

Collecting and analyzing data over multiple services under local differential privacy

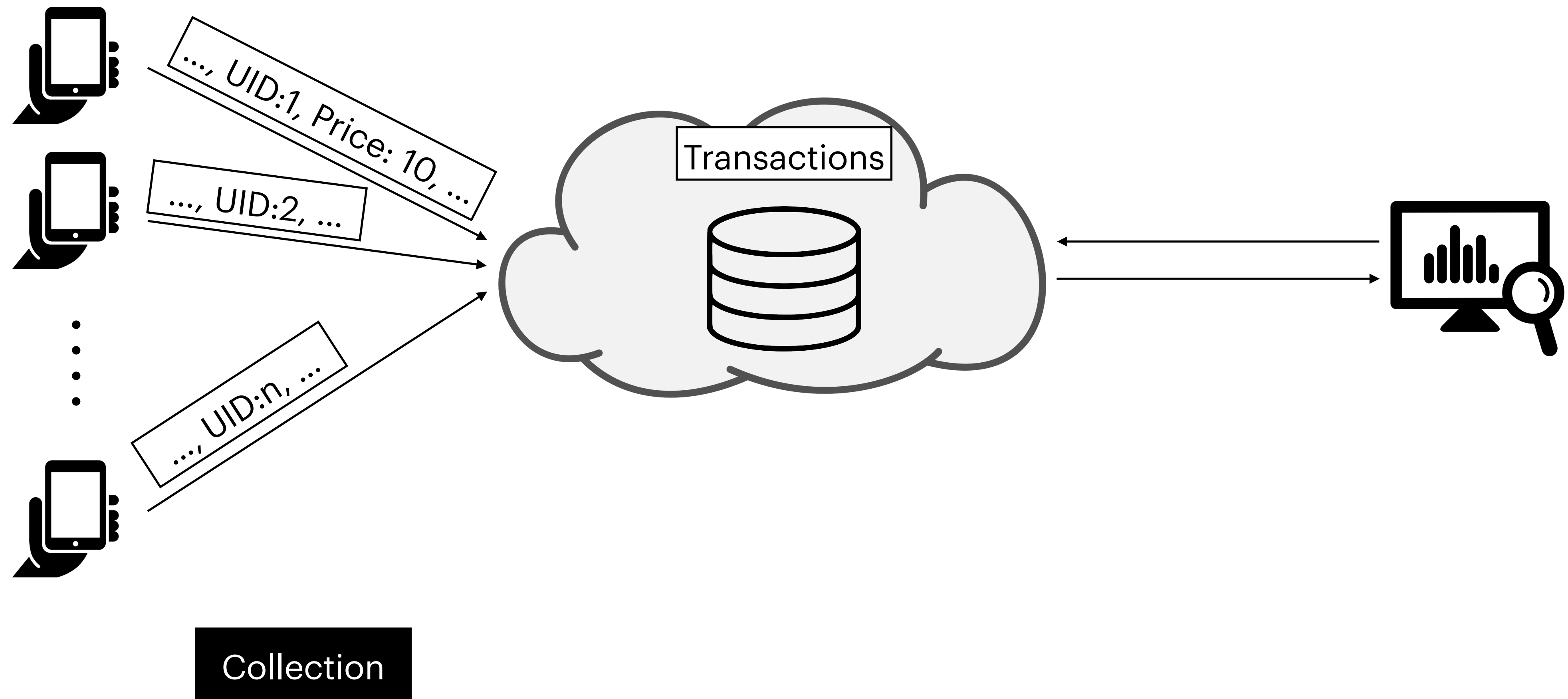
Min Xu, Bolin Ding, Tianhao Wang, Jingren Zhou



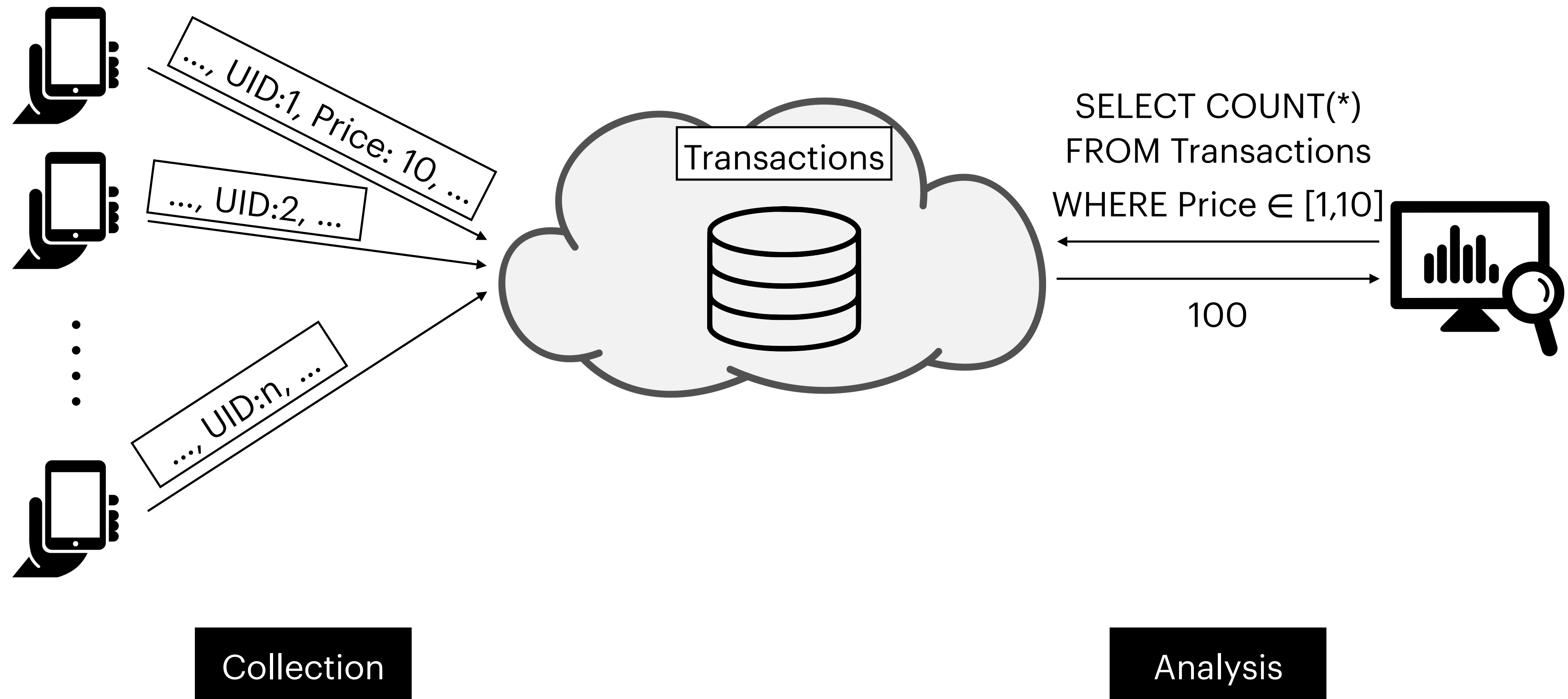
Data Collection and Analysis in the Cloud



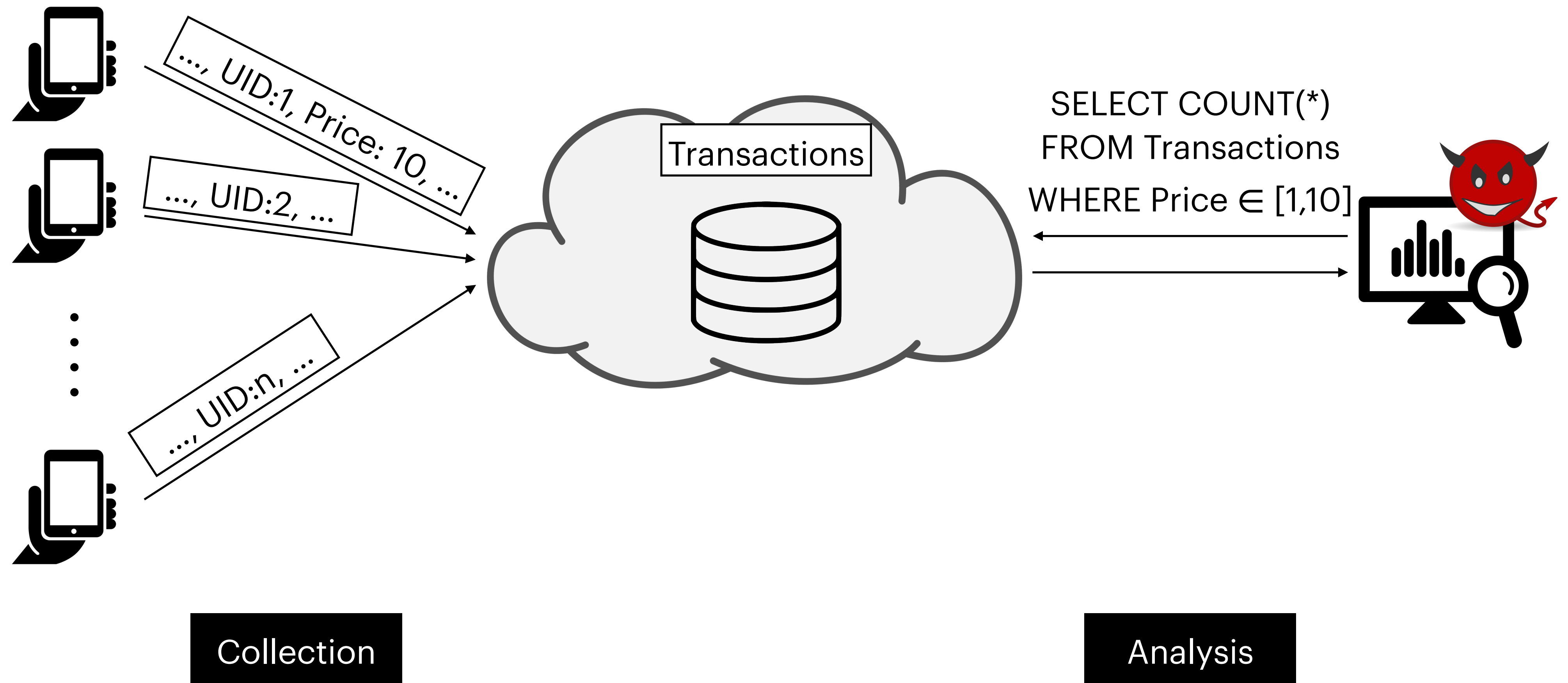
Data Collection and Analysis in the Cloud



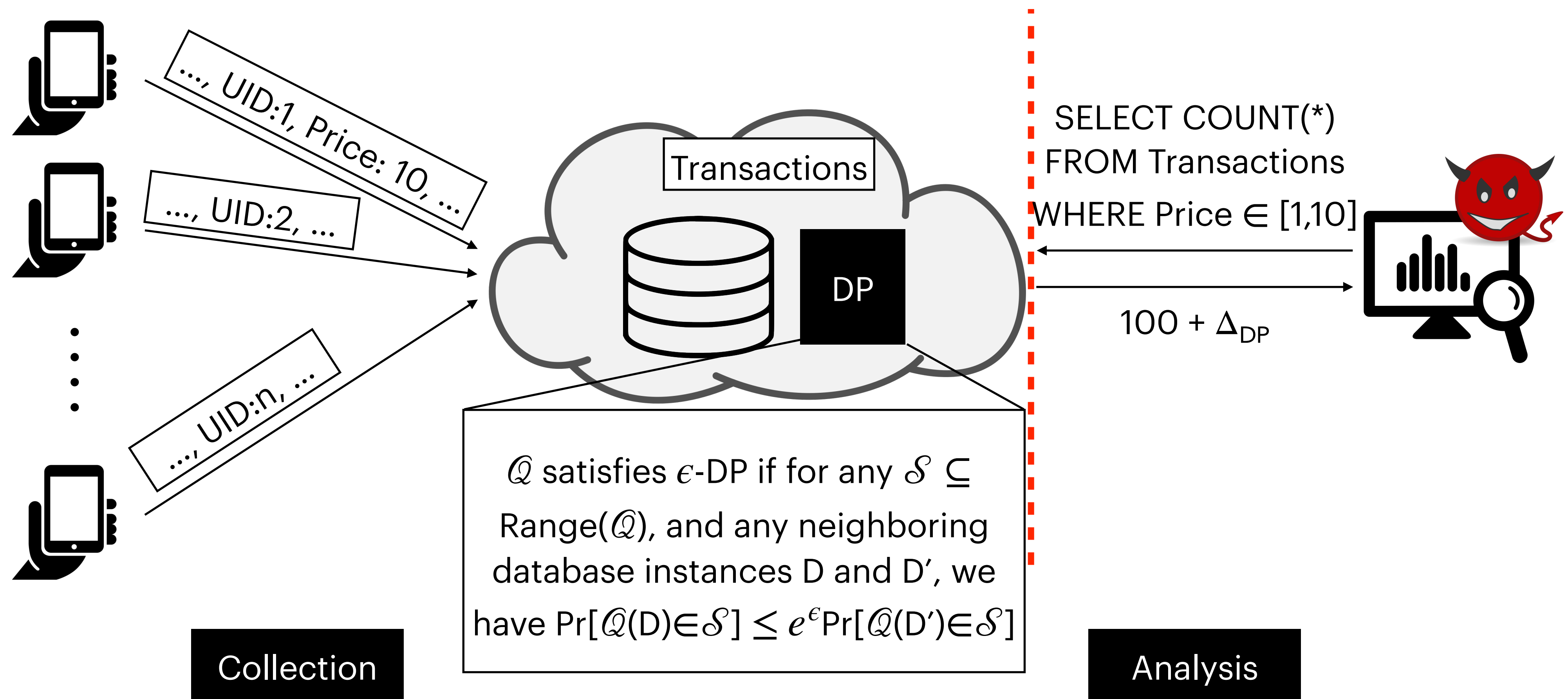
Data Collection and Analysis in the Cloud



