

Differentially-Private Batch Query Answering

Exploiting the Workload vs. Exploiting the Data

Gerome Miklau

University of Massachusetts, Amherst



Batch (non-interactive) query answering

- **Given:** a fixed set of queries
 - complex data analysis task into simpler queries.
 - multiple users each issuing one or more queries.
 - uncertainty about the eventual query answers needed--design workload to include all queries possibly of interest.

Batch (non-interactive) query answering

- **Given:** a fixed set of queries
 - complex data analysis task into simpler queries.
 - multiple users each issuing one or more queries.
 - uncertainty about the eventual query answers needed--design workload to include all queries possibly of interest.



Batch (non-interactive) query answering

- **Given:** a fixed set of queries
 - complex data analysis task into simpler queries.
 - multiple users each issuing one or more queries.
 - uncertainty about the eventual query answers needed--design workload to include all queries possibly of interest.
- **Goal:** release answers to all queries under ϵ - or (ϵ, δ) -differential privacy.

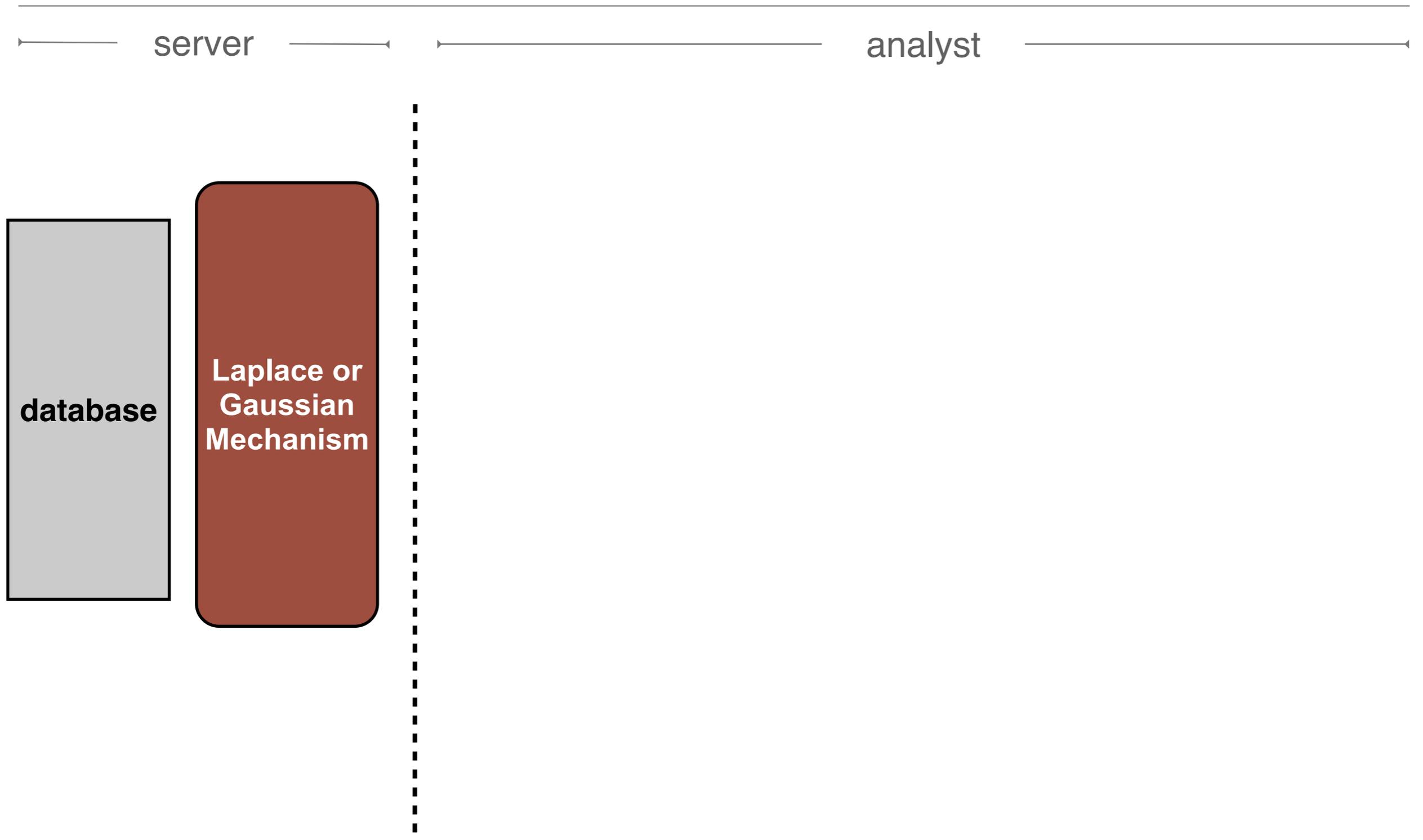


Batch (non-interactive) query answering

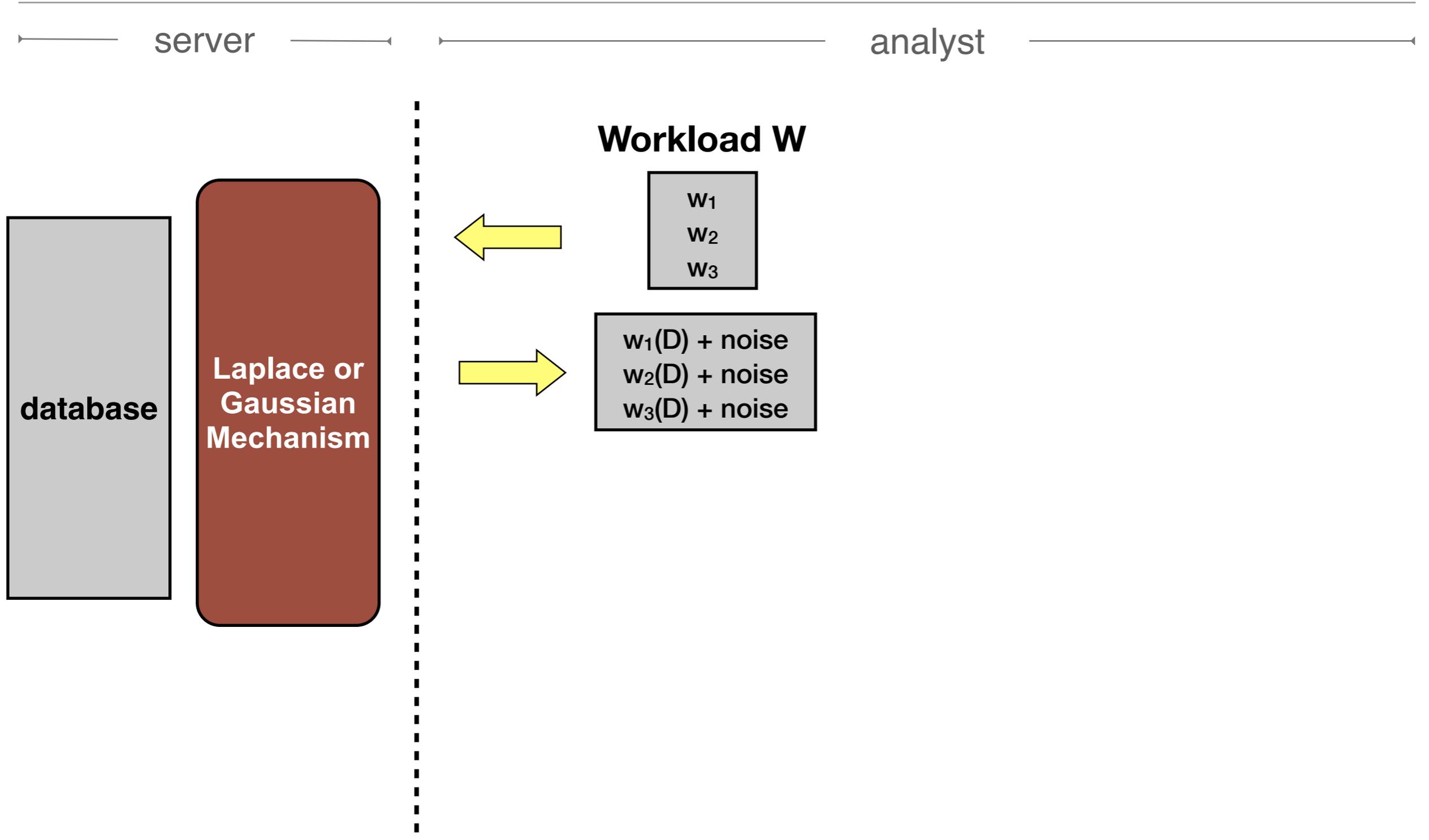
- **Given:** a fixed set of queries
 - complex data analysis task into simpler queries.
 - multiple users each issuing one or more queries.
 - uncertainty about the eventual query answers needed--design workload to include all queries possibly of interest.
- **Goal:** release answers to all queries under ϵ - or (ϵ, δ) -differential privacy.
- Linear counting queries
 - includes predicate counting queries, spatial queries, multi-dimensional range queries, marginals, data cubes, etc.



