# Partial Data Extraction via Noisy Histogram Query: The Information Theoretic Bounds

Wei-Ning Chen, joint work with Prof. I-Hsiang Wang
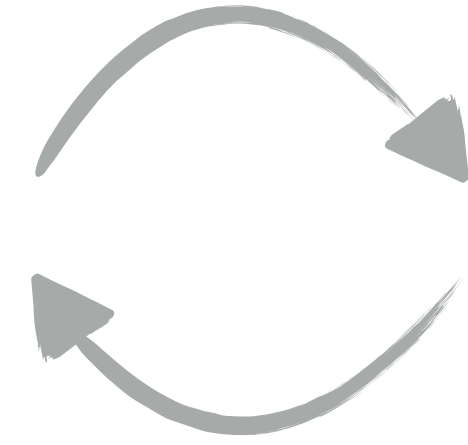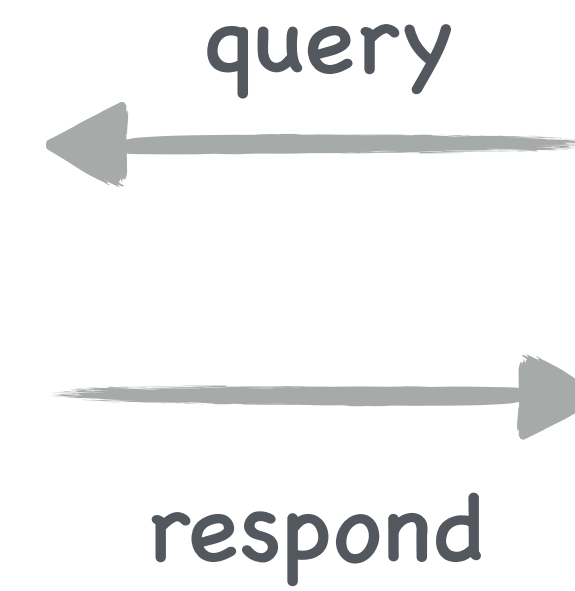National Taiwan University

Jun, 2017

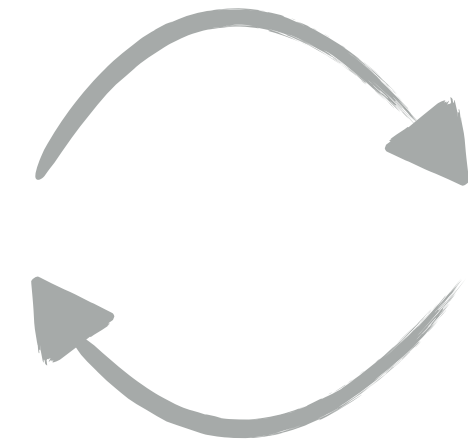# Query Model



Data set        Curator        Data Analyst

query

respond

- Query with the curator

# Query Model



**Data set**       **Curator**      **Data Analyst**

- Query with the curator
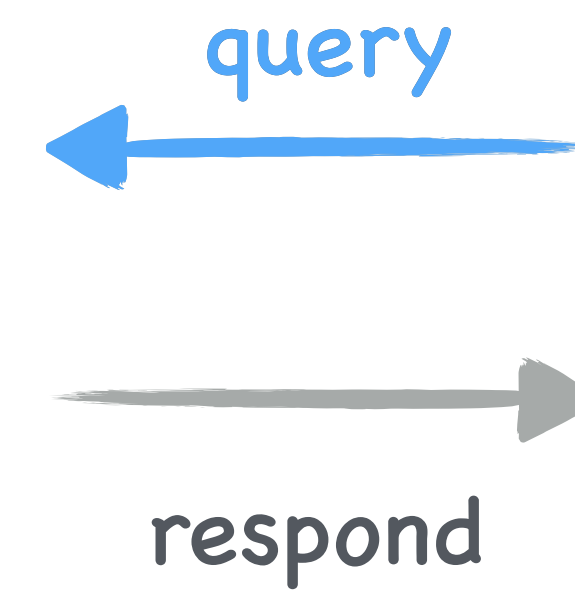
# Query Model



Data set                    Curator                    Data Analyst

- Query with the curator

# Query Model



Data set  ⟲  Curator  ←query  respond→  Data Analyst
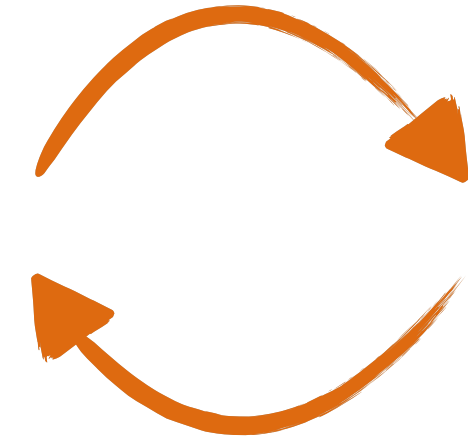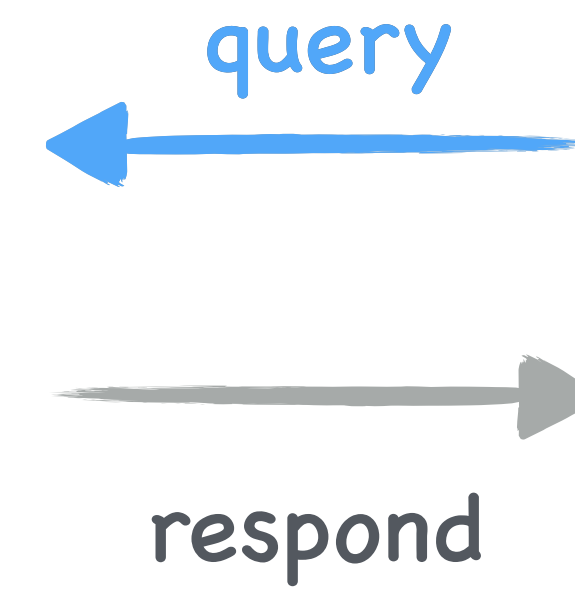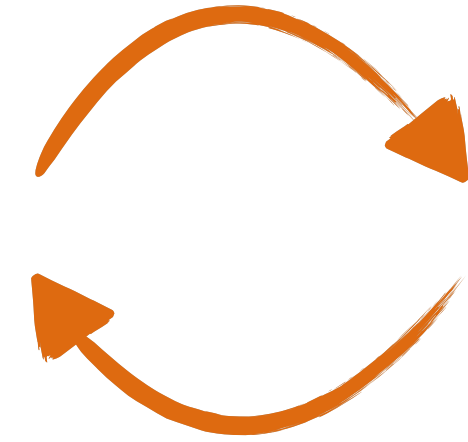
- Query with the curator

# Query Model



Data set             Curator             Data Analyst

- Query with the curator

- Certain types of queries are allowed
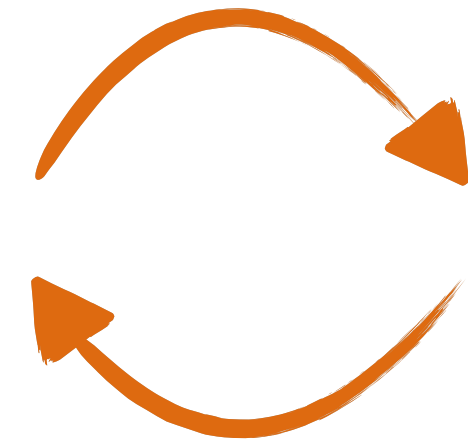
# Query Model



Data set             Curator             Data Analyst

- Query with the curator
- Certain types of queries are allowed
  - ▷ Subset query

