# Differential privacy without a central database

Uri Stemmer

## About this course

- The local model ✓
- The shuffle model ✓
- Streaming/online settings
- Differential privacy as a tool

# Streaming/online settings
## Today's Outline

1. **Private streaming algorithms**

2. **Privacy under continual observation**

# What is Streaming?

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

Alon, Matias, Szegedy 96

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

*Example: Can you count the number of <u>distinct</u> characters?*

# What is Streaming?

*Correct answer: 9*

# What is Streaming?

*Correct answer: 9*

- A stream of length $n$ over domain $X$ is a sequence of updates $(x_1, \ldots, x_n)$ where $x_i \in X$
- Let $g: X^* \to R$ be a function
- At every time $i \in [n]$ we obtain $x_i$
- At the end of the stream we need to output $z \approx g(x_1, \ldots, x_n)$
- **Requirement: small space**

# What is Streaming?

**Correct answer: 9**

- A stream of length $n$ over domain $X$ is a sequence of updates $(x_1, \ldots, x_n)$ where $x_i \in X$
- Let $g: X^* \to R$ be a function
- At every time $i \in [n]$ we obtain $x_i$
- At the end of the stream we need to output $z \approx g(x_1, \ldots, x_n)$
- **Requirement: small space**

- What does it mean for a streaming algorithm to be DP?

  ➤ A streaming algorithm $\mathcal{A}$ is $(\varepsilon, \delta)$-DP if for any two neighboring streams $\vec{x} = (x_1, \ldots, x_n)$ and $\vec{x}' = (x_1', \ldots, x_n')$ that differ on one update we have that $\mathcal{A}(\vec{x}) \approx_{(\varepsilon, \delta)} \mathcal{A}(x')$

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

0 0 0 1 1 0 0 1 1 0 0 →  → 0 1 0 0 1 1 0 1 0 1 1 →  → statistics

Private algorithm for "sanitizing streams"

Non-private streaming algorithm

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

0 0 0 1 1 0 0 1 1 0 0 →  → 0 1 0 0 1 1 0 1 0 1 1 →  → statistics

Private algorithm for "sanitizing streams"

Non-private streaming algorithm

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

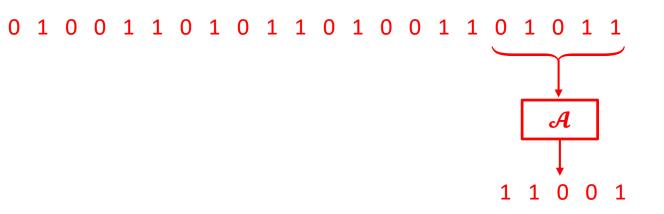0 0 0 1 1 0 0 1 1 0 0 → [Private algorithm for "sanitizing streams"] → 0 1 0 0 1 1 0 1 0 1 1 → [Non-private streaming algorithm] → statistics
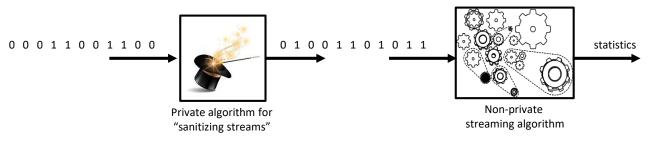
- For a predefined family of queries $Q$, a "streaming sanitizer" is a DP algorithm whose input is a stream $S$ and its output is a stream $\widehat{S}$ such that for every $q \in Q$ we have $\frac{1}{|S|} \sum_{x \in S} q(x) \approx \frac{1}{|\widehat{S}|} \sum_{x \in S} q(\widehat{x})$

- Recall the definition of standard "offline sanitizers" which is the same except that the input and outputs are specified all at once…