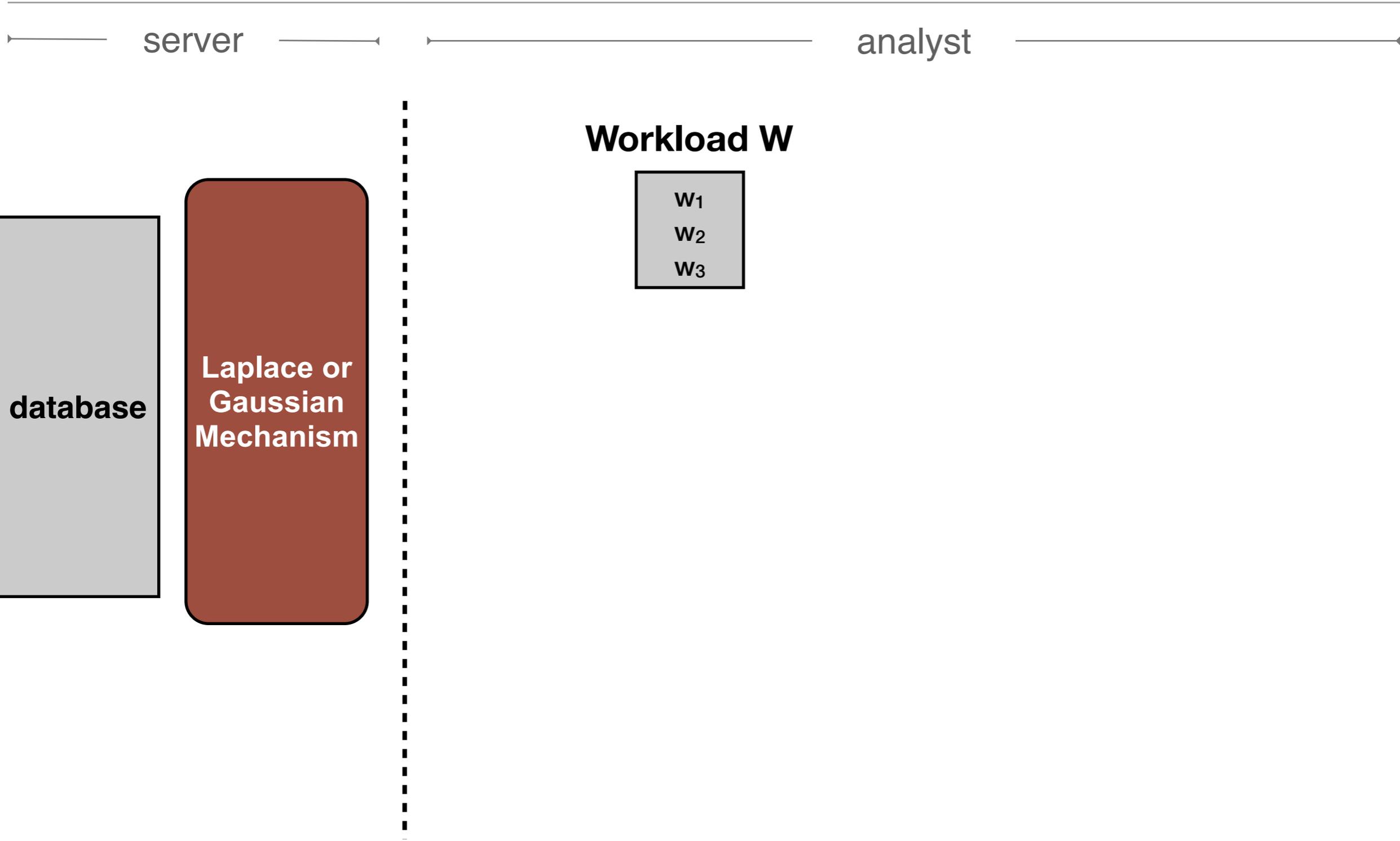
Approach 1: workload-aware mechanisms



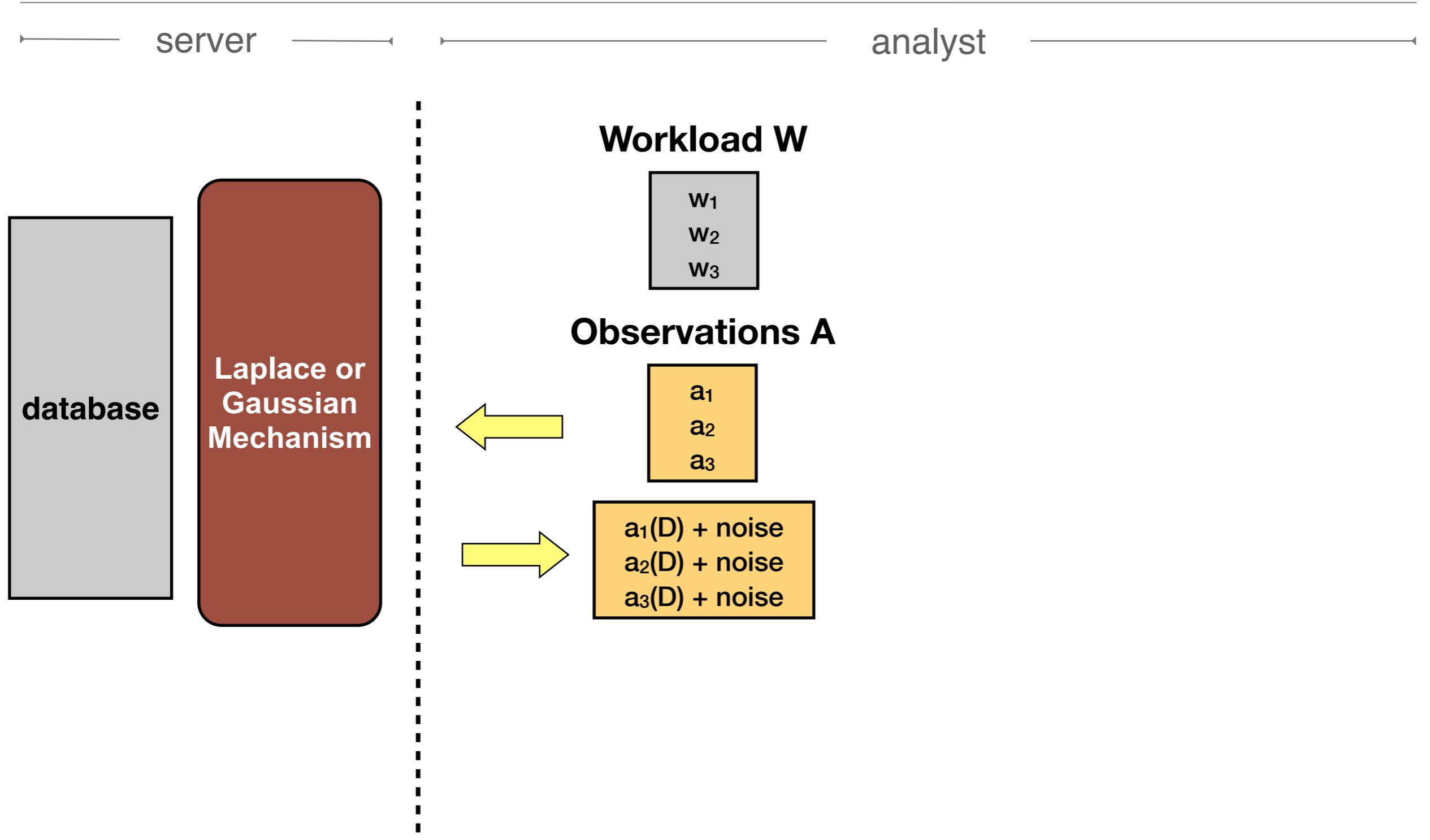
Approach 1: workload-aware mechanisms



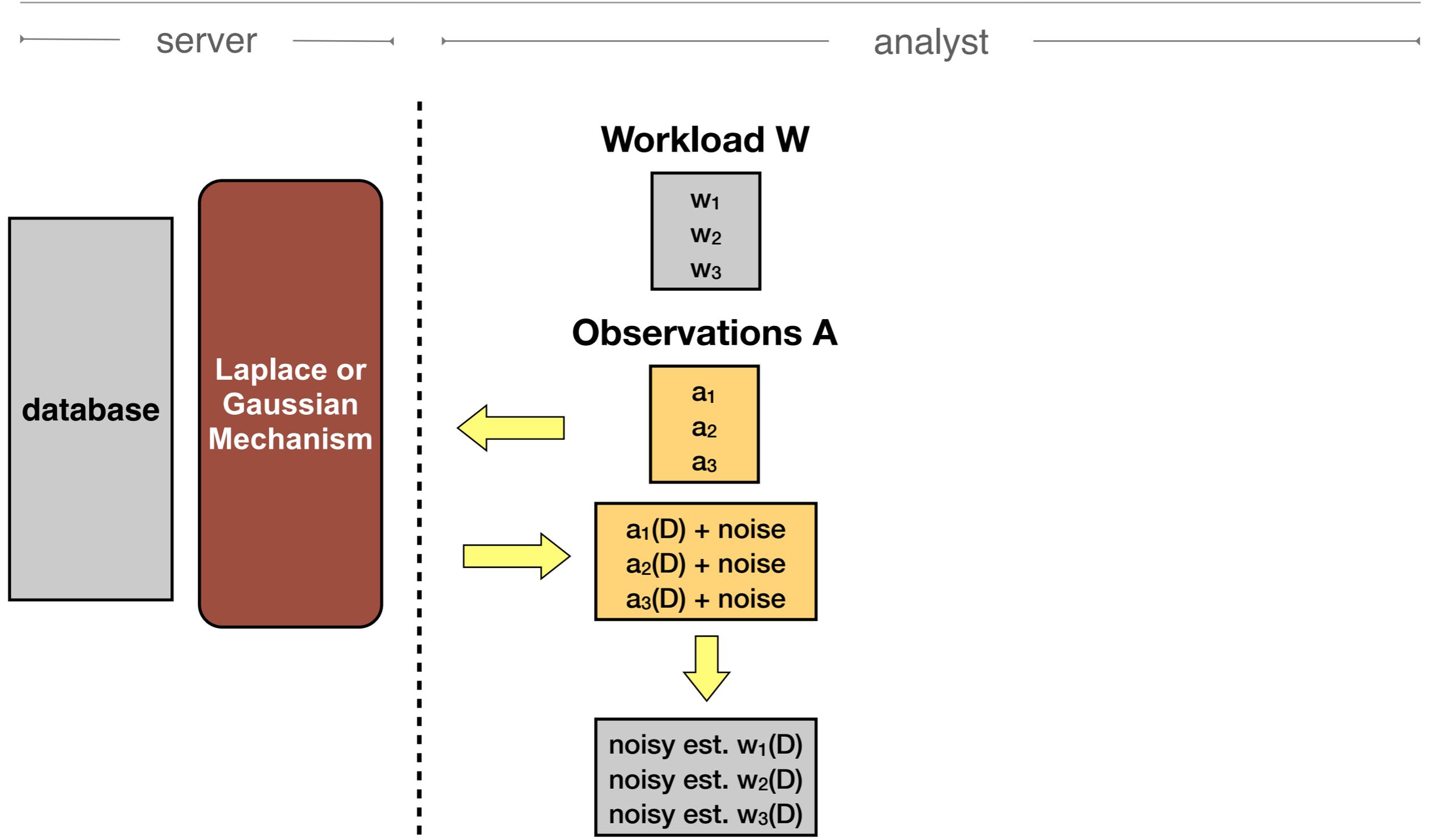
Approach 1: workload-aware mechanisms



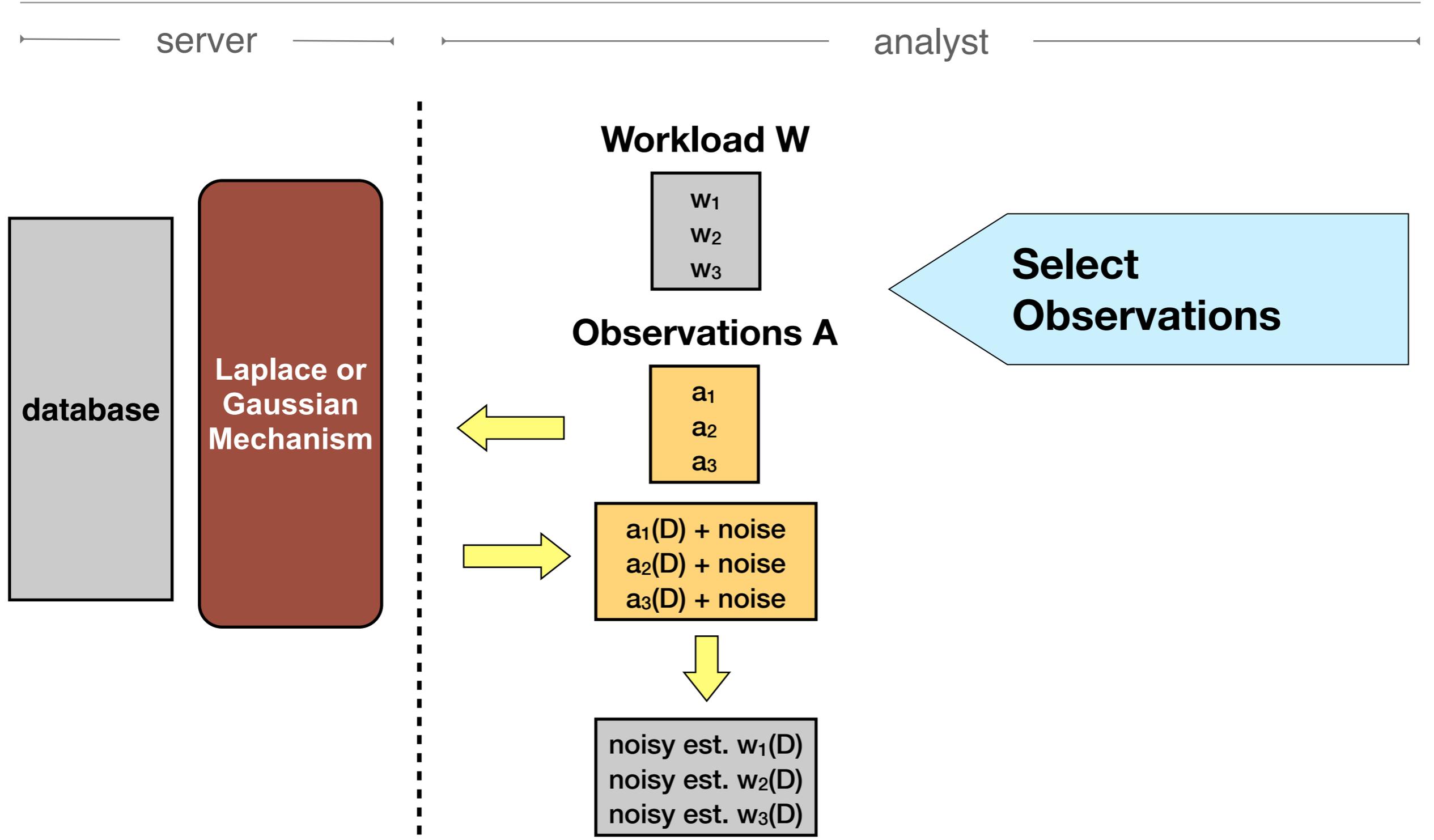
Approach 1: workload-aware mechanisms



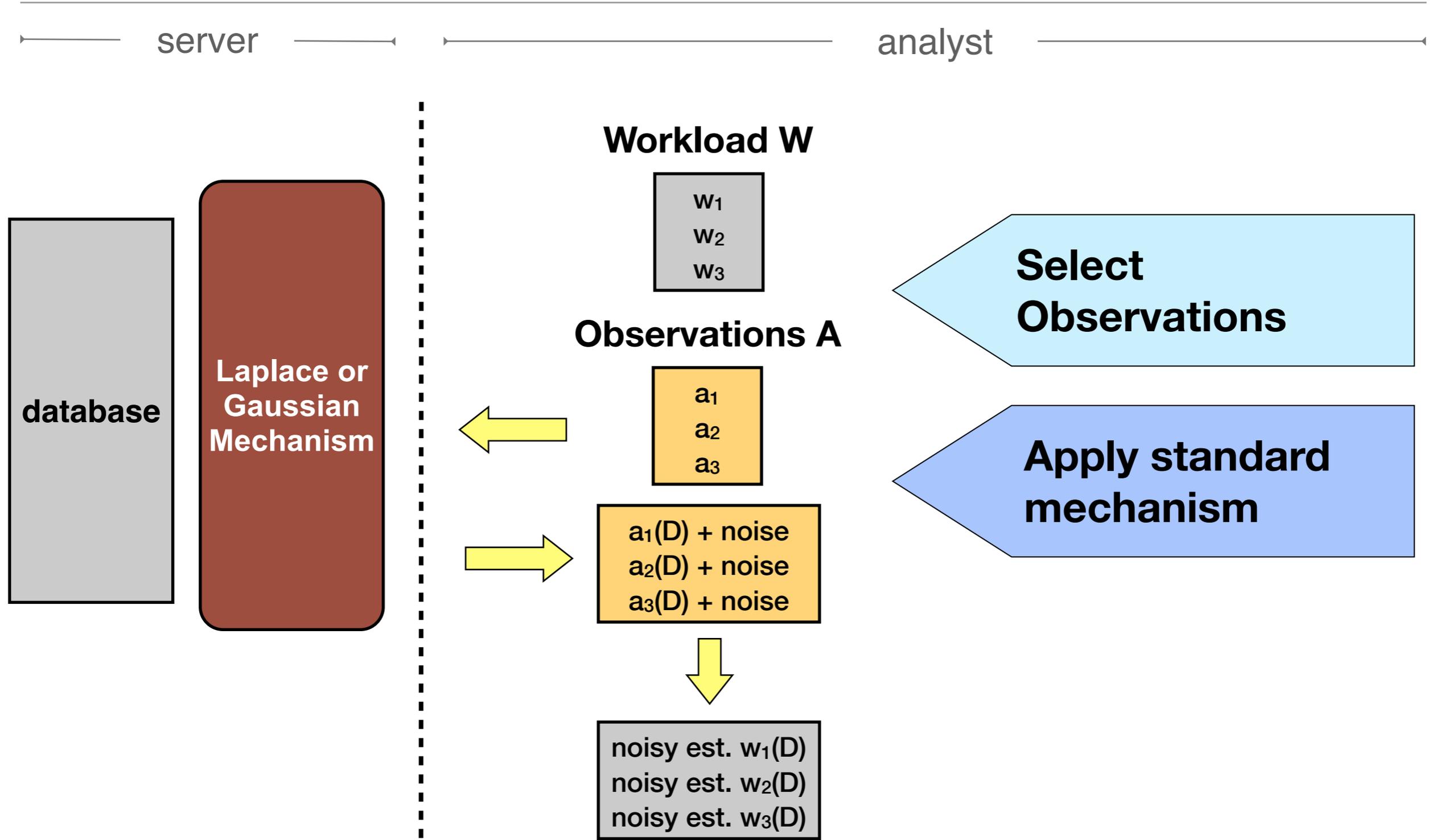
Approach 1: workload-aware mechanisms



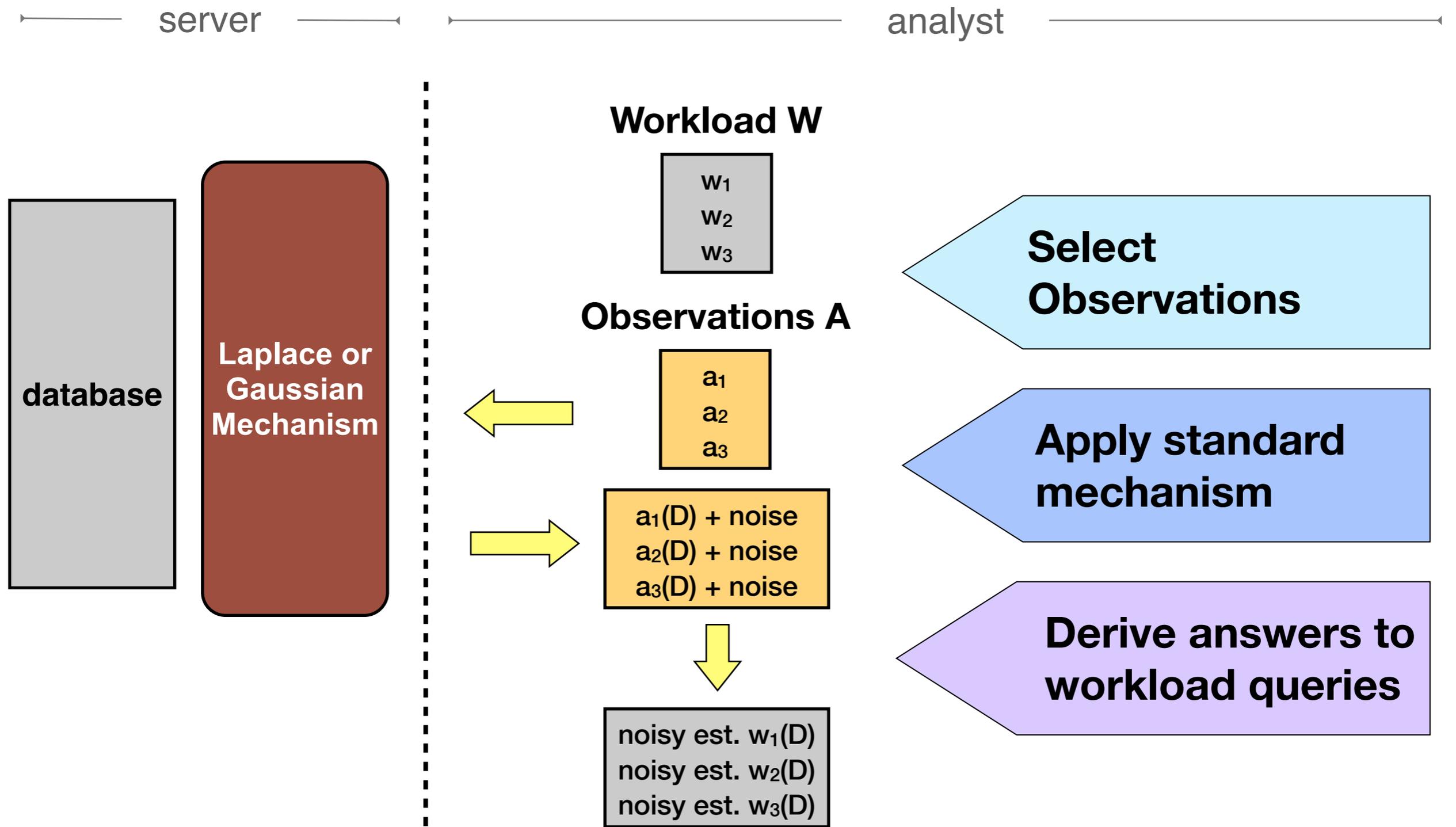
Approach 1: workload-aware mechanisms



Approach 1: workload-aware mechanisms



Approach 1: workload-aware mechanisms

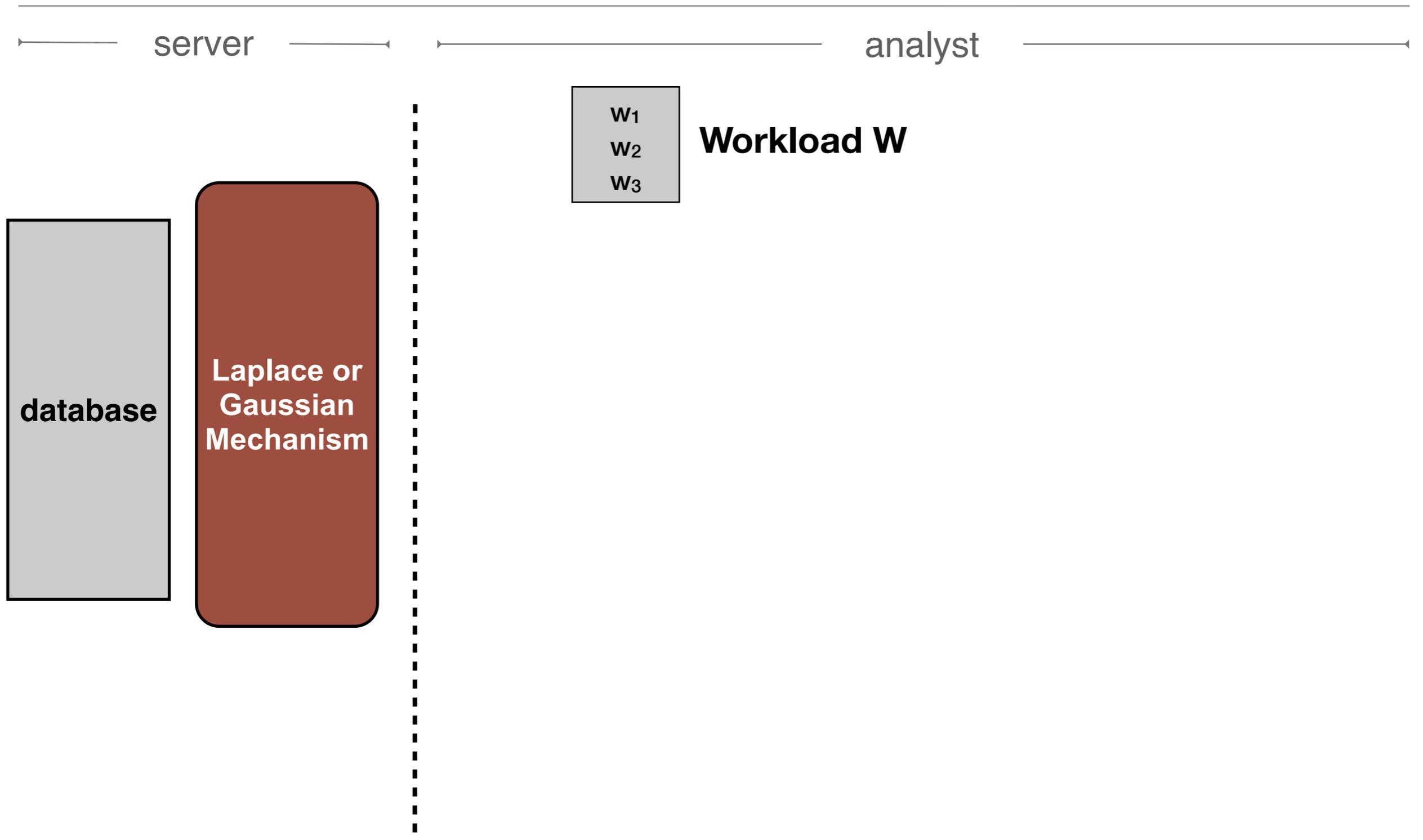


Workload-aware mechanisms

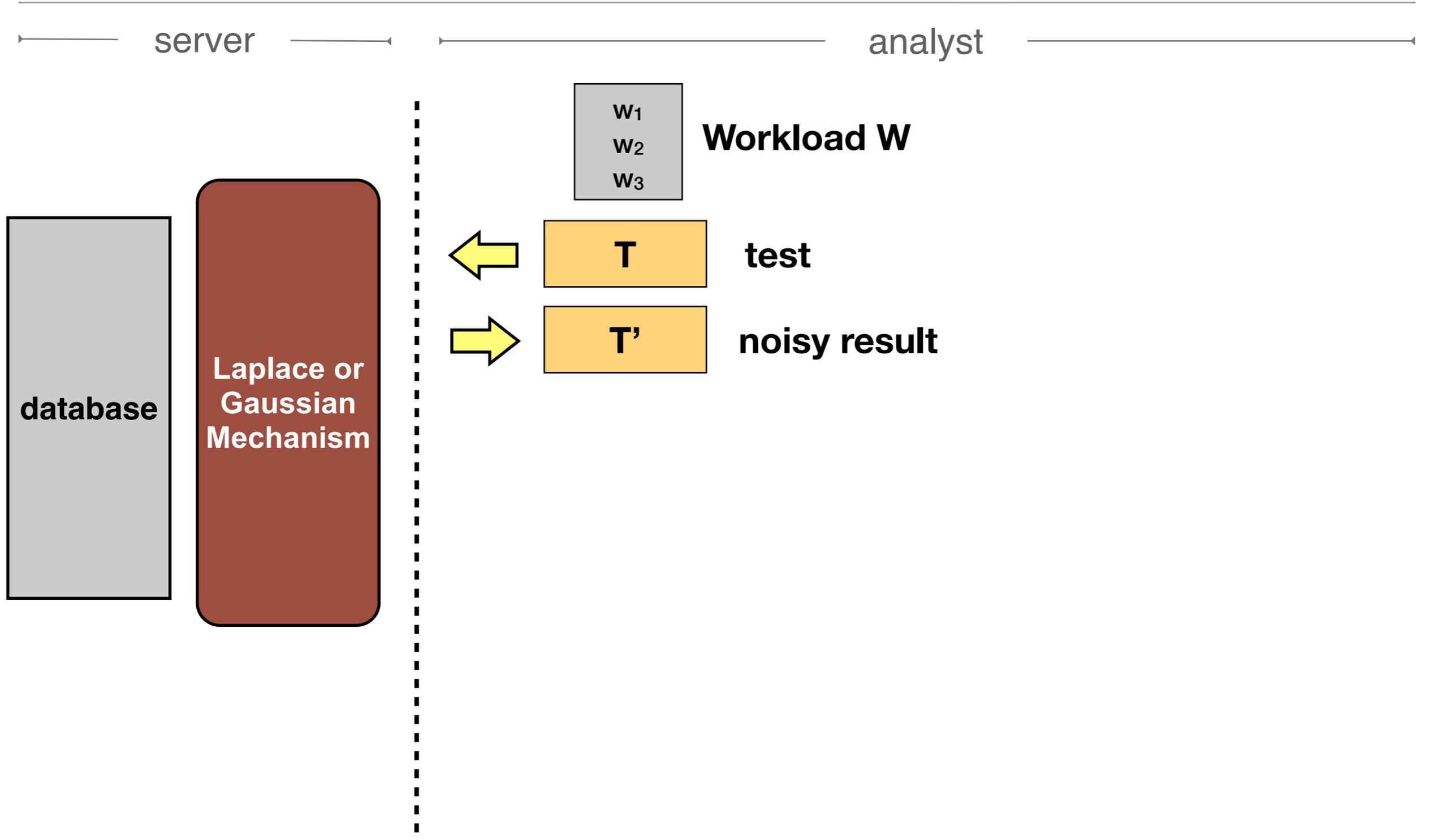
- Observations selected to match (only) the workload.

Workload	Observations		Citation
low-order marginals	Fixed	Fourier basis queries	[Barak, PODS '07]
all one-dim range queries		Hierarchical ranges	[Hay, PVLDB '10]
all (multi-dim) range queries		Haar wavelet queries	[Xiao, ICDE '10]
2-dim range queries		Quad-tree queries	[Cormode, ICDE '12]
sets of data cubes	Optimized	sets of data cubes	[Ding, SIGMOD '11]
set of linear queries		set of linear queries	[Li, PODS '10] [Li, PVLDB '12]
set of linear queries		set of linear queries	[Yuan, VLDB '12]