# Query Model



Data set          Curator          Data Analyst

- Query with the curator

- Certain types of queries are allowed

    ▷ Subset query

    ▷ Statistical information of subset

# Query Model



**Data set**  **Curator**  **Data Analyst**

query

respond

- Query with the curator

- Certain types of queries are allowed

  ▷ Subset query

  ▷ Statistical information of subset
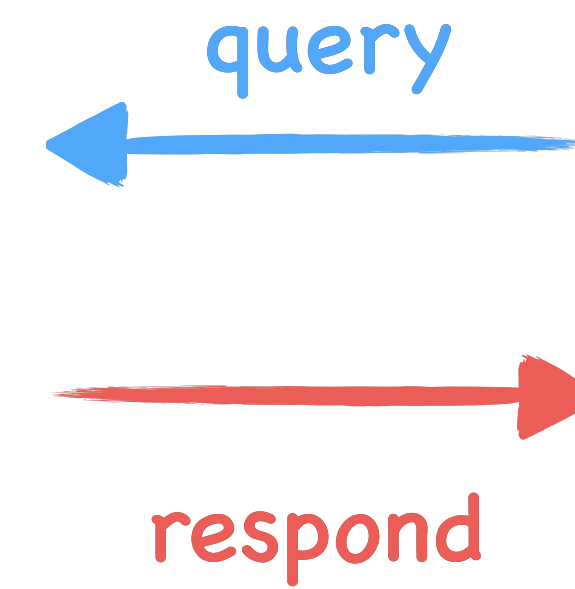
- Example :

# Query Model



Data set             Curator             Data Analyst

- Query with the curator
- Certain types of queries are allowed
  - ▷ Subset query
  - ▷ Statistical information of subset
- Example :

     A. Numerical data : statistical mean, variance etc.

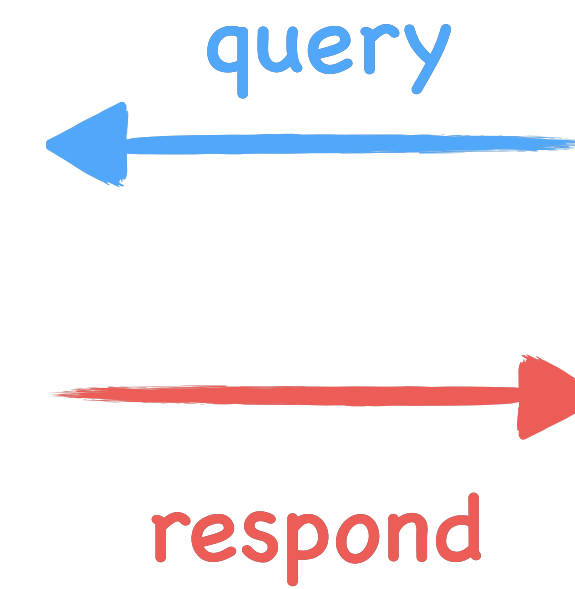     B. Categorical data : counting number, histogram etc.

# Query Model



Data set                    Curator                    Data Analyst

- Query with the curator

- Certain types of queries are allowed
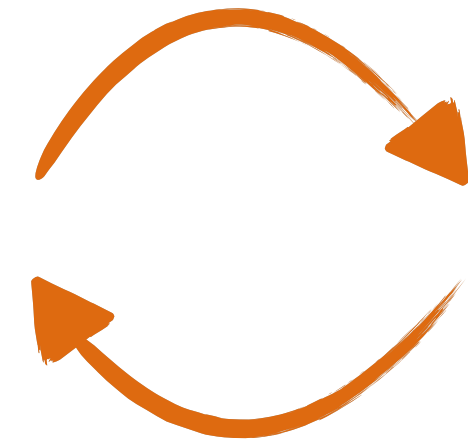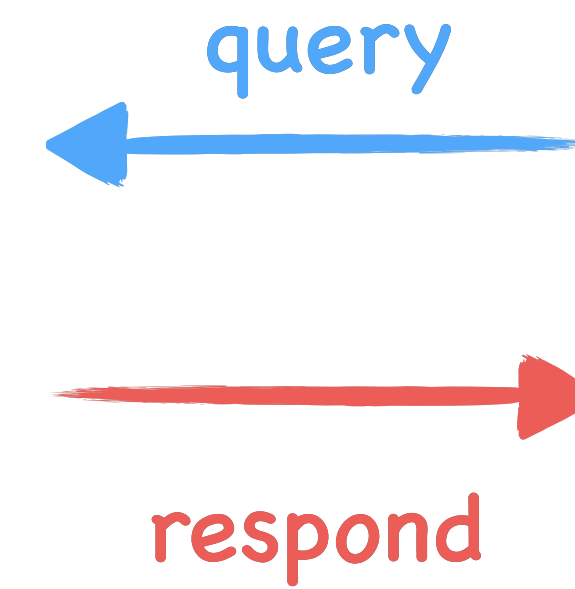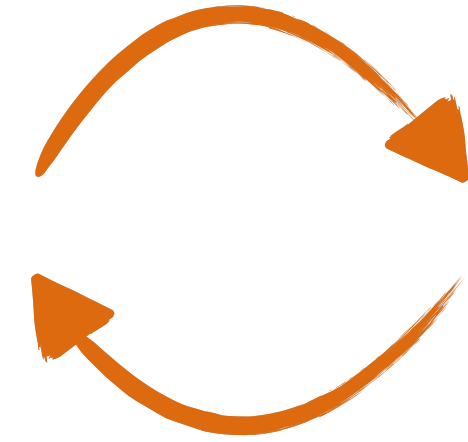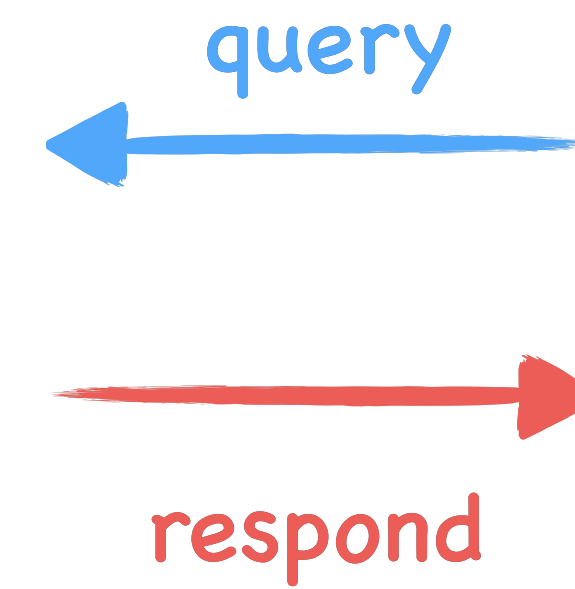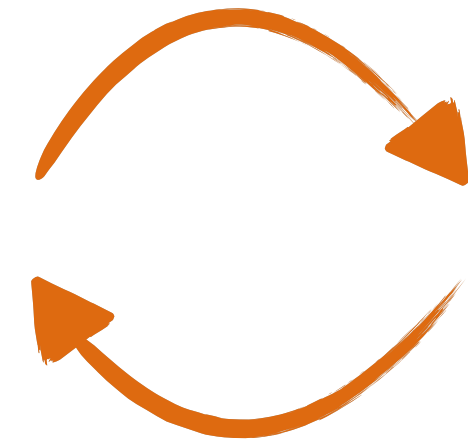
  ▷ Subset quer

  ▷ Statistical information of subset

- Example :

  A. ~~Numerical data : statistical mean, variance etc.~~

  B. Categorical data : counting number, histogram etc.

# Histogram Query

- Histogram Query

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| 4 | AB |
| 5 | O |
| 6 | O |

- Histogram Query

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| 4 | AB |
| 5 | O |
| 6 | O |

User{1,2,3,4}

⬅

- Histogram Query

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| 4 | AB |
| 5 | O |
| 6 | O |

User{1,2,3,4}

- Histogram Query

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| 4 | AB |
| 5 | O |
| 6 | O |

User{1,2,3,4}

A x 2, B x 1, AB x 1

# Histogram Query

- Histogram Query

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| 4 | AB |
| 5 | O |
| 6 | O |

User{1,2,3,4}

⟵

⟶

A x 2, B x 1, AB x 1

# Histogram Query

- Histogram Query

| Users | Blood |
|:-----:|:-----:|
| 1 | A |
| 2 | A |
| 3 | B |
| 4 | AB |
| 5 | O |
| 6 | O |

