2 \log N)$

KKNO16

This work

**Number of coupons (N)**

| record ID | matched query? | | | | | | |
|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 20 | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ |
| 23 | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | ✓ |
| 29 | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ |
| 89 | ✗ | ✓ | ✓ | ✗ | ✓ | ✓ | ✗ |
| 193 | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | ✓ |
| … | | | | | | | |

- Equality of matching defines a **partition** of records.
- Records in same class of partition cannot be distinguished.
- For complete reconstruction, we need $N$ classes – one class per value.

20
23
29
89
193

| record ID | matched query? | | | | | | |
|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 20 | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ |
| 23 | ✓ [1,100] | ✓ [18,82] | ✗ | ✗ | ✓ [16,96] | ✓ [16,30] | ✓ [21,61] |
| 29 | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ |
| 89 | ✗ | ✓ | ✓ | ✗ | ✓ | ✓ | ✗ |
| 193 | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | ✓ |
| … | | | | | | | |

Can also deduce from rank leakage that, e.g., records 23 and 193 have ranks in [21,30], by intersecting rank intervals.

20

23

29

89

193

1

2

3

4

5

6

Order partition into N classes by rank

Ranks [21,30]

records 23 and 193 (and more)

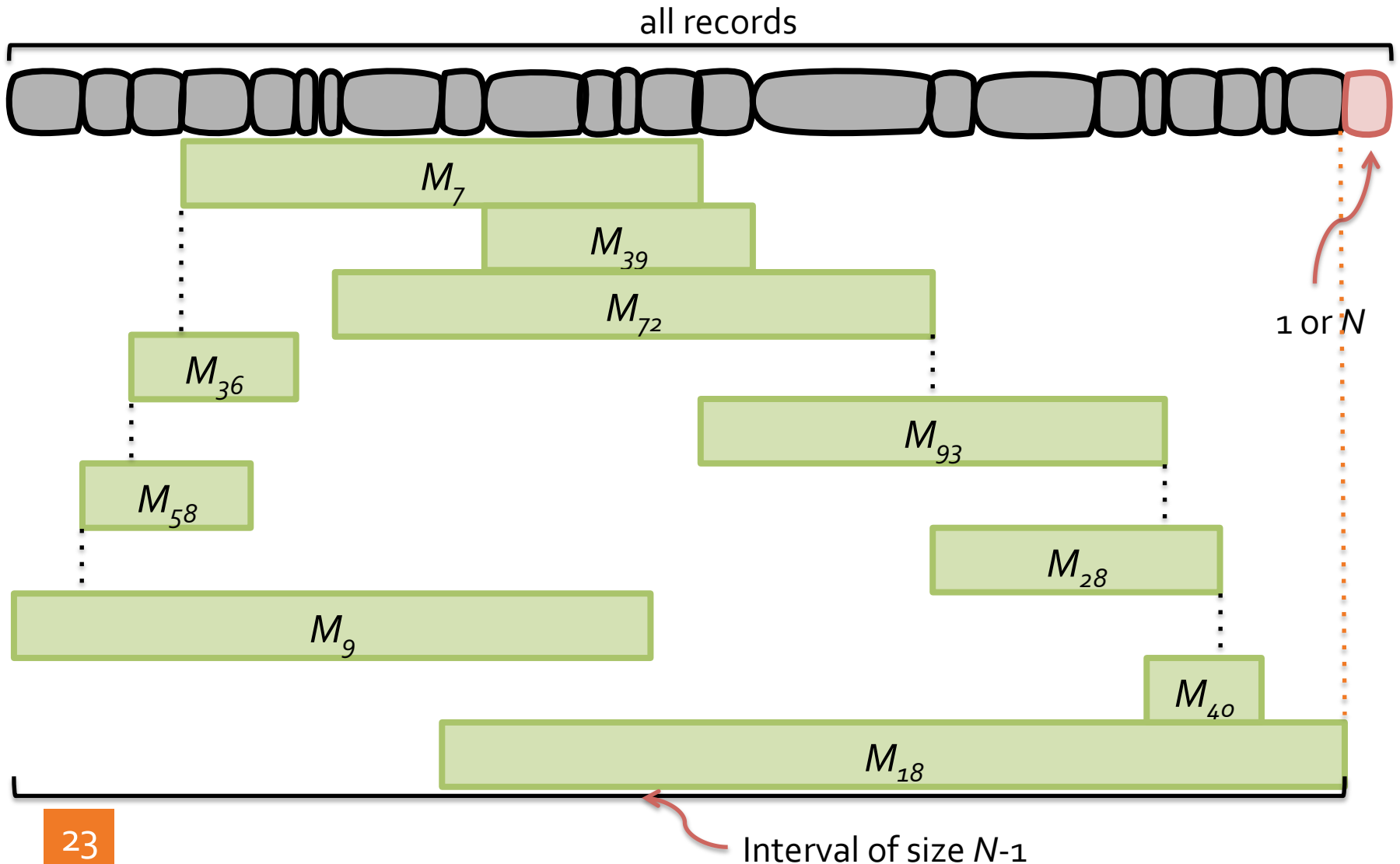# Full Reconstruction (with Rank Leakage): Proof Intuition