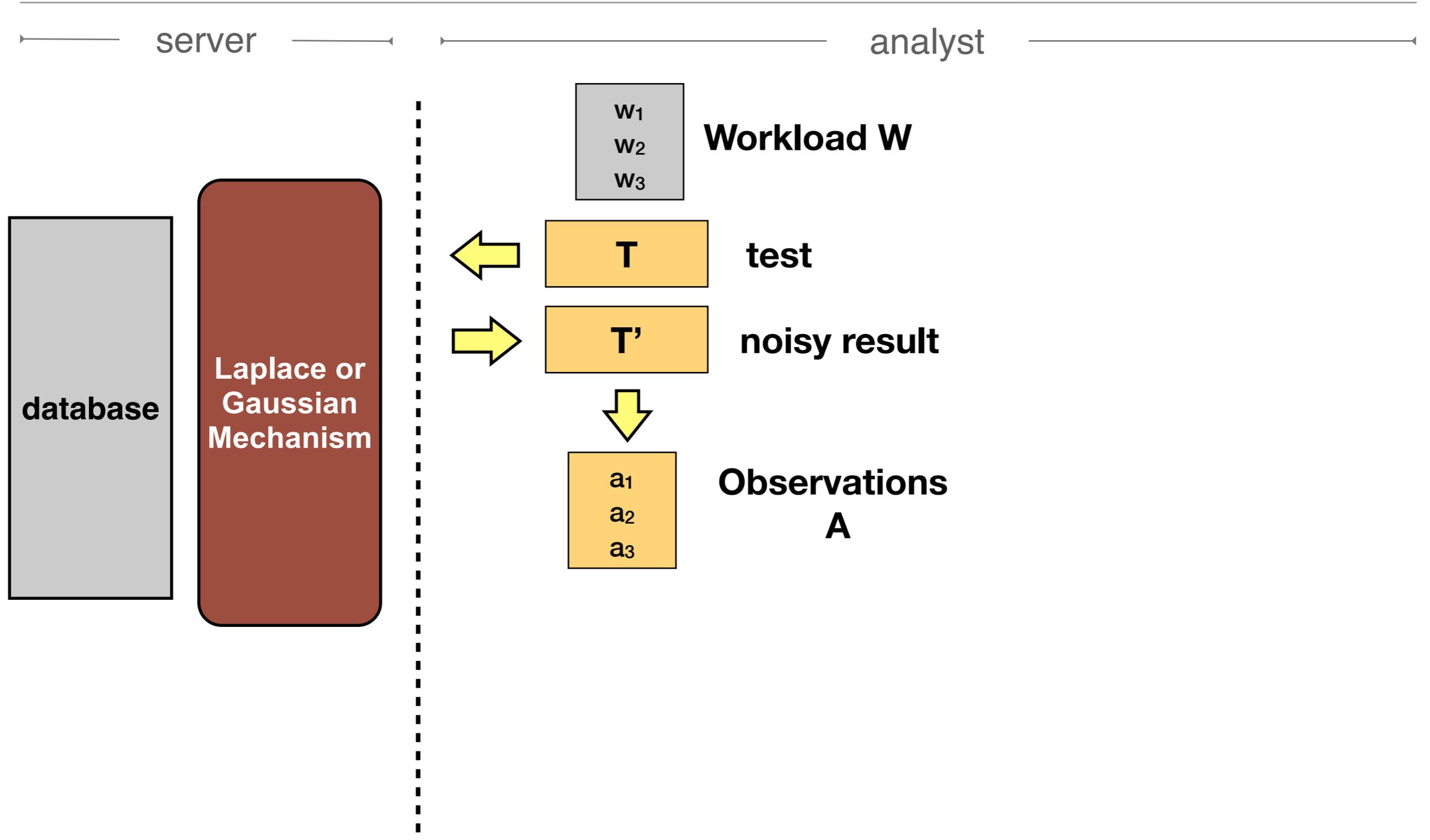
Approach 2: data-aware mechanisms



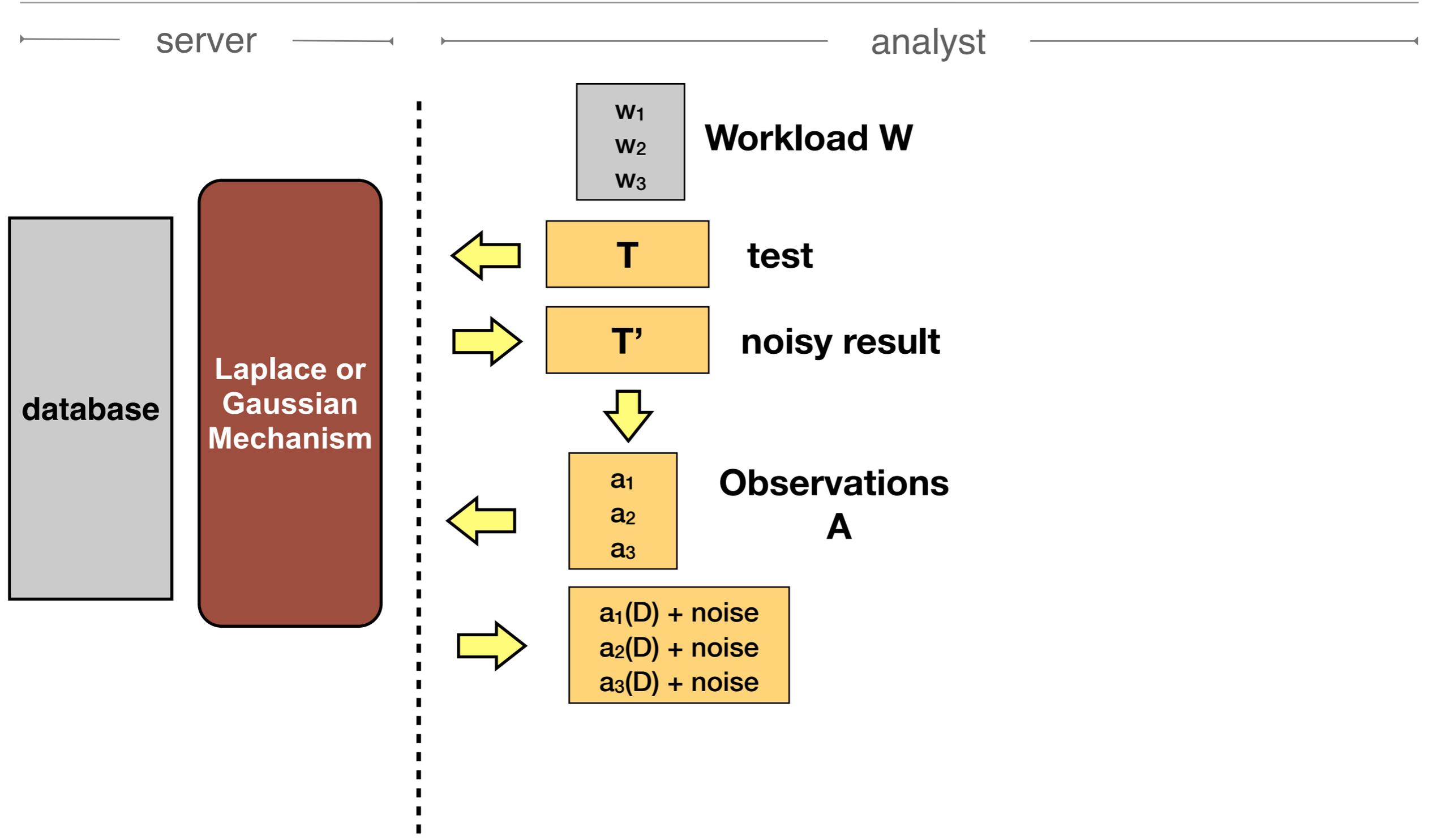
Approach 2: data-aware mechanisms



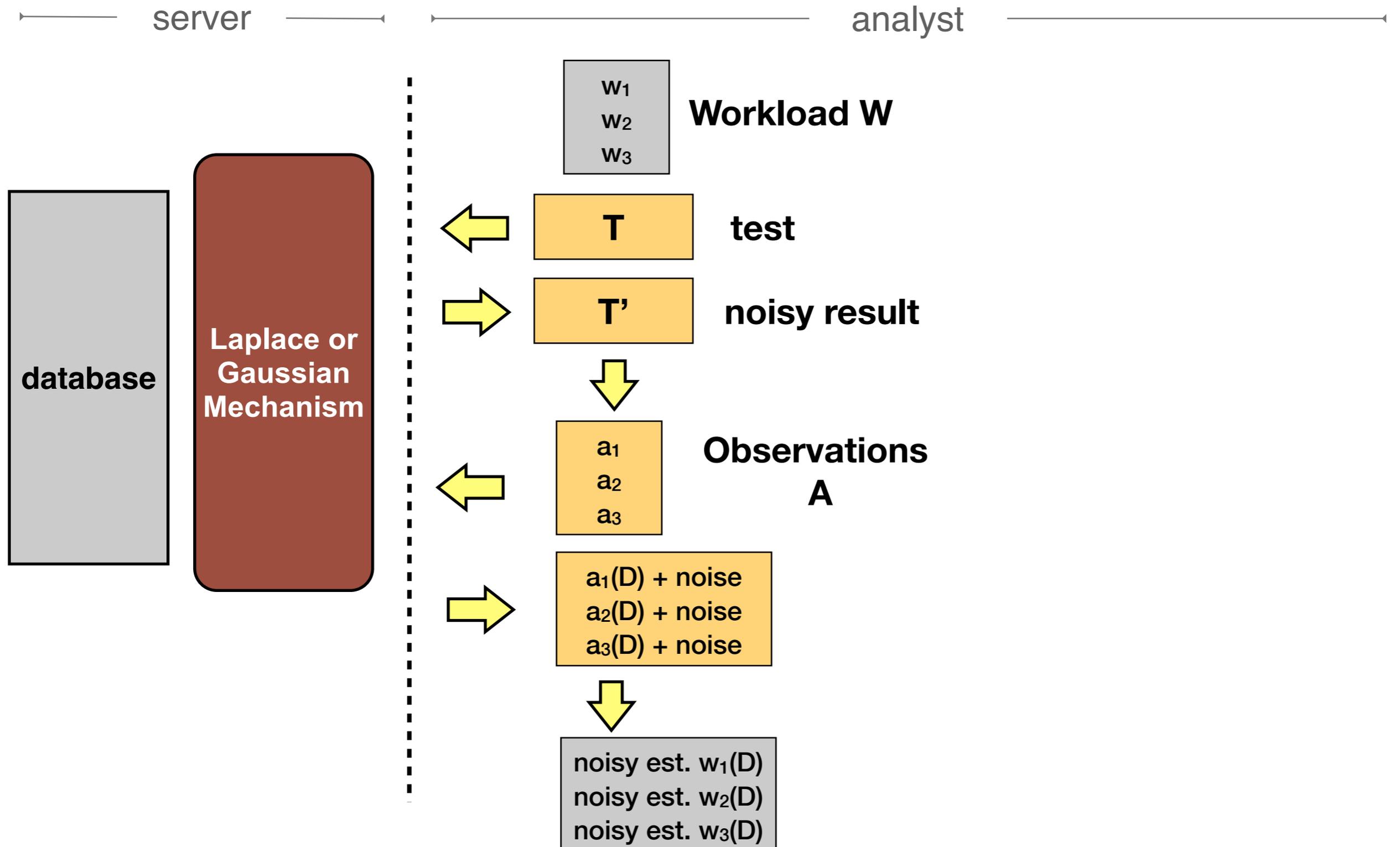
Approach 2: data-aware mechanisms



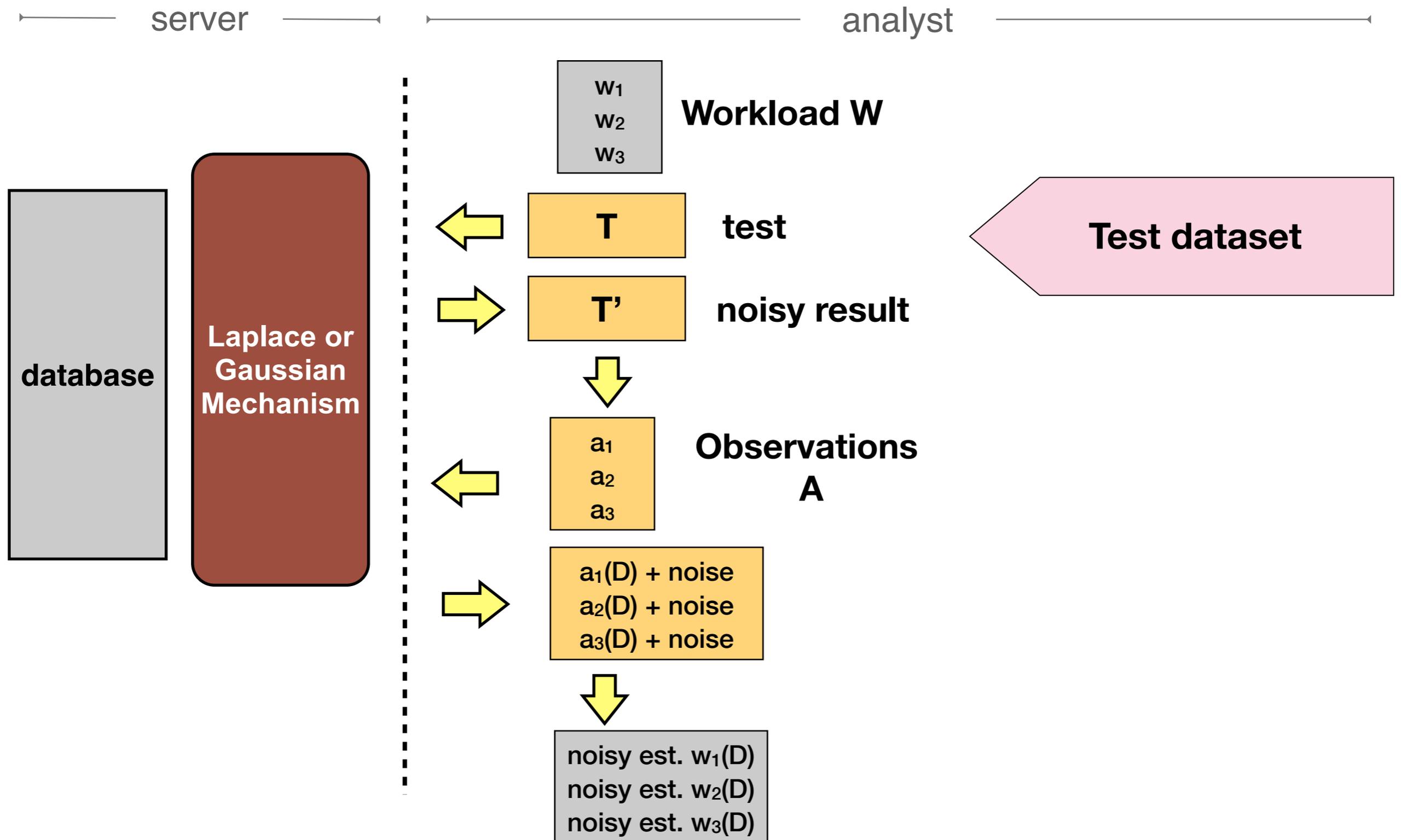
Approach 2: data-aware mechanisms



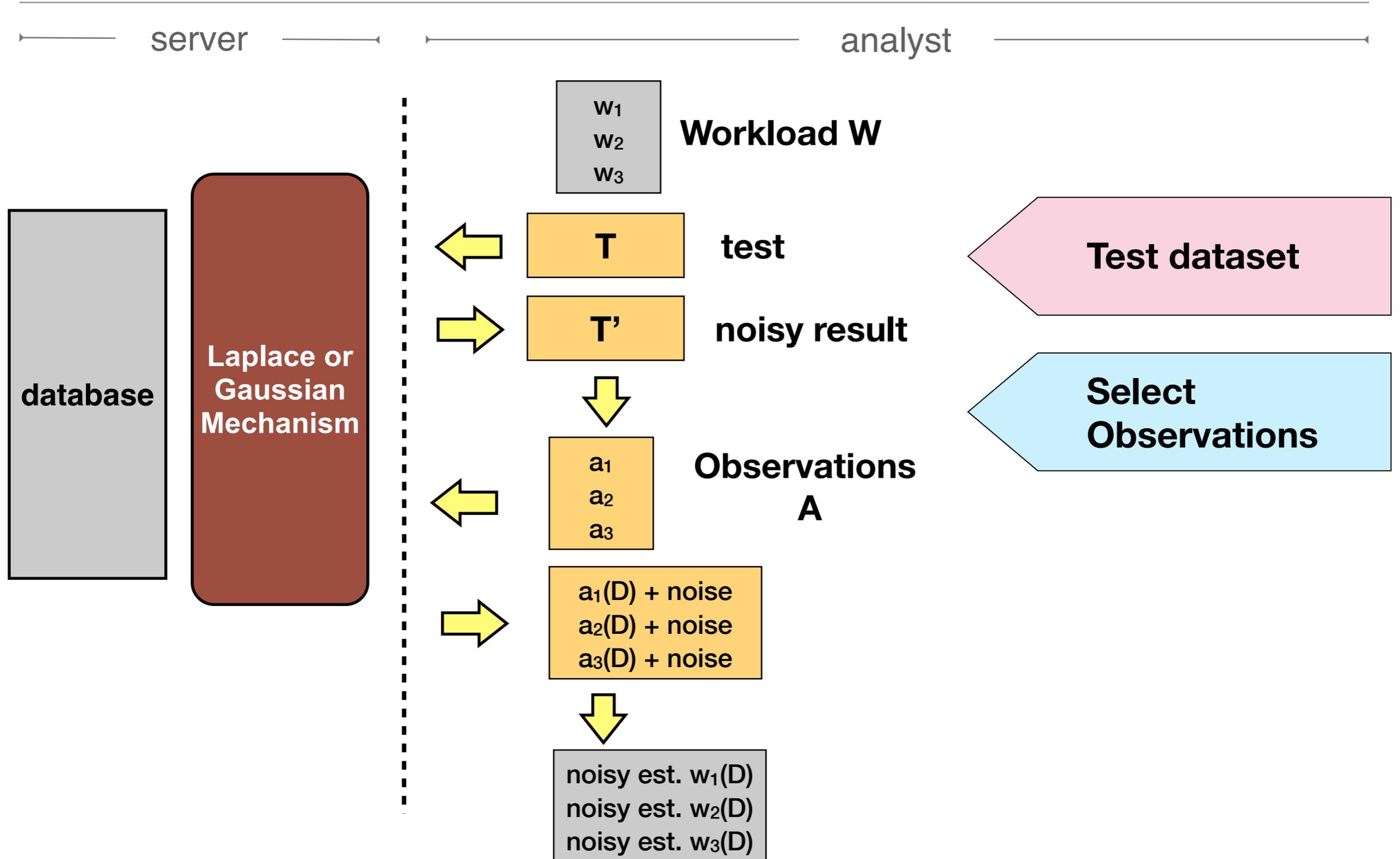
Approach 2: data-aware mechanisms



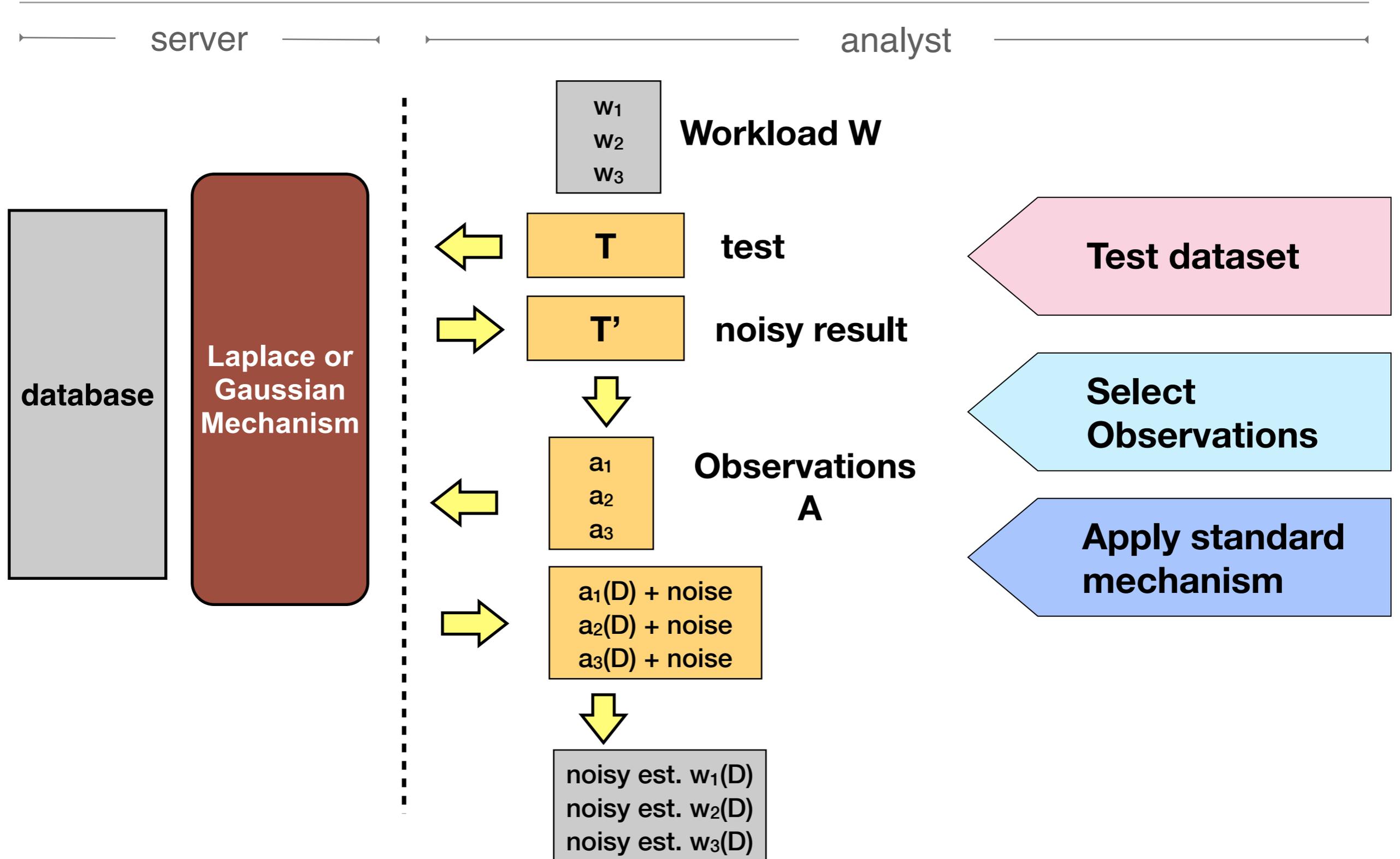
Approach 2: data-aware mechanisms



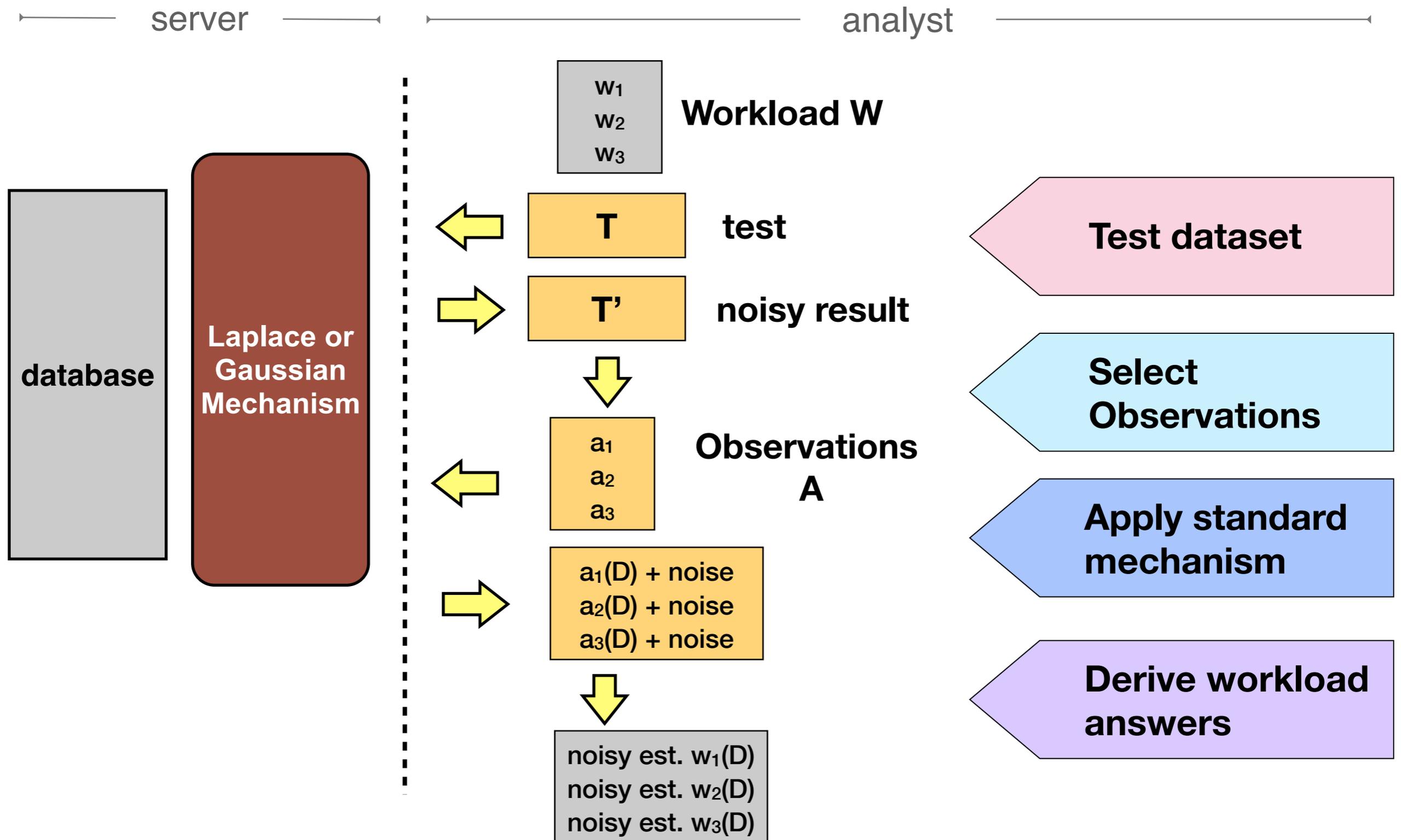
Approach 2: data-aware mechanisms



Approach 2: data-aware mechanisms



Approach 2: data-aware mechanisms



Data-aware mechanisms

- Observations selected to match properties of the database.

Workload	Observations	Citation
1D range queries	approx. v-optimal histogram	[Xu, ICDE '12]
2D range queries	kd-tree queries	[Xiao, SDM '10]
2D range queries	hybrid kd-tree queries	[Cormode, ICDE '12]
Marginals	scaled workload queries	[Xiao, SIGMOD '11]
Linear queries	subset of workload	[Hardt, NIPS '12]

Outline

1. Preliminaries

2. Approach 1: workload-aware

- Fixed Observations
- Optimized Observations

3. Approach 2: data-aware

4. Conclusions

Frequency representation of the database

name	gender	grade
Alice	Female	91
Bob	Male	84
Carl	Male	82
Dave	Male	97
Edwina	Female	88
Faith	Female	78
Ghita	Female	85
...

Relational database

{gender, grade}

gender	grade	count
Male	100	10
Male	99	13
Male	98	5
Male	97	7
...
Female	100	15
Female	99	21
Female	98	4
Female	97	14
Female	96	9

Frequency vector

x_1
x_2
x_3
x_4
x_5
x_6
x_7
x_8
...
x_n

x

Frequency representation of the database

name	gender	grade
Alice	Female	91
Bob	Male	84
Carl	Male	82
Dave	Male	97
Edwina	Female	88
Faith	Female	78
Ghita	Female	85
...

Relational database

Frequency vector

x

Frequency representation of the database

name	gender	grade
Alice	Female	91
Bob	Male	84
Carl	Male	82
Dave	Male	97
Edwina	Female	88
Faith	Female	78
Ghita	Female	85
...

Relational database

{grade}

grade	count
90-100	10
80-90	23
70-80	16
60-70	3

Frequency vector

x ₁
x ₂
x ₃
x ₄

x

Linear counting queries

A **linear counting query** w computes a linear combination of the frequency vector counts:

$$w(\mathbf{D}) = w_1x_1 + w_2x_2 + \dots + w_nx_n \quad \text{each } w_i \in \mathbb{R}$$

Linear counting queries

A **linear counting query** w computes a linear combination of the frequency vector counts:

$$\mathbf{w}(\mathbf{D}) = w_1x_1 + w_2x_2 + \dots + w_nx_n \quad \text{each } w_i \in \mathbb{R}$$

... as a length n row vector:

$$\mathbf{w} = [w_1, w_2, w_3 \dots w_n]$$

The query result is:

$$\mathbf{w}\mathbf{x}$$

Linear counting queries

A **linear counting query** w computes a linear combination of the frequency vector counts:

$$w(\mathbf{D}) = w_1x_1 + w_2x_2 + \dots + w_nx_n \quad \text{each } w_i \in \mathbb{R}$$

... as a length n row vector:

$$w = [w_1, w_2, w_3 \dots w_n]$$

The query result is:

$$wx$$

a set of linear counting queries is a matrix:

$$W = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 \end{bmatrix}$$

The query result is:

$$Wx$$

Queries and workloads

- **1-dimensional range queries:** intervals
- **Marginals / data cube queries / contingency tables:** aggregate over excluded dimensions.
- **k-dimensional range queries:** axis-aligned rectangles
- **Predicate counting queries:** only 0 or 1 coefficients
- **Linear counting queries:** arbitrary coefficients

Queries and workloads

1-dim ranges

- **1-dimensional range queries:** intervals
- **Marginals / data cube queries / contingency tables:** aggregate over excluded dimensions.
- **k-dimensional range queries:** axis-aligned rectangles
- **Predicate counting queries:** only 0 or 1 coefficients
- **Linear counting queries:** arbitrary coefficients

Queries and workloads

1-dim ranges

marginals

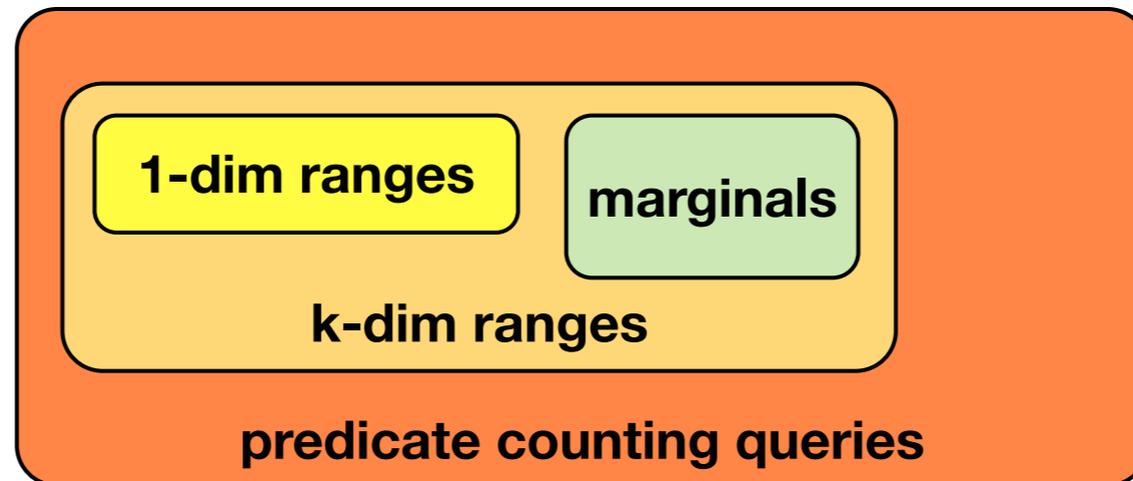
- **1-dimensional range queries:** intervals
- **Marginals / data cube queries / contingency tables:** aggregate over excluded dimensions.
- **k-dimensional range queries:** axis-aligned rectangles
- **Predicate counting queries:** only 0 or 1 coefficients
- **Linear counting queries:** arbitrary coefficients

Queries and workloads



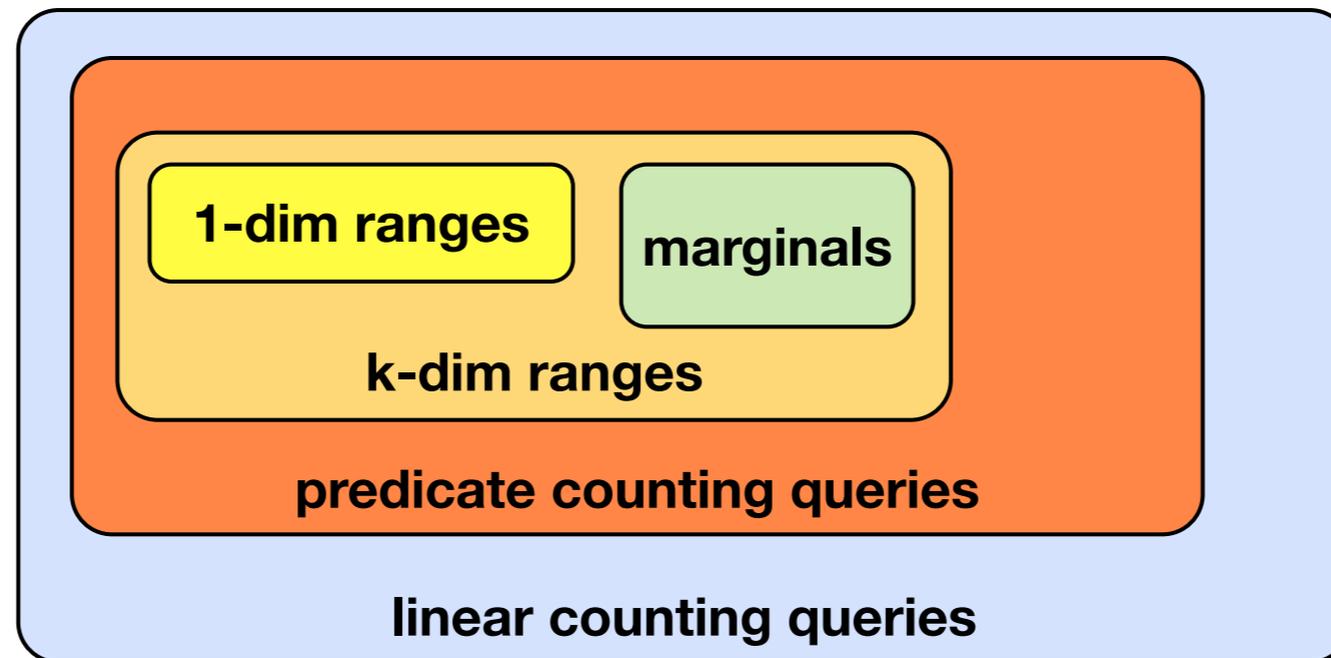
- **1-dimensional range queries:** intervals
- **Marginals / data cube queries / contingency tables:** aggregate over excluded dimensions.
- **k-dimensional range queries:** axis-aligned rectangles
- **Predicate counting queries:** only 0 or 1 coefficients
- **Linear counting queries:** arbitrary coefficients

Queries and workloads



- **1-dimensional range queries:** intervals
- **Marginals / data cube queries / contingency tables:** aggregate over excluded dimensions.
- **k-dimensional range queries:** axis-aligned rectangles
- **Predicate counting queries:** only 0 or 1 coefficients
- **Linear counting queries:** arbitrary coefficients

Queries and workloads



- **1-dimensional range queries:** intervals
- **Marginals / data cube queries / contingency tables:** aggregate over excluded dimensions.
- **k-dimensional range queries:** axis-aligned rectangles
- **Predicate counting queries:** only 0 or 1 coefficients
- **Linear counting queries:** arbitrary coefficients

Privacy definitions & mechanisms

- Differential privacy

A randomized algorithm \mathcal{A} provides (ϵ, δ) -**differential privacy** if:
for all neighboring databases D and D' , and
for any set of outputs S :

$$\Pr[\mathcal{A}(D) \in S] \leq e^\epsilon \Pr[\mathcal{A}(D') \in S] + \delta$$

- if $\delta=0$, standard ϵ -differential privacy:
 - **Laplace(0,b) noise where $b=\|q\|_1/\epsilon$**
- if $\delta>0$, approximate (ϵ, δ) -differential privacy:
 - **Gaussian(0, σ) noise where $\sigma= \|q\|_2 (2\ln(2/\delta))^{1/2}/\epsilon$**
- Multi-query Laplace/Gaussian mechanism adds independent noise to **each** query answer.
- Exponential mechanism

The sensitivity of a query matrix

- For two neighboring databases D and D' , their frequency vectors x and x' will differ in one position, by exactly 1.

$$\begin{array}{c} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \\ \text{answers} \end{array} = \begin{array}{c} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 \end{bmatrix} \\ \text{query matrix } \mathbf{W} \end{array} \times \begin{array}{c} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 + 1 \\ x_8 \\ x_9 \\ x_{10} \end{bmatrix} \\ \mathbf{x}' \end{array}$$

The sensitivity of a query matrix

- For two neighboring databases D and D', their frequency vectors x and x' will differ in one position, by exactly 1.

$$\begin{array}{c} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \\ \text{answers} \end{array} = \begin{array}{c} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 \end{bmatrix} \\ \text{query matrix } \mathbf{W} \end{array} \times \begin{array}{c} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 + 1 \\ x_8 \\ x_9 \\ x_{10} \end{bmatrix} \\ \mathbf{x}' \end{array}$$

The sensitivity of a query matrix

- For two neighboring databases D and D' , their frequency vectors x and x' will differ in one position, by exactly 1.

$$\begin{array}{c} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \\ \text{answers} \end{array} = \begin{array}{c} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 \end{bmatrix} \\ \text{query matrix } \mathbf{W} \\ \|\mathbf{W}\|_1 = 4 \end{array} \times \begin{array}{c} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 + 1 \\ x_8 \\ x_9 \\ x_{10} \end{bmatrix} \\ \mathbf{x}' \end{array}$$

The sensitivity of a query matrix

- For two neighboring databases D and D' , their frequency vectors x and x' will differ in one position, by exactly 1.

$$\begin{array}{c}
 \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \\
 \text{answers}
 \end{array}
 =
 \begin{array}{c}
 \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\
 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 \end{bmatrix} \\
 \text{query matrix } \mathbf{W} \\
 \|\mathbf{W}\|_1 = 4
 \end{array}
 \times
 \begin{array}{c}
 \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 + 1 \\ x_8 \\ x_9 \\ x_{10} \end{bmatrix} \\
 \mathbf{x}'
 \end{array}$$

The **L₁ sensitivity** of a query matrix is: the maximum L1 norm of the columns.

The sensitivity of a query matrix

- For two neighboring databases D and D' , their frequency vectors x and x' will differ in one position, by exactly 1.

$$\begin{array}{c}
 \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \\
 \text{answers}
 \end{array}
 =
 \begin{array}{c}
 \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 \end{bmatrix} \\
 \text{query matrix } \mathbf{W} \\
 \|\mathbf{W}\|_1 = 4
 \end{array}
 \times
 \begin{array}{c}
 \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 + 1 \\ x_8 \\ x_9 \\ x_{10} \end{bmatrix} \\
 \mathbf{x}'
 \end{array}$$

The **L_1 sensitivity** of a query matrix is: the maximum L_1 norm of the columns.