User{1,2,3,4}

A x 2, B x 1, AB x 1

A  B  O  AB

# The Noisy Response

- Histogram Query

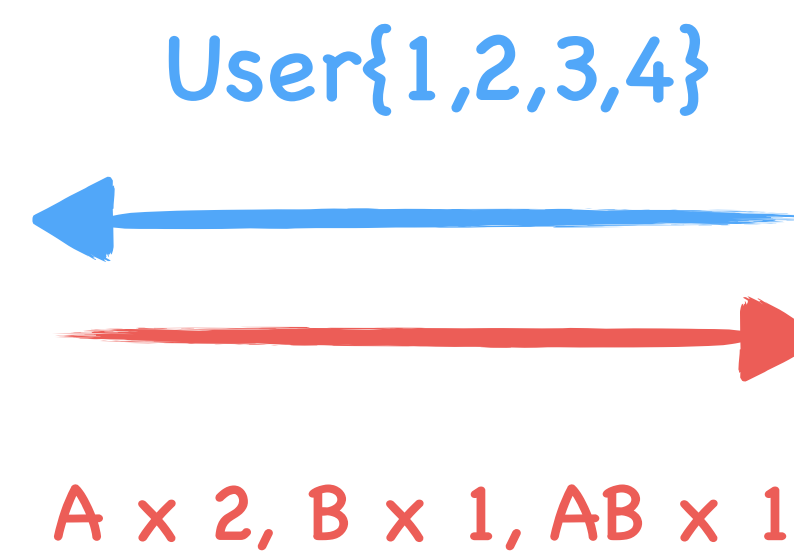| Users | Blood |
|:-----:|:-----:|
| 1 | A |
| 2 | A |
| 3 | B |
| ⋮ | ⋮ |
| n | O |

- Histogram Query

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| ⋮ | ⋮ |
| n | O |

# The Noisy Response

- Histogram Query

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| ⋮ | ⋮ |
| n | O |

Honest Response



A     B     O     AB

- Histogram Query

- Noisy response : ex. to guarantee stronger privacy

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| $\vdots$ | $\vdots$ |
| n | O |

Honest Response



A　　　B　　　O　　　AB

- Histogram Query

- Noisy response : ex. to guarantee stronger privacy

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| ⋮ | ⋮ |
| n | O |

Honest Response

Perturbed Response

A  B  O  AB

A  B  O  AB

- Histogram Query

- Noisy response : ex. to guarantee stronger privacy

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| ⋮ | ⋮ |
| n | O |

Honest Response

Perturbed Response

A    B    O    AB

A    B    O    AB

- Histogram Query

- Noisy response : ex. to guarantee stronger privacy

define the maximum difference as the noise level

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| ⋮ | ⋮ |
| n | O |

Honest Response

Perturbed Response

A   B   O   AB

A   B   O   AB

# The Noisy Response

- Histogram Query

- Noisy response : ex. to guarantee stronger privacy

- The added noise is at most $\delta_n$

| Users | Blood |
|-------|-------|
| 1 | A |
| 2 | A |
| 3 | B |
| ⋮ | ⋮ |
| n | O |

define the maximum difference as the noise level

Honest Response

Perturbed Response

A   B   O   AB

A   B   O   AB

# Problem Statement



Data set                Curator                Data Analyst

- Goal : to extract the data set partially

[1] I.-H. Wang, et. al "Data extraction via histogram and arithmetic mean queries: Fundamental limits and algorithms," Proceedings of IEEE International Symposium on Information Theory 2016

[2] Ahmed El Alaoui , et. al "Decoding from Pooled Data: Phase Transitions of Message Passing ," Proceedings of IEEE International Symposium on Information Theory 2017

# Problem Statement



Data set           Curator           Data Analyst

- Goal : to extract the data set partially
  - ▷ motivation: privacy, cost of data extraction, etc.

[1] I.-H. Wang, et. al "Data extraction via histogram and arithmetic mean queries: Fundamental limits and algorithms," Proceedings of IEEE International Symposium on Information Theory 2016

[2] Ahmed El Alaoui , et. al "Decoding from Pooled Data: Phase Transitions of Message Passing ," Proceedings of IEEE International Symposium on Information Theory 2017

Data set          Curator          Data Analyst

- Goal : to extract the data set partially
  ▷ motivation: privacy, cost of data extraction, etc.

- Key question : how many queries does the analyst required ?

[1] I.-H. Wang, et. al "Data extraction via histogram and arithmetic mean queries: Fundamental limits and algorithms," Proceedings of IEEE International Symposium on Information Theory 2016

[2] Ahmed El Alaoui , et. al "Decoding from Pooled Data: Phase Transitions of Message Passing ," Proceedings of IEEE International Symposium on Information Theory 2017
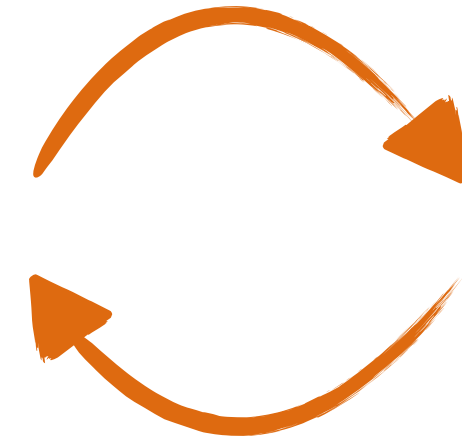
# Problem Statement



Data set          Curator          Data Analyst

- Goal : to extract the data set partially
  - motivation: privacy, cost of data extraction, etc.

- Key question : how many queries does the analyst required ?

- Query complexity : minimum number of queries to reconstruct the data set

---

[1] I.-H. Wang, et. al "Data extraction via histogram and arithmetic mean queries: Fundamental limits and algorithms," Proceedings of IEEE International Symposium on Information Theory 2016

[2] Ahmed El Alaoui , et. al "Decoding from Pooled Data: Phase Transitions of Message Passing ," Proceedings of IEEE International Symposium on Information Theory 2017
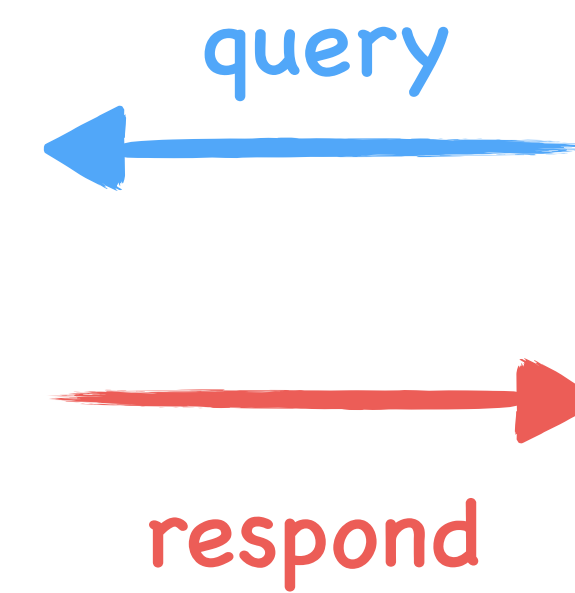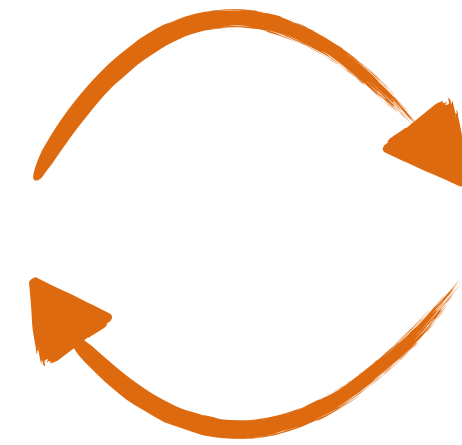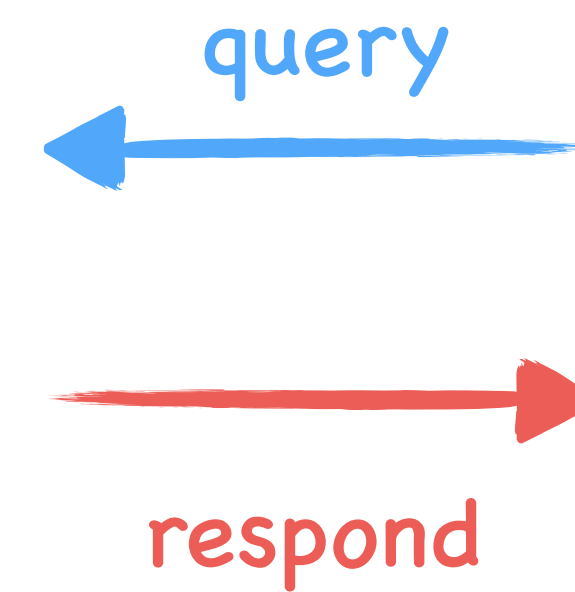
# Problem Statement



Data set                            Curator                          Data Analyst

- Goal : to extract the data set partially
  - ▷ motivation: privacy, cost of data extraction, etc.