- Hard part is to show that O($N$ log $N$) queries suffice with a small constant.

- Proof consists of showing that **if** certain favourable queries are made, then partitioning succeeds in constructing $N$ classes.

- Roughly speaking, for our proof we hope for queries on ranges:

  1. [x,*] for all $1 \leq x \leq N/2$  (left coupons)

  2. [*,y] for all $N/2+1 \leq y \leq N$  (right coupons)

  3. [$N/2+1$,y] and [x,$N$] for some y ≥ x.

- Assuming these all arise, then a combinatorial argument establishes the success of the partitioning step.

- First two cases are essentially a pair of coupon collector problems – success with high probability with O($N$ log $N$) queries.

- Third case is a high probability event: $1 - e^{-Q/(2N+2)}$ for $Q$ queries.
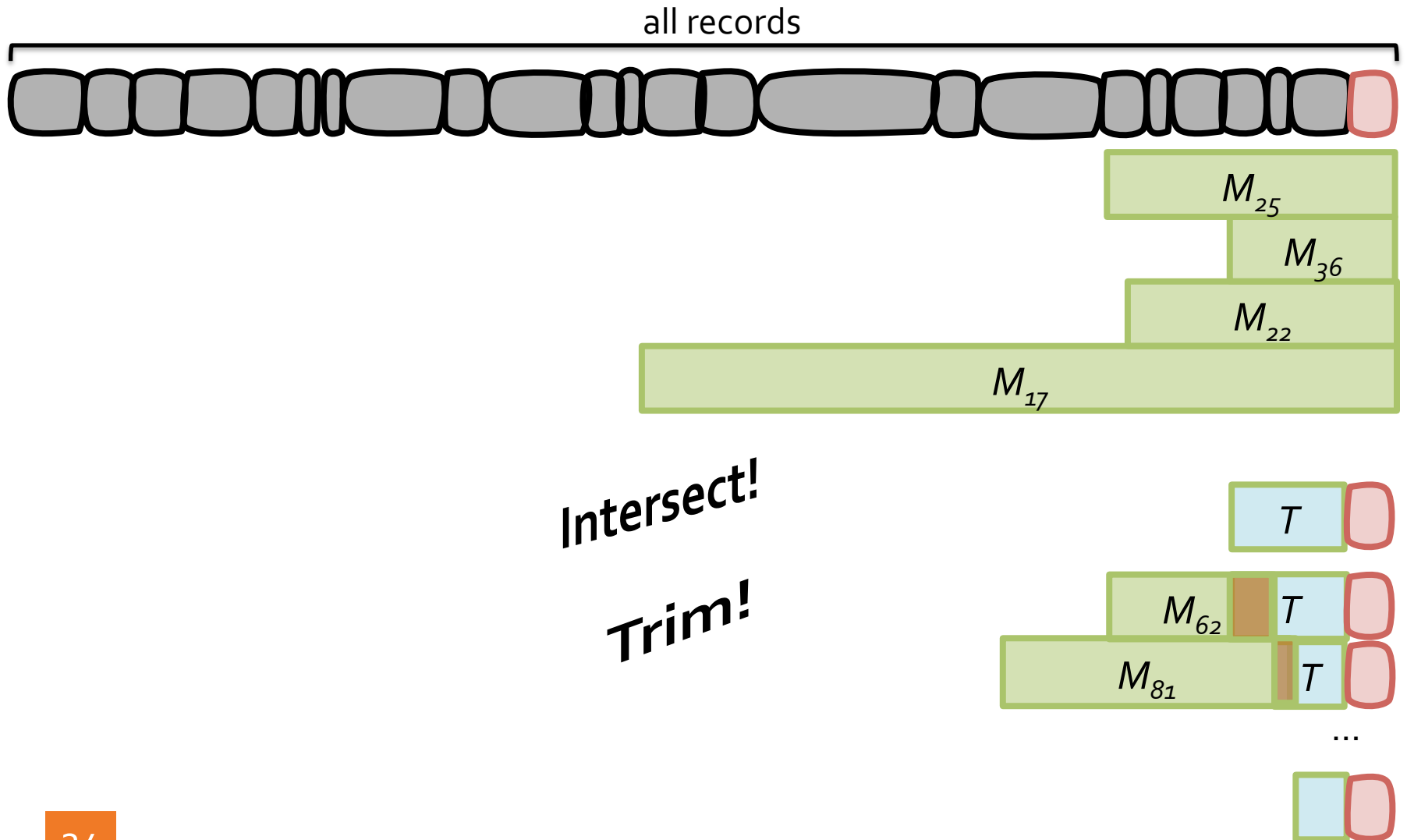
# Full Reconstruction (**without** Rank Leakage)