The **L_2 sensitivity** of a query matrix is: the maximum L_2 norm of the columns.

Outline

1. Preliminaries

2. Approach 1: workload-aware

- Fixed Observations
- Optimized Observations

3. Approach 2: data-aware

4. Conclusions

Outline

1. Preliminaries

2. Approach 1: workload-aware

- Fixed Observations
- Optimized Observations

3. Approach 2: data-aware

4. Conclusions

Answering all range queries

Goal: answer all **range-count queries** over x

$$\text{AllRange} = \{ w \mid w = x_i + \dots + x_j \text{ for } 1 \leq i \leq j \leq n \}$$

workload W

w_1	$\text{range}(x_1, x_4)$	x_1	+	x_2	+	x_3	+	x_4
w_2	$\text{range}(x_1, x_3)$	x_1	+	x_2	+	x_3		
w_3	$\text{range}(x_2, x_4)$			x_2	+	x_3	+	x_4
w_4	$\text{range}(x_1, x_2)$	x_1	+	x_2				
w_5	$\text{range}(x_2, x_3)$			x_2	+	x_3		
w_6	$\text{range}(x_3, x_4)$					x_3	+	x_4
w_7	$\text{range}(x_1, x_1)$	x_1						
w_8	$\text{range}(x_2, x_2)$			x_2				
w_9	$\text{range}(x_3, x_3)$					x_3		
w_{10}	$\text{range}(x_4, x_4)$							x_4

$X =$

10	23	16	3
----	----	----	---

Answering all range queries

Goal: answer all **range-count queries** over x

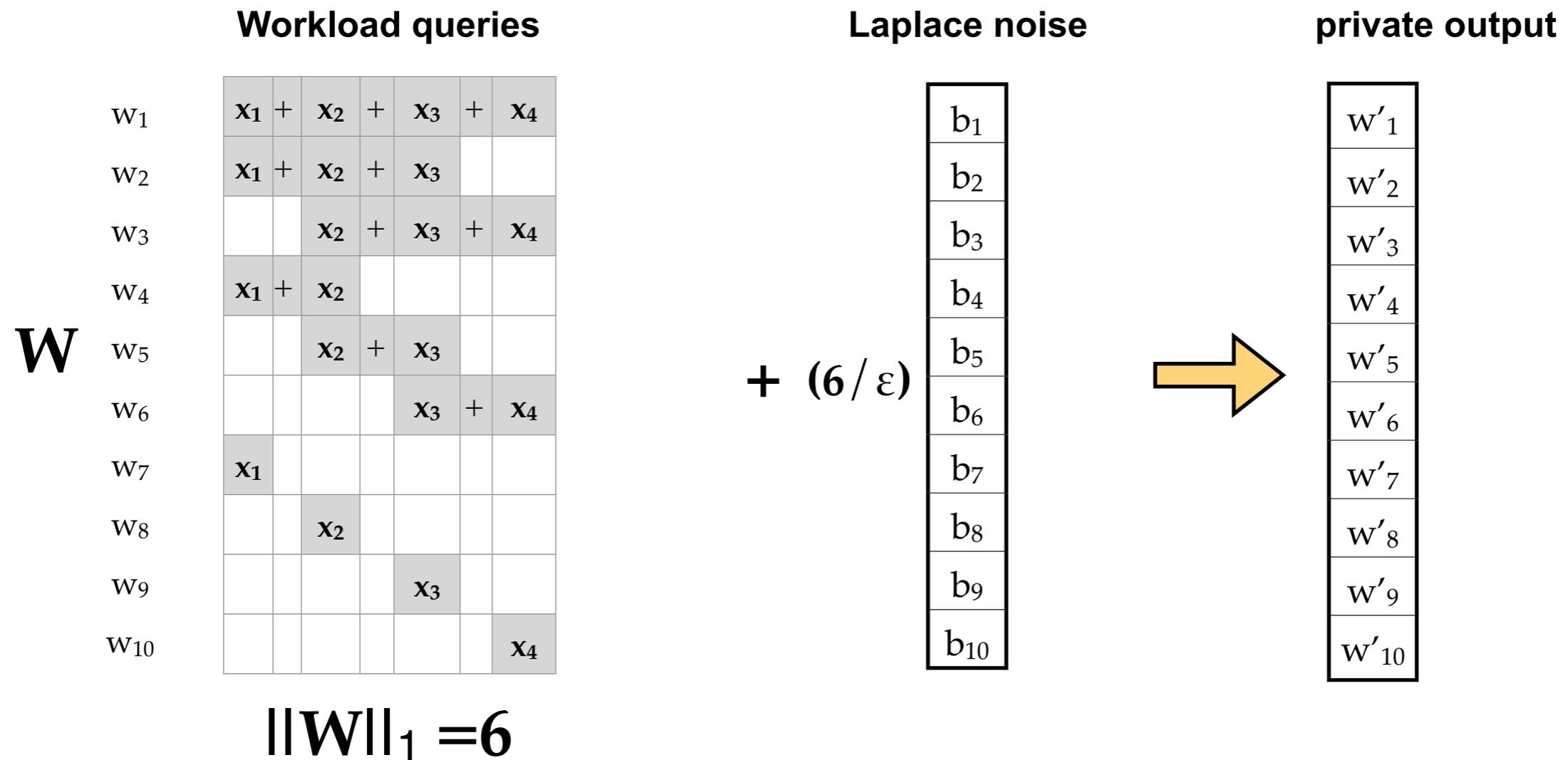
$$\text{AllRange} = \{ w \mid w = x_i + \dots + x_j \text{ for } 1 \leq i \leq j \leq n \}$$

workload W

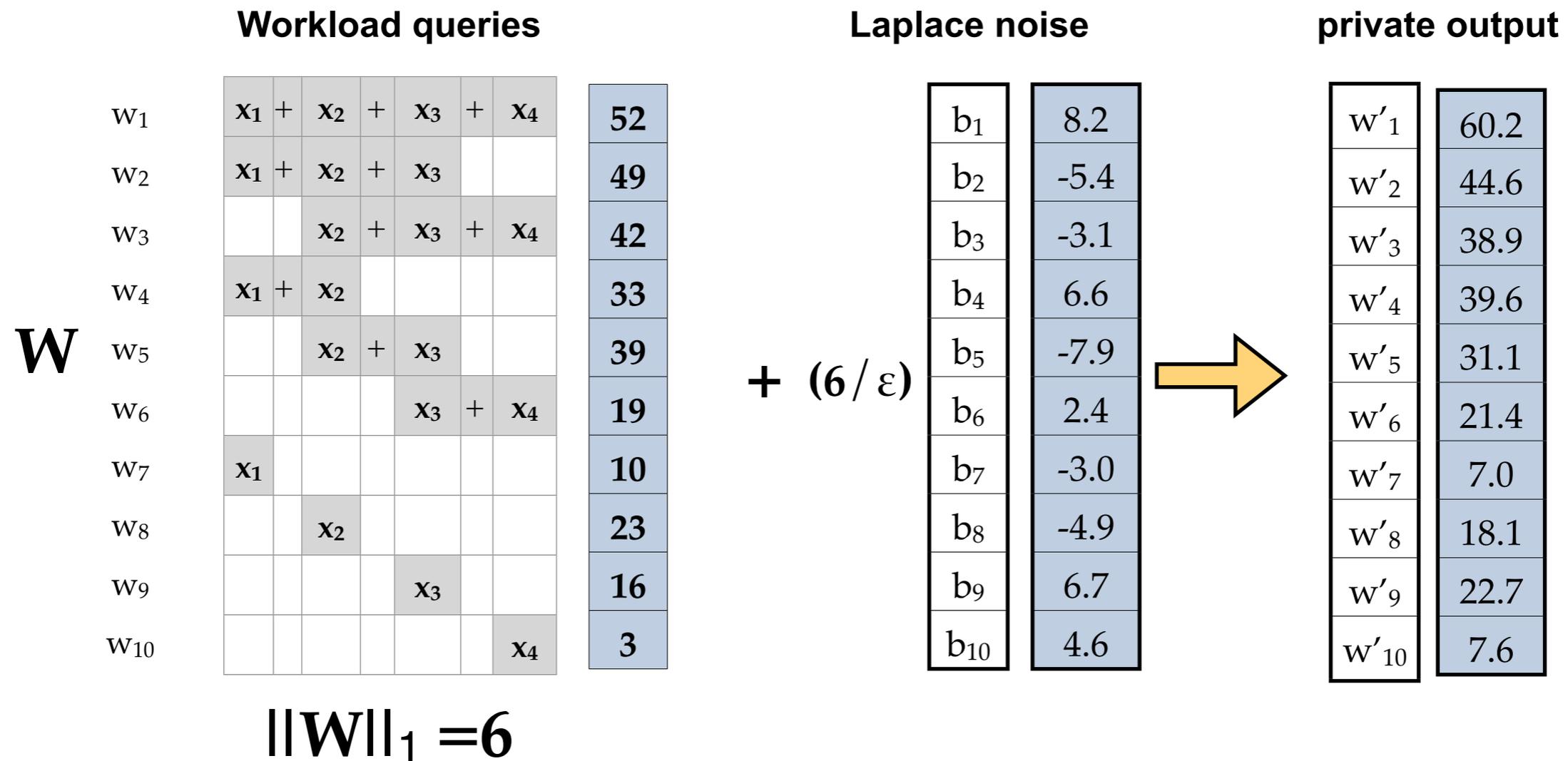
w ₁	range(x ₁ ,x ₄)	x ₁ + x ₂ + x ₃ + x ₄	w ₁	52
w ₂	range(x ₁ ,x ₃)	x ₁ + x ₂ + x ₃	w ₂	49
w ₃	range(x ₂ ,x ₄)	x ₂ + x ₃ + x ₄	w ₃	42
w ₄	range(x ₁ ,x ₂)	x ₁ + x ₂	w ₄	33
w ₅	range(x ₂ ,x ₃)	x ₂ + x ₃	w ₅	39
w ₆	range(x ₃ ,x ₄)	x ₃ + x ₄	w ₆	19
w ₇	range(x ₁ ,x ₁)	x ₁	w ₇	10
w ₈	range(x ₂ ,x ₂)	x ₂	w ₈	23
w ₉	range(x ₃ ,x ₃)	x ₃	w ₉	16
w ₁₀	range(x ₄ ,x ₄)	x ₄	w ₁₀	3

X= 10 23 16 3

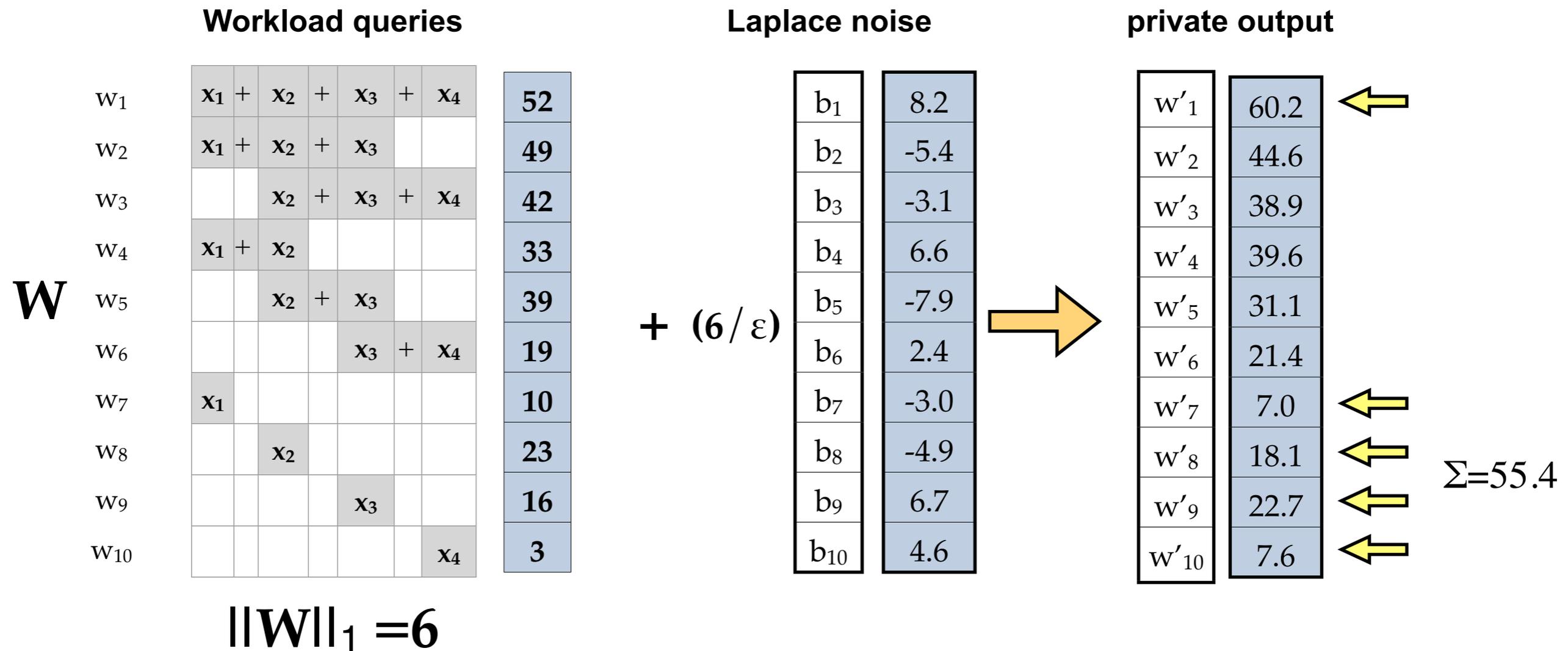
Method 1: basic Laplace mechanism



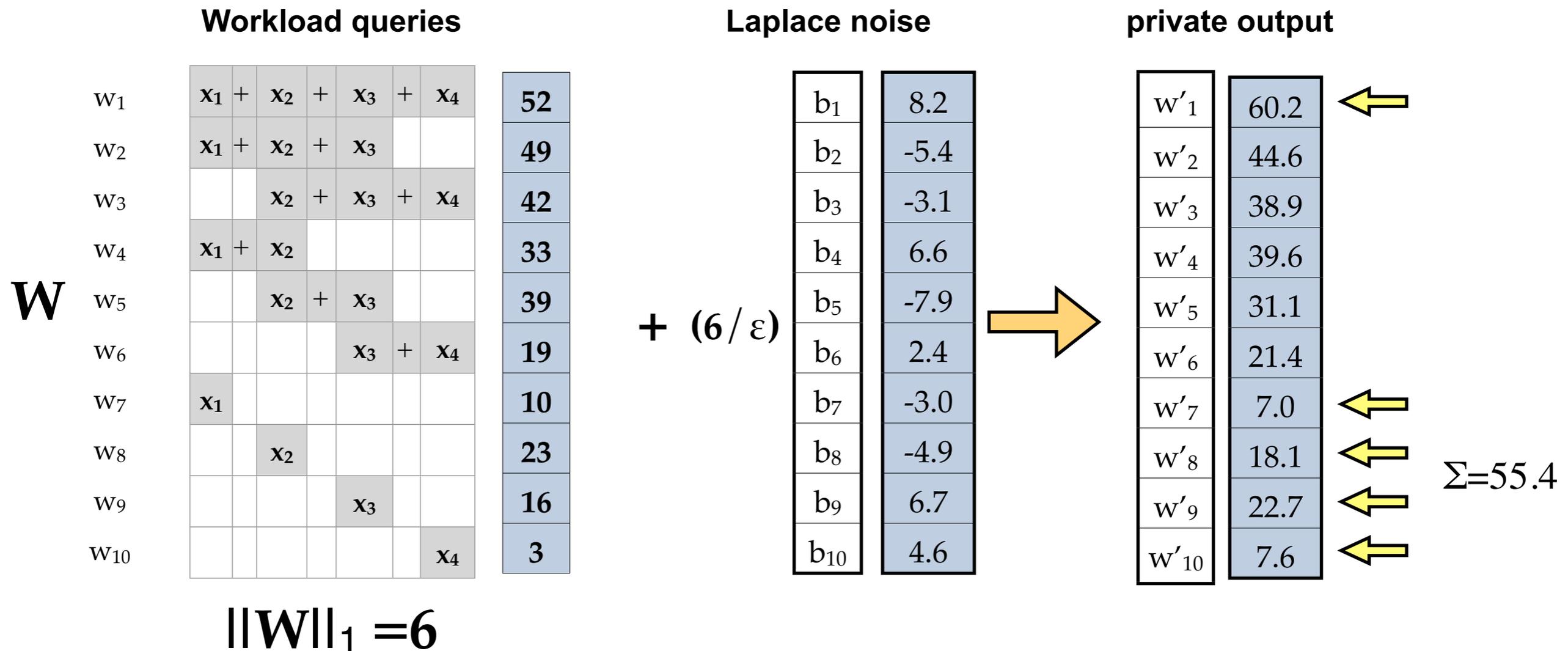
Method 1: basic Laplace mechanism



Method 1: basic Laplace mechanism



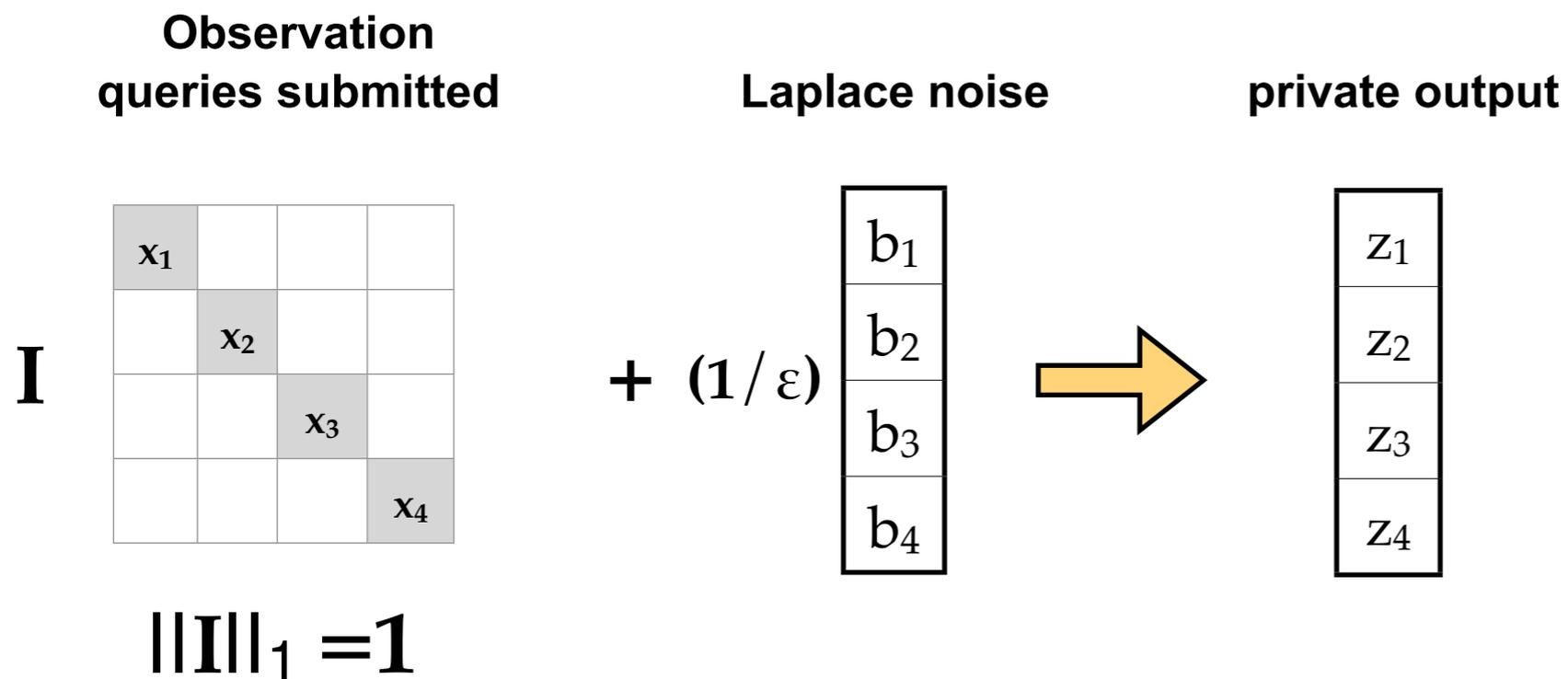
Method 1: basic Laplace mechanism



	n=4	n
Sensitivity $\ W\ _1$	6	$O(n^2)$
Error per query	$2(\ W\ _1/\epsilon)^2 = 72/\epsilon^2$	$2(\ W\ _1/\epsilon)^2 = O(n^4)/\epsilon^2$

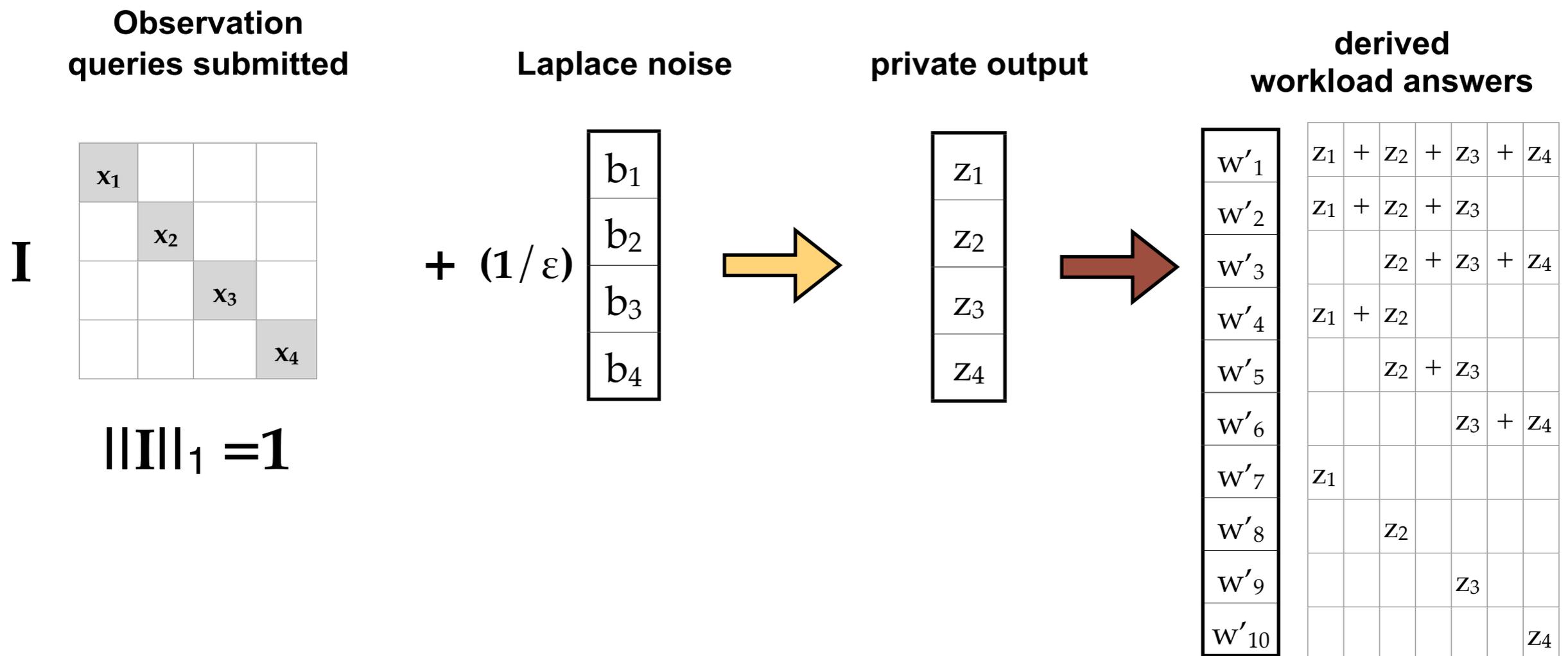
Method 2: noisy frequency counts

Use Laplace mechanism to get noisy estimates for each x_i .



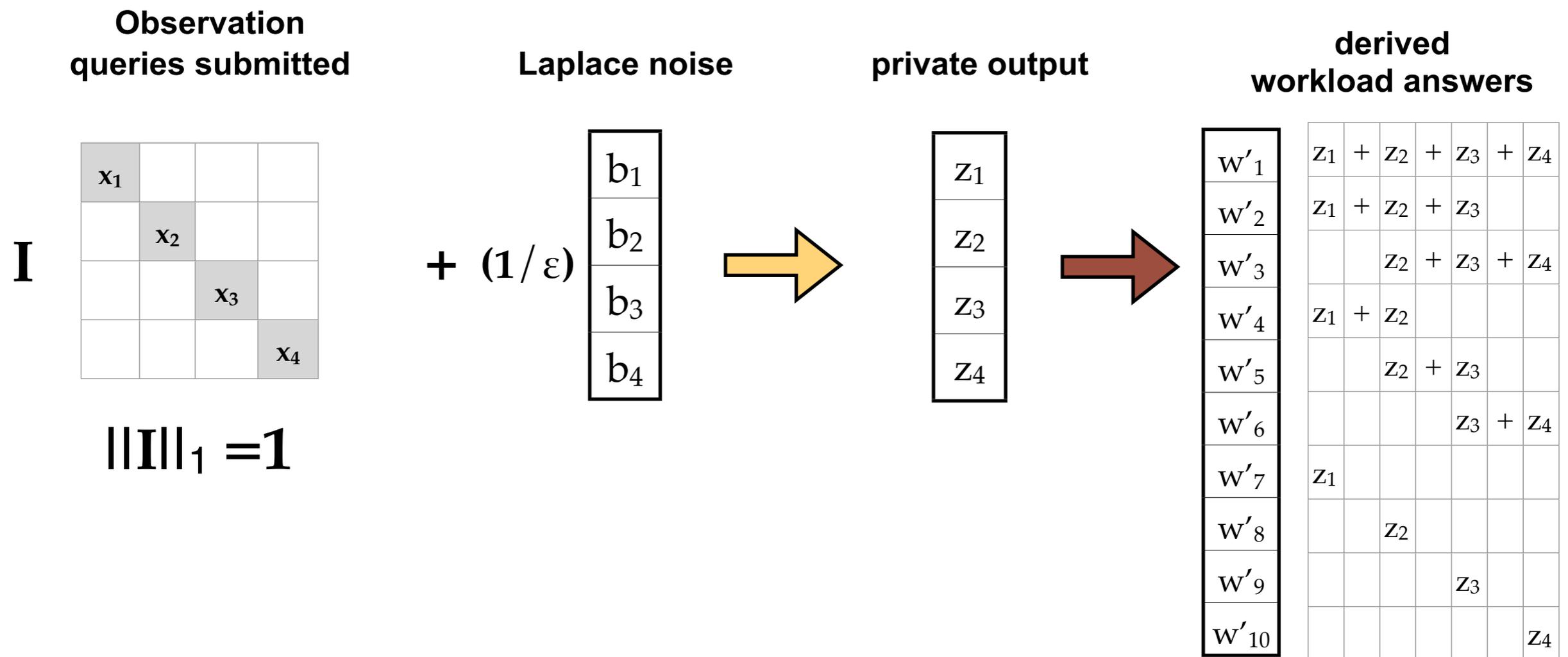
Method 2: noisy frequency counts

Use Laplace mechanism to get noisy estimates for each x_i .