- Key question : how many queries does the analyst required ?

- Query complexity : minimum number of queries to reconstruct the data set

- In noiseless case, i.e. $\delta_n = 0$, the query complexity in [1] is proven to be $\Theta\left(\dfrac{n}{\log n}\right)$
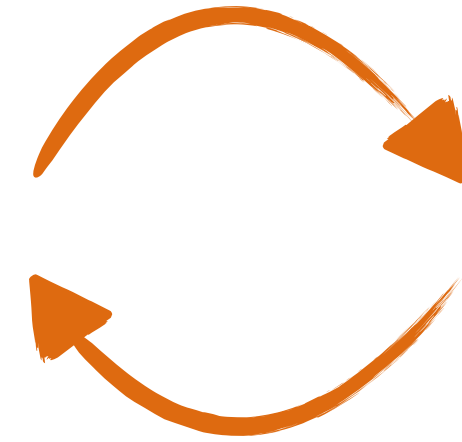  Also, in [2], an AMP algorithm is proposed to decode the data set

[1] I.-H. Wang, et. al "Data extraction via histogram and arithmetic mean queries: Fundamental limits and algorithms," Proceedings of IEEE International Symposium on Information Theory 2016

[2] Ahmed El Alaoui , et. al "Decoding from Pooled Data: Phase Transitions of Message Passing ," Proceedings of IEEE International Symposium on Information Theory 2017
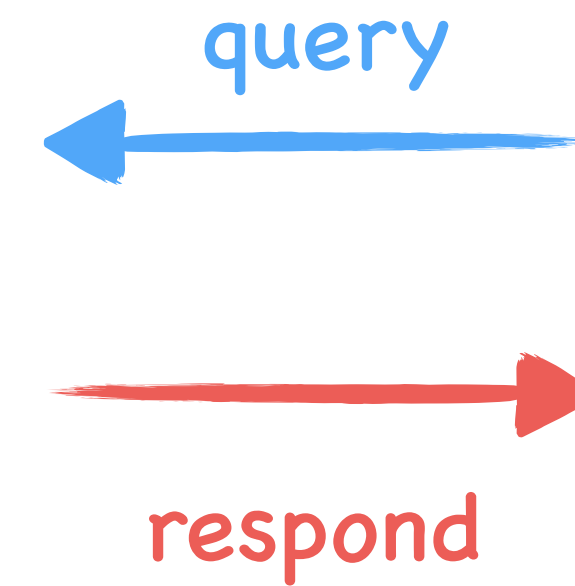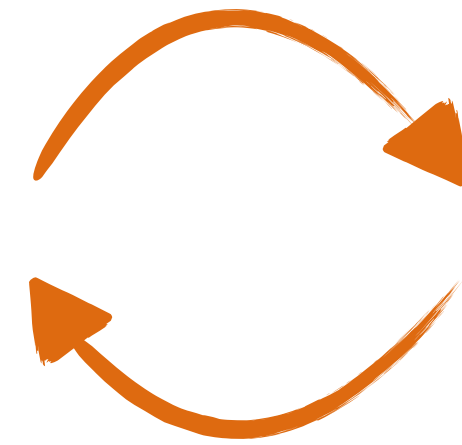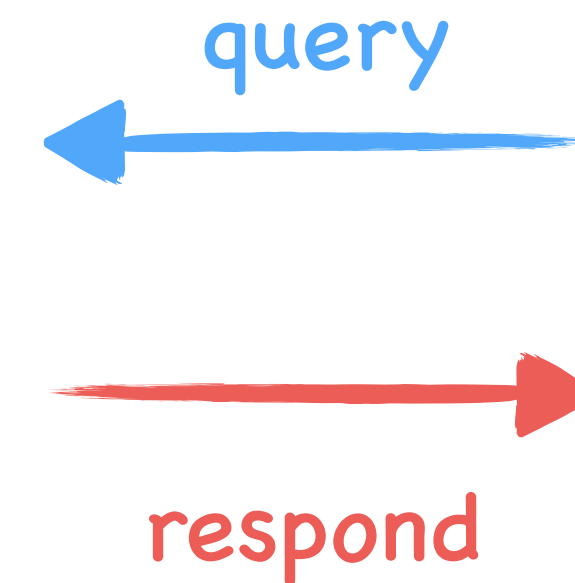
# Partial Data Reconstruction

- $\mathbf{x}$ : original data set

  $\hat{\mathbf{x}}$ : recovered data set

- $k_n$-distortion : $d_{\mathrm{Hamming}}(\mathbf{x}, \hat{\mathbf{x}}) \leq k_n$

$$
\mathbf{x} \qquad \hat{\mathbf{x}}
$$

$$
\begin{bmatrix} A \\ B \\ A \\ O \\ AB \\ \vdots \\ O \end{bmatrix}
\begin{bmatrix} {\color{red}B} \\ B \\ A \\ {\color{red}AB} \\ AB \\ \vdots \\ O \end{bmatrix}
$$

$$\delta_n = \Theta\left(n^d\right), \ k_n = \Theta\left(n^\kappa\right)$$

# Main Result

$$\delta_n = \Theta\left(n^d\right), \; k_n = \Theta\left(n^\kappa\right)$$

query complexity: non-polynomial



$$d > \left(\frac{1}{2} + \epsilon\right)\kappa$$

$$d < \frac{1}{2}\kappa$$

d

0.5

0.25

0

κ (distortion)

0.5

1

# Main Result

$$\boxed{\delta_n = \Theta\left(n^d\right), \, k_n = \Theta\left(n^\kappa\right)}$$

query complexity: non-polynomial



$$d > \left(\frac{1}{2} + \epsilon\right)\kappa$$

$$d < \frac{1}{2}\kappa$$

d

κ (distortion)

query complexity: sub-linear
(same as the noiseless case ! )

- Problem Formulation

  - Data extraction as a linear inverse problem

- Sketch of Proof :

  A. Regime 1 : Impossibility of Poly-n Query

  B. Regime 2 : The Fundamental Limit of Query Complexity

- Summary

# Histogram Query as Linear Multiplication



Data set

Curator

Data Analyst

query

respond

| Users | Blood |
|-------|-------|
| 1 | O |
| 2 | B |
| 3 | B |
| ⋮ | ⋮ |
| n | A |

User{1,2,n}

# Histogram Query as Linear Multiplication
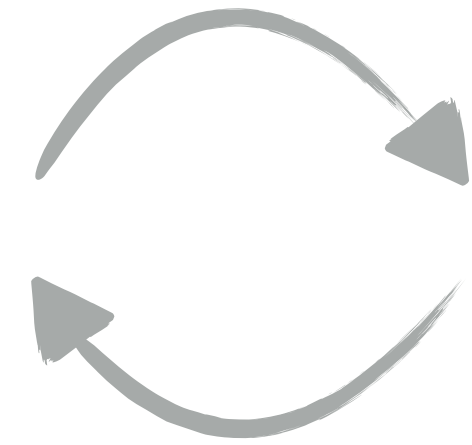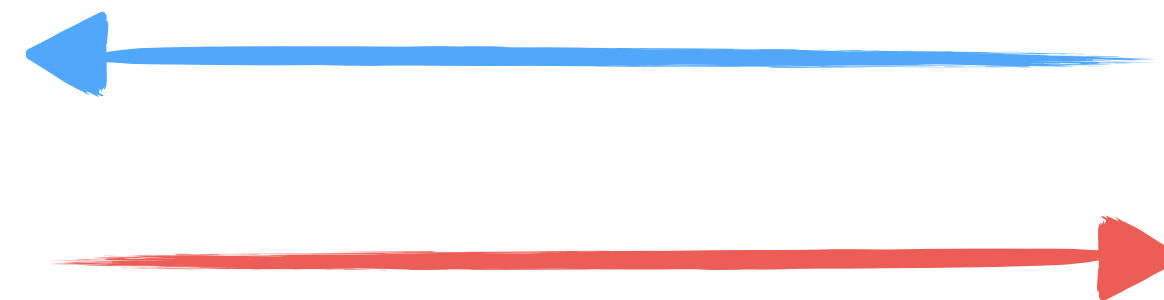