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

$$0\ 0\ 0\ 1\ 1\ 0\ 0\ 1\ 1\ 0\ 0 \rightarrow$$



Private algorithm for "sanitizing streams"

$$0\ 1\ 0\ 0\ 1\ 1\ 0\ 1\ 0\ 1\ 1 \rightarrow$$



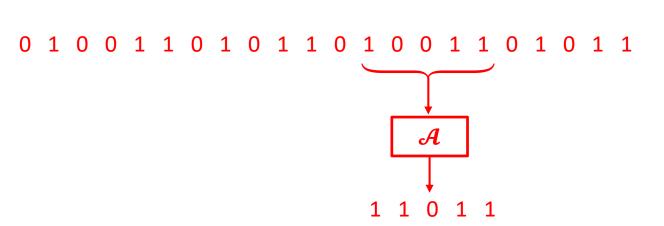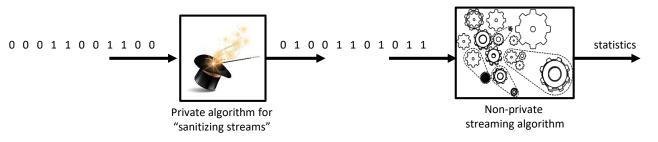Non-private streaming algorithm

$\rightarrow$ statistics

- For a predefined family of queries $Q$, a "streaming sanitizer" is a DP algorithm whose input is a stream $S$ and its output is a stream $\widehat{S}$ such that for every $q \in Q$ we have $\frac{1}{|S|}\sum_{x \in S} q(x) \approx \frac{1}{|\widehat{S}|}\sum_{x \in S} q(\widehat{x})$

- Recall the definition of standard "offline sanitizers" which is the same except that the input and outputs are specified all at once...

- Standard "offline sanitizers" can be transformed into "streaming sanitizers": Suppose we have an "offline sanitizer" $\mathcal{A}$ for $Q$ that operates on datasets of size $m$

- Construct a streaming sanitizer as follows:
    1) Let $D$ denote the next $m$ items
    2) Output $\mathcal{A}(D)$
    3) goto step (1)

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

0 0 0 1 1 0 0 1 1 0 0 → [Private algorithm for "sanitizing streams"] → 0 1 0 0 1 1 0 1 0 1 1 → [Non-private streaming algorithm] → statistics
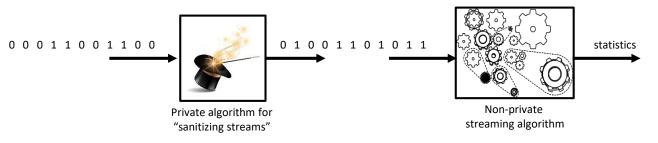
- For a predefined family of queries $Q$, a "streaming sanitizer" is a DP algorithm whose input is a stream $S$ and its output is a stream $\widehat{S}$ such that for every $q \in Q$ we have $\frac{1}{|S|}\sum_{x \in S} q(x) \approx \frac{1}{|\widehat{S}|}\sum_{x \in S} q(\widehat{x})$

- Recall the definition of standard "offline sanitizers" which is the same except that the input and outputs are specified all at once…

- Standard "offline sanitizers" can be transformed into "streaming sanitizers": Suppose we have an "offline sanitizer" $\mathcal{A}$ for $Q$ that operates on datasets of size $m$

- Construct a streaming sanitizer as follows:
  1) Let $D$ denote the next $m$ items
  2) Output $\mathcal{A}(D)$
  3) goto step (1)

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

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

0 0 0 1 1 0 0 1 1 0 0 → [Private algorithm for "sanitizing streams"] → 0 1 0 0 1 1 0 1 0 1 1 → [Non-private streaming algorithm] → statistics
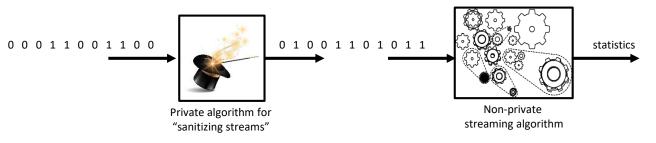
- For a predefined family of queries $Q$, a "streaming sanitizer" is a DP algorithm whose input is a stream $S$ and its output is a stream $\widehat{S}$ such that for every $q \in Q$ we have $\frac{1}{|S|}\sum_{x \in S} q(x) \approx \frac{1}{|\widehat{S}|}\sum_{x \in S} q(\widehat{x})$

- Recall the definition of standard "offline sanitizers" which is the same except that the input and outputs are specified all at once...