- Can only recover values up to **reflection**.

- Data-optimal algorithm (fails ⇨ full reconstruction is impossible).

- Expected number of sufficient queries is at most
$$N \cdot (3 + \log N) \text{ for } N \geq 26$$

- Partition (as before), then sort*.

- Expected number of necessary queries is at least
$$1/2 \cdot N \cdot \log N - O(N)$$

  **- for any algorithm**.

*Not quite.

all records

$M_7$

$M_{39}$

$M_{72}$

$M_{36}$

$M_{93}$

$M_{58}$

$M_{28}$

$M_9$

$M_{40}$

$M_{18}$

1 or $N$

Interval of size $N$-1

all records



$M_{25}$

$M_{36}$

$M_{22}$

$M_{17}$

Intersect!

Trim!

$T$

$M_{62}$ $T$

$M_{81}$ $T$

...

24

all records

all records

$M_3$

$M_{39}$

$M_{27}$

$M_{13}$

Intersect!

Trim!

$T$

$M_{52}T$

$M_{99}$    $T$

all records

...

## Full Reconstruction (without Rank Leakage): Proof Intuition

- Hard part is again to show that O($N$ log $N$) queries suffice, with a small constant.

- Proof again consists of showing that **if** certain favourable queries are made, then partitioning succeeds in constructing $N$ classes.

- Coupon collecting bounds then establish that O($N$ log $N$) queries are enough.

# Attack 2: Approximate Reconstruction

# Approximate Reconstruction Attack (without Rank Leakage)

- Recover values up to **reflection** and with relative error ε.

- Expected number of sufficient queries is
$$5/4 \cdot \ N \cdot (\log 1/\varepsilon) + O(N).$$

- Expected number of necessary queries is at least
$$1/2 \cdot N \cdot (\log 1/\varepsilon) - O(N)$$

  **for any algorithm.**

- Not data-optimal without rank leakage (but *is* with it)

Collecting *n* of 125 coupons

**Collecting fraction *(1-ε)* of 125 coupons**

Expected number of draws

Coupon number (*n*)