Method 2: noisy frequency counts

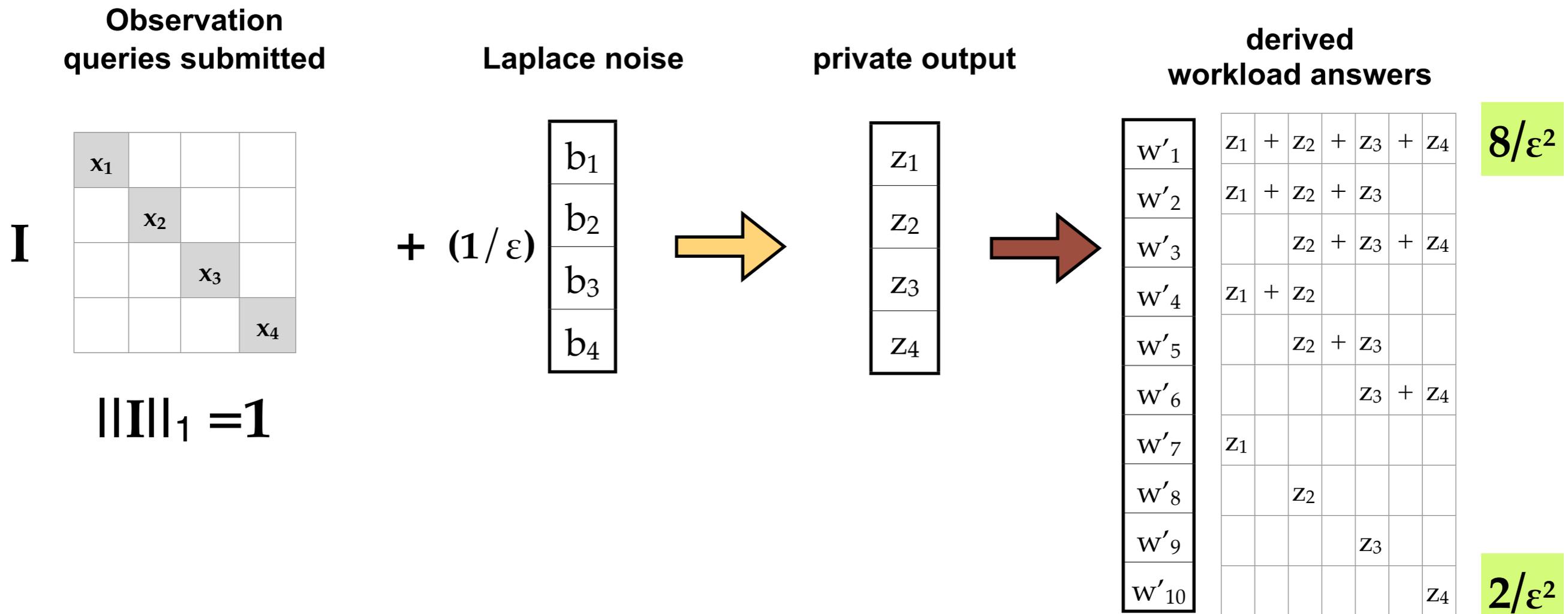
Use Laplace mechanism to get noisy estimates for each x_i .



For $w = \text{range}(x_i, x_j)$ $\text{Error}(w) = 2(j-i+1) / \epsilon^2$

Method 2: noisy frequency counts

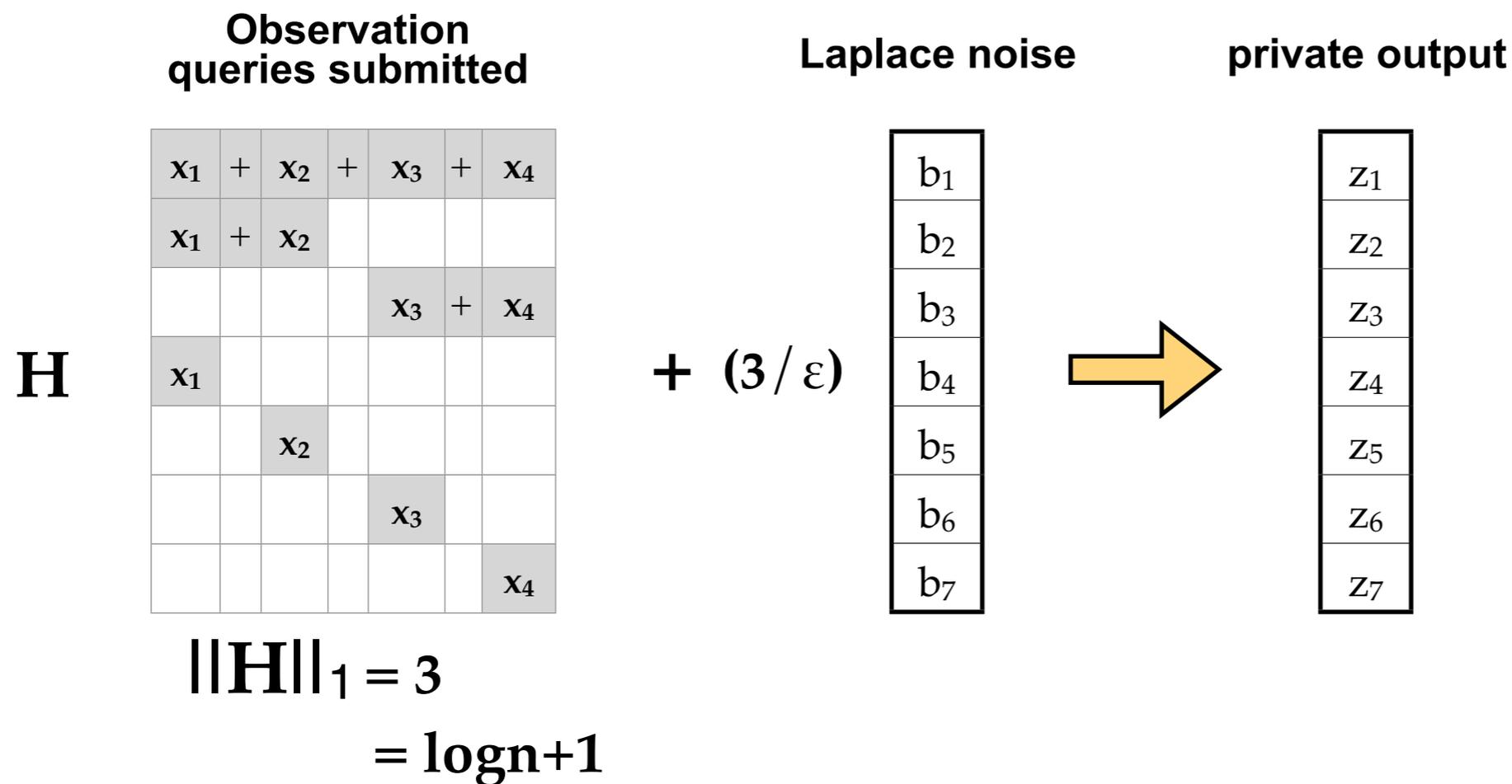
Use Laplace mechanism to get noisy estimates for each x_i .



For $w = \text{range}(x_i, x_j)$ $\text{Error}(w) = 2(j-i+1) / \epsilon^2$

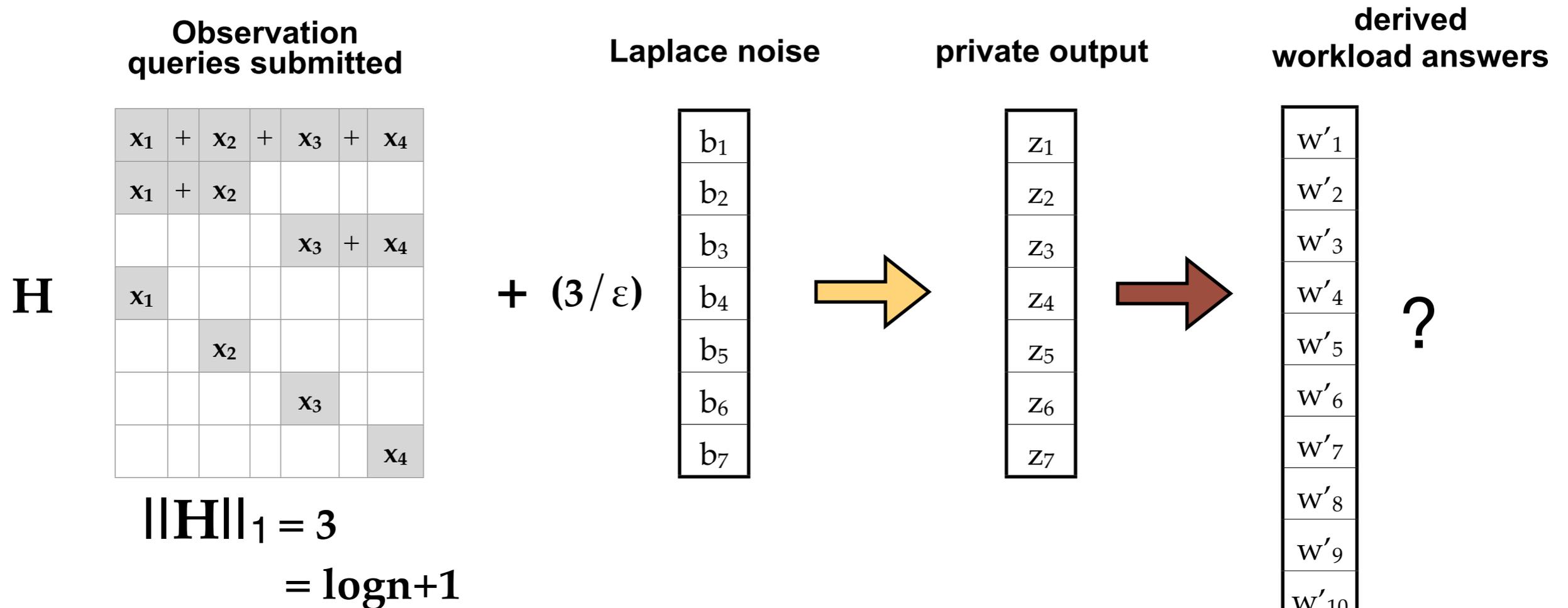
Method 3: hierarchical observations

Hierarchical queries: recursively partition the domain, computing sums of each interval.



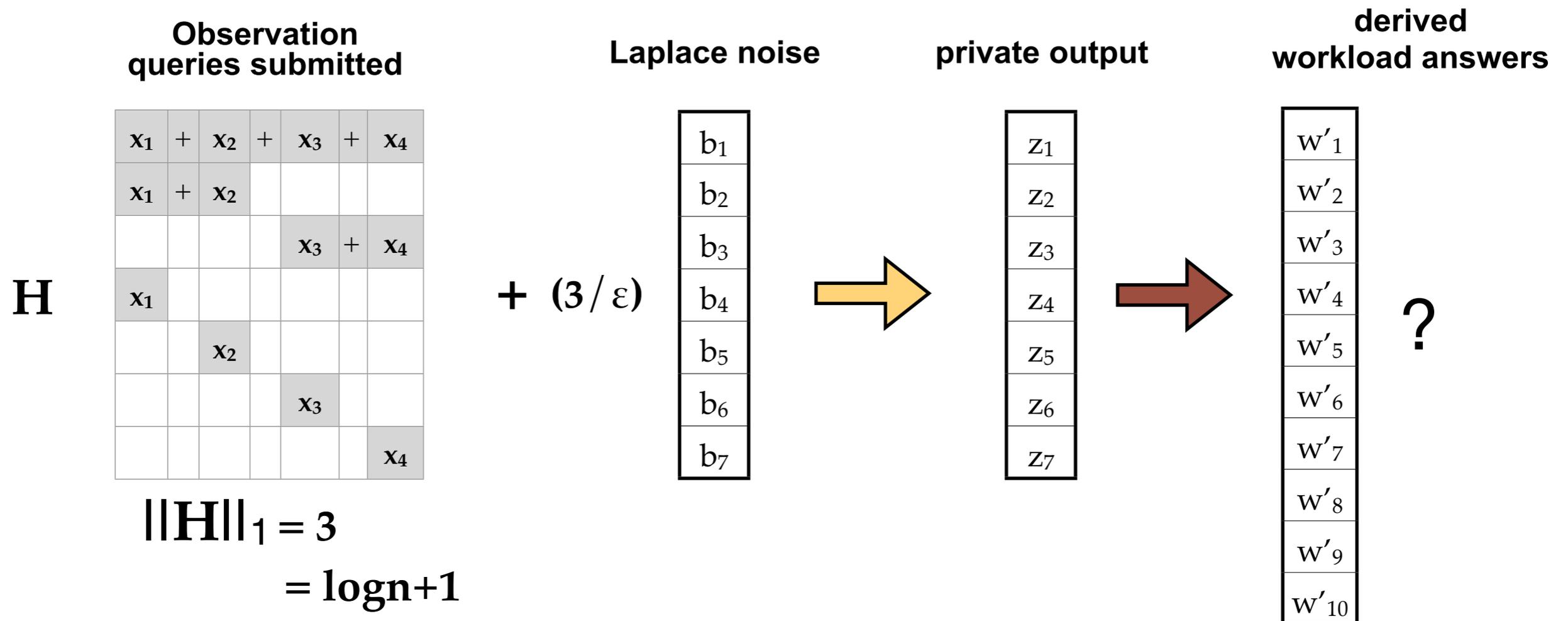
Method 3: hierarchical observations

Hierarchical queries: recursively partition the domain, computing sums of each interval.



Method 3: hierarchical observations

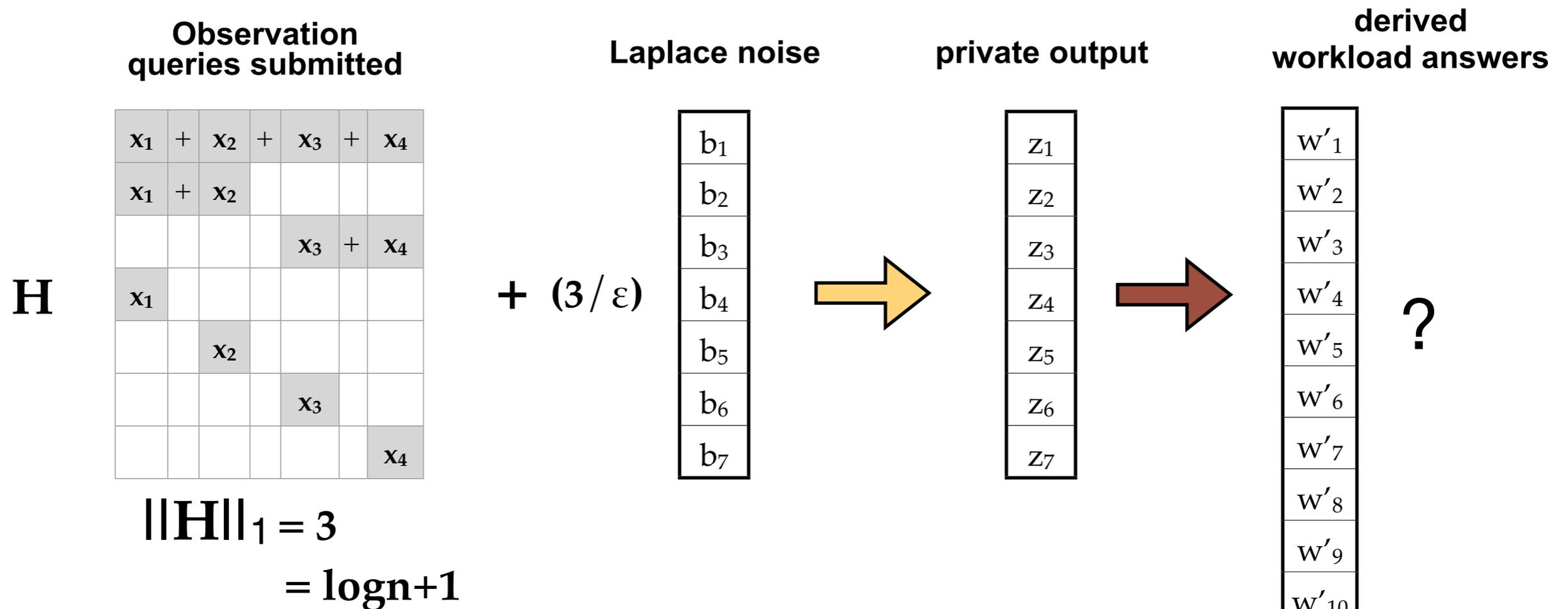
Hierarchical queries: recursively partition the domain, computing sums of each interval.



Possible estimates for query $\text{range}(x_2, x_3) = x_2 + x_3$

Method 3: hierarchical observations

Hierarchical queries: recursively partition the domain, computing sums of each interval.

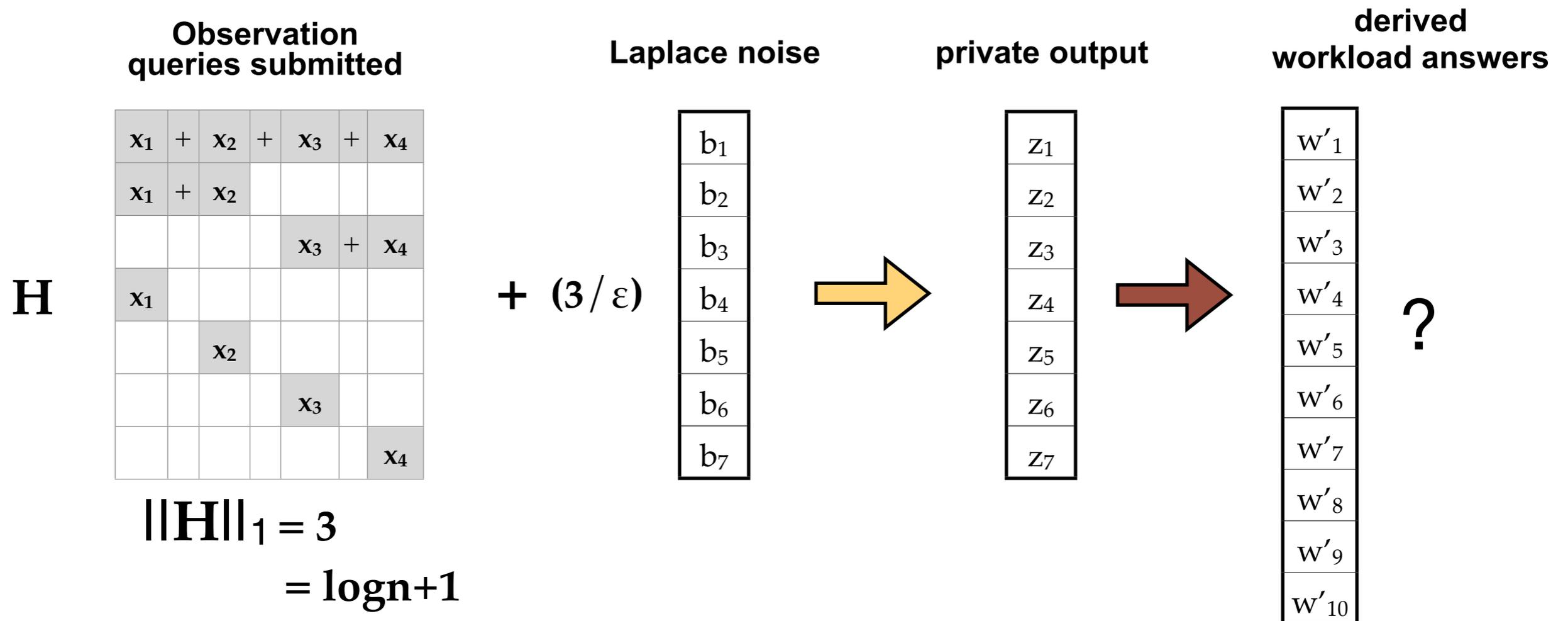


Possible estimates for query $\text{range}(x_2, x_3) = x_2 + x_3$

$$z_5 + z_6$$

Method 3: hierarchical observations

Hierarchical queries: recursively partition the domain, computing sums of each interval.



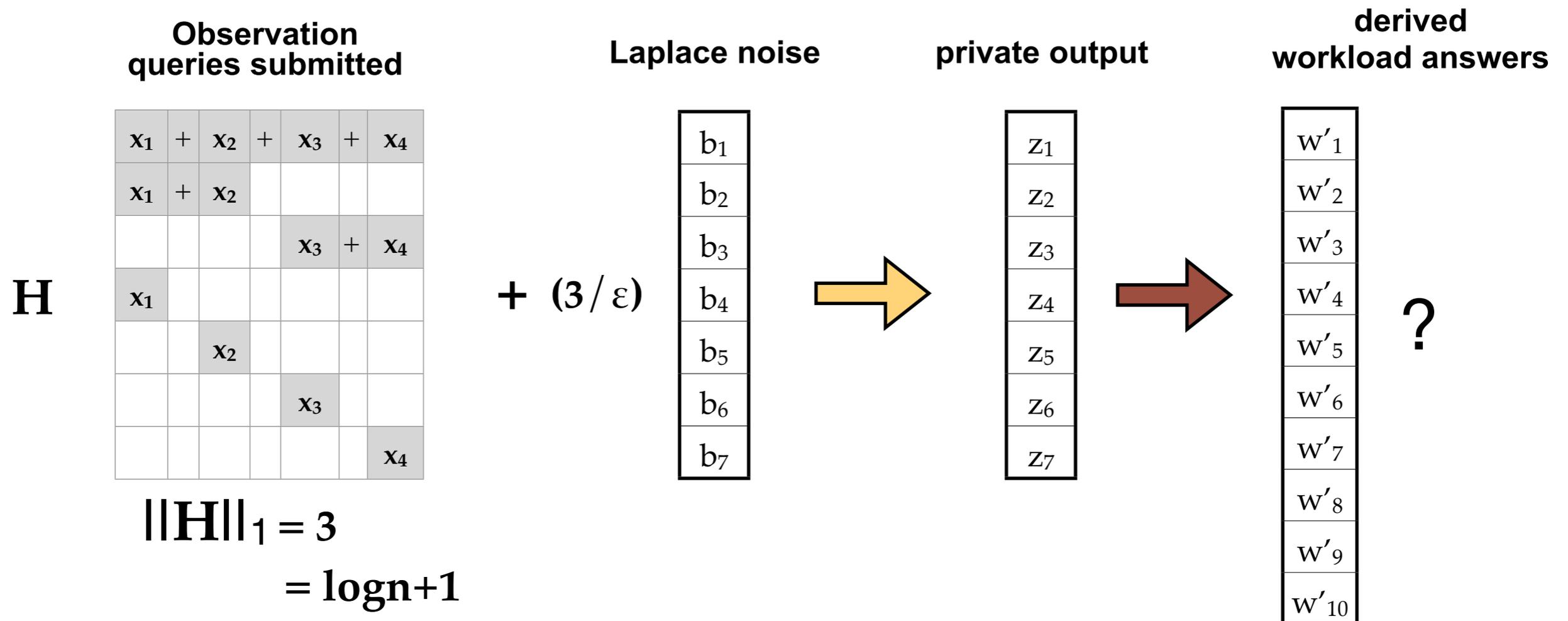
Possible estimates for query $\text{range}(x_2, x_3) = x_2 + x_3$

$$z_5 + z_6$$

$$z_2 - z_4 + z_6$$

Method 3: hierarchical observations

Hierarchical queries: recursively partition the domain, computing sums of each interval.



Possible estimates for query $\text{range}(x_2, x_3) = x_2 + x_3$

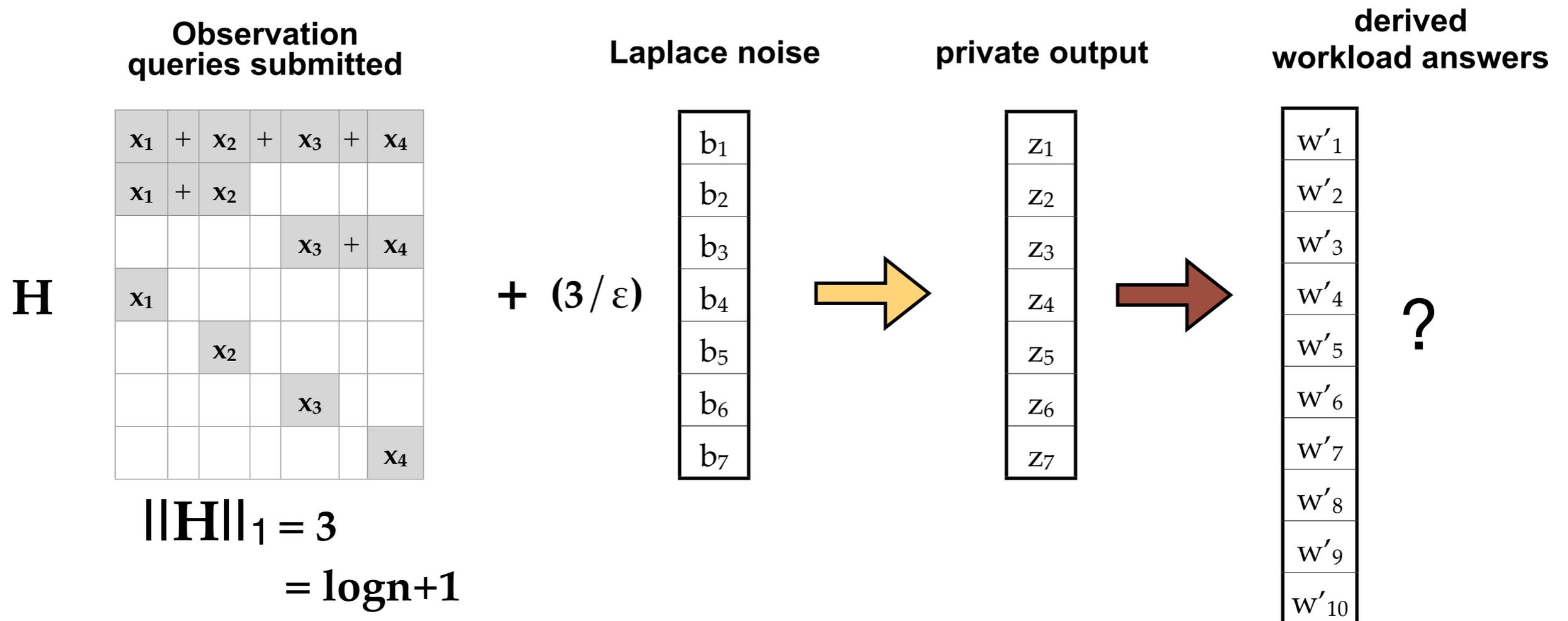
$$z_5 + z_6$$

$$z_2 - z_4 + z_6$$

$$z_1 - z_4 - z_7$$

Method 3: hierarchical observations

Hierarchical queries: recursively partition the domain, computing sums of each interval.