Data set   Curator   Data Analyst

query

respond

**A, B, AB, O**

User{1,2,n}

$$n \left\{ \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 0 \end{bmatrix} \right.$$

**X**

# Histogram Query as Linear Multiplication



**Data set**

**Curator**

**Data Analyst**

query

respond

A, B, AB, O

$$n \begin{cases} \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 0 \end{bmatrix} \end{cases}$$

**X**

User{1,2,n}

$$\boldsymbol{q}_i^{\mathsf{T}} = \underbrace{\begin{bmatrix} 1 & 1 & 0 & \cdots & 0 & 1 \end{bmatrix}}_{n}$$

# Histogram Query as Linear Multiplication



Data set

Curator

query

respond

Data Analyst

**A, B, AB, O**

$$n\left\{\begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 0 \end{bmatrix}\right.$$

$$\mathbf{X}$$

User{1,2,n}

$$\boldsymbol{q}_i^{\mathsf{T}} = \underbrace{\begin{bmatrix} 1 & 1 & 0 & \cdots & 0 & 1 \end{bmatrix}}_{n}$$

$$\boldsymbol{y}_i = \boldsymbol{q}_i^{\mathsf{T}}\mathbf{X}$$

# Histogram Query as Linear Multiplication
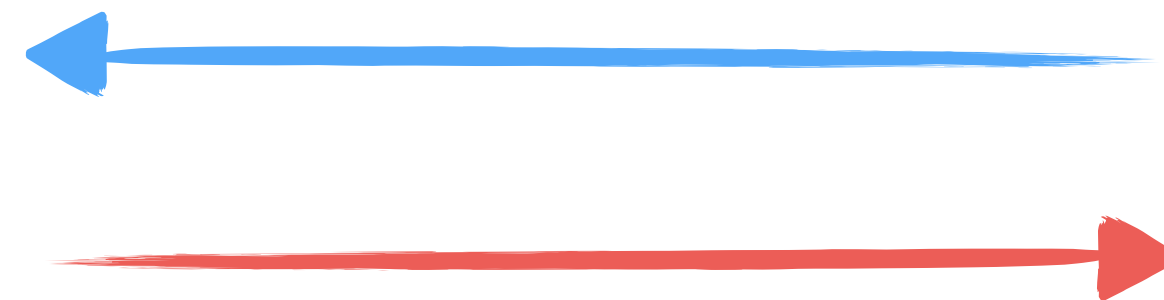
**Data set**

**Curator**

query

respond

**Data Analyst**

**A, B, AB, O**

$$n \left\{ \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 0 \end{bmatrix} \right.$$
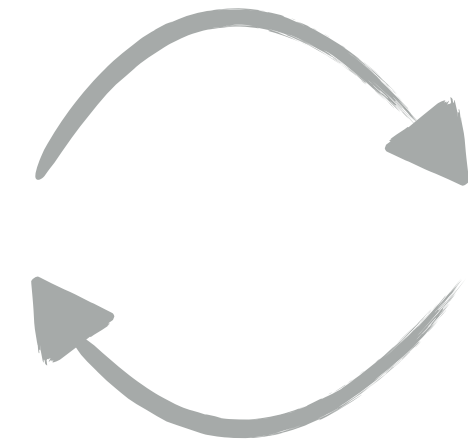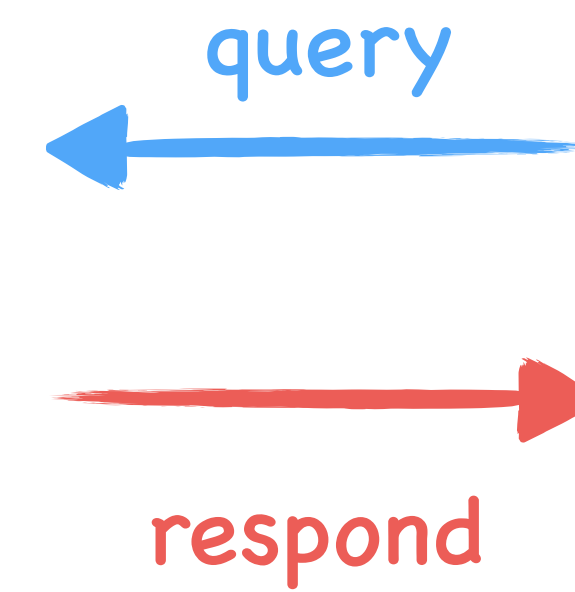
$$\mathbf{X}$$

User{1,2,n}

$$\boldsymbol{q}_i^{\mathsf{T}} = \underbrace{\begin{bmatrix} 1 & 1 & 0 & \cdots & 0 & 1 \end{bmatrix}}_{n}$$

$$\boldsymbol{y}_i = \boldsymbol{q}_i^{\mathsf{T}} \mathbf{X} + \Delta_i$$

# Histogram Query as Linear Multiplication



Data set

Curator

query

respond

Data Analyst
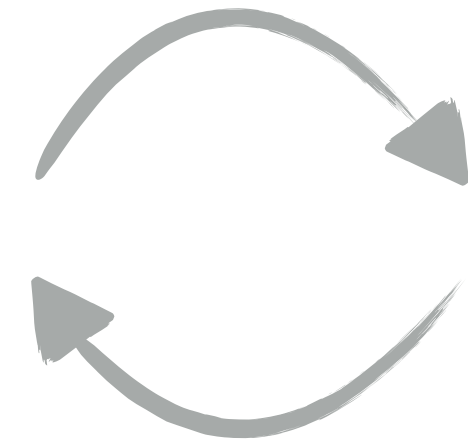
**A, B, AB, O**

$$n \left\{ \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 0 \end{bmatrix} \right.$$

$$\mathbf{X}$$

**User{1,2,n}**

$$\boldsymbol{q}_i^\mathsf{T} = \underbrace{\begin{bmatrix} 1 & 1 & 0 & \cdots & 0 & 1 \end{bmatrix}}_{n}$$

$$\boldsymbol{y}_i = \boldsymbol{q}_i^\mathsf{T} \mathbf{X} + \Delta_i$$

**A, B, O, AB**

$$\mathbf{y}_i = [10, 20, 18, 28]$$

A   B   O   AB

# Histogram Query as Linear Multiplication



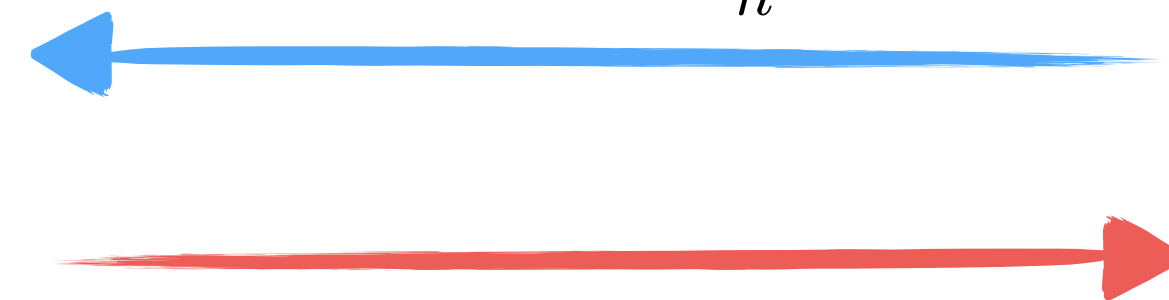**Data set**

**Curator**

**Data Analyst**

query

respond

A, B, AB, O

$$\mathbf{X}$$

$$n \begin{cases} \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 0 \end{bmatrix} \end{cases}$$

User{1,2,n}

$$\boldsymbol{q}_i^\mathsf{T} = \underbrace{\begin{bmatrix} 1 & 1 & 0 & \cdots & 0 & 1 \end{bmatrix}}_{n}$$

$$\boldsymbol{y}_i = \boldsymbol{q}_i^\mathsf{T} \mathbf{X} + \Delta_i$$

A, B, O, AB

$$\mathbf{y}_i = [10, 20, 18, 28]$$

A   B   O   AB

# Histogram Query as Linear Multiplication



Data set

Curator

Data Analyst

query

respond

**A, B, AB, O**

$$n \left\{ \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 0 \end{bmatrix} \right.$$