- Standard "offline sanitizers" can be transformed into "streaming sanitizers": Suppose we have an "offline sanitizer" $\mathcal{A}$ for $Q$ that operates on datasets of size $m$

- Construct a streaming sanitizer as follows:
  1) Let $D$ denote the next $m$ items
  2) Output $\mathcal{A}(D)$
  3) goto step (1)

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

$\mathcal{A}$

1 1 0 0 1

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

$$0\ 0\ 0\ 1\ 1\ 0\ 0\ 1\ 1\ 0\ 0 \longrightarrow \boxed{\text{magic}} \longrightarrow 0\ 1\ 0\ 0\ 1\ 1\ 0\ 1\ 0\ 1\ 1 \longrightarrow \boxed{\text{gears}} \longrightarrow \text{statistics}$$

Private algorithm for "sanitizing streams"     Non-private streaming algorithm
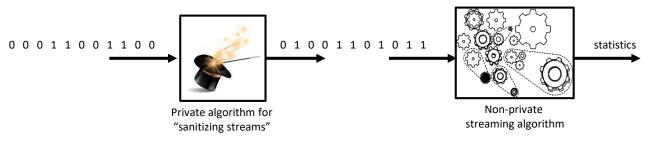
- For a predefined family of queries $Q$, a "streaming sanitizer" is a DP algorithm whose input is a stream $S$ and its output is a stream $\widehat{S}$ such that for every $q \in Q$ we have $\frac{1}{|S|}\sum_{x \in S} q(x) \approx \frac{1}{|\widehat{S}|}\sum_{x \in S} q(\widehat{x})$

- Recall the definition of standard "offline sanitizers" which is the same except that the input and outputs are specified all at once...

- Standard "offline sanitizers" can be transformed into "streaming sanitizers": Suppose we have an "offline sanitizer" $\mathcal{A}$ for $Q$ that operates on datasets of size $m$

- Construct a streaming sanitizer as follows:
  1) Let $D$ denote the next $m$ items
  2) Output $\mathcal{A}(D)$
  3) goto step (1)

$$0\ 1\ 0\ 0\ 1\ 1\ 0\ 1\ 0\ 1\ 1\ 0\ \underbrace{1\ 0\ 0\ 1\ 1}\ 0\ 1\ 0\ 1\ 1$$

$$\boxed{\mathcal{A}}$$

$$1\ 1\ 0\ 1\ 1$$

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

0 0 0 1 1 0 0 1 1 0 0 → [Private algorithm for "sanitizing streams"] → 0 1 0 0 1 1 0 1 0 1 1 → [Non-private streaming algorithm] → statistics

- For a predefined family of queries $Q$, a "streaming sanitizer" is a DP algorithm whose input is a stream $S$ and its output is a stream $\widehat{S}$ such that for every $q \in Q$ we have $\frac{1}{|S|} \sum_{x \in S} q(x) \approx \frac{1}{|\widehat{S}|} \sum_{x \in S} q(\hat{x})$

- Recall the definition of standard "offline sanitizers" which is the same except that the input and outputs are specified all at once…

- Standard "offline sanitizers" can be transformed into "streaming sanitizers": Suppose we have an "offline sanitizer" $\mathcal{A}$ for $Q$ that operates on datasets of size $m$

- Construct a streaming sanitizer as follows:
  1) Let $D$ denote the next $m$ items
  2) Output $\mathcal{A}(D)$
  3) goto step (1)