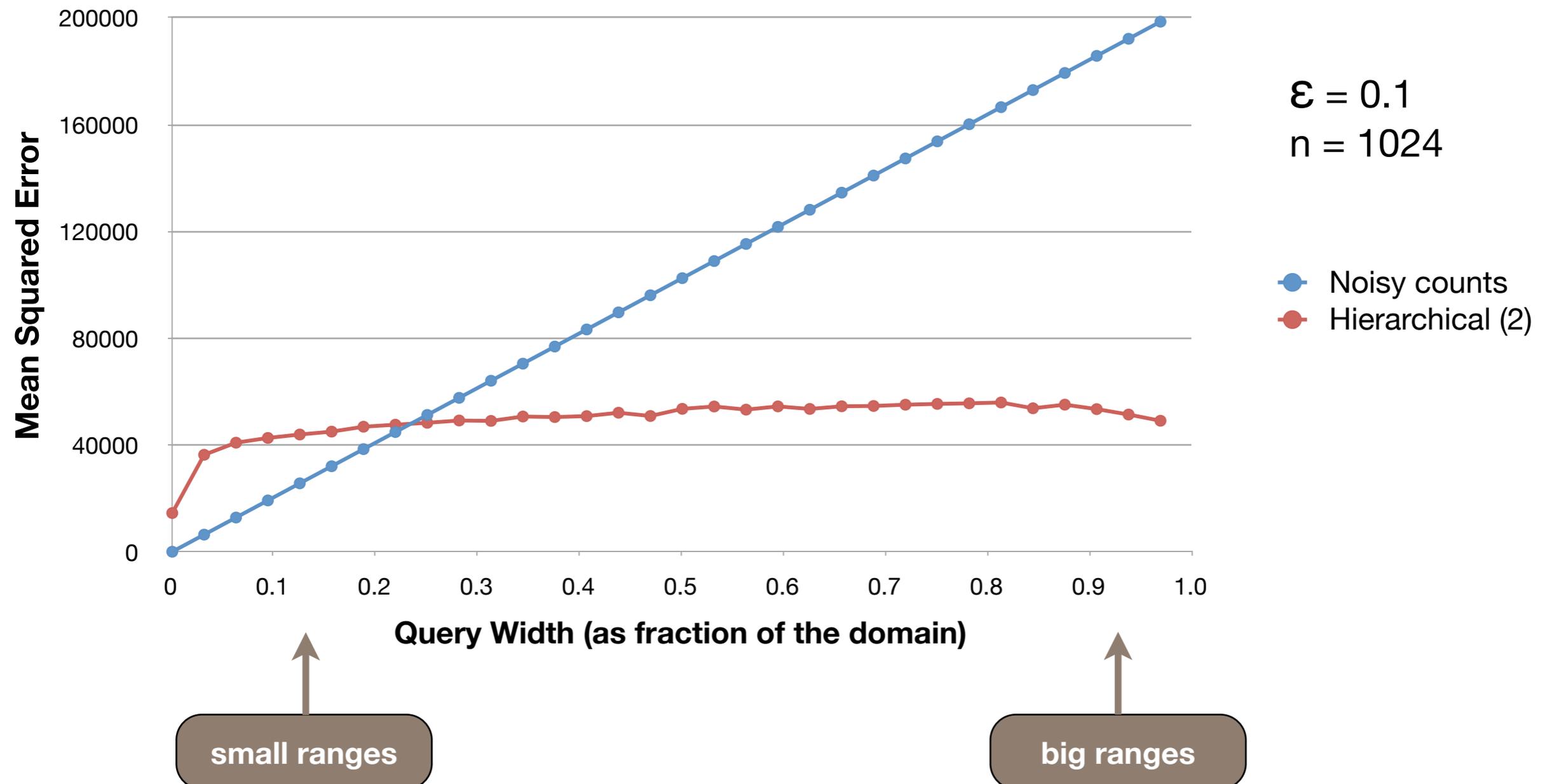
Possible estimates for query $\text{range}(x_2, x_3) = x_2 + x_3$

Least-squares
estimate

$$(6z_1 + 3z_2 + 3z_3 - 9z_4 + 12z_5 + 12z_6 - 9z_7) / 21$$

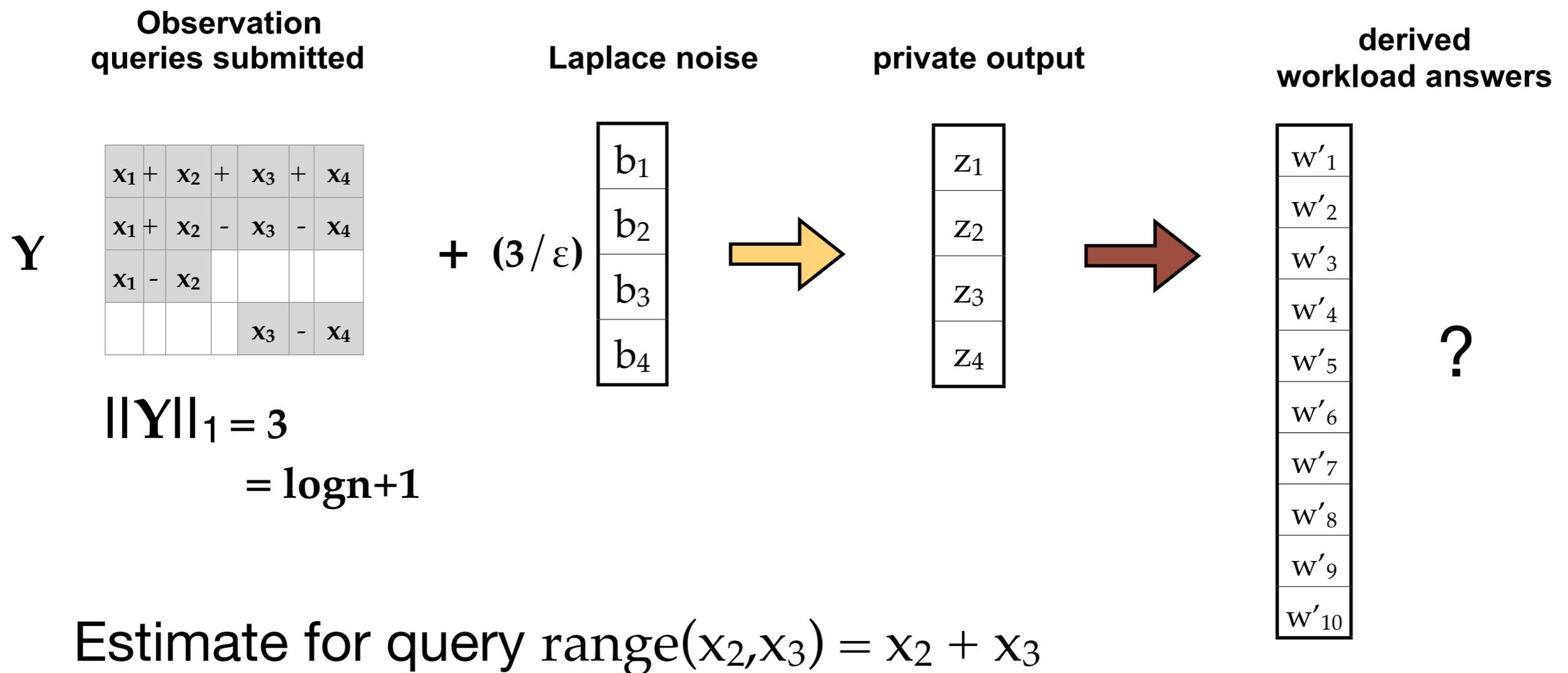
Error rates: workload of all range queries

ϵ -differential privacy



Method 4: wavelet queries

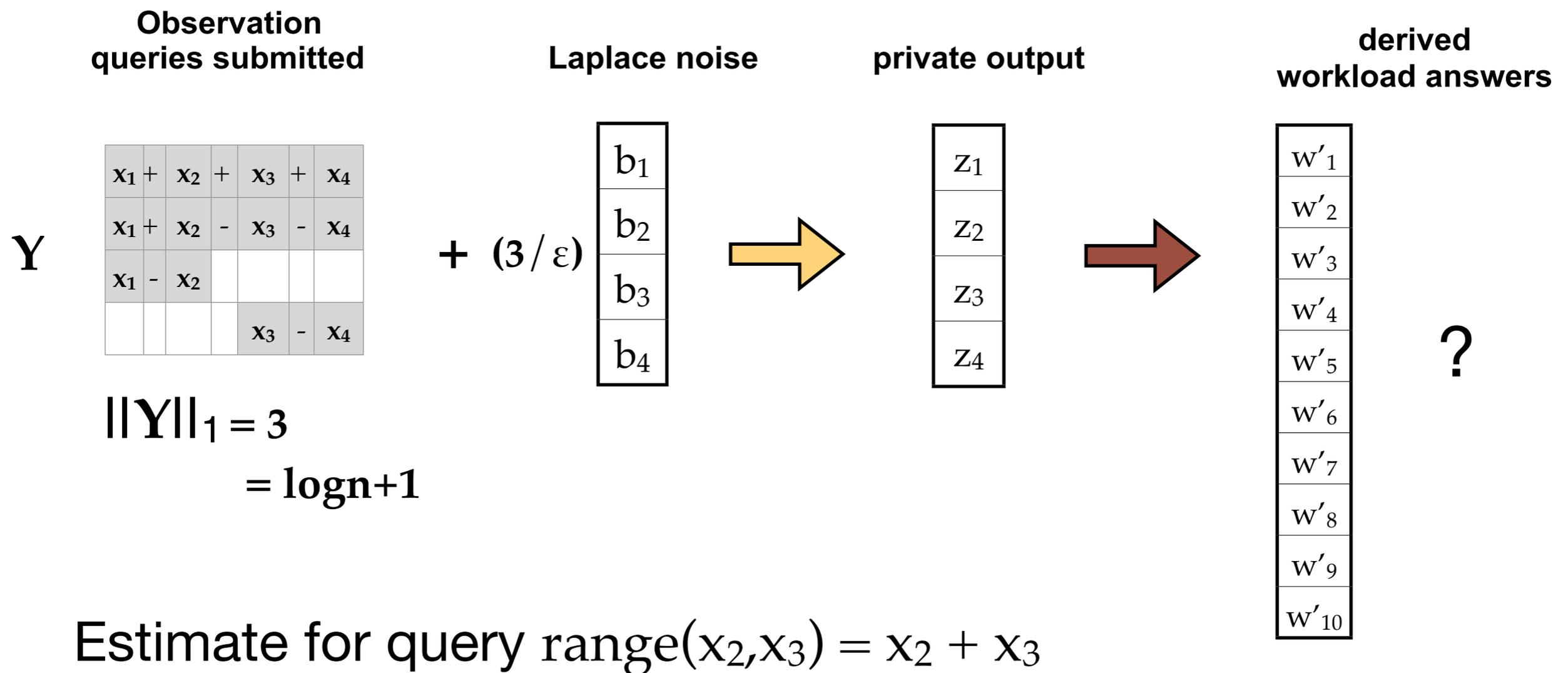
Wavelet: use Haar wavelet as observations.



$$.5z_1 + 0z_2 - .5z_3 + .5z_4$$

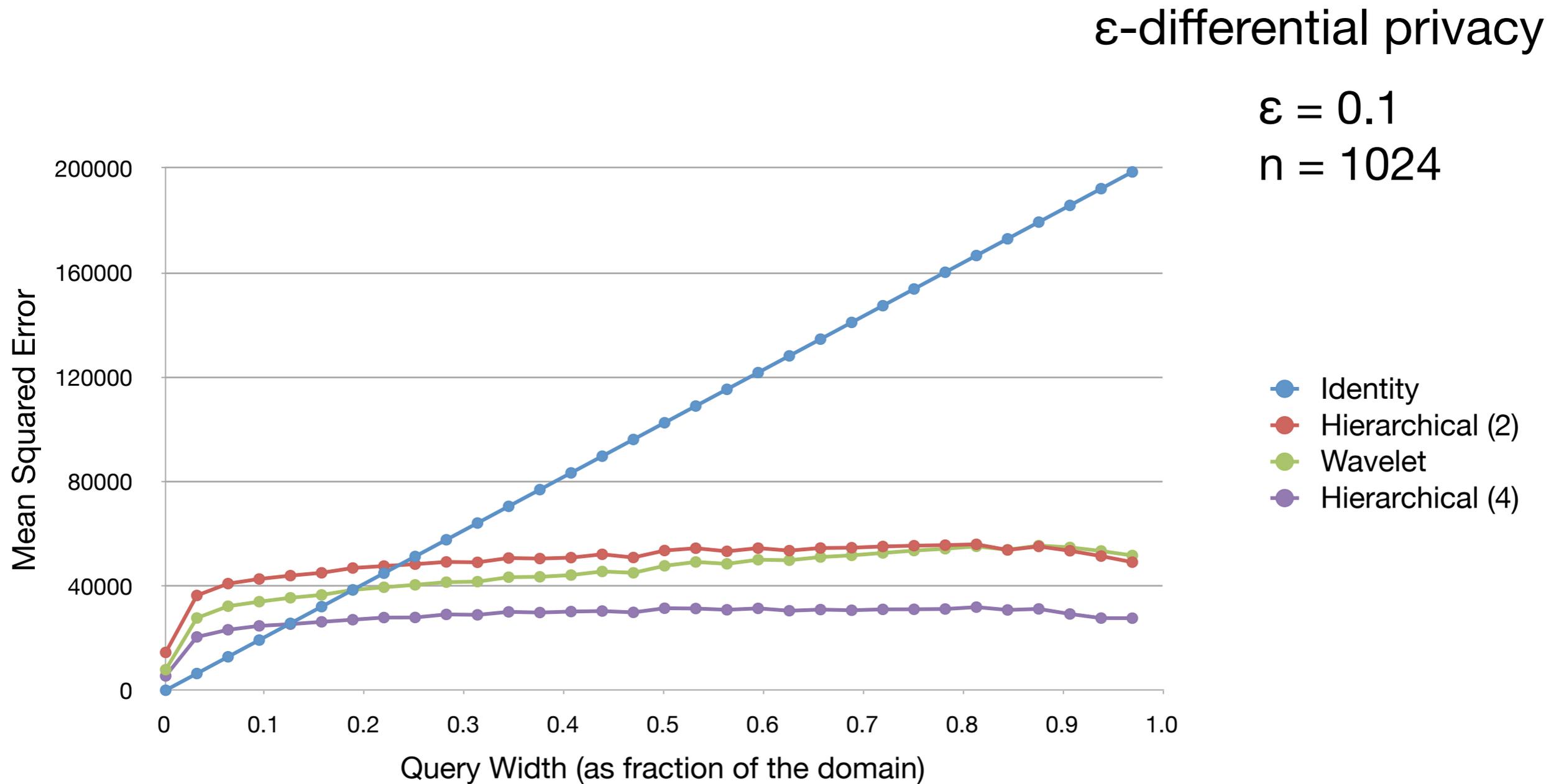
Method 4: wavelet queries

Wavelet: use Haar wavelet as observations.



$$.5z_1 + 0z_2 - .5z_3 + .5z_4$$

Error: workload of all range queries



Observations for the workload of all range queries

Noisy counts

x_1						
		x_2				
				x_3		
						x_4

I

Very low sensitivity, but large ranges estimated badly.

Max/Avg error $O(n/\epsilon^2)$

Hierarchical

x_1	+	x_2	+	x_3	+	x_4
x_1	+	x_2				
				x_3	+	x_4
x_1						
		x_2				
				x_3		
						x_4

H

Low sensitivity, and all range queries can be estimated using no more than $\log n$ output entries.

1-dim $O(\log^3 n / \epsilon^2)$

k-dim

Wavelet

x_1	+	x_2	+	x_3	+	x_4
x_1	+	x_2	-	x_3	-	x_4
x_1	-	x_2				
				x_3	-	x_4

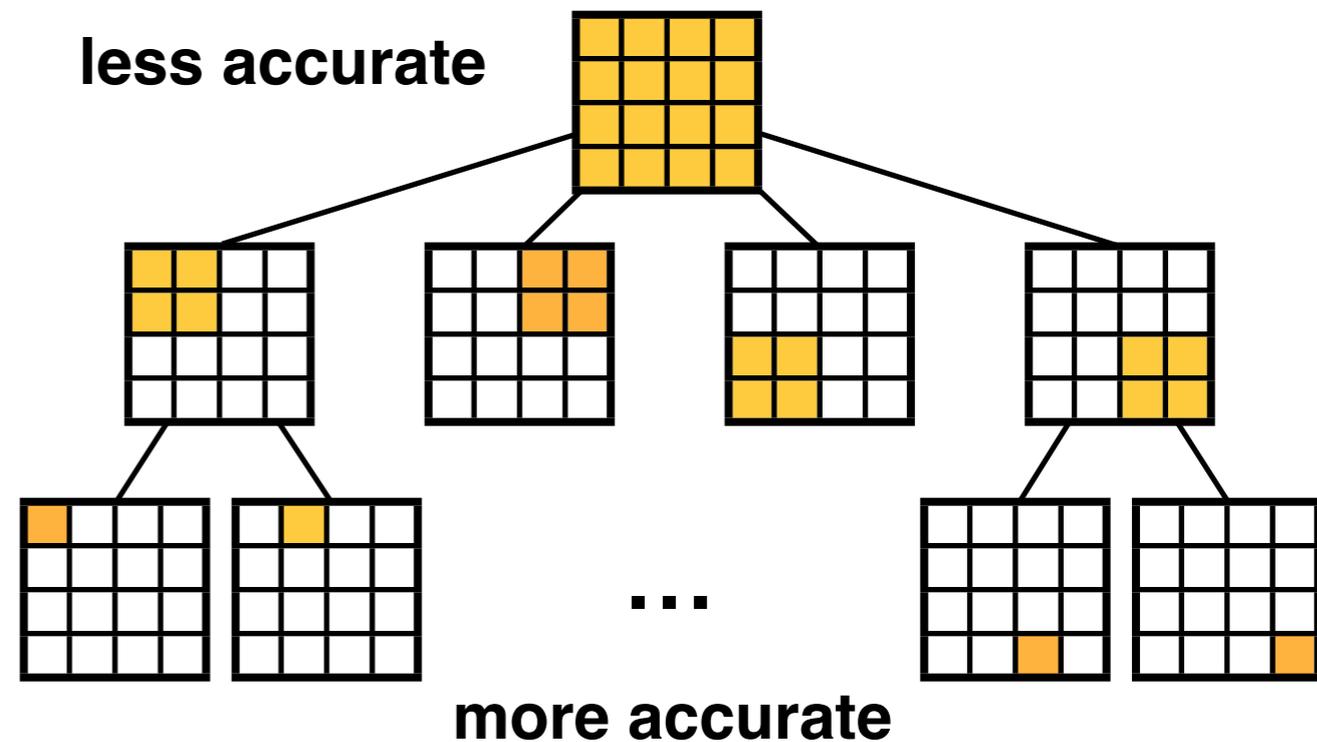
Y

$O(\log^3 n / \epsilon^2)$

$O(\log^{3k} n / \epsilon^2)$

Observations for alternative workloads

- **Workload:** sets of 2D range queries
- **Observations:** [Cormode, ICDE '12]
 - Quad-tree queries
 - Geometrically increasing ϵ by level



- **Workload:** sets of low-order marginals
- **Observations:** [Barak, PODS '07]
 - Fourier basis queries

$$H_i = \begin{bmatrix} H_{i-1} & H_{i-1} \\ H_{i-1} & -H_{i-1} \end{bmatrix}$$

Questions raised

	Workload	Observations	Citation
Non-adaptive	low-order marginals	Fourier basis queries	[Barak, PODS '07]
	all one-dim range queries	Hierarchical ranges	[Hay, PVLDB '10]
	all (multi-dim) range queries	Haar wavelet queries	[Xiao, ICDE '10]
	2-dim range queries	Quad-tree queries	[Cormode, ICDE '12]

- Are these observations optimal for the targeted workloads?
- Which observations should we use for other custom workloads?

Questions raised

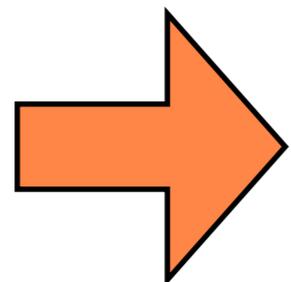
	Workload		Observations	Citation
Non-adaptive	low-order marginals		Fourier basis queries	[Barak, PODS '07]
	all one-dim range queries		Hierarchical ranges	[Hay, PVLDB '10]
	all (multi-dim) range queries		Haar wavelet queries	[Xiao, ICDE '10]
	2-dim range queries		Quad-tree queries	[Cormode, ICDE '12]

- Are these observations optimal for the targeted workloads?
- Which observations should we use for other custom workloads?

Questions raised

	Workload		Observations	Citation
Non-adaptive	low-order marginals	↔	Fourier basis queries	[Barak, PODS '07]
	all one-dim range queries	↔	Hierarchical ranges	[Hay, PVLDB '10]
	all (multi-dim) range queries	↔	Haar wavelet queries	[Xiao, ICDE '10]
	2-dim range queries	↔	Quad-tree queries	[Cormode, ICDE '12]

- Are these observations optimal for the targeted workloads?
- Which observations should we use for other custom workloads?



Adapt observations to workload

Outline

1. Preliminaries

2. Approach 1: workload-aware

- Fixed Observations
- Optimized Observations

3. Approach 2: data-aware

4. Conclusions

Outline

1. Preliminaries

2. Approach 1: workload-aware

- Fixed Observations

- Optimized Observations

3. Approach 2: data-aware

4. Conclusions

Laplace mechanism (matrix notation)

$$\text{Laplace}(W, x) = Wx + (\|W\|_1 / \epsilon) \mathbf{b}$$

W	m×n	workload
x	n×1	database
$\ W\ _1$	scalar	sensitivity
b	m×1	noise: independent samples from Laplace(1)

Laplace mechanism (matrix notation)

$$\text{Laplace}(W, x) = Wx + (\|W\|_1 / \epsilon) \mathbf{b}$$

W	$m \times n$	workload
x	$n \times 1$	database
$\ W\ _1$	scalar	sensitivity
\mathbf{b}	$m \times 1$	noise: independent samples from Laplace(1)

$$\text{Error}(w) = 2 (\|W\|_1 / \epsilon)^2$$

The matrix mechanism: justification

The matrix mechanism: justification

- 1 **(Select Observations)** Choose a (full rank) query matrix A

The matrix mechanism: justification

- ① **(Select Observations)** Choose a (full rank) query matrix A
- ② **(Apply Laplace)** Use the Laplace mechanism to answer A

The matrix mechanism: justification

- 1 **(Select Observations)** Choose a (full rank) query matrix \mathbf{A}
- 2 **(Apply Laplace)** Use the Laplace mechanism to answer \mathbf{A}

$$\mathbf{z} = \mathbf{Ax} + (\|\mathbf{A}\|_1 / \varepsilon) \mathbf{b}$$

The matrix mechanism: justification

① **(Select Observations)** Choose a (full rank) query matrix \mathbf{A}

② **(Apply Laplace)** Use the Laplace mechanism to answer \mathbf{A}

$$\mathbf{z} = \mathbf{Ax} + (\|\mathbf{A}\|_1 / \varepsilon) \mathbf{b}$$

③ **(Derive answers)** Compute estimate $\underline{\mathbf{x}}$ of \mathbf{x} using answers \mathbf{z} .

The matrix mechanism: justification

- 1 **(Select Observations)** Choose a (full rank) query matrix \mathbf{A}
- 2 **(Apply Laplace)** Use the Laplace mechanism to answer \mathbf{A}

$$\mathbf{z} = \mathbf{Ax} + (\|\mathbf{A}\|_1 / \varepsilon) \mathbf{b}$$

- 3 **(Derive answers)** Compute estimate $\underline{\mathbf{x}}$ of \mathbf{x} using answers \mathbf{z} .
 - compute estimate $\underline{\mathbf{x}}$ of \mathbf{x} that minimizes squared error:

$$\|\mathbf{A}\underline{\mathbf{x}} - \mathbf{z}\|_2^2$$

The matrix mechanism: justification

- 1 **(Select Observations)** Choose a (full rank) query matrix \mathbf{A}
- 2 **(Apply Laplace)** Use the Laplace mechanism to answer \mathbf{A}

$$\mathbf{z} = \mathbf{Ax} + (\|\mathbf{A}\|_1 / \varepsilon) \mathbf{b}$$

- 3 **(Derive answers)** Compute estimate $\underline{\mathbf{x}}$ of \mathbf{x} using answers \mathbf{z} .

- compute estimate $\underline{\mathbf{x}}$ of \mathbf{x} that minimizes squared error:

$$\|\mathbf{A}\underline{\mathbf{x}} - \mathbf{z}\|_2^2$$

- solution is the ordinary least squares estimator:

$$\underline{\mathbf{x}} = \mathbf{A}^+ \mathbf{z}$$

$$\text{where } \mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$$

The matrix mechanism: justification

- 1 **(Select Observations)** Choose a (full rank) query matrix \mathbf{A}
- 2 **(Apply Laplace)** Use the Laplace mechanism to answer \mathbf{A}

$$\mathbf{z} = \mathbf{Ax} + (\|\mathbf{A}\|_1 / \epsilon) \mathbf{b}$$

- 3 **(Derive answers)** Compute estimate $\underline{\mathbf{x}}$ of \mathbf{x} using answers \mathbf{z} .