$\mathbf{X}$

**User{1,2,n}**

$$\boldsymbol{q}_i^\mathsf{T} = \underbrace{\begin{bmatrix} 1 & 1 & 0 & \cdots & 0 & 1 \end{bmatrix}}_{n}$$

$$\boldsymbol{y}_i = \boldsymbol{q}_i^\mathsf{T} \mathbf{X} + \Delta_i$$

**A, B, O, AB**

$$\mathbf{y}_i = [10, 20, 18, 28]$$

A   B   O   AB

# Histogram Query as Linear Multiplication



**Data set**

Decode column by column

A, B, AB, O

$$n \left\{ \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 0 \end{bmatrix} \right.$$

$\mathbf{X}$

**Curator**
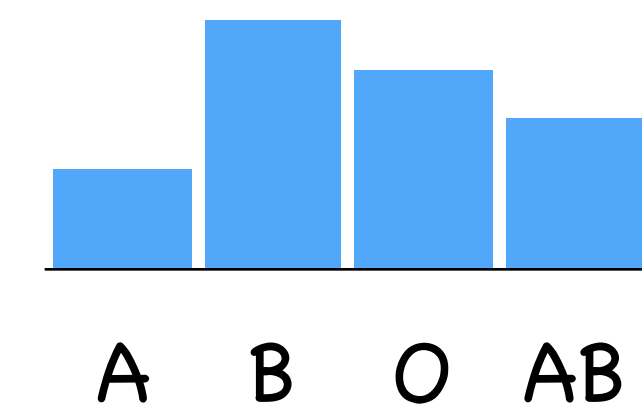
query

respond

User{1,2,n}

$$\boldsymbol{q}_i^{\mathsf{T}} = \underbrace{\begin{bmatrix} 1 & 1 & 0 & \cdots & 0 & 1 \end{bmatrix}}_{n}$$

$$\boldsymbol{y}_i = \boldsymbol{q}_i^{\mathsf{T}} \mathbf{X} + \Delta_i$$
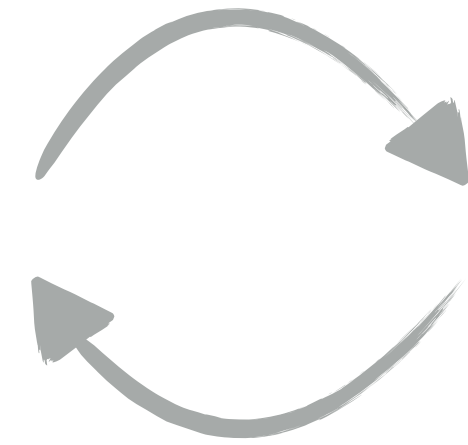
**Data Analyst**

A, B, O, AB

$$\mathbf{y}_i = [10, 20, 18, 28]$$

A  B  O  AB

# Histogram Query as Linear Multiplication

**Data set**

**Curator**

**Data Analyst**

Decode column by column

**A, B, AB, O**

$$n \begin{cases} \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 0 \end{bmatrix} \end{cases}$$

$$\mathbf{X}$$

**User{1,2,n}**

$$\boldsymbol{q}_i^\mathsf{T} = \underbrace{\begin{bmatrix} 1 & 1 & 0 & \cdots & 0 & 1 \end{bmatrix}}_{n}$$

$$\boldsymbol{y}_i = \boldsymbol{q}_i^\mathsf{T} \mathbf{X} + \Delta_i$$

**A, B, O, AB**

$$\mathbf{y}_i = [10, 20, 18, 28]$$

A  B  O  AB

**Data set**

**Curator**

query

respond

**Data Analyst**

$$\mathbf{x}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 1 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$\boldsymbol{q}_i^{\mathsf{T}} = \underbrace{\begin{bmatrix} 1 & 1 & 0 & \cdots & 0 & 1 \end{bmatrix}}_{n}$$

$$y_i : \# \text{ of } 1 \text{ in } \mathbf{x}$$

$$y_i = \boldsymbol{q}_i^{\mathsf{T}} \mathbf{x} + \boldsymbol{\Delta}$$

# Histogram Query as Linear Multiplication



**Data set**

$\mathbf{x}$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 1 \\ \vdots \\ \vdots \\ 0 \end{bmatrix}$$

**Curator**

$$\mathbf{Q} = \begin{bmatrix} \boldsymbol{q}_1^\mathsf{T} \\ \boldsymbol{q}_2^\mathsf{T} \\ \vdots \\ \boldsymbol{q}_{T_n}^\mathsf{T} \end{bmatrix}$$

$$\boldsymbol{y} = \mathbf{Q}\mathbf{x} + \boldsymbol{\Delta}$$

query

respond

**Data Analyst**

$$\boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{T_n} \end{bmatrix} + \begin{bmatrix} \Delta_1 \\ \Delta_2 \\ \vdots \\ \Delta_{T_n} \end{bmatrix}$$

- Equivalent linear inverse problem :

- Equivalent linear inverse problem :

  ▷ Given an output $y = \mathbf{Q}\mathbf{x} + \mathbf{\Delta}$

  ▷ Find the corresponding data set $\hat{\mathbf{x}} : \|\mathbf{Q}\hat{\mathbf{x}} - y\|_\infty \leq \delta_n$, and $d_{\mathrm{Hamming}}(\mathbf{x}, \hat{\mathbf{x}}) \leq k_n$

- Equivalent linear inverse problem :

  ▷ Given an output $y = \mathbf{Q}\mathbf{x} + \boldsymbol{\Delta}$

  ▷ Find the corresponding data set $\hat{\mathbf{x}} : \|\mathbf{Q}\hat{\mathbf{x}} - y\|_\infty \leq \delta_n$, and $d_{\mathrm{Hamming}}(\mathbf{x}, \hat{\mathbf{x}}) \leq k_n$

- The noise level is $\delta_n$, if the difference in each single query is at most $\delta_n$

- Equivalent linear inverse problem :

  ▷ Given an output $y = \mathbf{Q}\mathbf{x} + \boldsymbol{\Delta}$

  ▷ Find the corresponding data set $\hat{\mathbf{x}} : \|\mathbf{Q}\hat{\mathbf{x}} - \boldsymbol{y}\|_\infty \leq \delta_n$, and $d_{\mathrm{Hamming}}(\mathbf{x}, \hat{\mathbf{x}}) \leq k_n$

- The noise level is $\delta_n$, if the difference in each single query is at most $\delta_n \iff \|\boldsymbol{\Delta}\|_\infty \leq \delta_n (\iff \forall i, \ \Delta_i \leq \delta_n)$

- Equivalent linear inverse problem :

  ▷ Given an output $y = \mathbf{Q}\mathbf{x} + \boldsymbol{\Delta}$

  ▷ Find the corresponding data set $\hat{\mathbf{x}} : \|\mathbf{Q}\hat{\mathbf{x}} - y\|_\infty \leq \delta_n$, and $d_{\mathrm{Hamming}}(\mathbf{x}, \hat{\mathbf{x}}) \leq k_n$

- The noise level is $\delta_n$, if the difference in each single query is at most $\delta_n \iff \|\boldsymbol{\Delta}\|_\infty \leq \delta_n (\iff \forall i, \Delta_i \leq \delta_n)$

- Query complexity $T_n^*(k_n, \delta_n)$ : minimum number of queries required to extract data set within distortion $k_n$, under noise level $\delta_n$

$$\delta_n = \Theta\left(n^d\right),\ k_n = \Theta\left(n^\kappa\right)$$



$$d > \left(\frac{1}{2} + \epsilon\right)\kappa$$

$$d < \frac{1}{2}\kappa$$

d

κ (distortion)

$$\delta_n = \Theta\left(n^d\right),\ k_n = \Theta\left(n^\kappa\right)$$

query complexity: non-polynomial $\Omega\left(\exp\left(n^\epsilon\right)\right)$

$$d > \left(\frac{1}{2} + \epsilon\right)\kappa$$

$$d < \frac{1}{2}\kappa$$

d

κ (distortion)

# Main Result



$\delta_n = \Theta\left(n^d\right),\ k_n = \Theta\left(n^\kappa\right)$

query complexity: non-polynomial $\Omega\left(\exp\left(n^\epsilon\right)\right)$

$d > \left(\dfrac{1}{2} + \epsilon\right)\kappa$

$d < \dfrac{1}{2}\kappa$

query complexity: sub-linear $\Theta\left(\dfrac{n}{\log n}\right)$

d

κ (distortion)

0.5

0.25

0

0.5

1

- Regime 1: $d > (\frac{1}{2} + \epsilon)\kappa$, for any $\epsilon > 0$ (the noise is too large)

$$T_n^*(k_n, \delta_n) = \Omega\left(\exp\left(n^\epsilon\right)\right)$$

- Proof idea : without sufficient number of queries, there exists more than one possible data set which are consistent with the response.