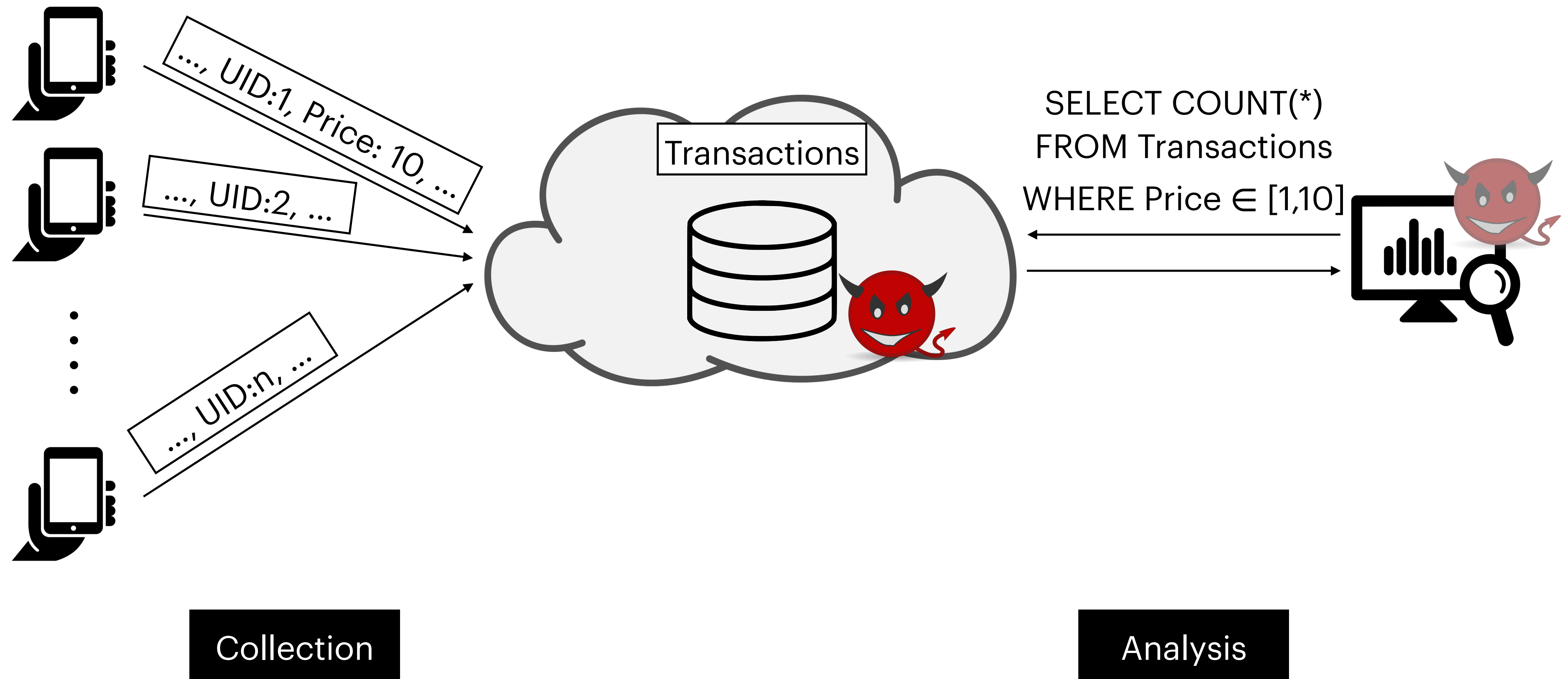
Data Collection and Analysis in the Cloud



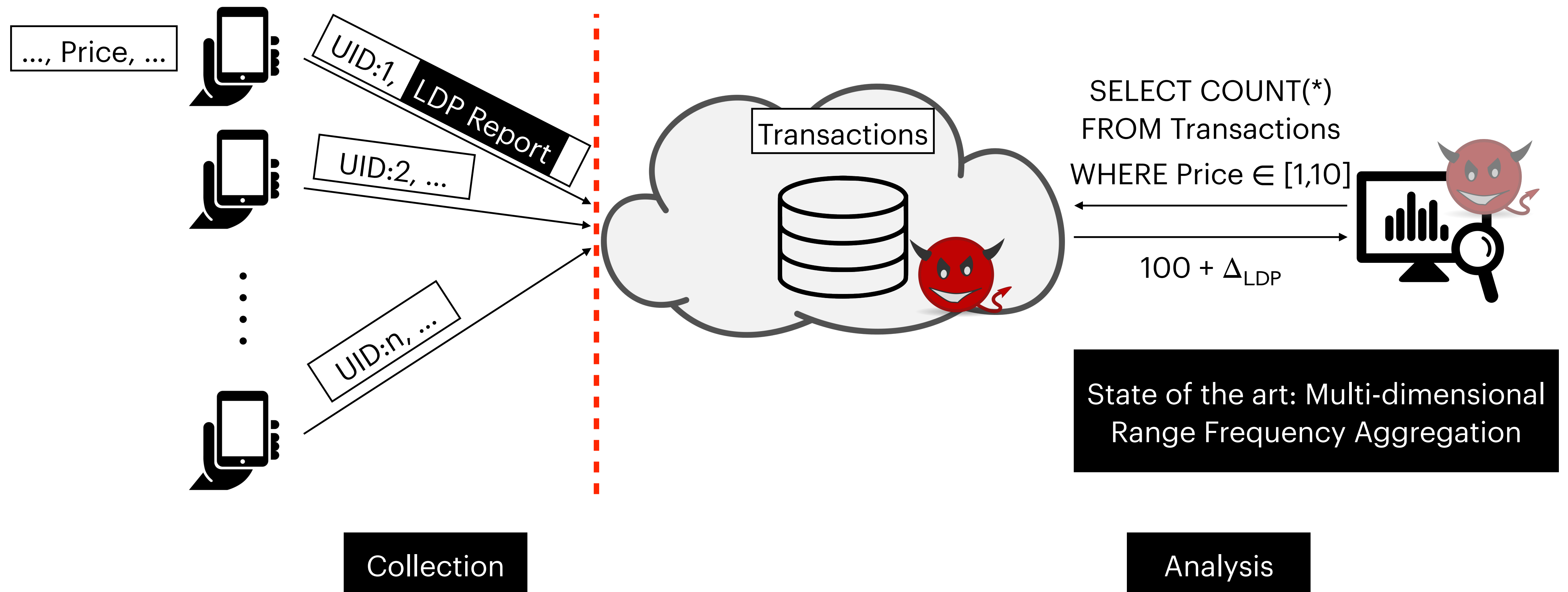
Data Collection and Analysis in the Cloud



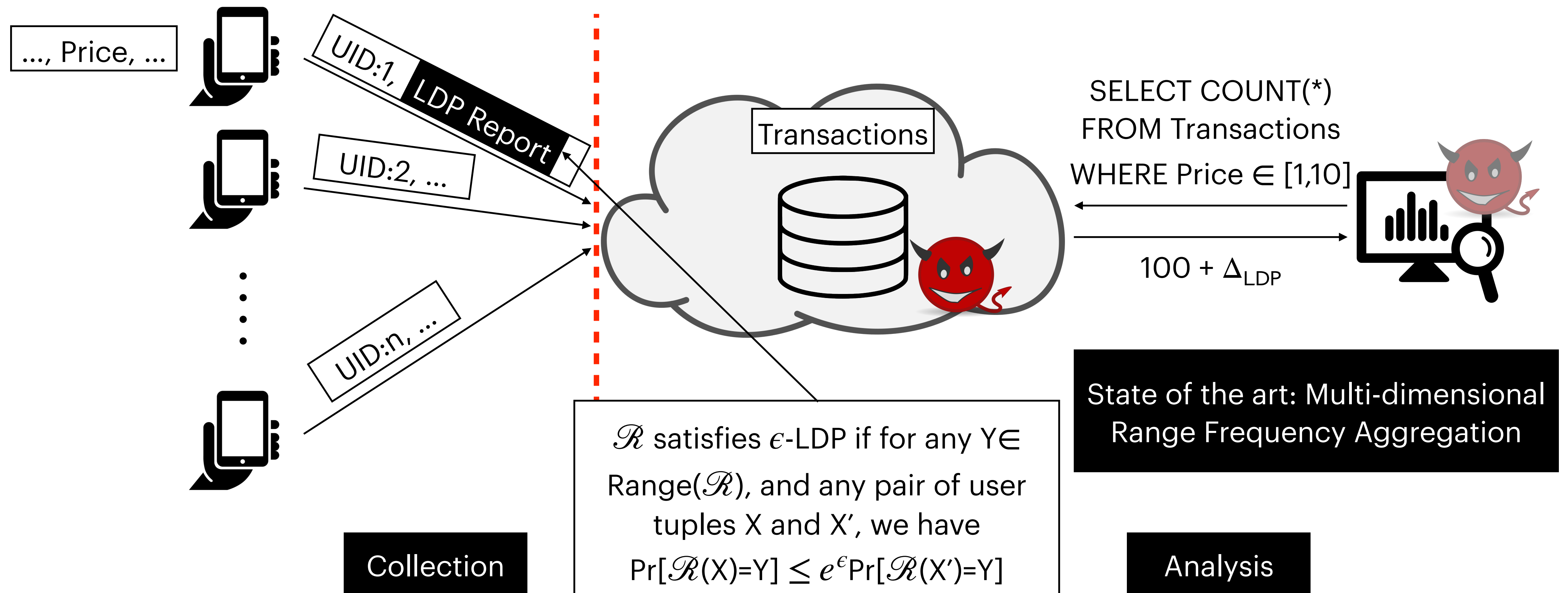
Data Collection and Analysis in the Cloud



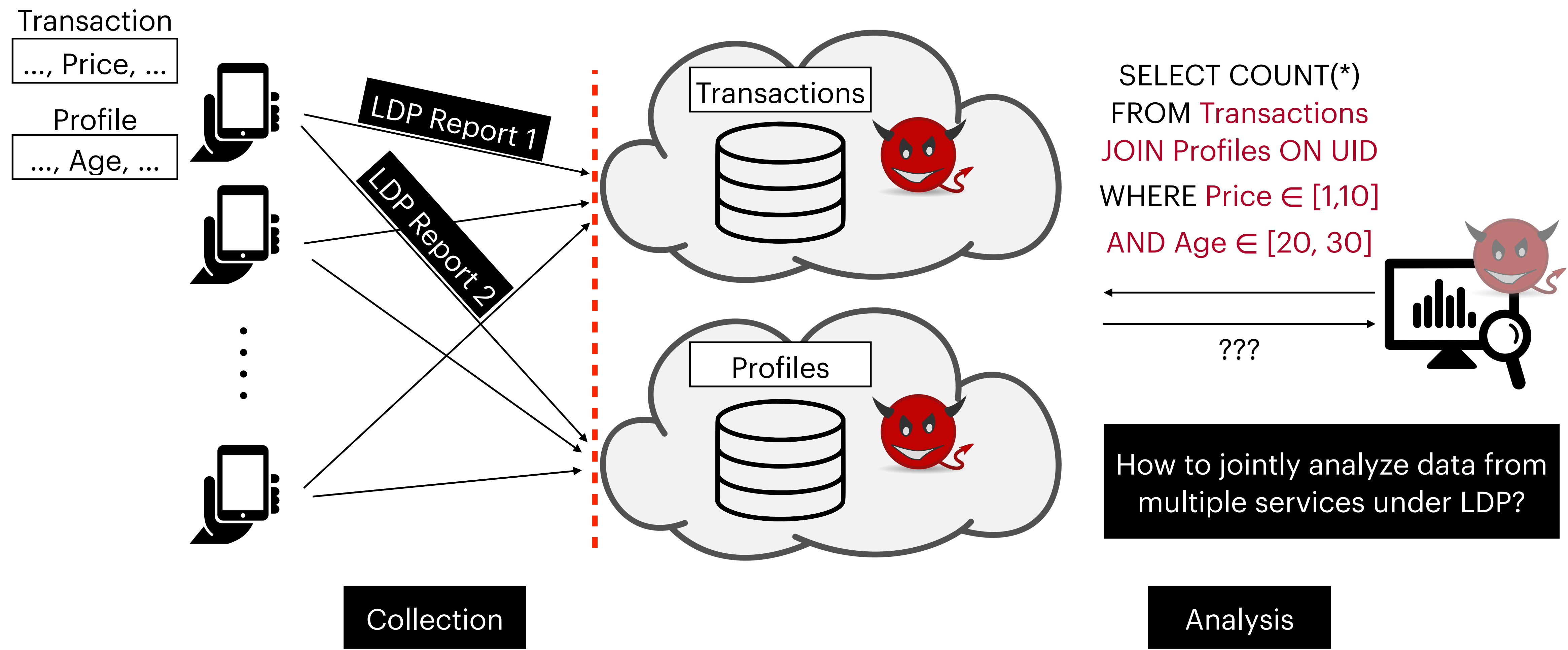
Data Collection and Analysis in the Cloud