- compute estimate $\underline{\mathbf{x}}$ of \mathbf{x} that minimizes squared error:

$$\|\mathbf{Ax} - \mathbf{z}\|_2^2$$

- solution is the ordinary least squares estimator:

$$\underline{\mathbf{x}} = \mathbf{A}^+ \mathbf{z}$$

$$\text{where } \mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$$

Thm: $\underline{\mathbf{x}}$ is unbiased and has the least variance among all linear unbiased estimators.

The matrix mechanism: justification

- 1 **(Select Observations)** Choose a (full rank) query matrix \mathbf{A}
- 2 **(Apply Laplace)** Use the Laplace mechanism to answer \mathbf{A}

$$\mathbf{z} = \mathbf{Ax} + (\|\mathbf{A}\|_1 / \epsilon) \mathbf{b}$$

- 3 **(Derive answers)** Compute estimate $\underline{\mathbf{x}}$ of \mathbf{x} using answers \mathbf{z} .

- compute estimate $\underline{\mathbf{x}}$ of \mathbf{x} that minimizes squared error:

$$\|\mathbf{A}\underline{\mathbf{x}} - \mathbf{z}\|_2^2$$

- solution is the ordinary least squares estimator:

$$\underline{\mathbf{x}} = \mathbf{A}^+ \mathbf{z}$$

$$\text{where } \mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$$

- Compute workload queries using estimate $\underline{\mathbf{x}}$:

$$W_{\underline{\mathbf{x}}}$$

Thm: $\underline{\mathbf{x}}$ is unbiased and has the least variance among all linear unbiased estimators.

The matrix mechanism

Given a workload W , and any full-rank strategy matrix A , the following randomized algorithm is ϵ -differentially private:

$$\text{Matrix}_A(W, x) = Wx + (\|A\|_1 / \epsilon) WA^+ \mathbf{b} \quad \mathbf{b} = \text{Lap}(1)$$

The matrix mechanism

Given a workload W , and any full-rank strategy matrix A , the following randomized algorithm is ϵ -differentially private:

$$\text{Matrix}_A(W, x) = Wx + (\|A\|_1 / \epsilon) WA^+ \mathbf{b} \quad \mathbf{b} = \text{Lap}(1)$$

instantiated with
observations A

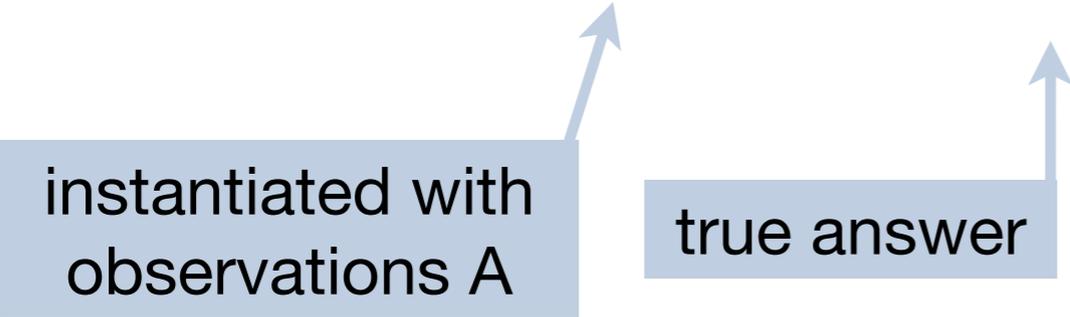


The matrix mechanism

Given a workload W , and any full-rank strategy matrix A , the following randomized algorithm is ϵ -differentially private:

$$\text{Matrix}_A(W, x) = Wx + (\|A\|_1 / \epsilon) WA^+ \mathbf{b} \quad \mathbf{b} = \text{Lap}(1)$$

instantiated with
observations A



true answer

The matrix mechanism

Given a workload W , and any full-rank strategy matrix A , the following randomized algorithm is ϵ -differentially private:

$$\text{Matrix}_A(W, x) = Wx + (\|A\|_1 / \epsilon) WA^+ \mathbf{b} \quad \mathbf{b} = \text{Lap}(1)$$

instantiated with
observations A

true answer

scaling by
 $\|A\|_1$

The matrix mechanism

Given a workload W , and any full-rank strategy matrix A , the following randomized algorithm is ϵ -differentially private:

$$\text{Matrix}_A(W, x) = Wx + (\|A\|_1 / \epsilon) WA^+ \mathbf{b} \quad \mathbf{b} = \text{Lap}(1)$$

instantiated with
observations A

true answer

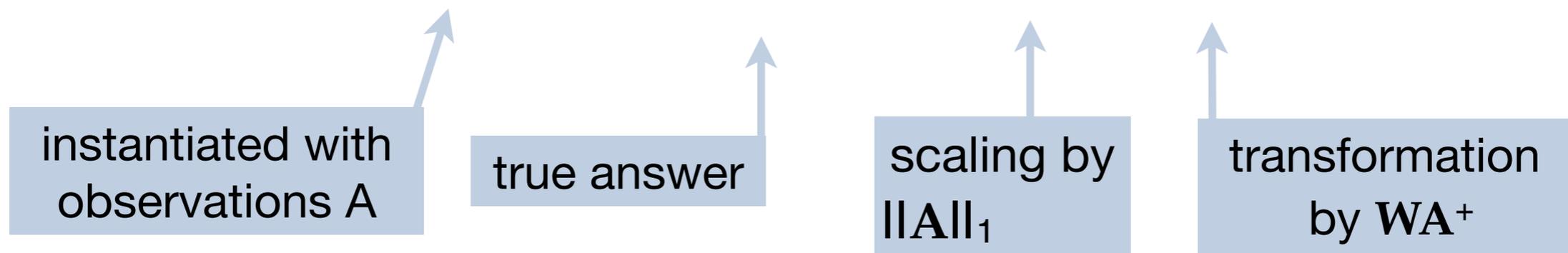
scaling by
 $\|A\|_1$

transformation
by WA^+

The matrix mechanism

Given a workload W , and any full-rank strategy matrix A , the following randomized algorithm is ϵ -differentially private:

$$\text{Matrix}_A(W, x) = Wx + (\|A\|_1 / \epsilon) WA^+ \mathbf{b} \quad \mathbf{b} = \text{Lap}(1)$$



Compare with the Laplace mechanism:

$$\text{Laplace}(W, x) = Wx + (\|W\|_1 / \epsilon) \mathbf{b}$$

Instances of the matrix mechanism

Given workload W of linear queries:

Observation Matrix A	Resulting mechanism
$A = W$	Never worse than Laplace -- sometimes better
$A =$ Identity matrix	a common baseline
$A =$ Haar wavelet	[Xiao, ICDE '10]
$A =$ tree based	[Hay, PVLDB '10] [Cormode, ICDE '12]
$A =$ fourier basis	[Barak, PODS '07]

Observation matrices equivalent to wavelet

1	1	1	1
1	1	-1	-1
1	-1	0	0
0	0	1	-1

Wavelet Y

$$\|Y\|_1 = 3$$

\equiv

1	1	0	0
0	0	1	1
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

Y'

$$\|Y'\|_1 = 3$$

\succ

1	1	0	0
0	0	1	1
$\sqrt{2}$	0	0	0
0	$\sqrt{2}$	0	0
0	0	$\sqrt{2}$	0
0	0	0	$\sqrt{2}$

Y''

$$\|Y''\|_1 = 2.414$$

Observation matrices equivalent to wavelet

1	1	1	1
1	1	-1	-1
1	-1	0	0
0	0	1	-1

≡

Equivalent error for all queries

1	1	0	0
0	0	1	1
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

>

1	1	0	0
0	0	1	1
$\sqrt{2}$	0	0	0
0	$\sqrt{2}$	0	0
0	0	$\sqrt{2}$	0
0	0	0	$\sqrt{2}$

Wavelet Y

$$\|Y\|_1 = 3$$

Y'

$$\|Y'\|_1 = 3$$

Y''

$$\|Y''\|_1 = 2.414$$

Observation matrices equivalent to wavelet

1	1	1	1
1	1	-1	-1
1	-1	0	0
0	0	1	-1

Wavelet Y

$$\|Y\|_1 = 3$$

\equiv

Equivalent
error for all
queries

1	1	0	0
0	0	1	1
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

Y'

$$\|Y'\|_1 = 3$$

$>$

Lower
error for all
queries

1	1	0	0
0	0	1	1
$\sqrt{2}$	0	0	0
0	$\sqrt{2}$	0	0
0	0	$\sqrt{2}$	0
0	0	0	$\sqrt{2}$

Y''

$$\|Y''\|_1 = 2.414$$

Observation matrices equivalent to wavelet

1	1	1	1
1	1	-1	-1
1	-1	0	0
0	0	1	-1

≡
Equivalent error for all queries

1	1	0	0
0	0	1	1
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

>
Lower error for all queries

1	1	0	0
0	0	1	1
$\sqrt{2}$	0	0	0
0	$\sqrt{2}$	0	0
0	0	$\sqrt{2}$	0
0	0	0	$\sqrt{2}$

Wavelet Y

$$\|Y\|_1 = 3$$

Y'

$$\|Y'\|_1 = 3$$

Y''

$$\|Y''\|_1 = 2.414$$

The haar wavelet observation matrix Y is **dominated** by alternative matrix Y'' .

Error of matrix mechanism

Given an observation matrix \mathbf{A} and workload \mathbf{W} , the error under the mechanism \mathbf{Matrix}_A is:

For a single query \mathbf{w} in \mathbf{W} :

$$\mathbf{Error}_A(\mathbf{w}) = (2 / \varepsilon^2)(\|\mathbf{A}\|_1)^2 \mathbf{w}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{w}^T$$

Total error for workload \mathbf{W} :

$$\mathbf{TotalError}_A(\mathbf{w}) = (2 / \varepsilon^2)(\|\mathbf{A}\|_1)^2 \text{trace}(\mathbf{W}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{W}^T)$$

Error independent of the input data

Optimal selection of observations

Objective: given workload \mathbf{W} , find the observation matrix \mathbf{A} that minimizes the **total** error.

Optimal selection of observations

Objective: given workload \mathbf{W} , find the observation matrix \mathbf{A} that minimizes the **total** error.

Privacy	Optimization Objective	Problem Type	Runtime
---------	------------------------	--------------	---------

Optimal selection of observations

Objective: given workload \mathbf{W} , find the observation matrix \mathbf{A} that minimizes the **total** error.

Privacy	Optimization Objective	Problem Type	Runtime
ϵ DP	Given \mathbf{W} consisting of data cube queries, choose \mathbf{A} consisting of data cube queries to minimize simplified error measure. [Ding, SIGMOD '11]	set-cover approx	$O(n)$

Optimal selection of observations

Objective: given workload W , find the observation matrix A that minimizes the **total** error.

Privacy	Optimization Objective	Problem Type	Runtime
ϵ DP	Given W consisting of data cube queries, choose A consisting of data cube queries to minimize simplified error measure. [Ding, SIGMOD '11]	set-cover approx	$O(n)$
ϵ DP	Given W , choose A to minimize $\text{TotalError}_A(W)$ [Li, PODS '10]	SDP w/ rank constraints	$O(n^8)$

Optimal selection of observations

Objective: given workload W , find the observation matrix A that minimizes the **total** error.

Privacy	Optimization Objective	Problem Type	Runtime
ϵ DP	Given W consisting of data cube queries, choose A consisting of data cube queries to minimize simplified error measure. [Ding, SIGMOD '11]	set-cover approx	$O(n)$
ϵ DP	Given W , choose A to minimize $\text{TotalError}_A(W)$ [Li, PODS '10]	SDP w/ rank constraints	$O(n^8)$
(ϵ, δ) DP	Given W , choose A to minimize $\text{TotalError}_A(W)$ [Li, PODS '10]	SDP	$O(n^8)$

Optimal selection of observations

Objective: given workload W , find the observation matrix A that minimizes the **total** error.

Privacy	Optimization Objective	Problem Type	Runtime
ϵ DP	Given W consisting of data cube queries, choose A consisting of data cube queries to minimize simplified error measure. [Ding, SIGMOD '11]	set-cover approx	$O(n)$
ϵ DP	Given W , choose A to minimize $\text{TotalError}_A(W)$ [Li, PODS '10]	SDP w/ rank constraints	$O(n^8)$
(ϵ, δ) DP	Given W , choose A to minimize $\text{TotalError}_A(W)$ [Li, PODS '10]	SDP	$O(n^8)$
ϵ DP	Given W , choose $AB \approx W$ to minimize $\text{TotalError}_A(AB)$ [Yuan, VLDB '12]	bi-convex opt	$O(n^4)$

Optimal selection of observations

Objective: given workload W , find the observation matrix A that minimizes the **total** error.

Privacy	Optimization Objective	Problem Type	Runtime
ϵ DP	Given W consisting of data cube queries, choose A consisting of data cube queries to minimize simplified error measure. [Ding, SIGMOD '11]	set-cover approx	$O(n)$
ϵ DP	Given W , choose A to minimize $\text{TotalError}_A(W)$ [Li, PODS '10]	SDP w/ rank constraints	$O(n^8)$
(ϵ, δ) DP	Given W , choose A to minimize $\text{TotalError}_A(W)$ [Li, PODS '10]	SDP	$O(n^8)$
ϵ DP	Given W , choose $AB \approx W$ to minimize $\text{TotalError}_A(AB)$ [Yuan, VLDB '12]	bi-convex opt	$O(n^4)$
(ϵ, δ) DP	Given W , choose optimal scaling of eigenvectors of W to minimize $\text{TotalError}_A(W)$ [Li, PVLDB '12]	convex opt	$O(n^4)$

Approximately optimal selection of observations

Matrix Mechanism under (ϵ, δ) -Differential Privacy

- Given \mathbf{W} , choose a set of **basis queries** for the observations:

- $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ (the eigenvectors of \mathbf{W})

- compute optimal scalars to minimize error c_1, c_2, \dots, c_n

- Resulting observation matrix is:

$$\mathbf{A} = \begin{bmatrix} c_1 \mathbf{v}_1 \\ c_2 \mathbf{v}_2 \\ \dots \\ c_n \mathbf{v}_n \end{bmatrix}$$

- Algorithm running time: $O(n \text{rank}(\mathbf{W})^3)$

- Efficiently solvable and achieves optimal error rates in practice.

Representative experimental findings

- **Benefit of fixed observations:**

- $W=\{\text{All Range Queries}\}$ can be reduced by a factor of 2-4 by using wavelet or hierarchical observations. [Xiao, ICDE '10] [Hay, PVLDB '10]

- **Benefit of optimized observations:**

- ϵ -DP: Error reduced by 2-3 times compared with fixed observation methods. [Yuan, VLDB '12]
- (ϵ, δ) -DP: Error reduced by 2-6 times on range and marginal workloads for which fixed observation methods were designed; up to 10 times reduction for ad hoc workloads. [Li, PVLDB '12]

Note 1: comparisons don't depend on input data or privacy parameters.*

Note 2: ratios based on **root** mean squared error.

Lower bound on error

- Given workload \mathbf{W} with singular values $\lambda_1 > \dots > \lambda_n$, the minimum total error of the matrix mechanism is greater than or equal to:

Privacy	Error Lower Bound
ϵ -DP	$(2 / \epsilon^2)(1 / n)(\lambda_1 + \dots + \lambda_n)^2$
(ϵ, δ) -DP	$(2 \log(2 / \delta) / \epsilon^2)(1 / n)(\lambda_1 + \dots + \lambda_n)^2$ (tight)

Runtime complexity

- Answering \mathbf{W} using Laplace/Gaussian mechanism takes $O(|\mathbf{W}|n)$ time.

Costs	Fixed Observations	Optimized Observations
1. Select observations	-	$\sim O(n^4)$
2. Apply standard mechanism	$O(\mathbf{A} n)$	$O(\mathbf{A} n)$
3. Derive answers	$O(\mathbf{W} n)$	$O(\mathbf{W} n^2)$

- Because of data-independence, observation matrix can be preprocessed:
 - Given fixed workload \mathbf{W} and observation matrix \mathbf{A} , runtime is $O(|\mathbf{W}|n)$ after pre-computation of \mathbf{WA}^+ : no worse than standard mechanisms

Summary: workload-aware mechanisms

- Methods can be seen as a generalization of Laplace/Gaussian mechanism, with error rates significantly reduced and independent of data.

Workload		Observations	Citation
low-order marginals	Fixed	Fourier basis queries	[Barak, PODS '07]
all one-dim range queries		Hierarchical ranges	[Hay, PVLDB '10]
all (multi-dim) range queries		Haar wavelet queries	[Xiao, ICDE '10]
2-dim range queries		Quad-tree queries	[Cormode, ICDE '12]
sets of data cubes	Optimized	sets of data cubes	[Ding, SIGMOD '11]
set of linear queries		set of linear queries	[Li, PODS '10] [Li, PVLDB '12]
set of linear queries		low-order set of linear queries	[Yuan, VLDB '12]

Summary: workload-aware mechanisms

- **Benefits**

- Independence of data makes error analysis easy, error rates publishable to analyst, and improves efficiency in some cases.

- **Limitations**

- Computational dependence on domain size, n .
- Error dependence on epsilon: $1/\epsilon^2$
- For some workloads, there is no set of observations that can help much.

- **Open questions**

- Alternative derivation methods: e.g. non-negative least squares
- Relationship with “universal” error lower bounds for DP.

Outline

1. Preliminaries

2. Approach 1: workload-aware

- Fixed Observations

- Optimized Observations

3. Approach 2: data-aware

4. Conclusions

Outline

1. Preliminaries

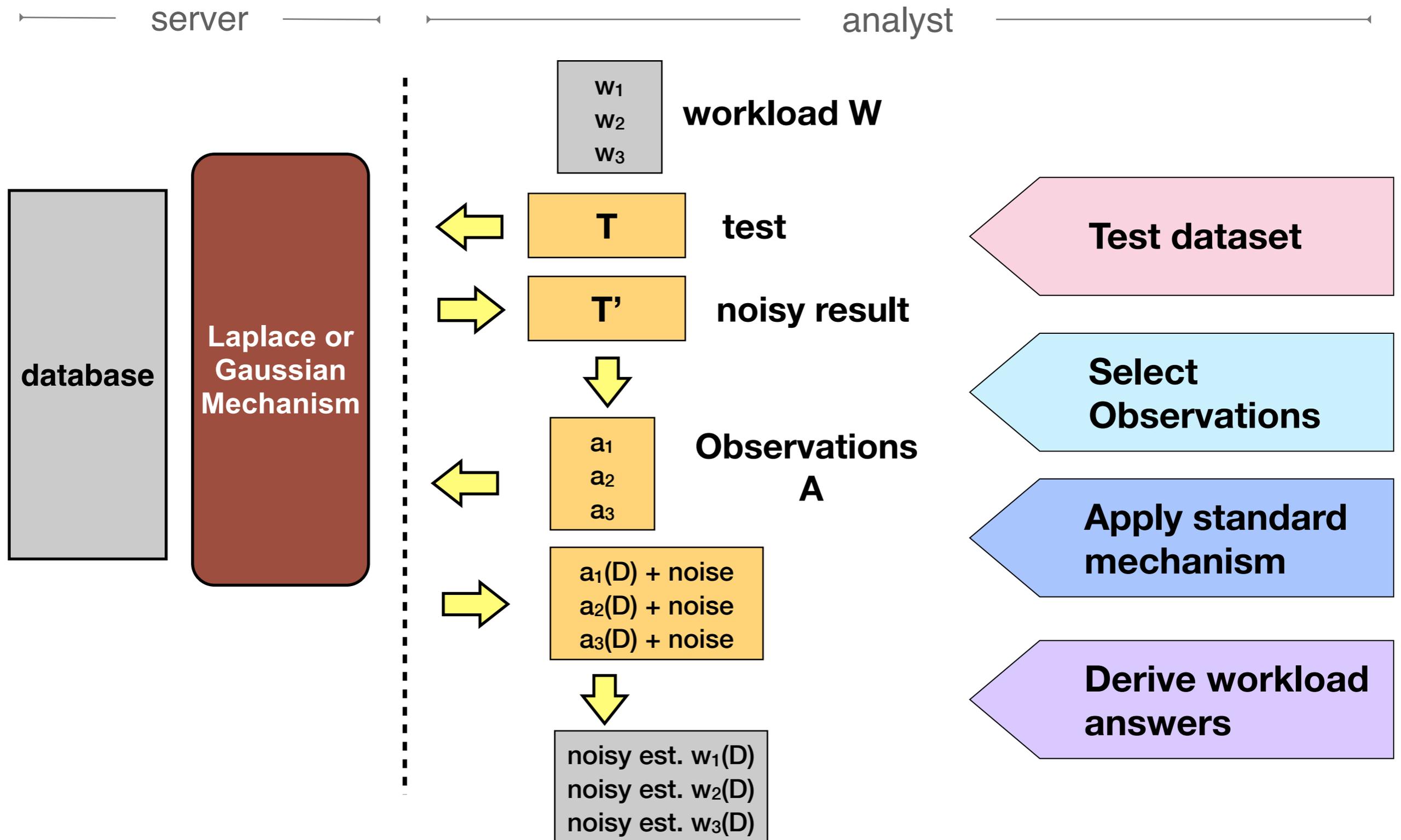
2. Approach 1: workload-aware

- Fixed Observations
- Optimized Observations

3. Approach 2: data-aware

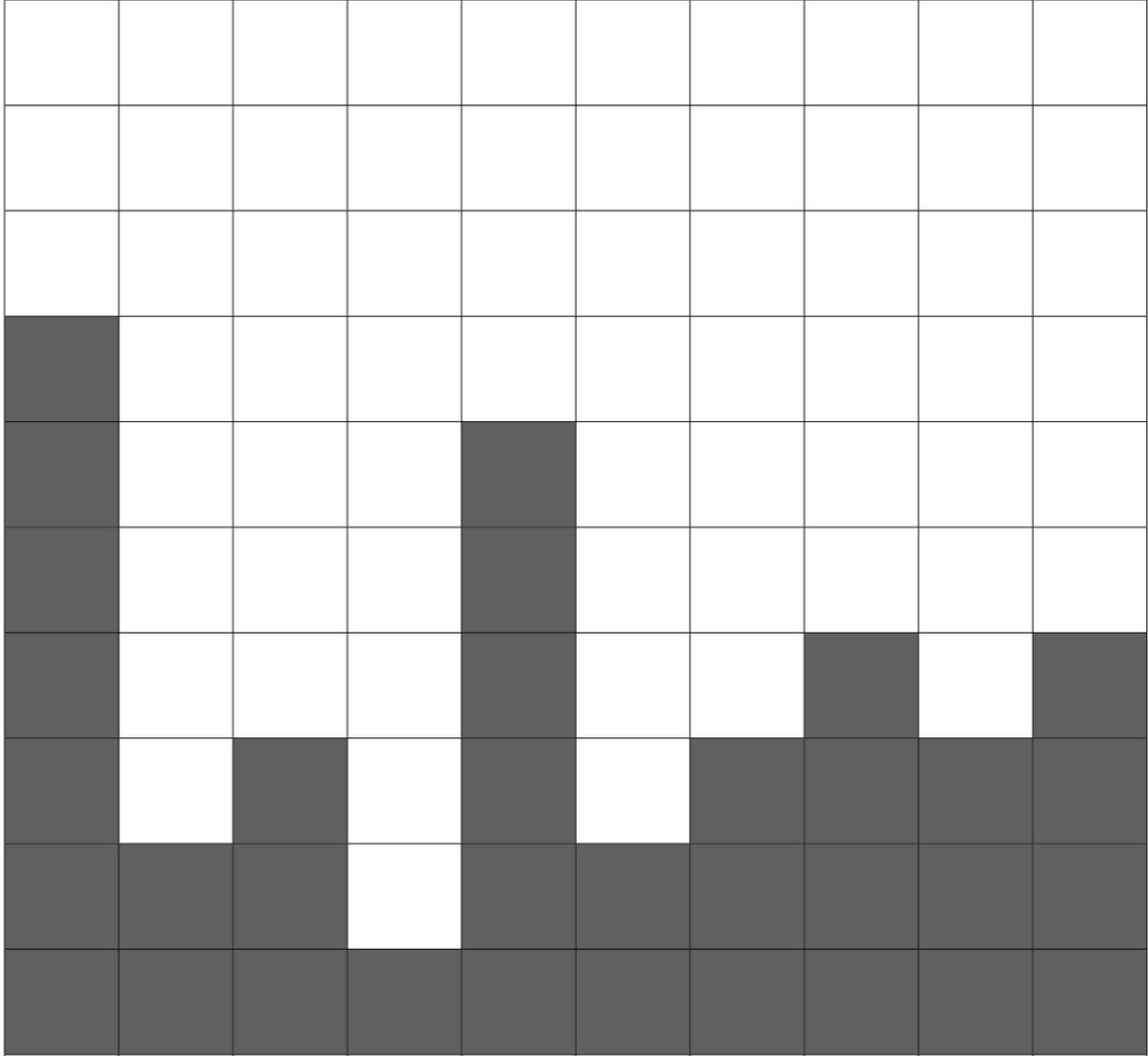
4. Conclusions

(Recall) Approach 2: data-aware mechanisms



A basic intuition

- Detect when additional observations won't help much.
- Challenges:
 - Balance privacy budget between testing data and usable observations.
 - When possible, incorporate test observations into query answers.
 - Perturbation error vs. approximation error.



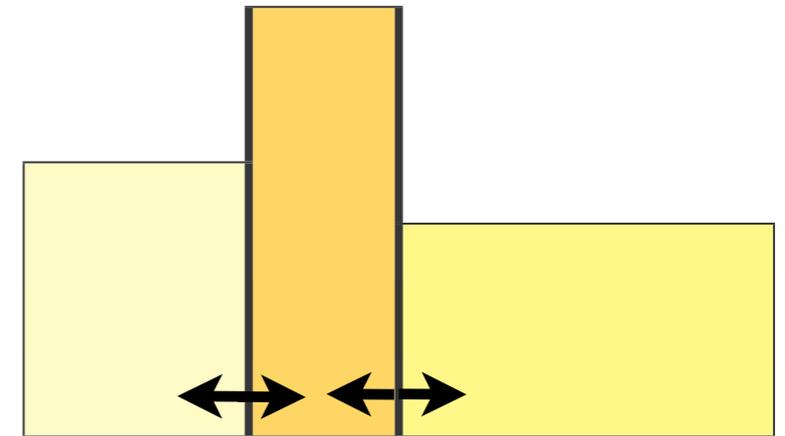
A 10x10 grid with columns labeled x_1 through x_{10} at the bottom. The grid contains several shaded gray cells, indicating a specific pattern of observations or data points. The shaded cells are located at the following coordinates (row, column): (1,1), (2,1), (3,1), (4,1), (4,5), (5,1), (5,3), (5,5), (5,7), (5,8), (5,9), (5,10), (6,1), (6,2), (6,3), (6,5), (6,6), (6,7), (6,8), (6,9), (6,10), (7,1), (7,2), (7,3), (7,4), (7,5), (7,6), (7,7), (7,8), (7,9), (7,10).