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

$\mathcal{A}$

0 1 0 1 0

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

0 0 0 1 1 0 0 1 1 0 0 → [Private algorithm for "sanitizing streams"] → 0 1 0 0 1 1 0 1 0 1 1 → [Non-private streaming algorithm] → statistics
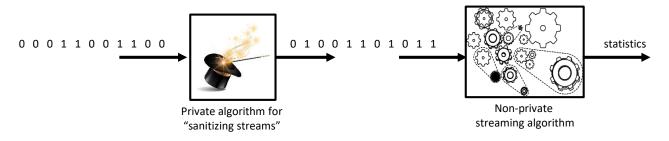
- For a predefined family of queries $Q$, a "streaming sanitizer" is a DP algorithm whose input is a stream $S$ and its output is a stream $\widehat{S}$ such that for every $q \in Q$ we have $\frac{1}{|S|}\sum_{x \in S} q(x) \approx \frac{1}{|\widehat{S}|}\sum_{x \in S} q(\widehat{x})$

- Recall the definition of standard "offline sanitizers" which is the same except that the input and outputs are specified all at once...

- Standard "offline sanitizers" can be transformed into "streaming sanitizers": Suppose we have an "offline sanitizer" $\mathcal{A}$ for $Q$ that operates on datasets of size $m$

- Construct a streaming sanitizer as follows:
  1) Let $D$ denote the next $m$ items
  2) Output $\mathcal{A}(D)$
  3) goto step (1)

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

$\mathcal{A}$

0 0 0 1 0

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

0 0 0 1 1 0 0 1 1 0 0 → [Private algorithm for "sanitizing streams"] → 0 1 0 0 1 1 0 1 0 1 1 → [Non-private streaming algorithm] → statistics

- For a predefined family of queries $Q$, a "streaming sanitizer" is a DP algorithm whose input is a stream $S$ and its output is a stream $\widehat{S}$ such that for every $q \in Q$ we have $\frac{1}{|S|}\sum_{x \in S} q(x) \approx \frac{1}{|\widehat{S}|}\sum_{x \in S} q(\widehat{x})$

- Recall the definition of standard "offline sanitizers" which is the same except that the input and outputs are specified all at once...

- Standard "offline sanitizers" can be transformed into "streaming sanitizers": Suppose we have an "offline sanitizer" $\mathcal{A}$ for $Q$ that operates on datasets of size $m$

- Construct a streaming sanitizer as follows:
  1) Let $D$ denote the next $m$ items
  2) Output $\mathcal{A}(D)$
  3) goto step (1)

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

0 0 0 1 1 0 0 1 1 0 0 → Private algorithm for "sanitizing streams" → 0 1 0 0 1 1 0 1 0 1 1 → Non-private streaming algorithm → statistics

- For a predefined family of queries $Q$, a "streaming sanitizer" is a DP algorithm whose input is a stream $S$ and its output is a stream $\widehat{S}$ such that for every $q \in Q$ we have $\frac{1}{|S|}\sum_{x \in S} q(x) \approx \frac{1}{|\widehat{S}|}\sum_{x \in S} q(\widehat{x})$

- Recall the definition of standard "offline sanitizers" which is the same except that the input and outputs are specified all at once…

- Standard "offline sanitizers" can be transformed into "streaming sanitizers": Suppose we have an "offline sanitizer" $\mathcal{A}$ for $Q$ that operates on datasets of size $m$

- Construct a streaming sanitizer as follows:
    1) Let $D$ denote the next $m$ items
    2) Output $\mathcal{A}(D)$
    3) goto step (1)

**Privacy analysis:**
- We apply $\mathcal{A}$ on disjoint portions of the input
- So no need for composition and privacy follows from $\mathcal{A}$

**Utility analysis:**
- Since we aim for <u>relative</u> error, the error do not accumulate

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

0 0 0 1 1 0 0 1 1 0 0

Private algorithm for
"sanitizing streams"