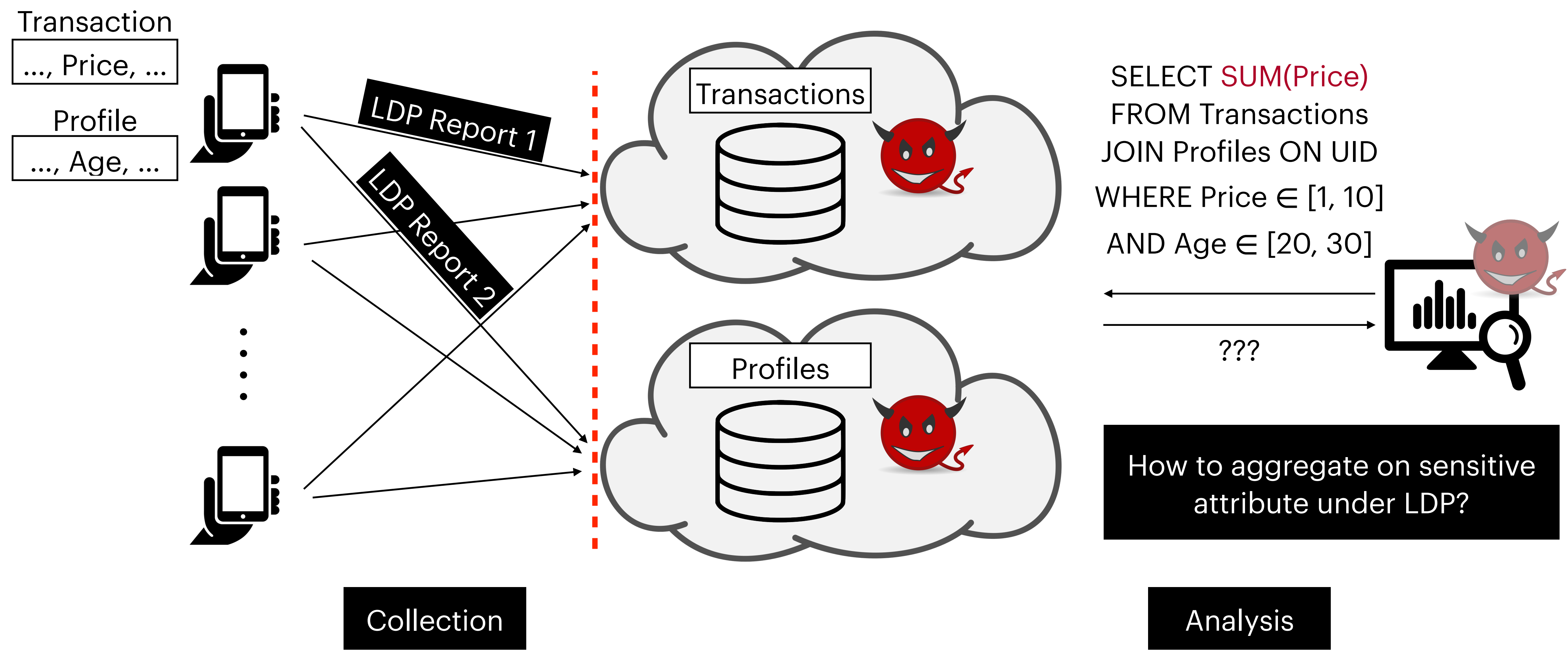
Data Collection and Analysis in the Cloud



Data Collection and Analysis in the Cloud - This Work



Data Collection and Analysis in the Cloud - This Work



Recap: Hierarchical Interval Optimized (HIO)

Collection

Hierarchy For Age $\in [1, 8]$

1-8							
1-4				5-8			
1-2		3-4		5-6		7-8	
1	2	3	4	5	6	7	8

Layer 1

Layer 2

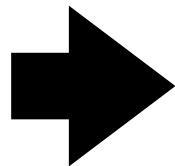
Layer 3

Layer 4

Recap: Multi-dimensional Range Frequency Oracles using Hierarchical Interval Optimized (HIO)

Collection

Age: 5

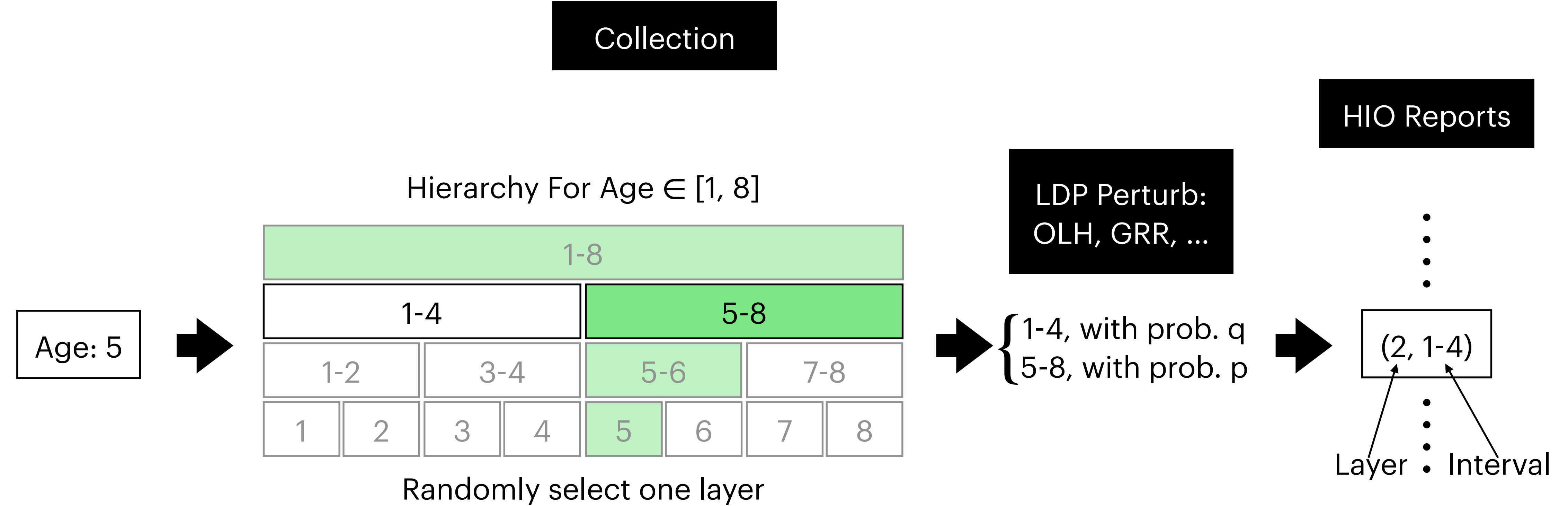


Hierarchy For Age $\in [1, 8]$

1-8							
1-4				5-8			
1-2		3-4		5-6		7-8	
1	2	3	4	5	6	7	8

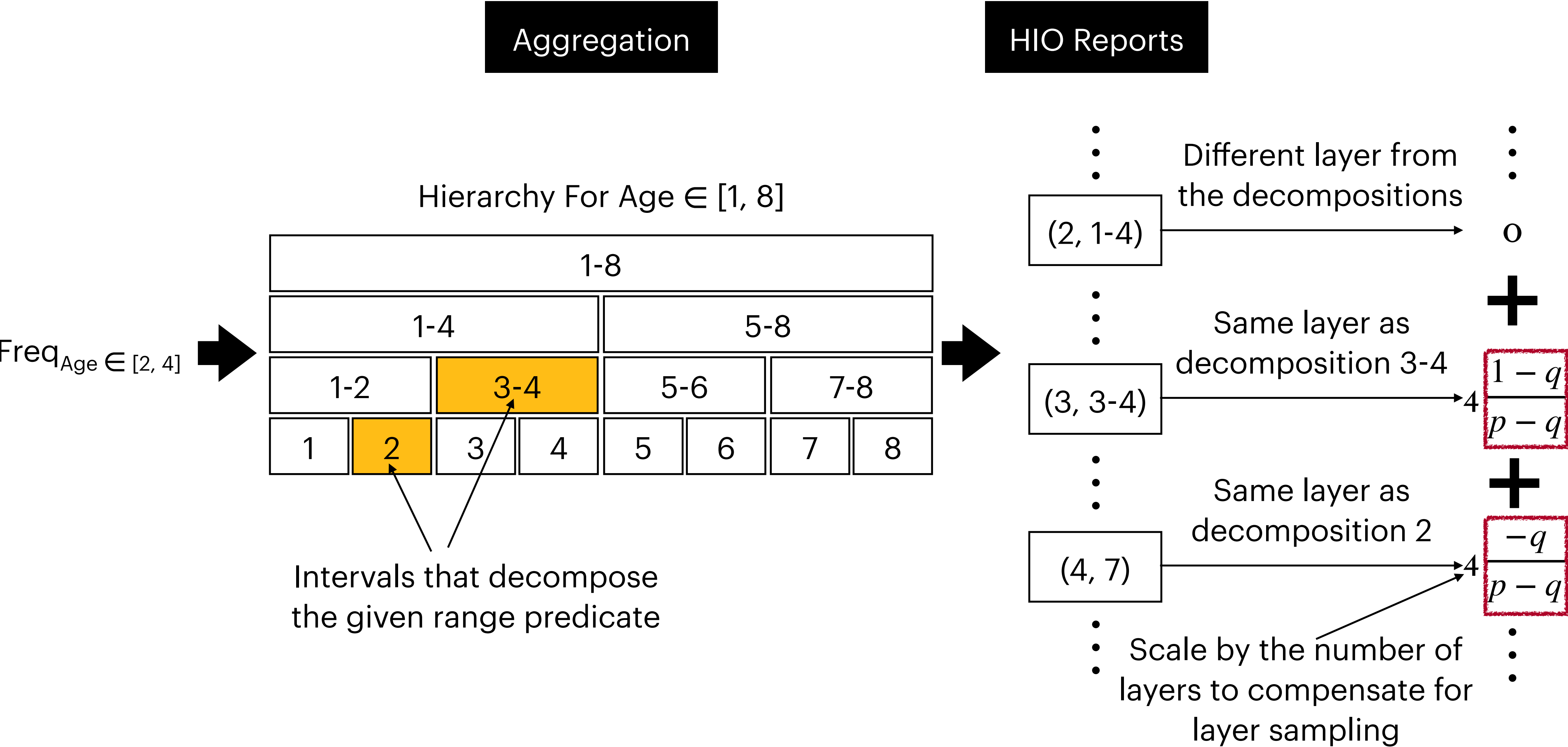
Intervals for
value 5

Recap: Multi-dimensional Range Frequency Oracles using Hierarchical Interval Optimized (HIO)

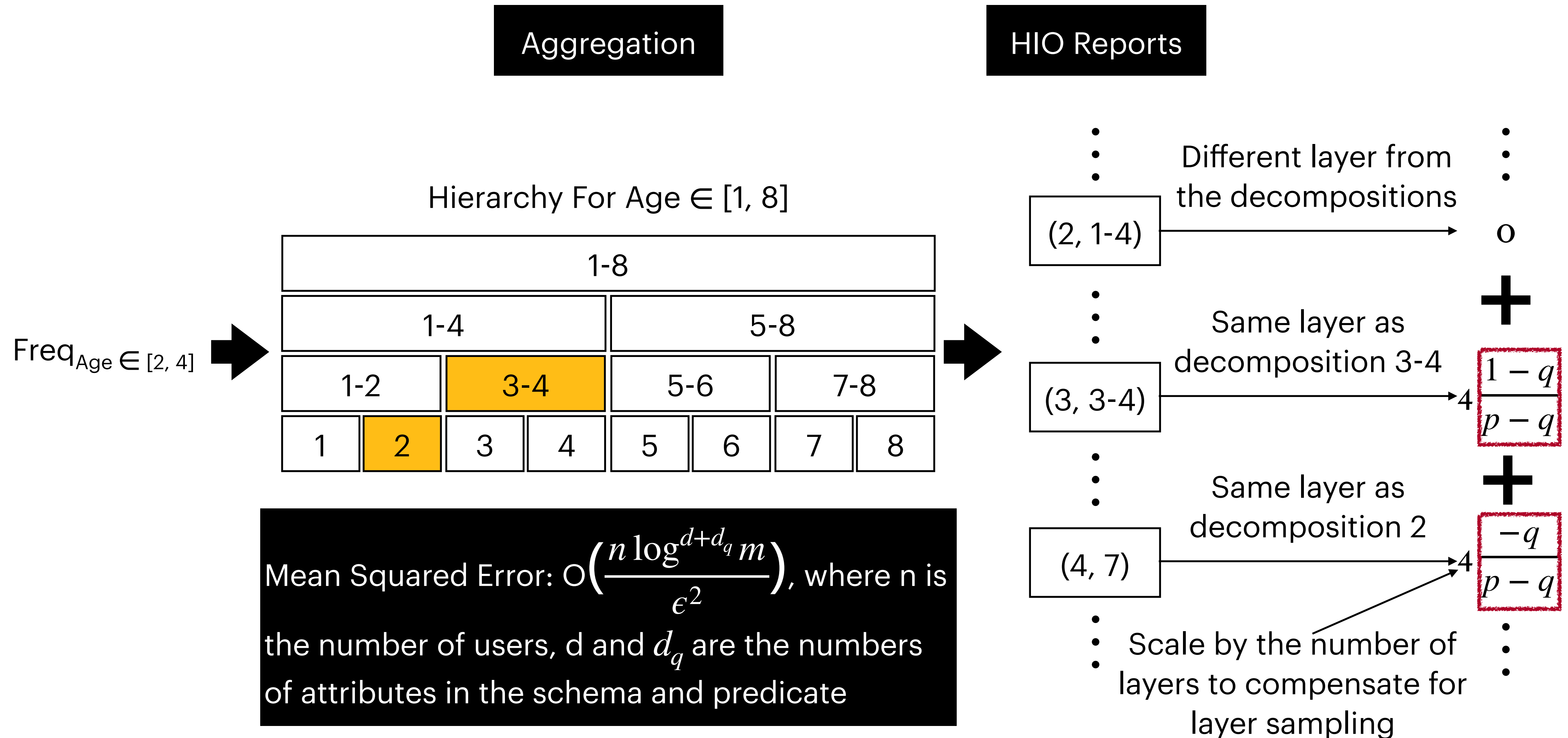


For GRR, $p = \frac{e^\epsilon}{e^\epsilon + m - 1}$, $q = \frac{1}{e^\epsilon + m - 1}$, where m is the cardinality of the domain

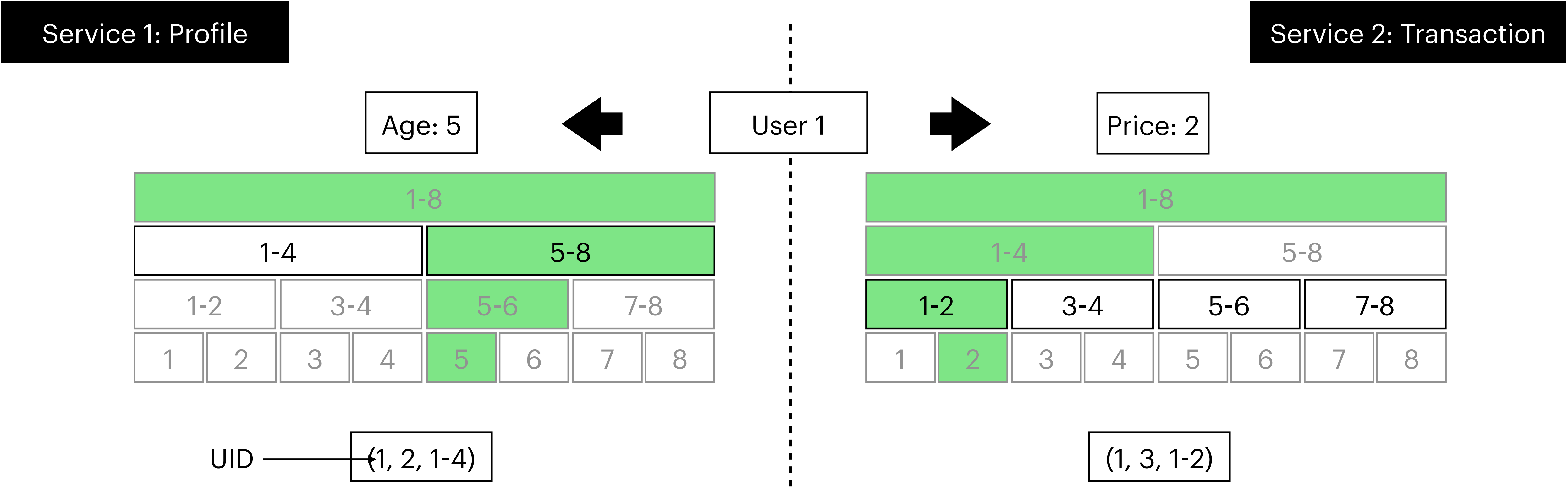
Recap: Multi-dimensional Range Frequency Oracles using Hierarchical Interval Optimized (HIO)