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}

Data-aware histogram

[Xu, ICDE '12]

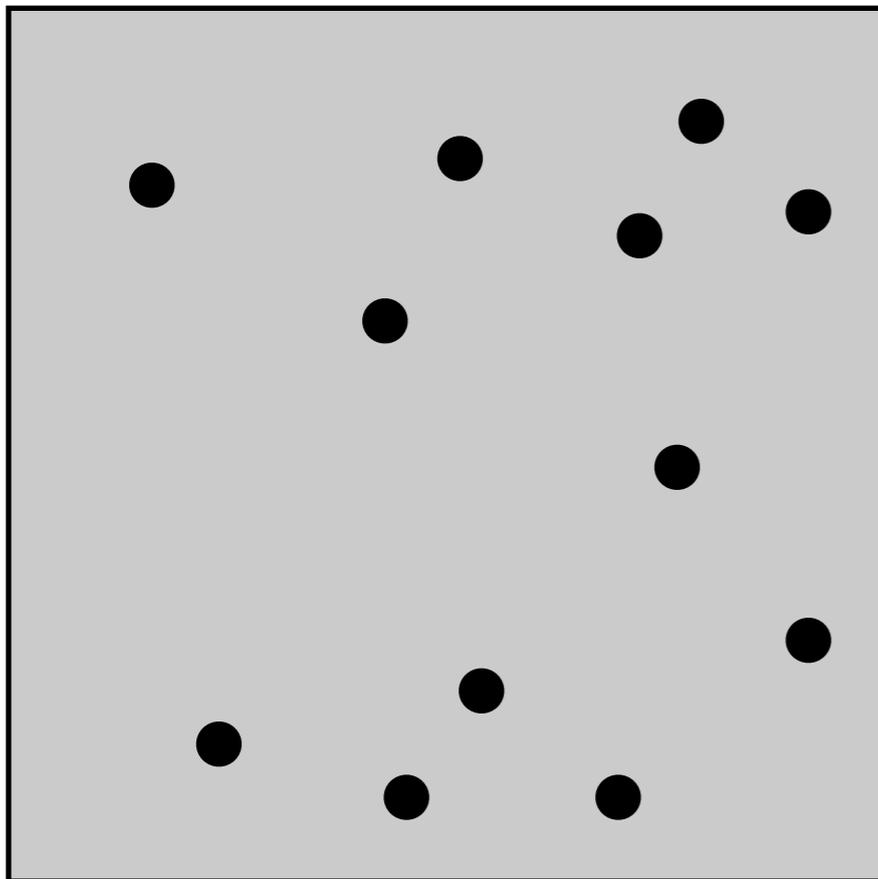
Workload	1D Range Queries
Parameters	$k, \epsilon_1, \epsilon_2$ s.t. $\epsilon_1 + \epsilon_2 = \epsilon$



1. Compute a private estimate of the k -bin, variance-optimal histogram using the exponential mechanism. ϵ_1
2. Use Laplace mechanism to get bin counts **and** all individual counts. ϵ_2
3. Derive answers to workload queries using least squares.

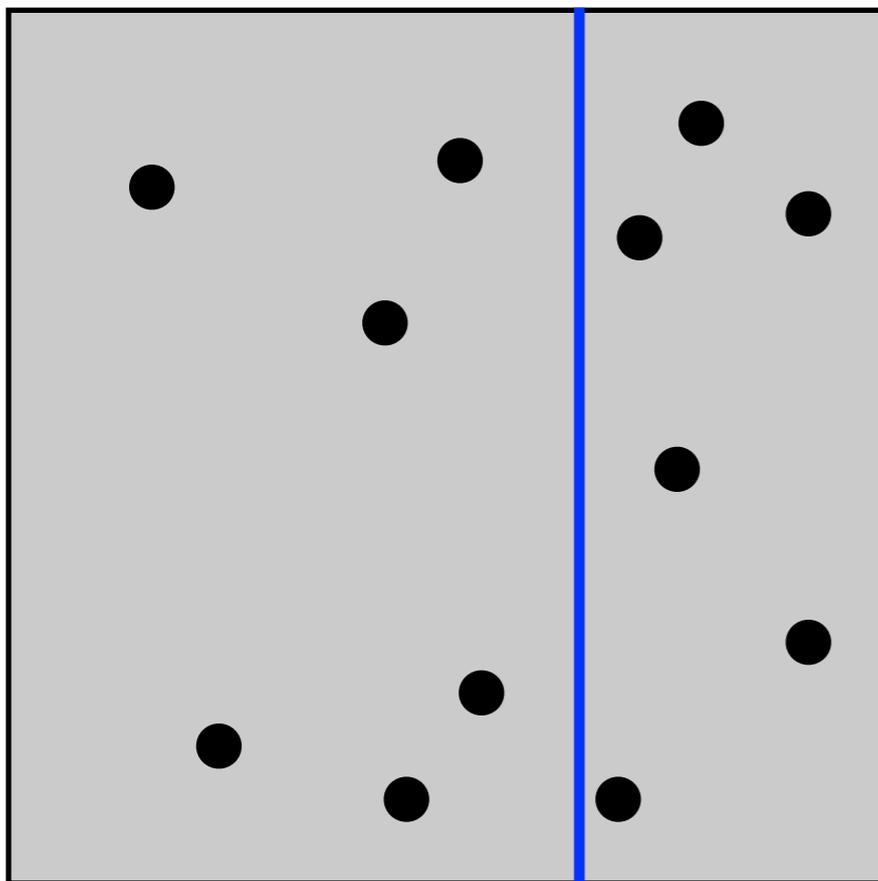
Techniques for spatial queries

- Spatial queries are 2 dimensional counting queries (typically range queries)
- kd-tree: a data-aware hierarchical space partitioning data structure.



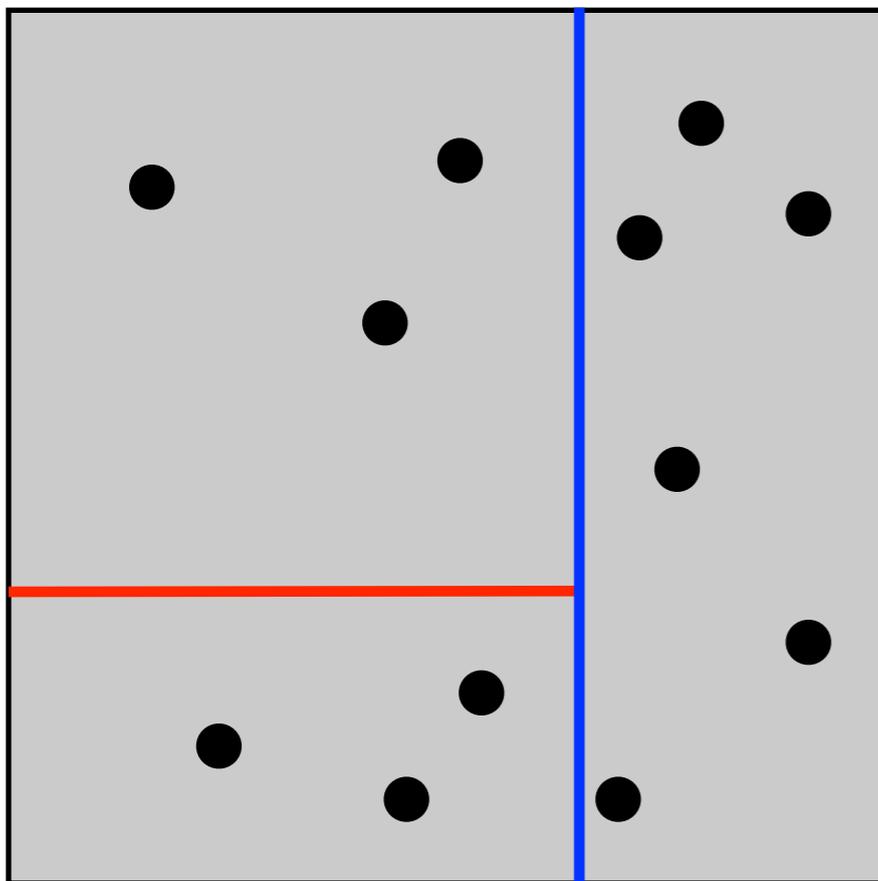
Techniques for spatial queries

- Spatial queries are 2 dimensional counting queries (typically range queries)
- kd-tree: a data-aware hierarchical space partitioning data structure.



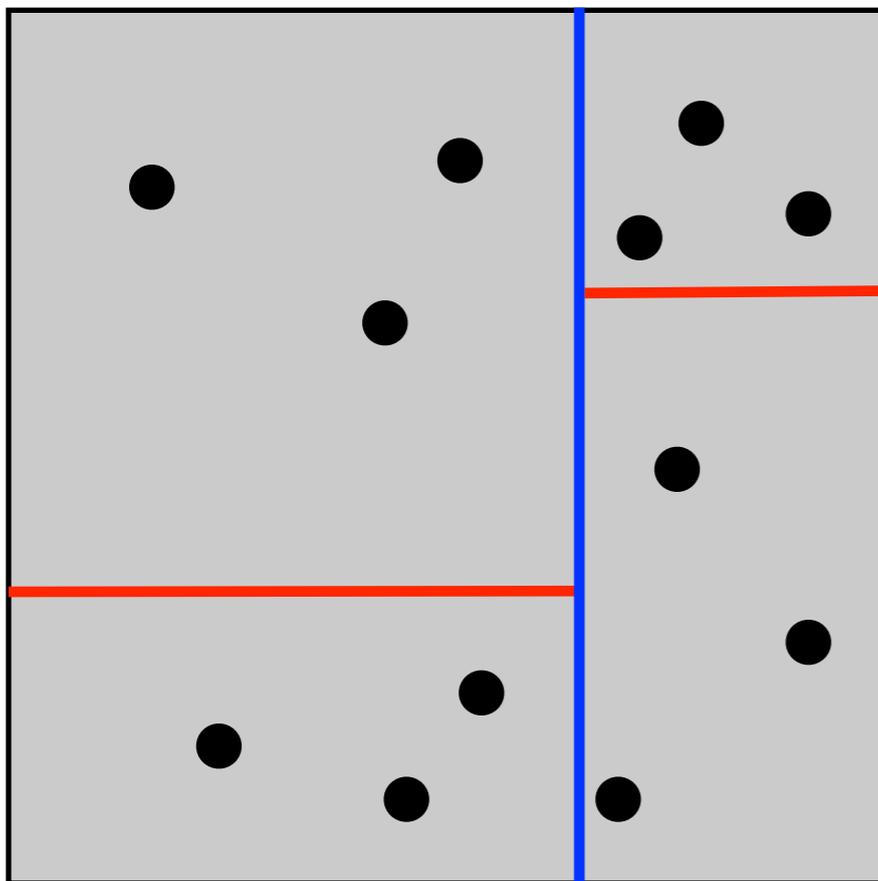
Techniques for spatial queries

- Spatial queries are 2 dimensional counting queries (typically range queries)
- kd-tree: a data-aware hierarchical space partitioning data structure.



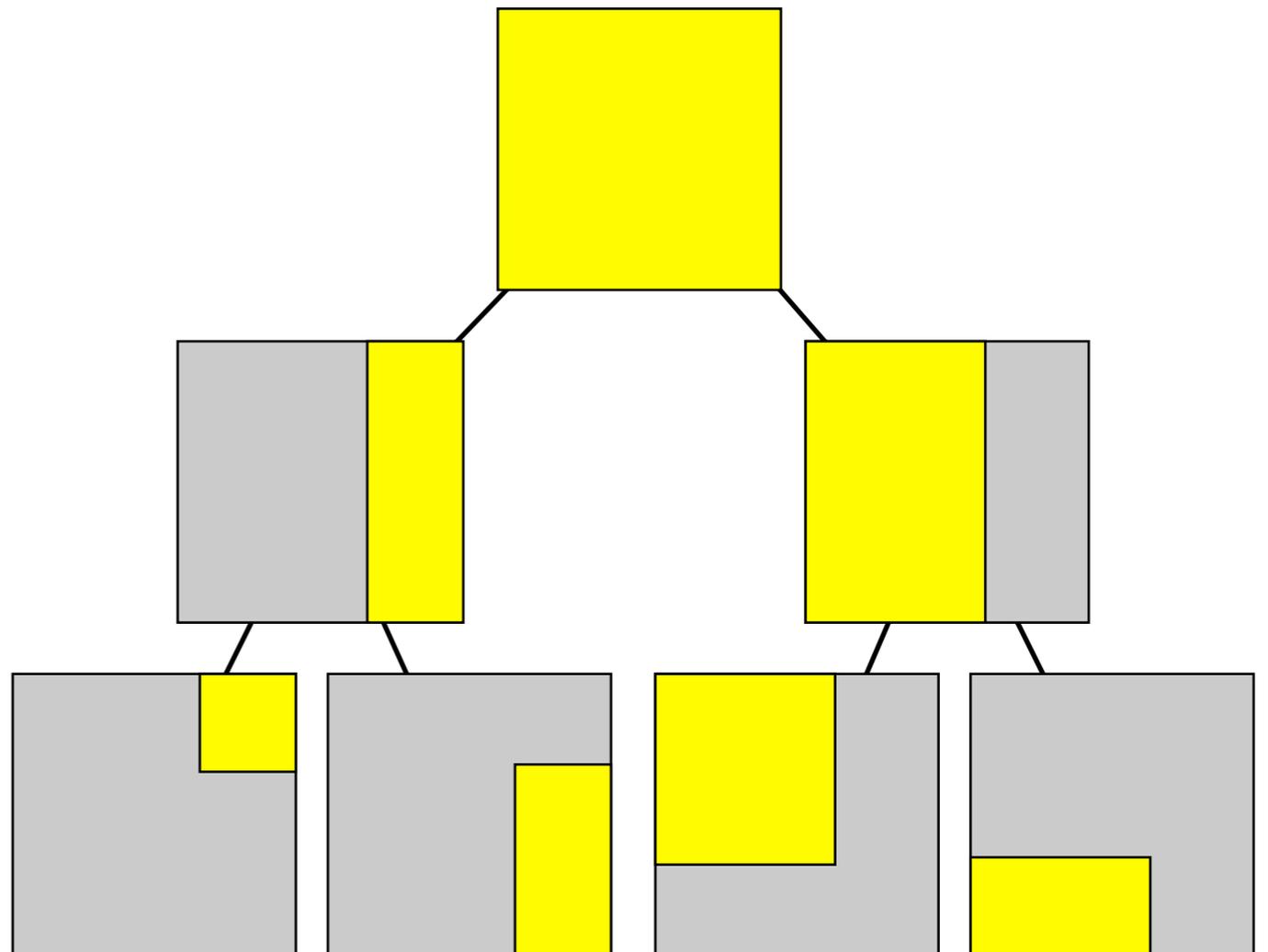
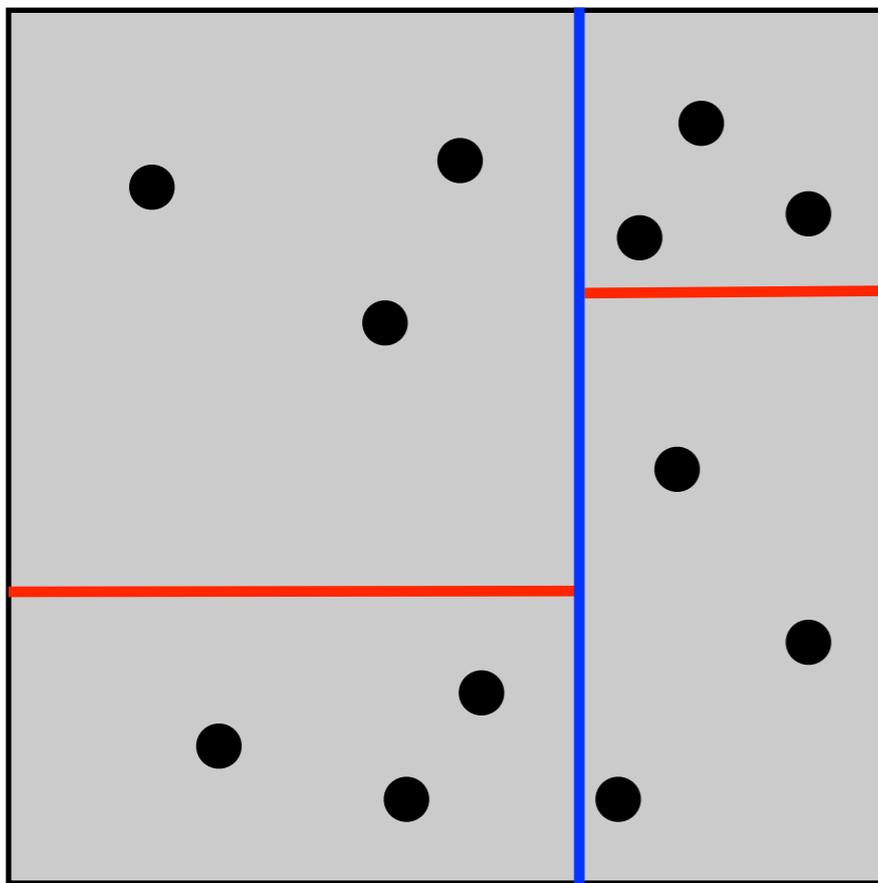
Techniques for spatial queries

- Spatial queries are 2 dimensional counting queries (typically range queries)
- kd-tree: a data-aware hierarchical space partitioning data structure.



Techniques for spatial queries

- Spatial queries are 2 dimensional counting queries (typically range queries)
- kd-tree: a data-aware hierarchical space partitioning data structure.



Data-aware kd-tree (1)

[Xiao, SDM '10]

Workload	2D Range Queries
Parameters	$p_1, p_2, \epsilon_1, \epsilon_2$ s.t. $\epsilon_1 + \epsilon_2 = \epsilon$

1. Use Laplace mechanism to get noisy counts: x'

2. Build kd-tree K from x' , but stop splitting if:

- sum of counts in current region is too small (p_1), or
- counts in current region are close to uniform (p_2)

3. Use Laplace mechanism to get noisy counts K' for all regions in K .

4. Compute workload answers from K' using least squares.

$$\epsilon_1 = \epsilon / 2$$

$$\epsilon_2 = \epsilon / 2$$

Data-aware kd-tree (2)

[Cormode, ICDE '12]

Workload	2D Range Queries
Parameters	$l, k, \epsilon_1, \epsilon_2$ s.t. $\epsilon_1 + \epsilon_2 = \epsilon$

1. Build hybrid hierarchical structure:

- l -levels of kd-tree using exponential mechanism to compute median. $\epsilon_1 = .3\epsilon$
- remaining $(k-l)$ levels uniform quad-tree.

2. Use Laplace mechanism to get noisy counts. $\epsilon_2 = .7\epsilon$

3. Derive workload query answers using least squares.

Optimizing for relative error

[Xiao, SIGMOD '11]

Workload	marginals
Parameters	T, ϵ

1. Answer all workload queries using Laplace mechanism with budget ϵ/T
2. Repeat $T-1$ times:
 - Refine query answers, by resampling queries with small values.
 - Final query answers have same privacy cost as single Laplace random variable with resulting error.

Multiplicative weights

[Hardt, NIPS '12]

Workload	linear queries
Parameters	$T, \epsilon_1, \epsilon_2$ s.t. $T(\epsilon_1 + \epsilon_2) = \epsilon$

- Begin with uniform estimate x_0 of database x
- For $i = 1 \dots T$:
 - Evaluate all workload queries using current estimate x_{i-1} . Select inaccurate q_i with exponential mechanism. $\epsilon_1 = \epsilon / 2T$
 - Laplace mechanism: get noisy estimate m_i of q_i .
 - Update $x_{i-1} \rightarrow x_i$ using m_i : multiplicative weights. $\epsilon_2 = \epsilon / 2T$

Multiplicative weights

[Hardt, NIPS '12]

- Provably better dependence on ε than workload-aware techniques: squared error $O(1/\varepsilon^2)$ vs. $O(1/\varepsilon^{2/3})$
- Observations customized to workload.
- Very good accuracy for sparse datasets.
- Output satisfies non-negativity constraints.
- Must compute all workload queries T times.

Representative experimental findings

- Building a data-aware histogram reduces error on range queries by 20-40% compared with fixed workload-aware methods like wavelet or tree-based. [Xu, ICDE '12]
- Neither of the data-aware kd-trees consistently outperform workload-aware quad-tree (on random sets of 2D range queries). [Cormode, ICDE '12]
- For reasonable privacy parameters, small workloads of random range queries on sparse data, multiplicative weights can reduce error by a factor of 10 over matrix mechanism. [Hardt, NIPS '12]
 - (But for other datasets, it can be outperformed by a factor of 10 by a fixed workload-aware method like wavelet.)

Note: ratios based on **root** mean squared error.

Data-aware mechanisms

- Observations selected to match properties of the database; generally efficient, but spending privacy budget on testing doesn't always pay off.

Workload	Observations	Citation
1D range queries	approx. v-optimal histogram	[Xu, ICDE '12]
2D range queries	kd-tree queries	[Xiao, SDM '10]
2D range queries	hybrid kd-tree queries	[Cormode, ICDE '12]
Marginals	scaled workload queries	[Xiao, SIGMOD '11]
Linear queries	subset of workload	[Hardt, NIPS '12]

Summary: data-aware mechanisms

- **Benefits:**

- Lower error than Approach 1 in some cases.

- **Limitations:**

- Parameters for algorithms must be selected carefully.
- Public error rates not available to analyst.
- Techniques are data-aware, but are they workload-aware?

- **Open questions:**

- Evaluation highly dependent on workload, dataset, epsilon. What are “real” data and workloads like? What properties of data determine error?

Outline

1. Preliminaries

2. Approach 1: workload-aware

- Fixed Observations
- Optimized Observations

3. Approach 2: data-aware

4. Conclusions

Outline

1. Preliminaries

2. Approach 1: workload-aware

- Fixed Observations
- Optimized Observations

3. Approach 2: data-aware

4. Conclusions

Summary and conclusions

- Two approaches to batch query answering, each of which provide significant error improvements by building on standard Laplace/Gaussian mechanisms, but using alternative observations.
 - Workload-aware methods ignore the input data, and choose observations solely by analyzing the workload.
 - Data-aware methods carefully (i.e. privately) exploit properties of the input data.
- Both approaches are efficient for modestly sized domains.

Workload-aware

- **Benefits**

- Independence of data makes error analysis easy, error rates publishable to analyst, and improves efficiency in some cases.

- **Limitations**

- Computational dependence on domain size, n .
- Error dependence on epsilon: $1/\epsilon^2$
- For some workloads, there is no set of observations that can help much.

- **Open questions**

- Alternative derivation methods: e.g. non-negative least squares
- Relationship with “universal” error lower bounds for DP.

Data-aware

- **Benefits:**

- Lower error than Approach 1 in some cases.

- **Limitations:**

- Parameters for algorithms must be selected carefully.
- Public error rates not available to analyst.
- Techniques are data-aware, but are they workload-aware?

- **Open questions:**

- Evaluation highly dependent on workload, dataset, epsilon. What are “real” data and workloads like? What properties of data determine error?

Open issues

- What makes one workload “harder” to answer than another?
- What makes one database “harder” to support accurately?
- Can we avoid the computational dependence on the domain size n ?
- How do we analyze the error resulting from non-negative least squares if applied in derivation of matrix mechanism?
- Methods for more expressive queries.

References

[Barak, PODS '07]	B. Barak, K. Chaudhuri, C. Dwork, S. Kale, F. McSherry, and K. Talwar. Privacy, accuracy, and consistency too: a holistic solution to contingency table release. In Principles of Database Systems (PODS) 2007.
[McSherry, FOCS '07]	F. McSherry and K. Talwar. Mechanism design via differential privacy. In FOCS '07
[McSherry, SIGMOD '09]	F. D. McSherry. Privacy integrated queries: an extensible platform for privacy-preserving data analysis. SIGMOD 2009.
[Hay, PVLDB '10]	M. Hay, V. Rastogi, G. Miklau, and D. Suciu. Boosting the accuracy of differentially-private queries through consistency. PVLDB, 2010.
[Xiao, ICDE '10]	X. Xiao, G. Wang, and J. Gehrke. Differential privacy via wavelet transforms. In International Conference on Data Engineering (ICDE), 2010.
[Li, PODS '10]	C. Li, M. Hay, V. Rastogi, G. Miklau, and A. McGregor. Optimizing Linear Counting Queries Under Differential Privacy. Principles of Database Systems (PODS) 2010.
[Xiao, SDM '10]	Y. Xiao, L. Xiong, and C. Yuan. Differentially private data release through multidimensional partitioning. Secure Data Management (SDM) 2010.
[Xiao, SIGMOD '11]	Xiaokui Xiao, Gabriel Bender, Michael Hay, and Johannes Gehrke. iReduct: Differential privacy with reduced relative errors. SIGMOD, 2011.
[Ding, SIGMOD '11]	B. Ding, M. Winslett, J. Han, and Z. Li. Differentially private data cubes: optimizing noise sources and consistency. In SIGMOD 2011.
[Xiao, SIGMOD '11]	X. Xiao, G. Bender, M. Hay, and J. Gehrke. iReduct: Differential privacy with reduced relative errors. In SIGMOD, 2011.
[Cormode, ICDE '12]	G. Cormode, M. Procopiuc, D. Srivastava, E. Shen, and T. Yu. Differentially private spatial decompositions. International Conference on Data Engineering (ICDE), 2012.

References (con't)

[Xu, ICDE '12]	J. Xu, Z. Zhang, X. Xiao, Y. Yang, and G. Yu. Differentially private histogram publication. In ICDE, 2012.
[Li, PVLDB '12]	C. Li and G. Miklau. An adaptive mechanism for accurate query answering under differential privacy. Proceedings of the VLDB Endowment (PVLDB) 2012.
[Yuan, VLDB '12]	G. Yuan, Z. Zhang, M. Winslett, X. Xiao, Y. Yang, and Z. Hao. Low-rank mechanism: Optimizing batch queries under differential privacy. VLDB, 2012.
[Hardt, NIPS '12]	M. Hardt, K. Ligett, and F. McSherry. A simple and practical algorithm for differentially private data release. In NIPS, 2012.