$\mathbf{x}$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$S_{k_n}$

$$S_{k_n} \triangleq \{(\mathbf{x}, \tilde{\mathbf{x}}) \mid \mathbf{x}, \tilde{\mathbf{x}} \in \{0,1\}^n, \|\mathbf{x} - \tilde{\mathbf{x}}\|_1 = k_n, \|\mathbf{x}\|_1 = \|\tilde{\mathbf{x}}\|_1\}$$



$$S_{k_n}$$

confused set

$$\mathbf{x}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$\mathbf{x}$$

$$S_{k_n}$$

$$\boxed{\text{confused set}}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$\mathbf{q}_i^\mathsf{T} = \underbrace{\begin{bmatrix} 1 & 0 & 1 & \cdots & 0 & 0 \end{bmatrix}}_{n}$$

$$\mathbf{y}_i = \mathbf{q}_i^\mathsf{T}\mathbf{X} + \Delta$$

$$S_{k_n}$$

confused set

$$\mathbf{x}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$\mathbf{q}_i^\mathsf{T} = \begin{bmatrix} 1 & 0 & 1 & \cdots & 0 & 0 \end{bmatrix}$$

$$\underbrace{\qquad}_{n}$$

$$\mathbf{y}_i = \mathbf{q}_i^\mathsf{T}\mathbf{X} + \Delta$$

$$V_i \triangleq \left\{ (\mathbf{x}, \tilde{\mathbf{x}}) \in S_{k_n} \mid |\mathbf{q}_i \cdot (\mathbf{x} - \tilde{\mathbf{x}})| > \delta_n \right\}.$$

$$\mathbf{q}_i^{\mathsf{T}} = \underbrace{\begin{bmatrix} 1 & 0 & 1 & \cdots & 0 & 0 \end{bmatrix}}_{n}$$

$$\mathbf{y}_i = \mathbf{q}_i^{\mathsf{T}} \mathbf{X} + \Delta$$

$$V_i$$

$$S_{k_n}$$

confused set

$$\mathbf{x}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$\mathbf{q}_i^\intercal = \begin{bmatrix} 1 & 0 & 1 & \cdots & 0 & 0 \end{bmatrix}$$

$$\underbrace{\phantom{1 \quad 0 \quad 1 \quad \cdots \quad 0 \quad 0}}_{n}$$

$$\mathbf{y}_i = \mathbf{q}_i^\intercal \mathbf{X} + \Delta$$

$V_i$

$S_{k_n}$

confused set

$$\mathbf{x}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$\mathbf{q}_i^\mathsf{T} = \begin{bmatrix} 1 & 0 & 1 & \cdots & 0 & 0 \end{bmatrix}$$

$$n$$

$$\mathbf{y}_i = \mathbf{q}_i^\mathsf{T} \mathbf{X} + \Delta$$

$$V_i$$

$$S_{k_n}$$

confused set

$$\mathbf{x}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$\mathbf{q}_i^\intercal = \begin{bmatrix} 1 & 0 & 1 & \cdots & 0 & 0 \end{bmatrix}$$

$$n$$

$$\mathbf{y}_i = \mathbf{q}_i^\intercal \mathbf{X} + \Delta$$

$$S_{k_n}$$

$$V_i$$

confused set

$$\mathbf{x}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$\mathbf{q}_i^{\mathsf{T}} = \begin{bmatrix} 1 & 0 & 1 & \cdots & 0 & 0 \end{bmatrix}$$
$$\underbrace{\qquad\qquad}_{n}$$

$$\mathbf{y}_i = \mathbf{q}_i^{\mathsf{T}} \mathbf{X} + \Delta$$

$$\mathbf{x}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$V_i$

$S_{k_n}$

confused set

$$\mathbf{q}_i^\mathsf{T} = \underbrace{\begin{bmatrix} 1 & 0 & 1 & \cdots & 0 & 0 \end{bmatrix}}_{n}$$

$$\mathbf{y}_i = \mathbf{q}_i^\mathsf{T}\mathbf{X} + \Delta$$

$$\mathbf{x} \quad \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$V_i$

$S_{k_n}$

confused set

at least $\dfrac{|S_{k_n}|}{\max_i |V_i|}$ queries are required

$V_i$

$S_{k_n}$

confused set

$\mathbf{q}_i^\mathsf{T} = \begin{bmatrix} 1 & 0 & 1 & \cdots & 0 & 0 \end{bmatrix}$
$\underbrace{\phantom{1 \quad 0 \quad 1 \quad \cdots \quad 0 \quad 0}}_{n}$

$\mathbf{y}_i = \mathbf{q}_i^\mathsf{T} \mathbf{X} + \Delta$

$\mathbf{x}$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

- Therefore, we have the following lower bound on $T_n^*(k_n, \delta_n)$ :

$$T_n^*(k_n, \delta_n) \geq \frac{|S_{k_n}|}{\max_{i \in \{1, \dots, T_n\}} |V_i|}$$

- Therefore, we have the following lower bound on $T_n^*(k_n, \delta_n)$ :

$$T_n^*(k_n, \delta_n) \geq \frac{|S_{k_n}|}{\max_{i \in \{1, \ldots, T_n\}} |V_i|}$$

$$\underbrace{\phantom{\text{solve the optimization over V,}}}_{\substack{\text{solve the optimization over V,} \\ \text{and apply Chernoff ineq.}}} \geq C \exp\left(\frac{\delta_n^2}{k_n}\right) = C \exp\left(n^{2d - \kappa}\right)$$

- **Regime 2:** $d < \dfrac{1}{2}\kappa$ (the noise is small enough)

$$T_n^*(k_n, \delta_n) = \Theta\left(n/\log n\right)$$

- Random sampling is considered

- Random sampling is considered

- Evert item is included to the queried subset with probability $1/2$

- **Random sampling** is considered

- Evert item is included to the queried subset with probability $1/2$

$$\Longleftrightarrow (\mathbf{Q})_{i,j} \sim \mathrm{Ber}\left(\frac{1}{2}\right)$$

- **Random sampling** is considered

- Evert item is included to the queried subset with probability $1/2$

$$\Longleftrightarrow (\mathbf{Q})_{i,j} \sim \mathrm{Ber}\left(\frac{1}{2}\right)$$

- The probability of failure :

$$P_f\left(\mathbf{x}; k_n, \delta_n\right) \triangleq$$

$$P\left\{\exists \text{ a confused } \tilde{\mathbf{x}} \text{ which is consistent with the query output}\right\}$$

- <span style="color:red">Random sampling</span> is considered

- Evert item is included to the queried subset with probability $1/2$

$$\Longleftrightarrow (\mathbf{Q})_{i,j} \sim \mathrm{Ber}\left(\frac{1}{2}\right)$$

- The probability of failure :

$$P_f\left(\mathbf{x}; k_n, \delta_n\right) \triangleq$$

$$P\left\{\exists \text{ a confused } \tilde{\mathbf{x}} \text{ which is consistent with the query output}\right\}$$

- If number of queries is $\Omega(n/\log n)$ then $P_f(\mathbf{x}; k_n, \delta_n) \to 0$ as $n \to \infty$

- **Random sampling** is considered

- Evert item is included to the queried subset with probability $1/2$

$$\Longleftrightarrow (\mathbf{Q})_{i,j} \sim \mathrm{Ber}\left(\frac{1}{2}\right)$$

- The probability of failure :

$$P_f\left(\mathbf{x}; k_n, \delta_n\right) \triangleq$$

$$P\left\{\exists \text{ a confused } \tilde{\mathbf{x}} \text{ which is consistent with the query output}\right\}$$

- If number of queries is $\Omega(n/\log n)$ then $P_f(\mathbf{x}; k_n, \delta_n) \to 0$ as $n \to \infty$