0 1 0 0 1 1 0 1 0 1 1

statistics

Non-private
streaming algorithm

- The point here is that our space does not depend directly on the length of the stream $n$
- Our space equals to the space of the sanitizer $\mathcal{A}$, which is independent of $n$, and the space of the non-private streaming algorithm

Since we aim for relative error, the error do not accumulate

# How can we design private streaming algorithms?

**Idea 1: Synthetic streams**

0 0 0 1 1 0 0 1 1 0 0 → *Private algorithm for "sanitizing streams"* → 0 1 0 0 1 1 0 1 0 1 1 → *Non-private streaming algorithm* → statistics

- F
  i

- F
  s

- S
  s

- C

- The point here is that our space does not depend directly on the length of the stream $n$
- Our space equals to the space of the sanitizer $\mathcal{A}$, which is independent of $n$, and the space of the non-private streaming algorithm

**Example where this is useful: Quantile estimation**
- Items in the stream are numbers $x_1, x_2, \ldots, x_n \in [0, 1]$
- The goal is, at the end of the stream, to get approximations for all quantiles of the data
- E.g., at the end of the stream we want to learn 9 numbers $y_1, y_2, \ldots, y_9 \in [0, 1]$ such that for every $\ell$ we have $|\{i : y_\ell \leq x_i \leq y_{\ell+1}\}| \approx \frac{n}{10}$

(this works well because we have very efficient "offline sanitizers" for this problem)

- Since we aim for <u>relative</u> error, the error do not accumulate

# The Counter Problem

**Updates are bits and we want to estimate their sum**

- Simple solution: Store the sum in memory using $\log n$ bits
  - Can we maintain a counter using smaller space?

# The Counter Problem

**Updates are bits and we want to estimate their sum**
- Simple solution: Store the sum in memory using $\log n$ bits
  - Can we maintain a counter using smaller space?

**Step 1:**

---
— Initialize estimate $\widehat{C} = 0$

— Given an update $x_i = 1$ flip a coin and increment $\widehat{C}$ only if coin is heads
---

- We expect that $\widehat{C} \approx C/2$ where $C$ is the true value
- The good: We still know $C$ (approximately) while storing only a smaller number
- The bad: We saved only 1 bit (and also this has large variance, but let's ignore it…)

# The Counter Problem

**Updates are bits and we want to estimate their sum**
- Simple solution: Store the sum in memory using $\log n$ bits
  - Can we maintain a counter using smaller space?

**Step 1:**

> — Initialize estimate $\widehat{C} = 0$
>
> — Given an update $x_i = 1$ flip a coin and increment $\widehat{C}$ only if coin is heads

- We expect that $\widehat{C} \approx C/2$ where $C$ is the true value
- The good: We still know $C$ (approximately) while storing only a smaller number
- The bad: We saved only 1 bit (and also this has large variance, but let's ignore it…)

**Step 2:**

> — Given an update $x_i = 1$ flip $\widehat{C}$ **coins** and increment $\widehat{C}$ only if **all** coin are heads

- Can show that in expectation $\widehat{C} \approx \log C$
- Thus if $C$ takes $\log n$ then $\widehat{C}$ takes $\approx \log\log n$ bits
- So we gained exponentially in storage! (again, let's ignore the variance…)

# The Counter Problem

**Updates are bits and we want to estimate their sum**
- Simple solution: Store the sum in memory using $\log n$ bits
  - Can we maintain a counter using smaller space?