ε = 0.04
ε = 0.08
ε = 0.12
ε = 0.16
ε = 0.2

32

all records

$M_7$

$M_{39}$

$M_{72}$

$M_{36}$

$M_{93}$

$M_{58}$

$M_{28}$

$M_9$

$M_{40}$

$M_{18}$

1.  Pick any record *r.*

2. Intersect all queries matching *r* to get *M.*

2. Intersect all queries matching *r* to get *M*.

3. Find $q_L$ and $q_R$ : $q_L \cap q_R = M$ and $|q_L \cup q_R|$ maximised.

4. Find $q'_L : q_L \cap q'_L \neq \varnothing$, $q'_L \cap q_R \subseteq M$, $|q_L \cup q'_L|$ maximised.

5. Find $q'_R : q_R \cap q'_R \neq \varnothing$, $q'_R \cap q_L \subseteq M$, $|q_R \cup q'_R|$ maximised.

6. Start over if not every record is in $q_L \cup q'_L \cup q_R \cup q'_R$.

7. Split into $half_L = q_L \cup q'_L$, $half_R = q_L \cup q'_L$, and $M$.

7.  Split into $half_L = q_L \cup q'_L$, $half_R = q_L \cup q'_L$, and $M$.

8. Form **left & right coupons** with queries containing $M$.

$half_L \setminus M$          $M$          $half_R \setminus M$

1      $N$

$n_R$ right coupons

$n_L$ left coupons

9. Use left & right coupons to **sort $half_L \setminus M$ & $half_R \setminus M$.**

9. Use left & right coupons to **sort $half_L$ \ $M$ & $half_R$ \ $M$.**

$half_L$ \ $M$        $M$        $half_R$ \ $M$

1           $N$

$$n_L + 1 + n_R = (1\text{-}\varepsilon) \cdot N$$

$$\Downarrow$$

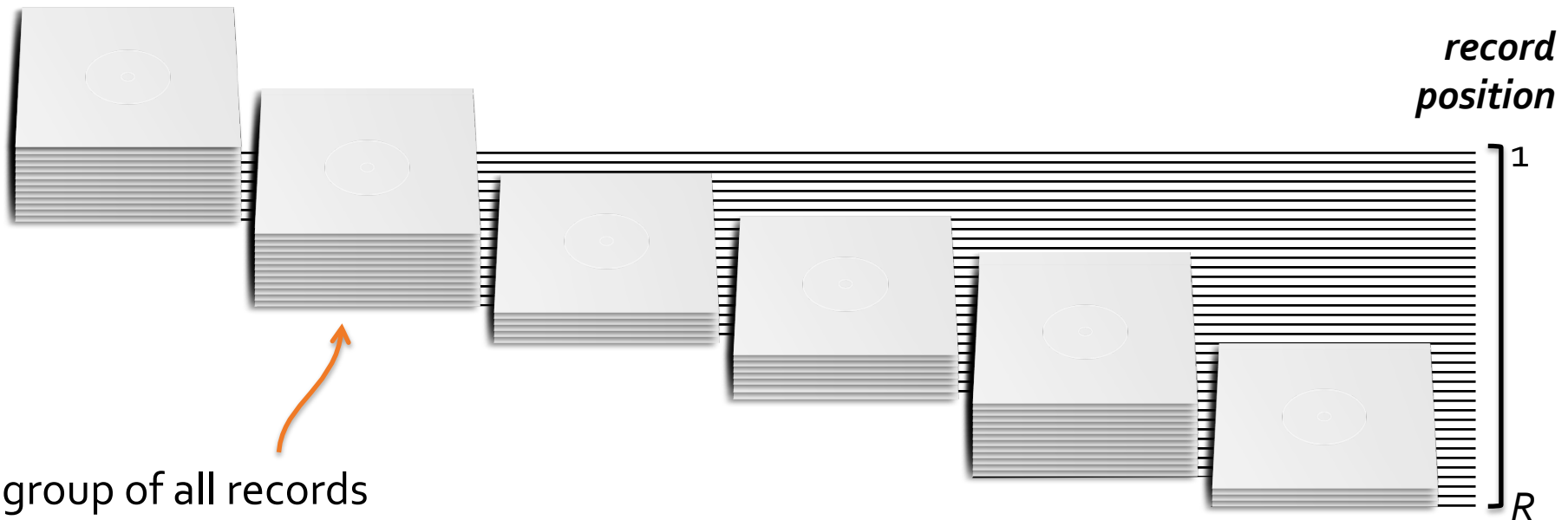reconstruction with precision $\varepsilon \cdot N$

# Attack 3: Reconstruction with Auxiliary Data

# Reconstruction with Auxiliary Data and Rank Leakage

- As before, queries have **ranges** chosen **uniformly at random**.

- Assume access pattern and rank are leaked.

- We now also assume that an **approximation to the distribution on values** is known.

  - "Auxiliary data".