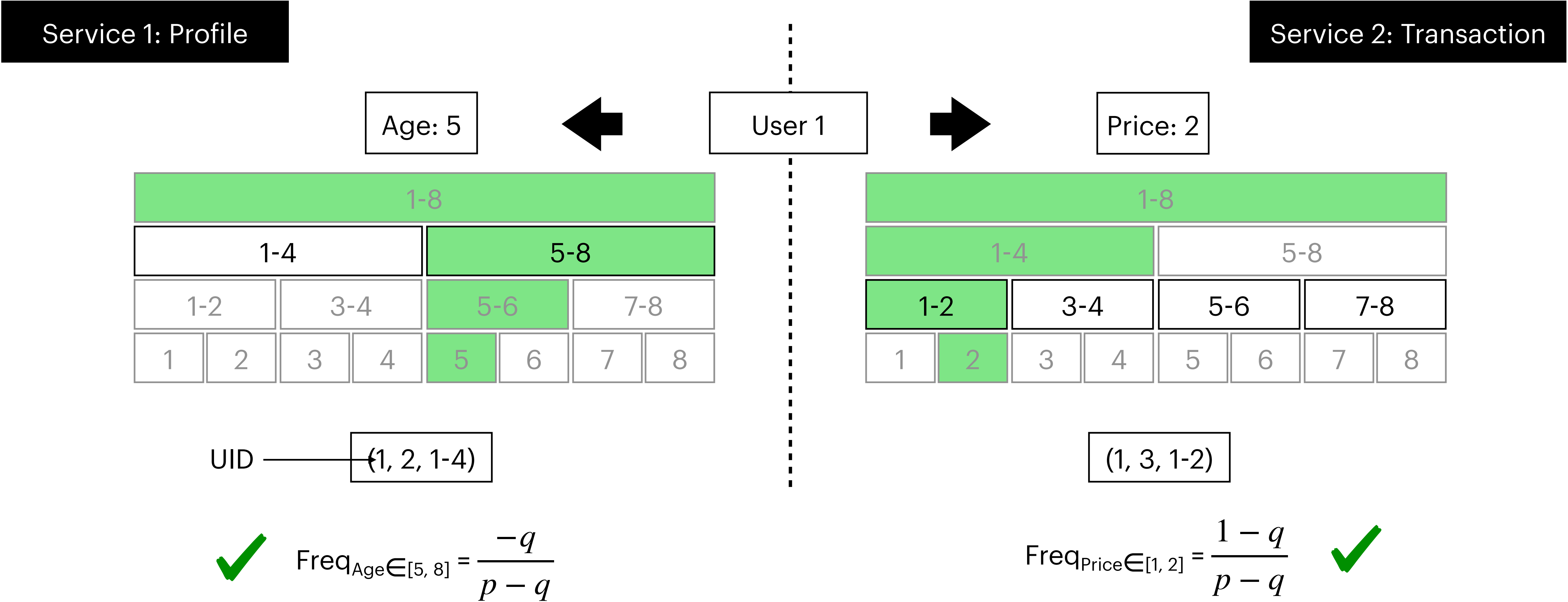
Recap: Multi-dimensional Range Frequency Oracles using Hierarchical Interval Optimized (HIO)



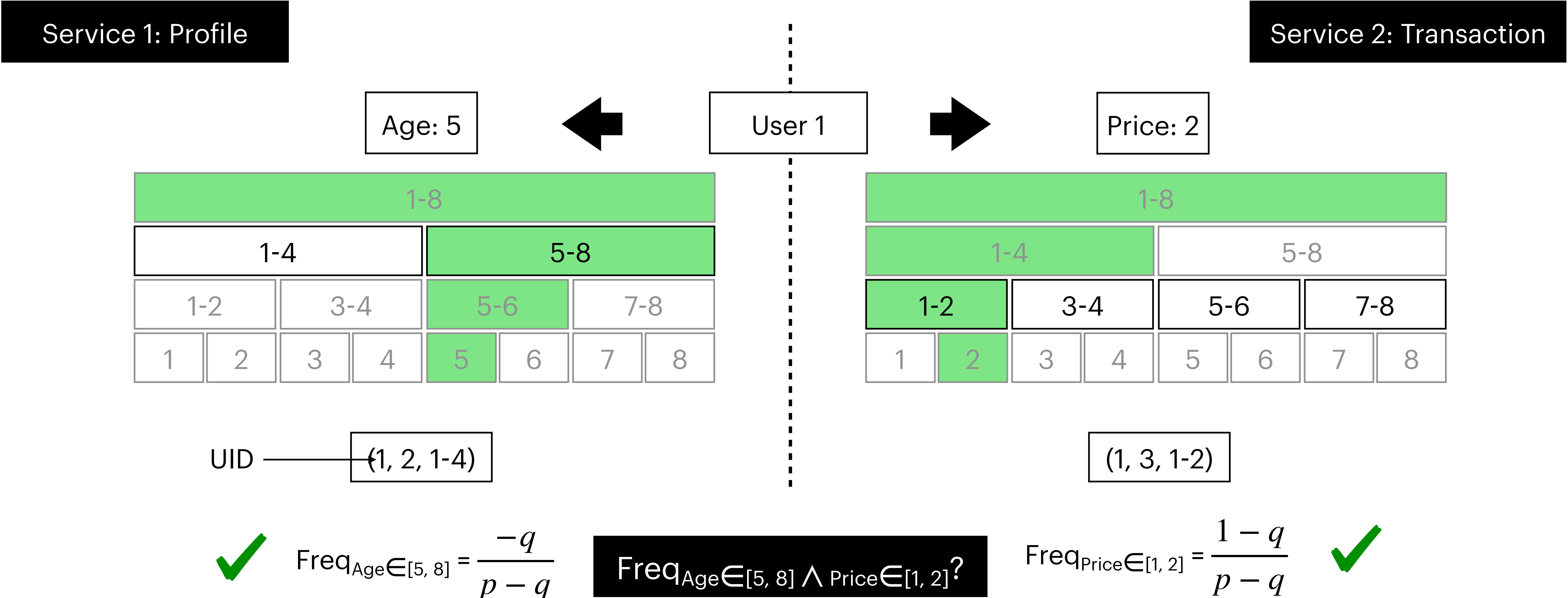
Data Collection with Two Independent Services using HIO



Data Collection with Two Independent Services using HIO



Data Collection with Two Independent Services



State Inversion with Two Independent Hierarchies

Target frequency
estimation

X[11] := Freq_{x1} ∈ [5, 8] ∧ x2 ∈ [1, 2]

X[10] := Freq_{x1} ∈ [5, 8] ∧ x2 ∉ [1, 2]

X[01] := Freq_{x1} ∉ [5, 8] ∧ x2 ∈ [1, 2]

X[00] := Freq_{x1} ∉ [5, 8] ∧ x2 ∉ [1, 2]

Y[11] := Freq_{y1} ∈ [5, 8] ∧ y2 ∈ [1, 2]

Y[10] := Freq_{y1} ∈ [5, 8] ∧ y2 ∉ [1, 2]

Y[01] := Freq_{y1} ∉ [5, 8] ∧ y2 ∈ [1, 2]

Y[00] := Freq_{y1} ∉ [5, 8] ∧ y2 ∉ [1, 2]

State Inversion with Two Independent Hierarchies

$$\mathbf{X}[11] := \text{Freq}_{\mathbf{x}_1 \in [5, 8] \wedge \mathbf{x}_2 \in [1, 2]}$$

$$\mathbf{X}[10] := \text{Freq}_{\mathbf{x}_1 \in [5, 8] \wedge \mathbf{x}_2 \notin [1, 2]}$$

$$\mathbf{X}[01] := \text{Freq}_{\mathbf{x}_1 \notin [5, 8] \wedge \mathbf{x}_2 \in [1, 2]}$$

$$\mathbf{X}[00] := \text{Freq}_{\mathbf{x}_1 \notin [5, 8] \wedge \mathbf{x}_2 \notin [1, 2]}$$

$$\mathbf{Y}[11] := \text{Freq}_{\mathbf{y}_1 \in [5, 8] \wedge \mathbf{y}_2 \in [1, 2]}$$

$$\mathbf{Y}[10] := \text{Freq}_{\mathbf{y}_1 \in [5, 8] \wedge \mathbf{y}_2 \notin [1, 2]}$$

$$\mathbf{Y}[01] := \text{Freq}_{\mathbf{y}_1 \notin [5, 8] \wedge \mathbf{y}_2 \in [1, 2]}$$

$$\mathbf{Y}[00] := \text{Freq}_{\mathbf{y}_1 \notin [5, 8] \wedge \mathbf{y}_2 \notin [1, 2]}$$

$$\mathbf{E} \begin{bmatrix} \mathbf{Y}[11] \\ \mathbf{Y}[10] \\ \mathbf{Y}[01] \\ \mathbf{Y}[00] \end{bmatrix} = \begin{bmatrix} \text{Pr}[\mathbf{Y}[11]|\mathbf{X}[11]] & \text{Pr}[\mathbf{Y}[11]|\mathbf{X}[10]] & \text{Pr}[\mathbf{Y}[11]|\mathbf{X}[01]] & \text{Pr}[\mathbf{Y}[11]|\mathbf{X}[00]] \\ \text{Pr}[\mathbf{Y}[10]|\mathbf{X}[11]] & \text{Pr}[\mathbf{Y}[10]|\mathbf{X}[10]] & \text{Pr}[\mathbf{Y}[10]|\mathbf{X}[01]] & \text{Pr}[\mathbf{Y}[10]|\mathbf{X}[00]] \\ \text{Pr}[\mathbf{Y}[01]|\mathbf{X}[11]] & \text{Pr}[\mathbf{Y}[01]|\mathbf{X}[10]] & \text{Pr}[\mathbf{Y}[01]|\mathbf{X}[01]] & \text{Pr}[\mathbf{Y}[01]|\mathbf{X}[00]] \\ \text{Pr}[\mathbf{Y}[00]|\mathbf{X}[11]] & \text{Pr}[\mathbf{Y}[00]|\mathbf{X}[10]] & \text{Pr}[\mathbf{Y}[00]|\mathbf{X}[01]] & \text{Pr}[\mathbf{Y}[00]|\mathbf{X}[00]] \end{bmatrix} * \begin{bmatrix} \mathbf{X}[11] \\ \mathbf{X}[10] \\ \mathbf{X}[01] \\ \mathbf{X}[00] \end{bmatrix}$$

State Inversion with Two Independent Hierarchies

X[11] := Freq $x_1 \in [5, 8] \wedge x_2 \in [1, 2]$

X[10] := Freq $x_1 \in [5, 8] \wedge x_2 \notin [1, 2]$

X[01] := Freq $x_1 \notin [5, 8] \wedge x_2 \in [1, 2]$

X[00] := Freq $x_1 \notin [5, 8] \wedge x_2 \notin [1, 2]$

Y[11] := Freq $y_1 \in [5, 8] \wedge y_2 \in [1, 2]$

Y[10] := Freq $y_1 \in [5, 8] \wedge y_2 \notin [1, 2]$

Y[01] := Freq $y_1 \notin [5, 8] \wedge y_2 \in [1, 2]$

Y[00] := Freq $y_1 \notin [5, 8] \wedge y_2 \notin [1, 2]$

$$\begin{array}{c} \mathbf{E} \end{array} \begin{bmatrix} Y[11] \\ Y[10] \\ Y[01] \\ Y[00] \end{bmatrix} = \begin{array}{c} \begin{bmatrix} \Pr[Y11|X11] & \Pr[Y11|X10] & \Pr[Y11|X01] & \Pr[Y11|X00] \\ \Pr[Y10|X11] & \Pr[Y10|X10] & \Pr[Y10|X01] & \Pr[Y10|X00] \\ \Pr[Y01|X11] & \Pr[Y01|X10] & \Pr[Y01|X01] & \Pr[Y01|X00] \\ \Pr[Y00|X11] & \Pr[Y00|X10] & \Pr[Y00|X01] & \Pr[Y00|X00] \end{bmatrix} \end{array} * \begin{array}{c} \begin{bmatrix} \mathbf{X[11]} \\ X[10] \\ X[01] \\ X[00] \end{bmatrix} \end{array} \Rightarrow \mathbf{X} \leftarrow \mathbf{M}^{-1} \mathbf{Y}$$

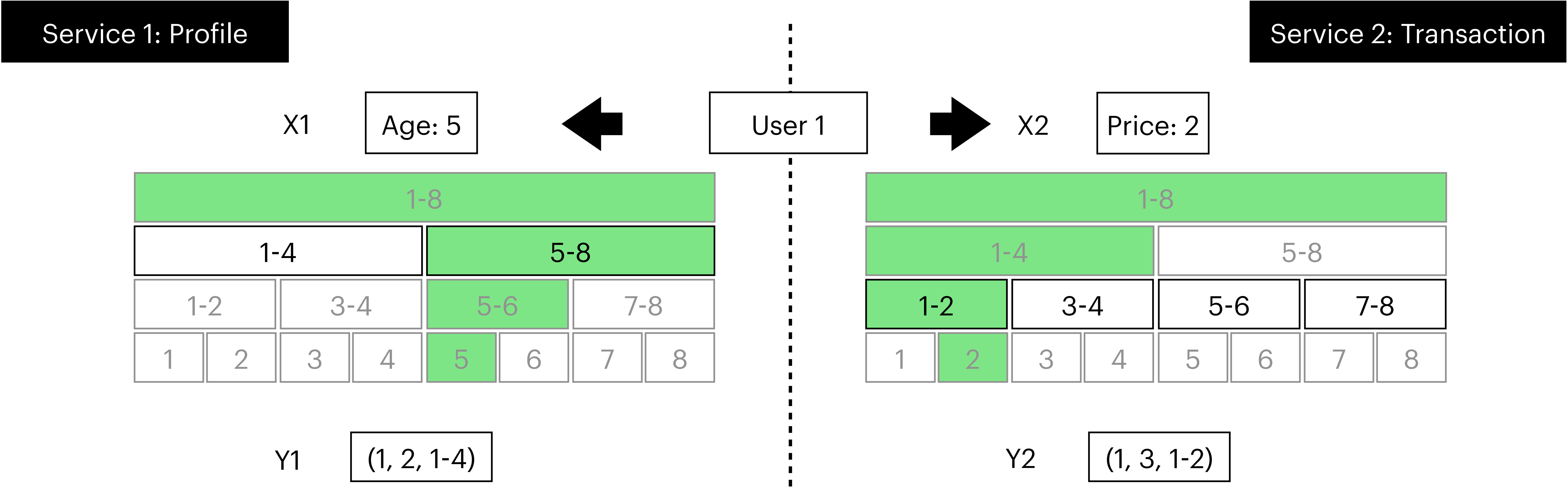
$\begin{array}{ccc} \uparrow & \uparrow & \uparrow \\ \mathbf{Y} & \mathbf{M} & \mathbf{X} \end{array}$

Unbiased estimation

- How to calculate $\Pr[Y11|X11]$?