  ▷ Applying Chernoff bound on failure event

- Necessary condition :

$$\forall \mathbf{x}, \tilde{\mathbf{x}} \in \mathcal{X}, \ \|\mathbf{x} - \tilde{\mathbf{x}}\|_1 > k_n \implies \|\mathbf{Q}\mathbf{x} - \mathbf{Q}\tilde{\mathbf{x}}\|_\infty > 2\delta_n$$

- Packing inequality :

$$2\delta_n\text{-packing number on } \mathcal{Y} \geq \tfrac{1}{2}k_n\text{-packing number on } \mathcal{X}$$

- Necessary condition :

$$\forall \mathbf{x}, \tilde{\mathbf{x}} \in \mathcal{X},\ \|\mathbf{x} - \tilde{\mathbf{x}}\|_1 > k_n \implies \|\mathbf{Q}\mathbf{x} - \mathbf{Q}\tilde{\mathbf{x}}\|_\infty > 2\delta_n$$

- Packing inequality :

$$2\delta_n\text{-packing number on } \mathcal{Y} \geq \tfrac{1}{2}k_n\text{-packing number on } \mathcal{X}$$
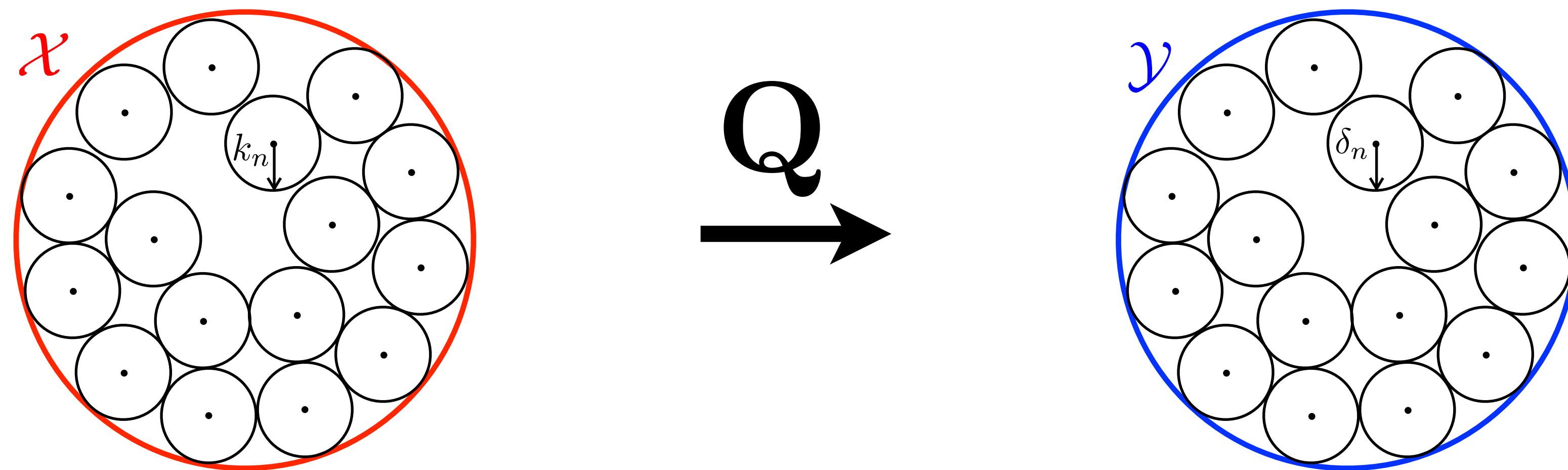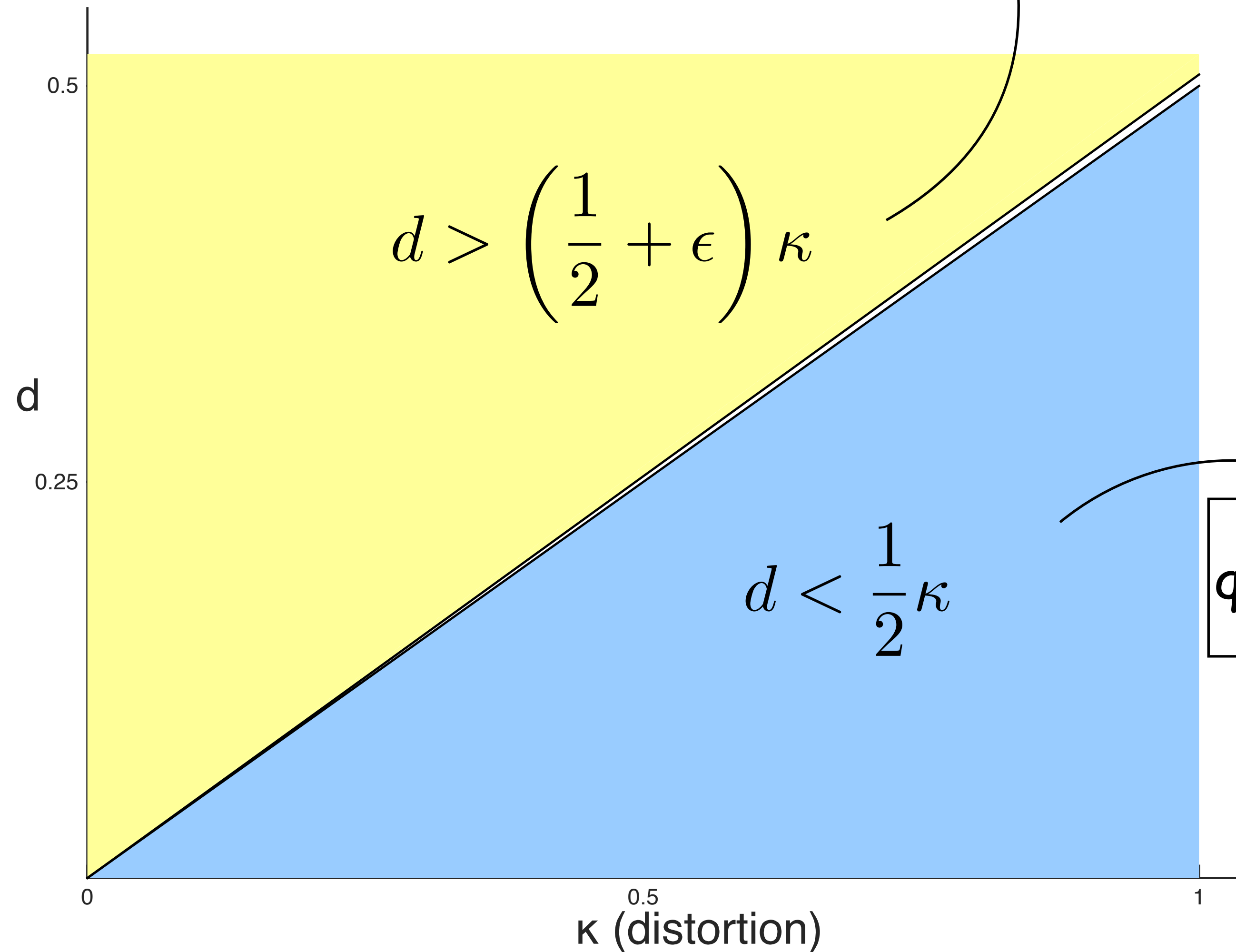
$$\delta_n = \Theta\left(n^d\right), \; k_n = \Theta\left(n^\kappa\right)$$

query complexity: non-polynomial $\Omega\left(\exp\left(n^\epsilon\right)\right)$

$$d > \left(\frac{1}{2} + \epsilon\right)\kappa$$

$$d < \frac{1}{2}\kappa$$

query complexity: sub-linear $\Theta\left(\dfrac{n}{\log n}\right)$

d

0.5

0.25

0

0.5

1

κ (distortion)

# Reference

[1] I.-H. Wang, S.-L. Huang et. al. "Data extraction via histogram and arithmetic mean queries: Fundamental limits and algorithms," *Proceedings of IEEE International Symposium on Information Theory*, July 2016.

[2] Ahmed El Alaoui , et. al "Decoding from Pooled Data: Phase Transitions of Message Passing ," *Proceedings of IEEE International Symposium on Information Theory,* June 2017

[3] C. Dwork, A. Roth, "The algorithmic foundations of differential privacy," *Theoretical Computer Science,* 2013

# Question ?

Thank you for your attention !