  - From aggregate data, or from another reference source.

- We show experimentally that, under these assumptions, far fewer queries are needed.

- Now no requirement on density, so interesting for OPE and ORE schemes too (OPE/ORE schemes are trivial to break in dense case).

1.  **Partition** records as in full reconstruction attack.



*record position*

1

R

group of all records appearing in exact same subset of queries

2. Assign a **position interval** to each partition.



*record position*

1

*a+1*

*b*

R

intersect leaked
rank intervals to get
**position interval**

3. Assign a **value** to each group's position interval

*record position*

*Inverse CDF of auxiliary distribution*

*value*

1

0

$rank^{-1}(a) + 1$

$a+1$

$x$

expected value restricted to *[x,y]*

**point guess *v***

$b$

$y$

$rank^{-1}(b)$

*R*

1

# Auxiliary Data Attack: Experimental Evaluation

- Ages, *N* = 125 (0 to 124).

- Health records from US hospitals (NIS HCUP 2009).

- Target data: individual hospitals' records.

- Auxiliary data: aggregate of 200 hospitals' records.

- Measure of success: proportion of records with value guessed within *ε.*

95% of records within 1.25 years!

All records off by > 20 years.

Fraction of records

Relative error $\varepsilon$

0.020   discrete Kolmogorov-Smirnov statistic   0.556

52

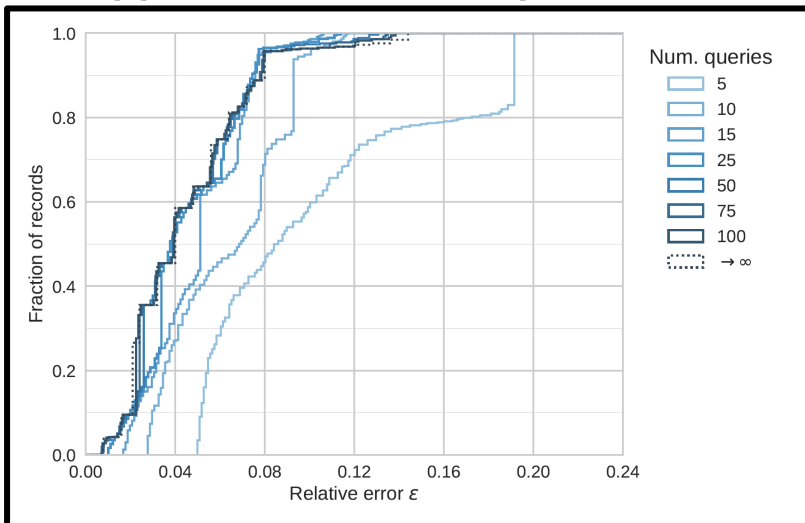# Auxiliary Data Attack: Removing Assumptions

- Estimating total number of records is **fast** if not known *a priori*

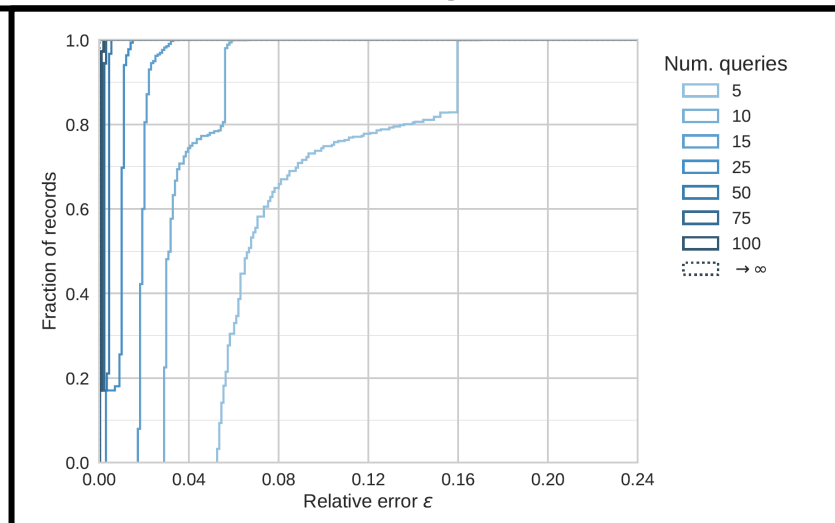- Learning set of record identifiers **can be slow** if not known *a priori:*