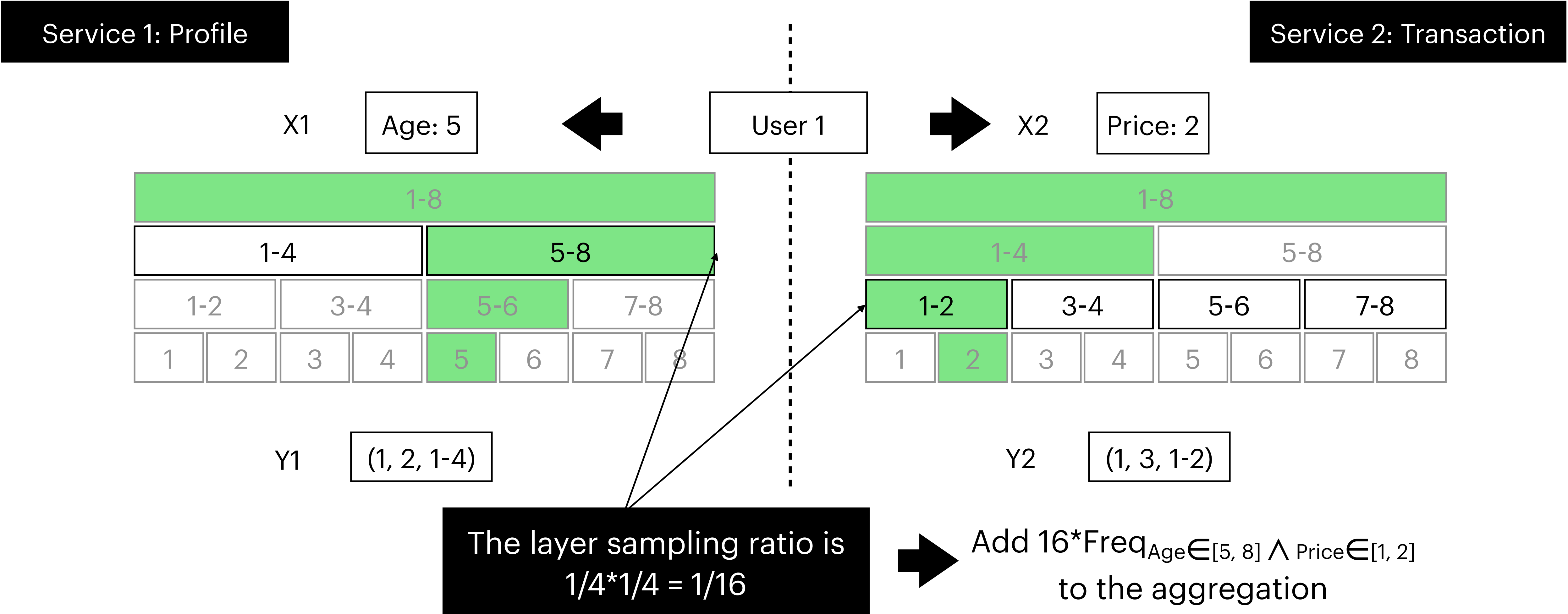
Because X_1 and X_2 are independently perturbed into Y_1 and Y_2 , $\Pr[Y11|X11] = \Pr[Y1 \in [2, 4] | X1 \in [2, 4]] * \Pr[Y2 \in [1, 4] | X2 \in [1, 4]] = p * p$. We can derive other conditional prob. similarly

Joint Frequency Oracles with Two Independent Services

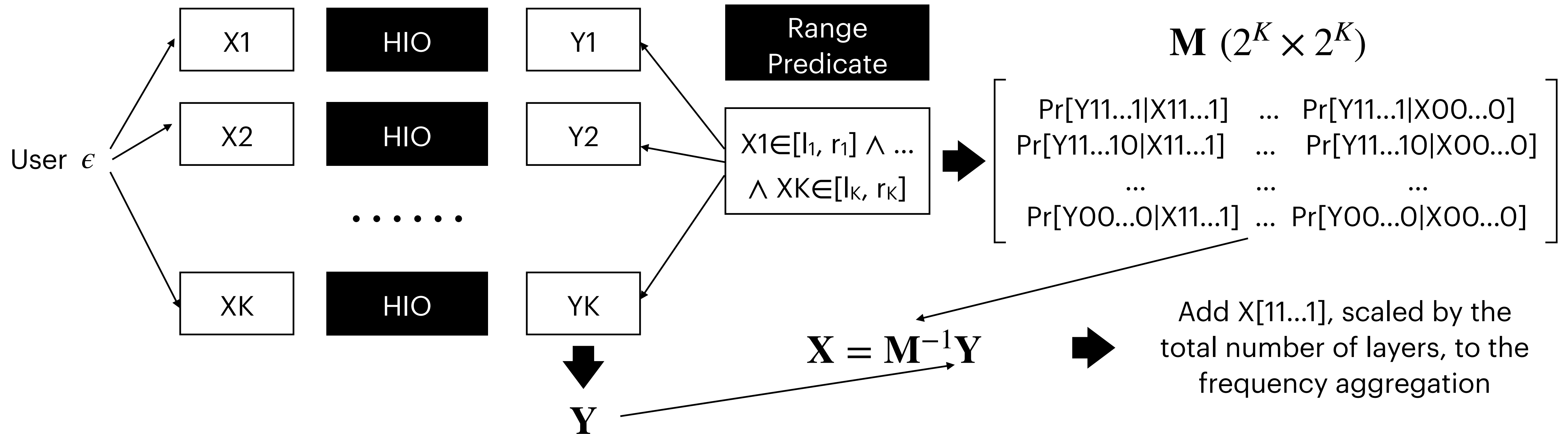


$$\text{Freq}_{\text{Age} \in [5, 8] \wedge \text{Price} \in [1, 2]} \leftarrow (\mathbf{M}^{-1} \mathbf{Y})[11 \dots 1]$$

Joint Frequency Oracles with Two Independent Services

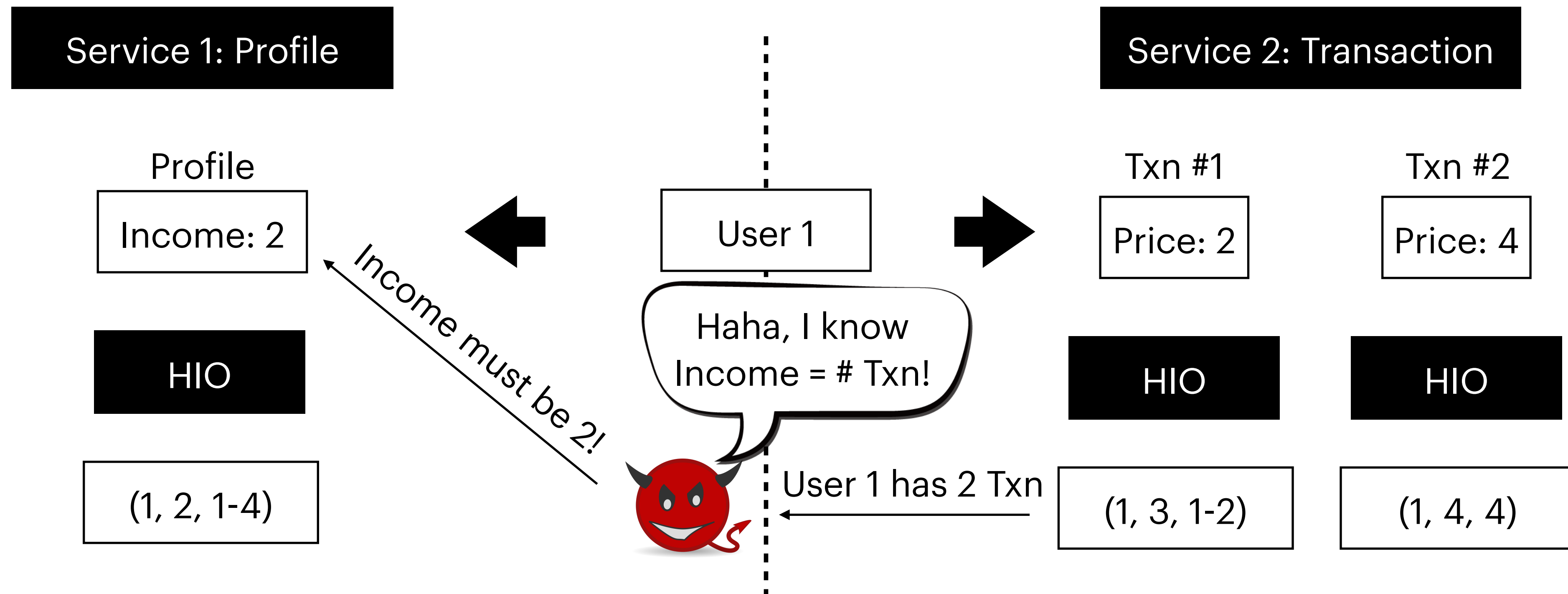


General Joint Frequency Oracles with K Services (HIO-JOIN)



- Mean Squared Error: $O\left(\frac{n \log^{K(d+d_q)} m}{(\epsilon/K)^{2K}}\right)$, where K is the number of services, d and d_q are the numbers of attributes in the schema and predicate for each service
- In fact, it can handle range predicate on attributes of arbitrary subset of the K services

1-Many join: Primary-Foreign-Key JOIN with Two Services



Frequency-Based Attack: Infer sensitive information based on the leaked number of tuples of a user collected by each service

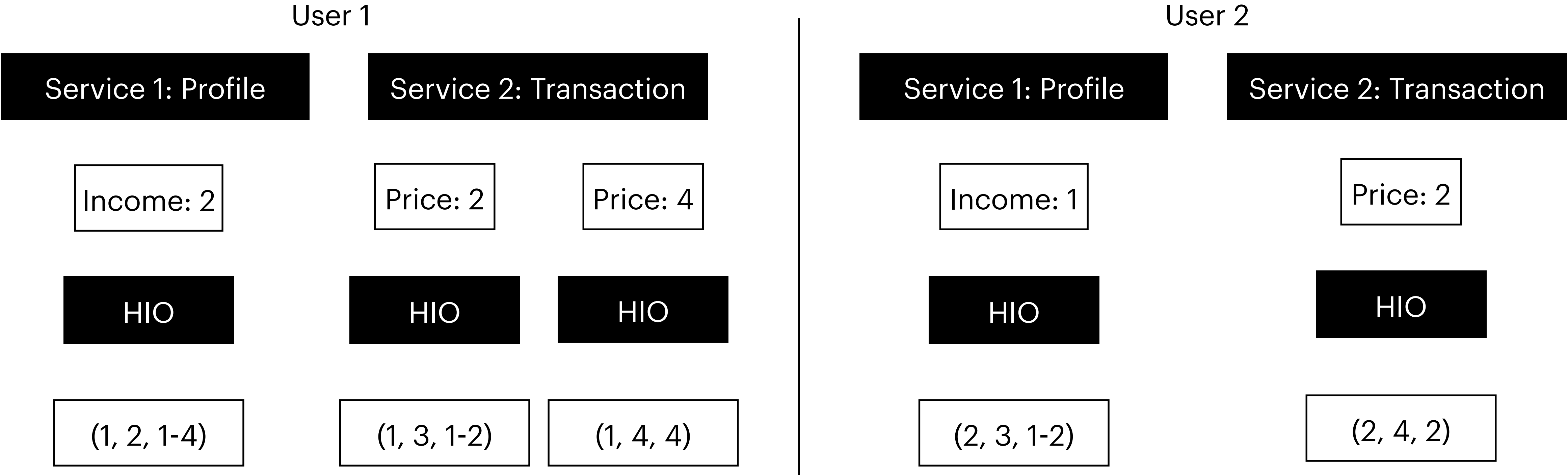
- The information that user 1 made 2 transactions with service 2 is leaked at join
- If such information is correlated with the sensitive information, say, the income, attacker can infer the income of the user based on the number of transactions she made

User-Level LDP (uLDP)

ϵ -uLDP: Intuitively, for any two users, the difference between the distributions of their collected data is bounded by e^ϵ