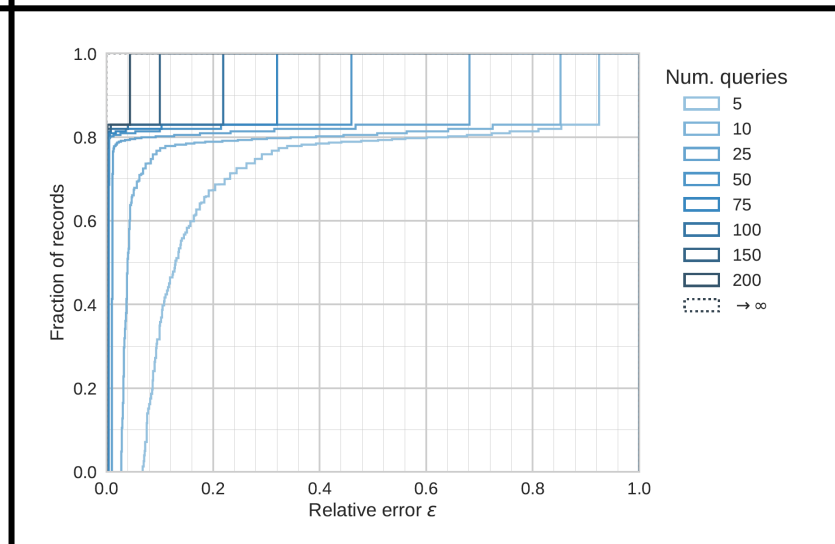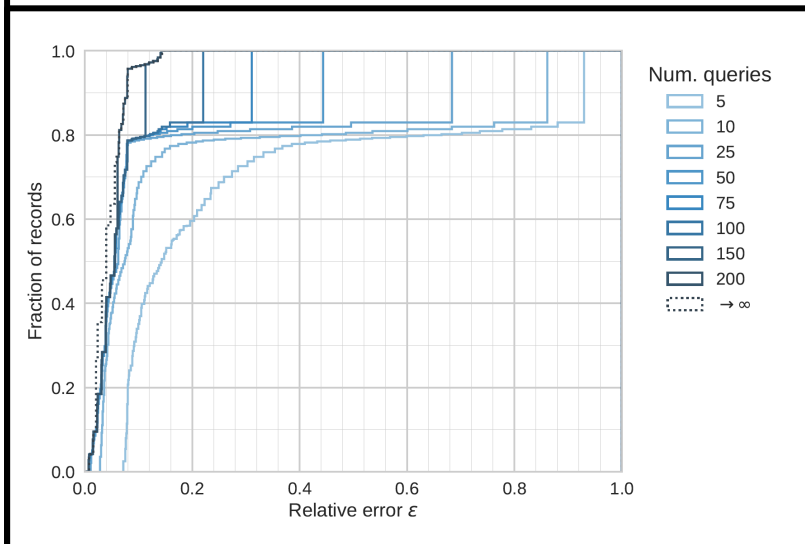# Auxiliary Data Attack: Removing Assumptions

# Summary and Conclusions

# Summary of Our Attacks

| Attack | Req'd leakage | Other req'ts | Suff. # queries |
|---|---|---|---|
| **Full** | AP + rank | Density | $N \cdot (\log N + 2)$ |
|  | AP | Density | $N \cdot (\log N + 3)$ |
| **ε-approximate** | AP | Density | $5/4\ N \cdot (\log 1/\varepsilon) + O(N)$ |
| **Auxiliary** | AP + rank | Auxiliary dist. | ??? |

# Conclusions

- Many clever schemes have been designed, enabling range queries on encrypted data:

    - OPE, ORE schemes.

    - POPE, [HK16],...

    - Blind seer, [Lu12], [FJKNRS15],...

    - FH-OPE, Lewi-Wu, Arx, Cipherbase, EncKV,...

- These schemes are surprisingly vulnerable to attack in realistic setting (density + uniform queries + access pattern leakage): **O(*N*log*N*) queries suffice!**

- Even more severe attacks are possible when auxiliary distribution + rank leakage is available.

- Read more at eprint 2017/701.