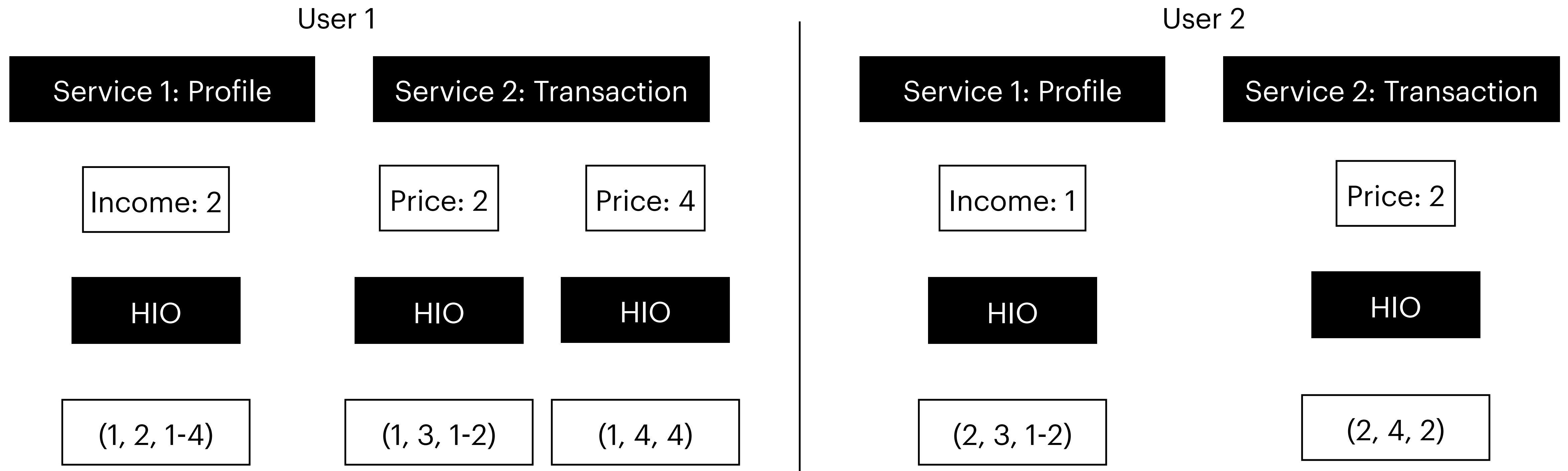
User-Level LDP (uLDP)

ϵ -uLDP: Intuitively, for any two users, the difference between the distributions of their collected data is bounded by e^ϵ



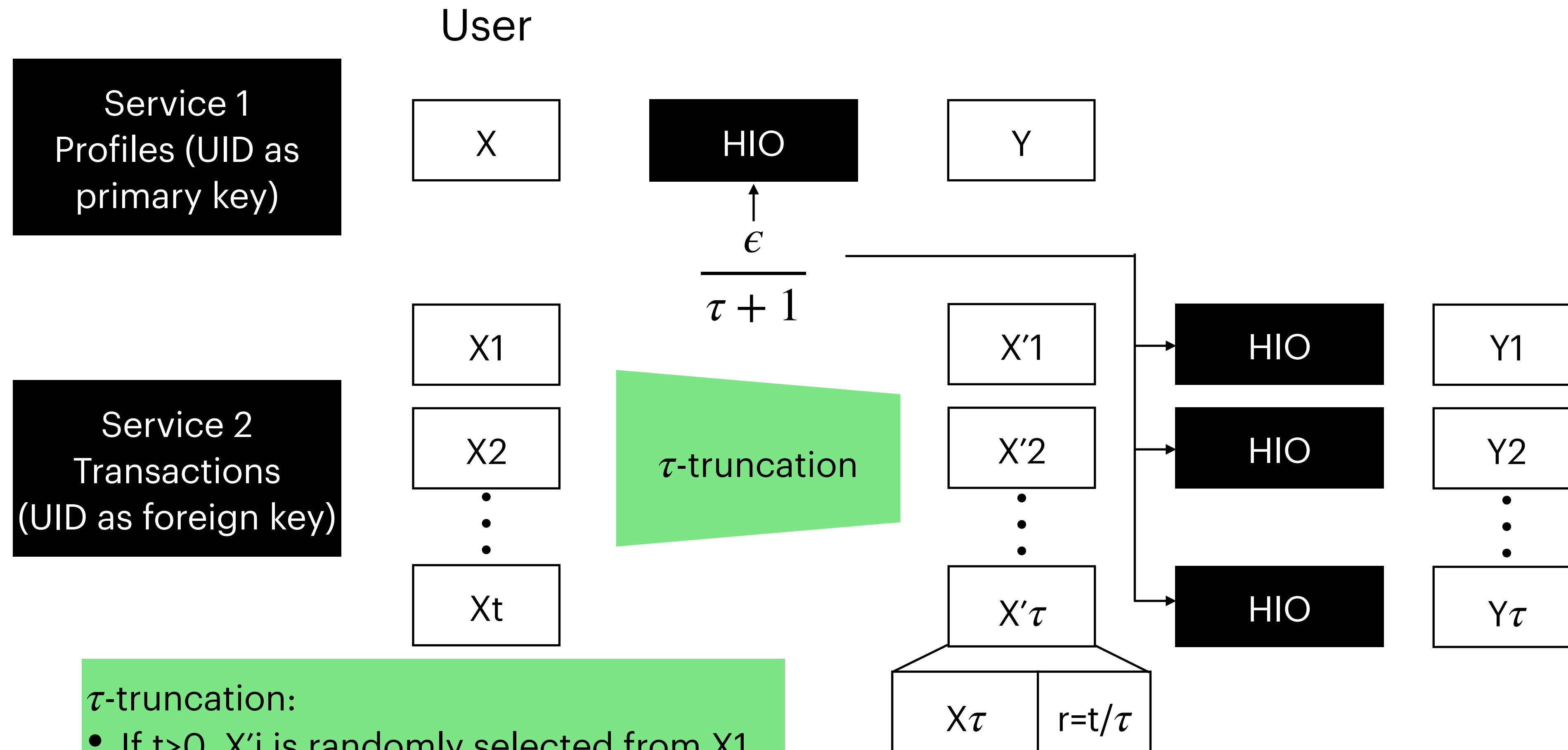
User-Level LDP (uLDP)

ϵ -uLDP: Intuitively, for any two users, the difference between the distributions of their collected data is bounded by e^ϵ



Impossible to have the same collected data for user 1 and 2! $\Rightarrow \infty$ -uLDP \Rightarrow no privacy at all!

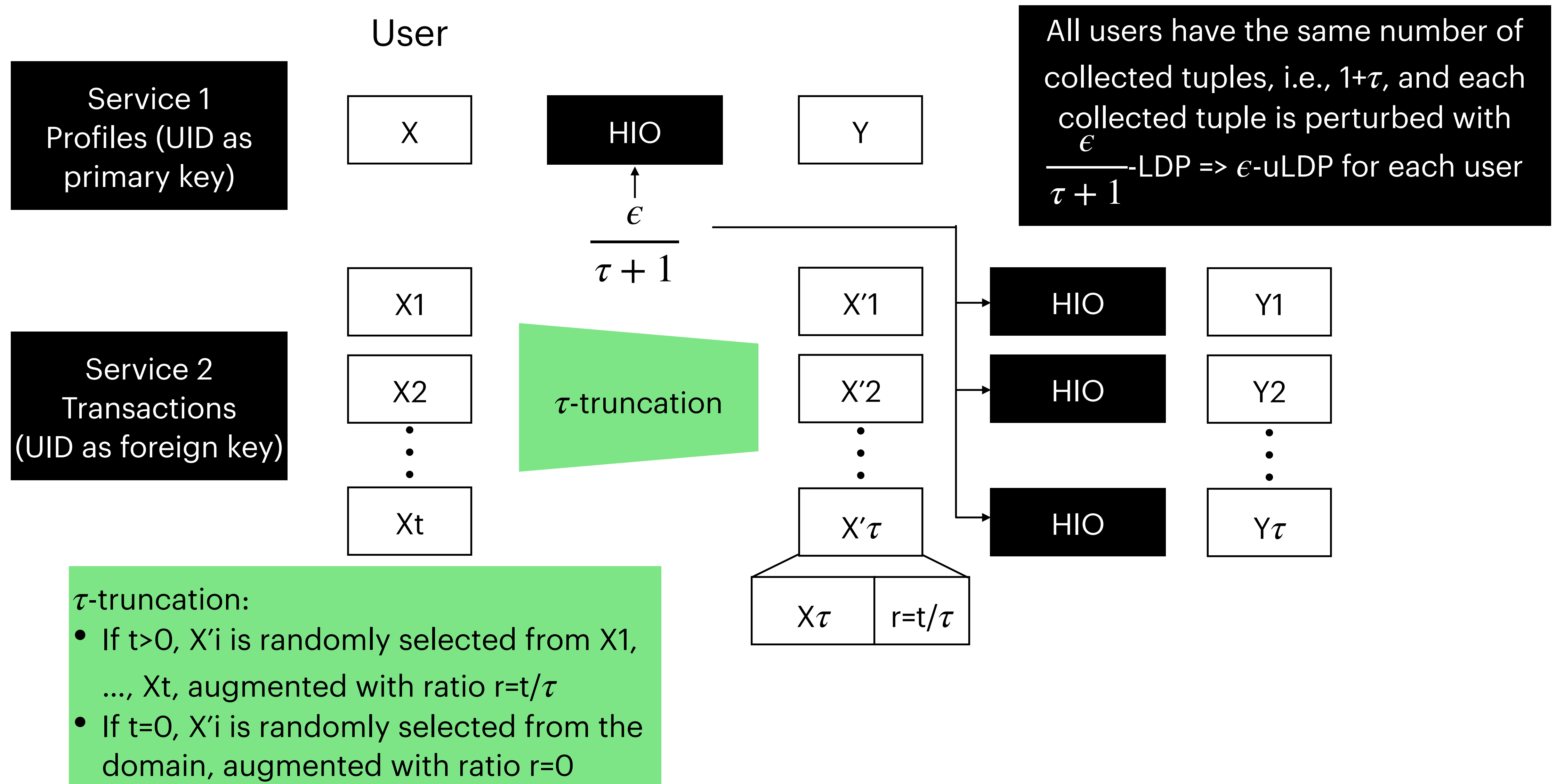
1-Many JOIN under ϵ -uLDP



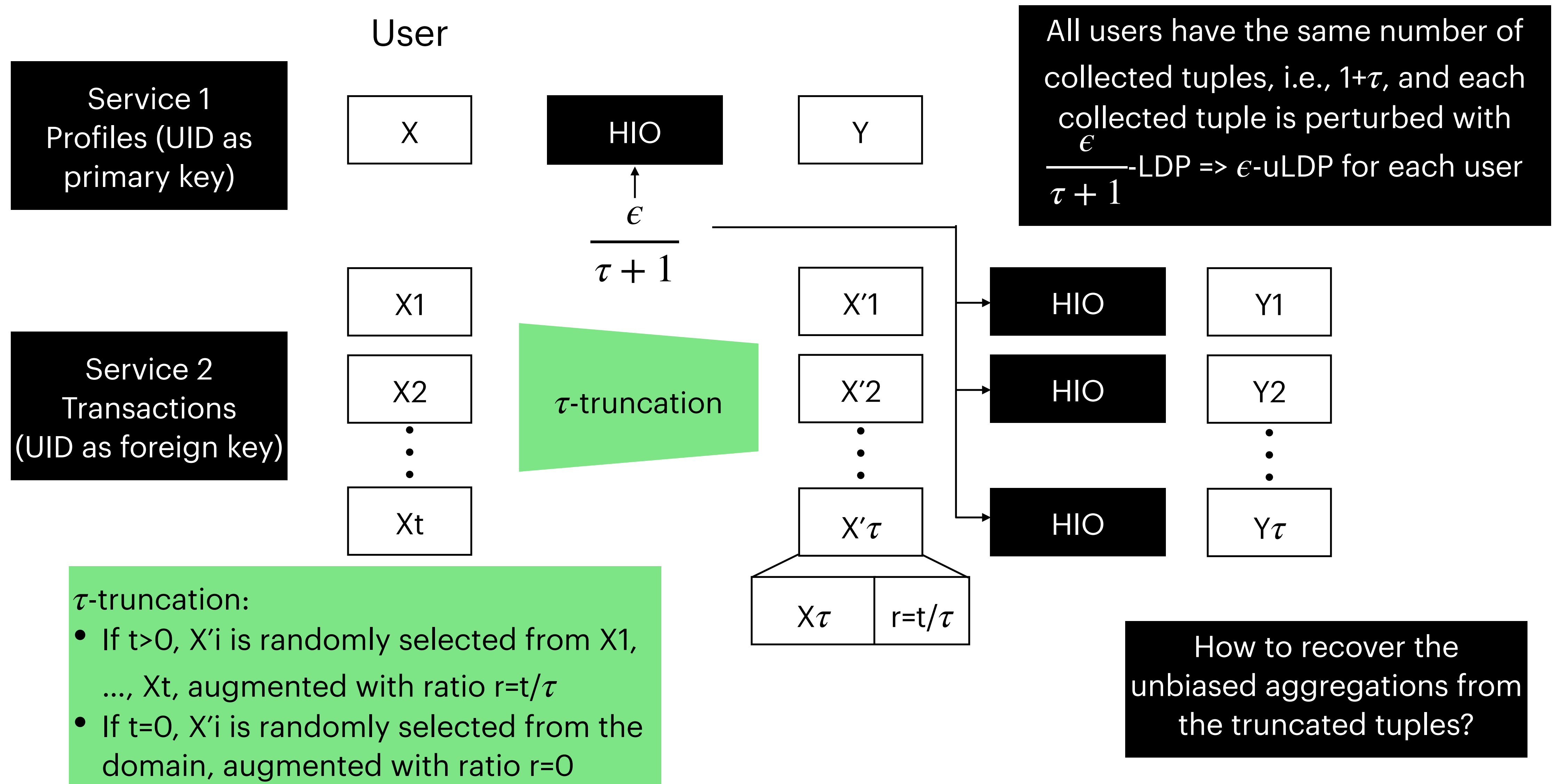
τ -truncation:

- If $t > 0$, X'_i is randomly selected from X_1, \dots, X_t , augmented with ratio $r = t/\tau$
- If $t = 0$, X'_i is randomly selected from the domain, augmented with ratio $r = 0$

1-Many JOIN under ϵ -uLDP



1-Many JOIN under ϵ -uLDP



Sensitive Weight Frequency Oracles (SWFO)

- So far, we can handle frequency oracles on the collected data with range predicate over arbitrary subset of attributes from the multiple services
- It is necessary to further support oracles for the aggregated frequency, weighted by some sensitive attribute of each tuple. We call such oracles the *sensitive weight frequency oracles*:

- **Attribute Aggregation**

SELECT SUM(Income) FROM Profiles WHERE Age \in [20,50]

- **Recover Unbiased Aggregation with τ -Truncation** (weighing each tuple by $r=t/\tau$)

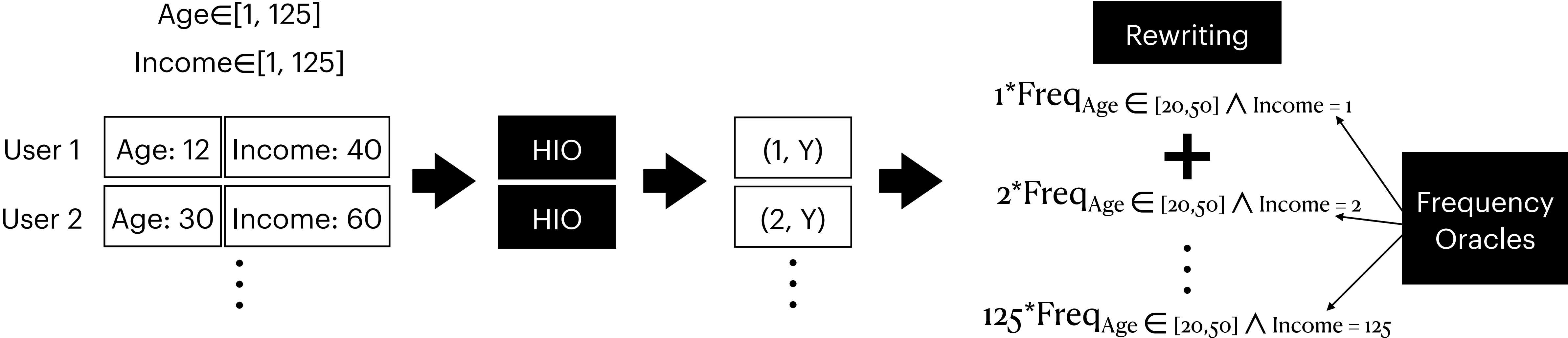
SELECT SUM(r) FROM Profiles JOIN Transactions ON UID WHERE Age \in [20,50]

SELECT SUM(r*Price) FROM Profiles JOIN Transactions ON UID WHERE Age \in [20,50]

- Previous works leave such aggregation analysis as open problems

(Joint) Frequency Oracles to SWFO (Straw-man)

SELECT SUM(Income) FROM Profiles WHERE Age $\in [20,50]$

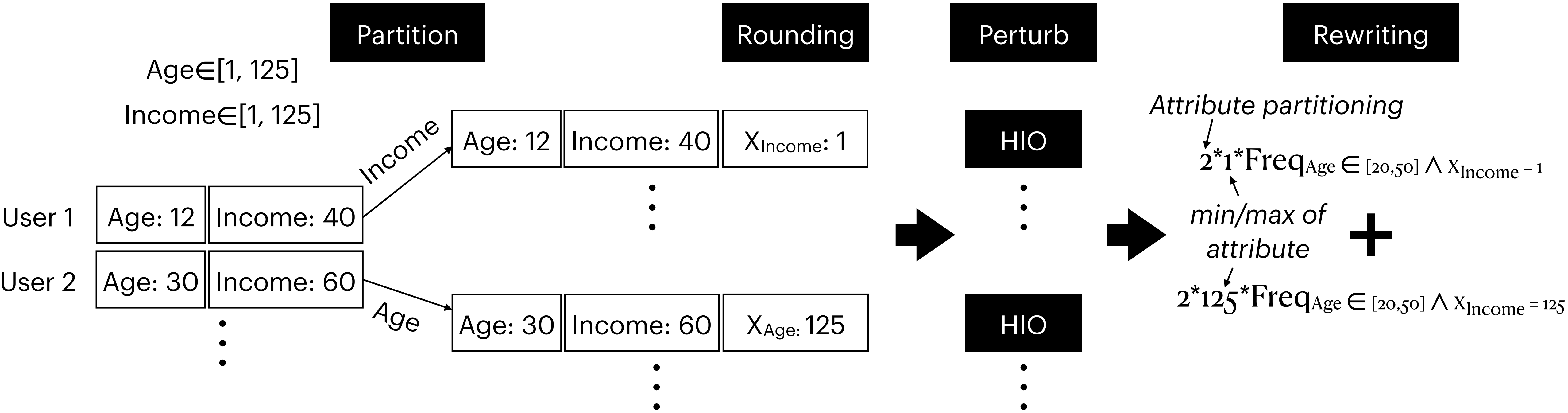


Mean Squared Error: $O\left(\sum_{v \in A} v^2 \cdot \frac{n \log^{d+d_q} m}{\epsilon^2}\right)$

The multiplicative factor of sum of squares of the domain values in the error bound is prohibitive for aggregation on large domains

Partition-Rounding-Perturb Framework for SWFO

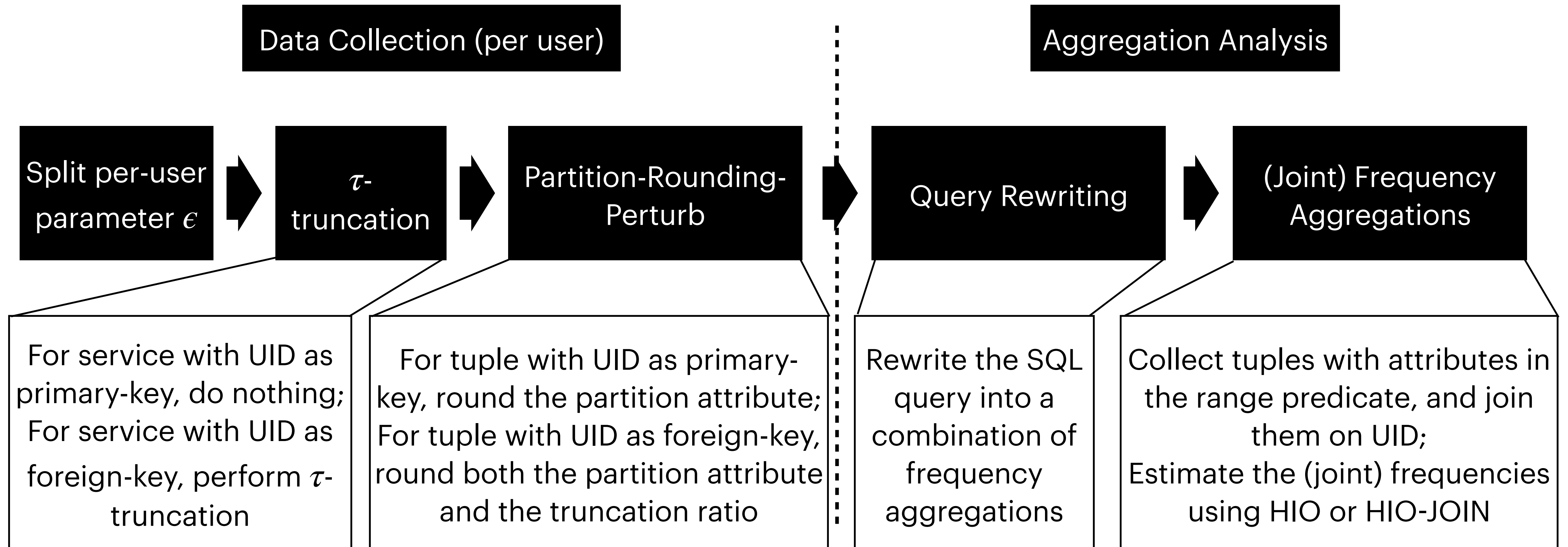
SELECT SUM(Income) FROM Profiles WHERE Age $\in [20,50]$



Mean Squared Error: $O\left(\frac{2(A_{\min}^2 + A_{\max}^2)d \cdot n \log^{d+d_q} m}{\epsilon^2}\right)$

The multiplicative factor is linear to the number of attributes times the sum of squares of the min/max

End-to-End Data Collection and Analysis Pipeline



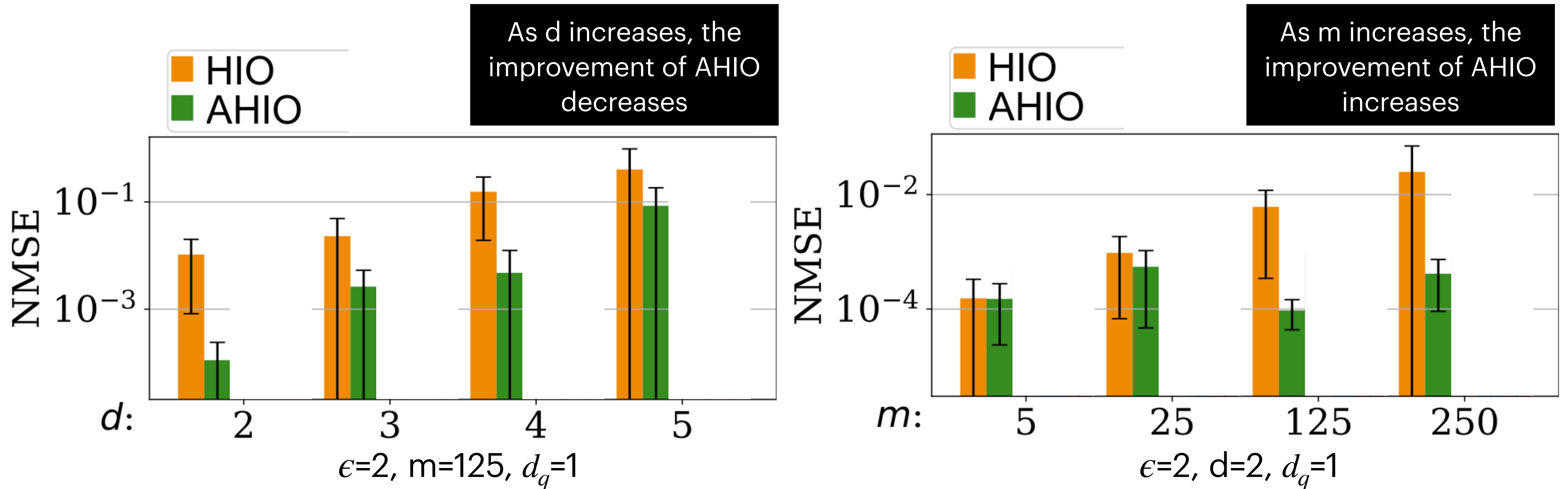
- Query independent - can handle arbitrary frequency and attribute aggregation with range predicate over any subset of attributes
- Minimal coordination among services: i) value of τ ; and ii) the total number of tuples to be collected for each user

Experimental Results

- We setup a single-node SparkSQL cluster for the experiments, and register the LDP perturb and aggregation estimation primitives as UDFs of SparkSQL
- Datasets
 - SYN-1: Single-table synthetic data , with 1M records and configurable number of attributes d and attribute cardinality m
 - SYN-2: Two-table synthetic data, with two join configurations: 1-1; and 1-many with join degree in range $[1, 10]$
- We focus on the utilities for COUNT, SUM and AVG aggregations of the following schemes
 - AHIO: Partition-Rounding-Perturb with the rounding value stored in separate attribute
 - HIO-JOIN: no truncation for 1-1 join; with τ -truncation for 1-many join

Attribute Aggregation

SELECT SUM(A) FROM SYN-1 WHERE C

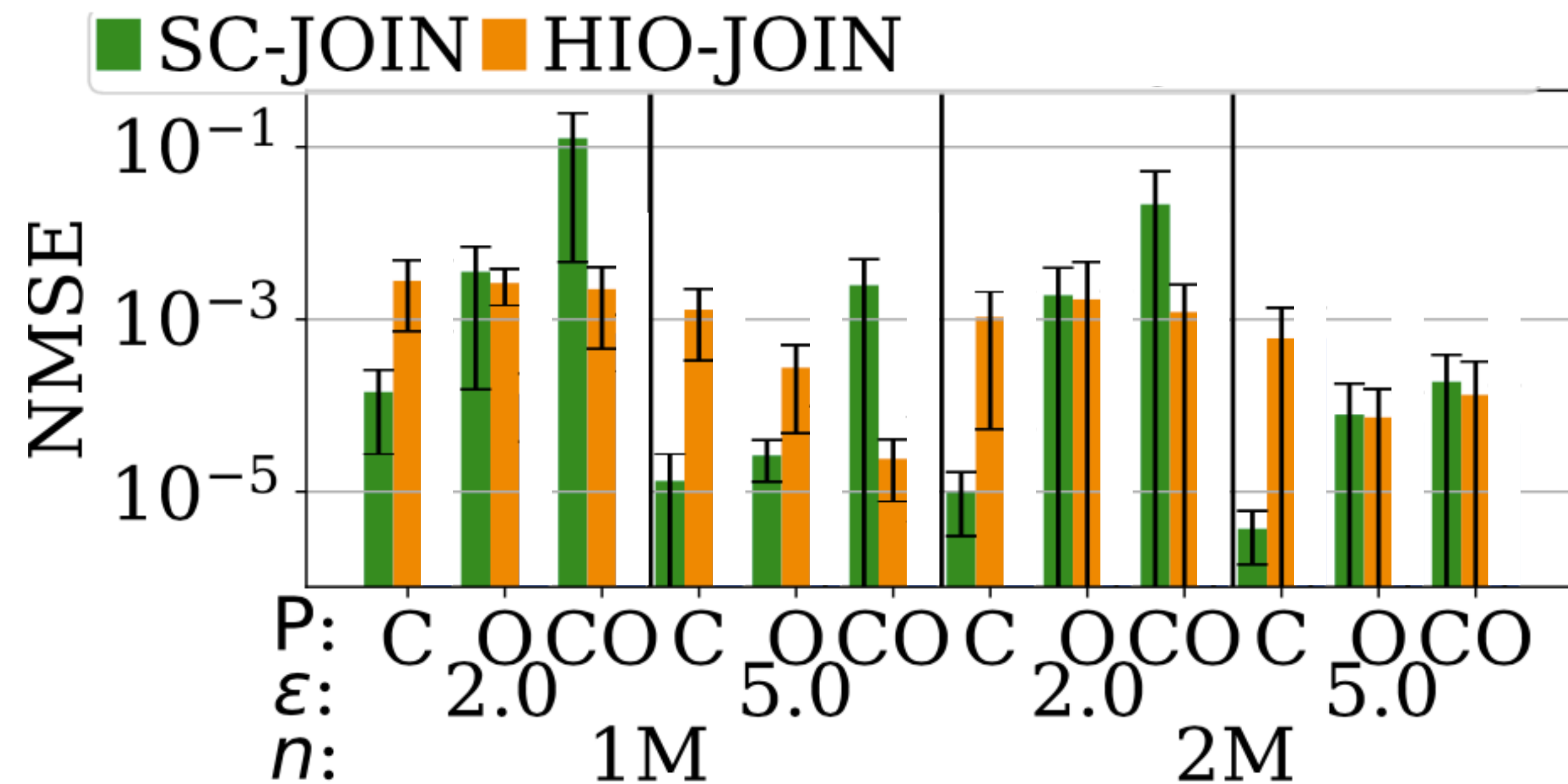


Asymptotic Mean Squared Error

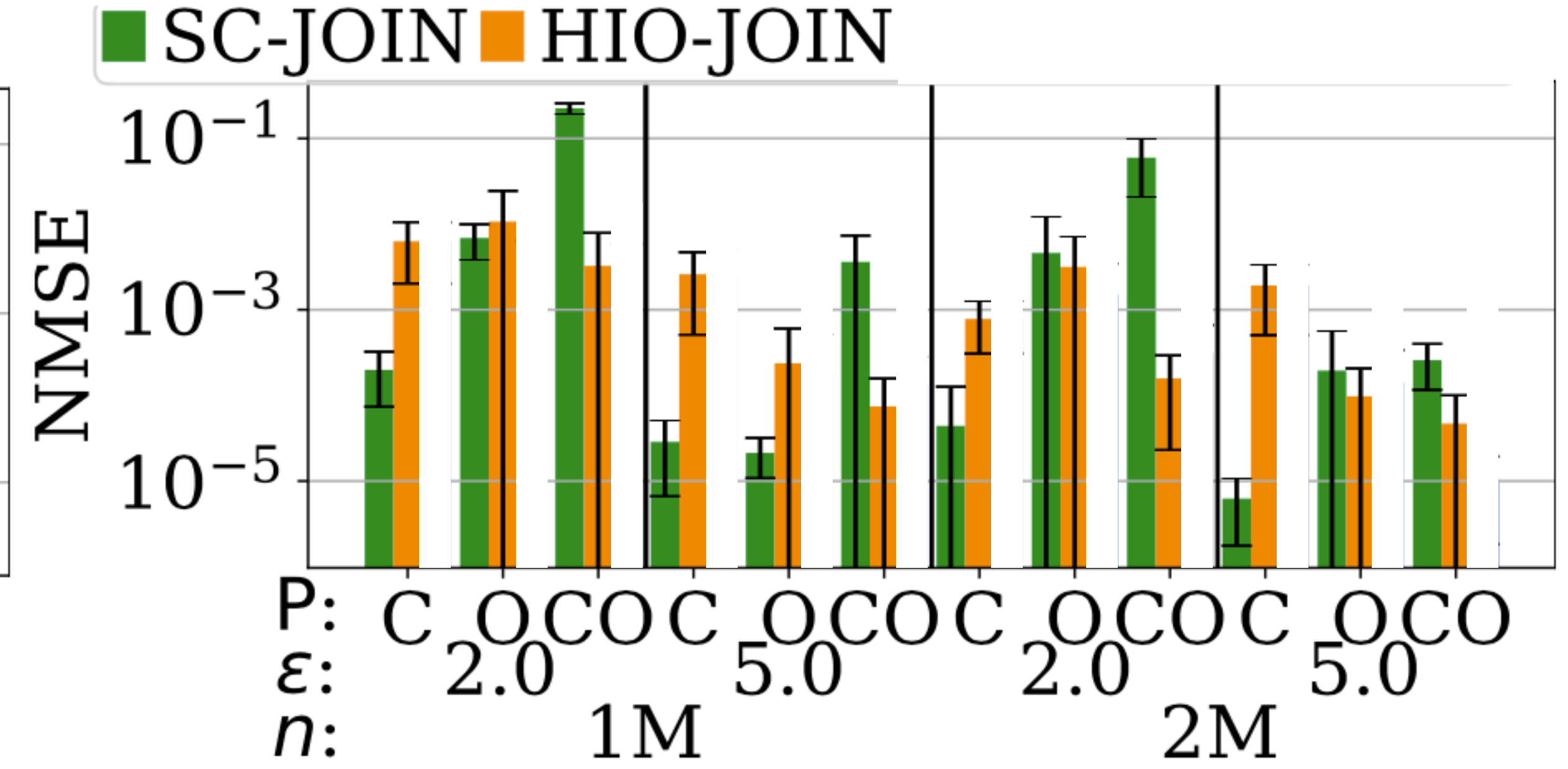
$$\text{HIO: } O\left(\sum_{v \in A} v^2 \frac{n \log^{d+d_q} m}{\epsilon^2}\right), \text{ AHIO: } O\left(\frac{2(A_{\min}^2 + A_{\max}^2)d \cdot n \log^{d+d_q} m}{\epsilon^2}\right)$$

1-1 Joint Aggregation

SELECT COUNT(*) FROM T1 JOIN T2 ON
UID WHERE C



SELECT SUM(A) FROM T1 JOIN T2 ON
UID WHERE C



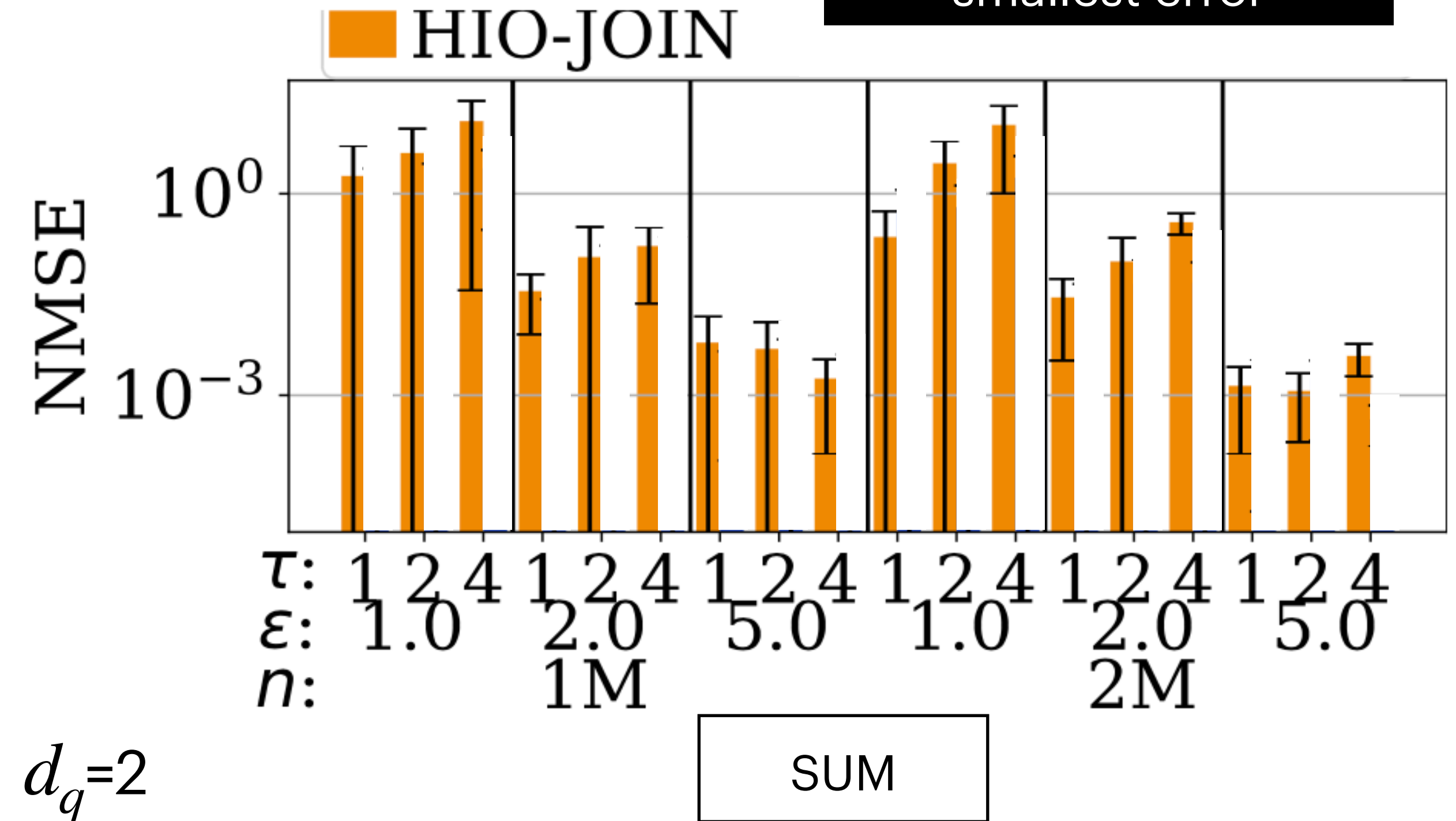
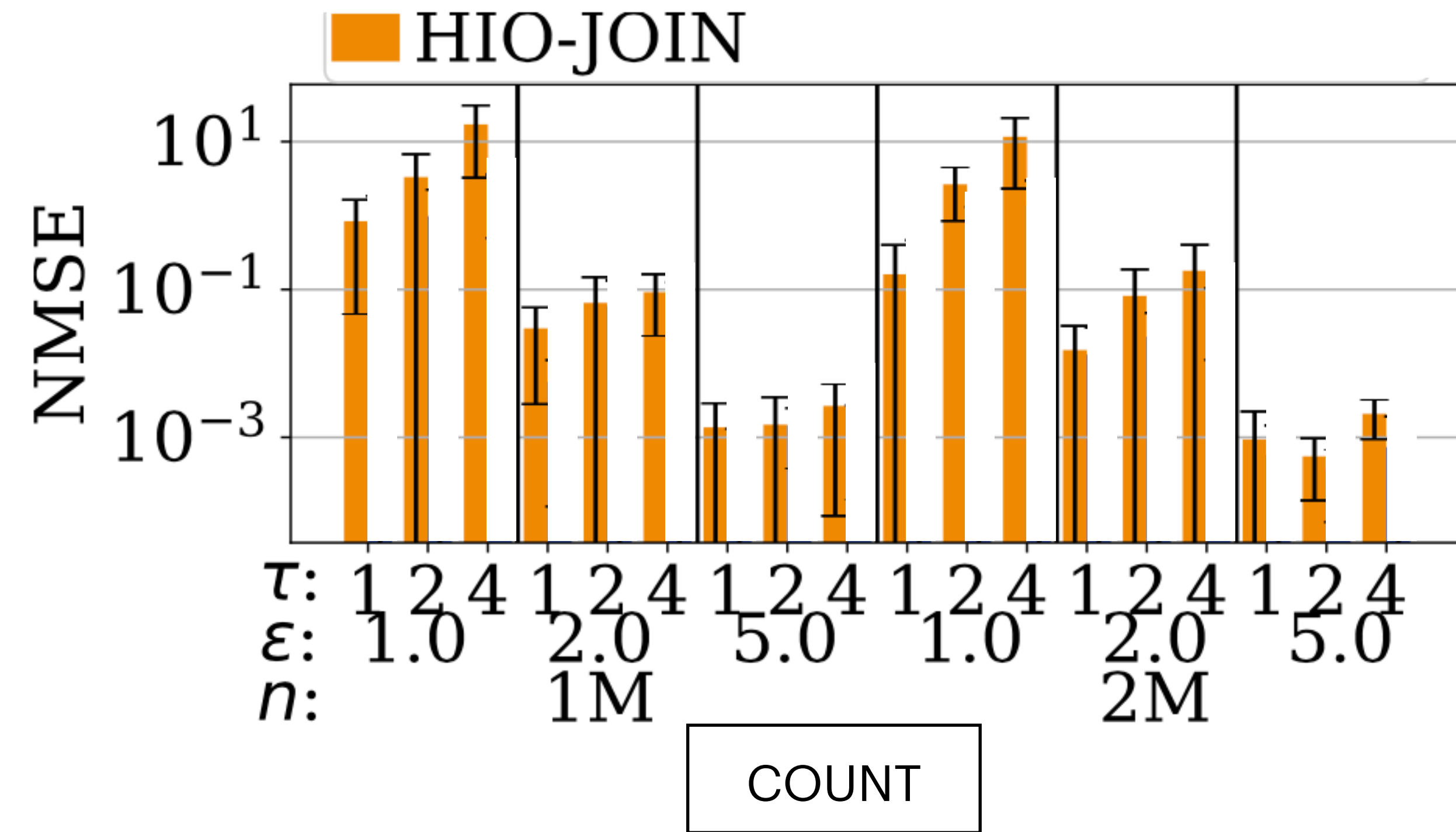
Asymptotic Mean Squared Error

$$\text{SC-JOIN: } O\left(\frac{n \log^{3Kd_q} m}{(\epsilon/Kd)^{2Kd_q}}\right), \text{ HIO-JOIN: } O\left(\frac{n \log^{K(d+d_q)} m}{(\epsilon/K)^{2K}}\right)$$

As range predicate gets complicated, i.e., d_q increases, error of SC-JOIN increases faster than that of HIO-JOIN

1-Many Joint Aggregation

$\tau = 1$ leads to the smallest error



$d_q=2$

Asymptotic Mean Squared Error of HIO-JOIN (with τ -truncation):

$$\text{COUNT: } O\left(\frac{nd\tau(1+\tau)^4 \log^{2(d+d_q)} m}{\epsilon^4} (r_{\min}^2 + r_{\max}^2)\right), \text{ SUM: } O\left(\frac{nd\tau(1+\tau)^4 \log^{2(d+d_q)} m}{\epsilon^4} (r_{\min}^2 + r_{\max}^2) (A_{\min}^2 + A_{\max}^2)\right)$$

More in the Paper

- Efficient range utility optimization based on range decomposition and consistency
- Support for group-by aggregations
- More experiments using real-world PUMS datasets and other tested schemes
- Alternative instantiation of the Partition-Rounding-Perturb framework by embedding the rounding value in the partition attribute

Thank you!
Q & A