**Step 1:**

---
— Initialize estimate $\widehat{C} = 0$
— Given an update $x_i = 1$ flip a coin and increment $\widehat{C}$ only if coin is heads
---

- We expect that $\widehat{C} \approx C/2$ where $C$ is the true value
- The good: We still know $C$ (approximately) while storing only a smaller number
- The bad: We saved only 1 bit (and also this has large variance, but let's ignore it…)

**Step 2:**

---
— Given an update $x_i = 1$ flip $\widehat{C}$ **coins** and increment $\widehat{C}$ only if **all** coin are heads
---

- Can show that in expectation $\widehat{C} \approx \log C$
- Thus if $C$ takes $\log n$ then $\widehat{C}$ takes $\approx \log \log n$ bits
- So we gained exponentially in storage! (again, let's ignore the variance…)

**Observe:** The outcome distribution of the algorithm depends only on $C$

# A Private Algorithm for the Counter Problem

[Dwork, Naor, Pitassi, Rothblum, Yekhanin]

- We can design a private variant as follows (informal):

---
— Sample $Y \sim \mathbf{Lap}\left(\frac{1}{\varepsilon}\right)$

— Run Morris' counter on a modified stream:
   ▪ If $Y < 0$ then ignore the first $|Y|$ ones in the stream
   ▪ If $Y \geq 0$ then add $Y$ ones before the stream begins

— The outcome distribution if a function of $(C + Y)$, satisfying privacy by post-processing
---

# A Private Algorithm for the Counter Problem
[Dwork, Naor, Pitassi, Rothblum, Yekhanin]

- We can design a private variant as follows (informal):

---

— Sample $Y \sim \mathbf{Lap}\left(\frac{1}{\varepsilon}\right)$

— Run Morris' counter on a modified stream:
  - If $Y < 0$ then ignore the first $|Y|$ ones in the stream
  - If $Y \geq 0$ then add $Y$ ones before the stream begins

— The outcome distribution if a function of $(C + Y)$, satisfying privacy by post-processing

---

- This idea is useful for other streaming problems

  ➤ Example in the context of the counter problem: Counting the number of people who viewed my YouTube video

# Streaming/online settings
## Today's Outline

✓ 1. **Private streaming algorithms**

⟹ 2. **Privacy under continual observation**

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

**Modified problem – Counter with continual reports:**

- On every time $t \in [n]$
    - We get a bit $x_t \in \{0, 1\}$
    - Need to respond with an approximation $\hat{c}_i$ for $c_i = \sum_{i=1}^{t} x_i$

# Private counter under continual observation

**Modified problem – Counter with continual reports:**

- On every time $t \in [n]$
    - We get a bit $x_t \in \{0, 1\}$
    - Need to respond with an approximation $\hat{c}_i$ for $c_i = \sum_{i=1}^{t} x_i$

> Algorithm $\mathcal{A}$ is $(\varepsilon, \delta)$-DP for this problem if for any two neighboring streams $\vec{x} = (x_1, \dots, x_n)$ and $\vec{x}' = (x_1', \dots, x_n')$ that differ on one update we have that $\mathcal{A}(\vec{u}) \approx_{(\varepsilon, \delta)} \mathcal{A}(\vec{u}')$

**Remarks:**

- Observe that now $\mathcal{A}(\vec{x})$ is a vector of length $m$
- This problem is interesting regardless of space, so let's forget about space from now on
- Sanity check: Is the previous algorithm private w.r.t. this definition?

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

**<u>Naïve attempts at solving the problem:</u>**

1) "LDP style": Every time $t \in [n]$ we release $\hat{x}_t = x_t + \textbf{Lap}\left(\frac{1}{\varepsilon}\right)$

   — This would maintain privacy, but sum of $n$ noises accumulates to $\approx \sqrt{n}/\varepsilon$

2) Using composition: Every time $t \in [n]$ we release $\hat{c}_t = \left(\sum_{i=1}^{t} x_t\right) + \textbf{Lap}(b)$

   — We would need $b \approx \sqrt{n}/\varepsilon$ due to composition

# Private counter under continual observation

**<u>Naïve attempts at solving the problem:</u>**

1) "LDP style": Every time $t \in [n]$ we release $\hat{x}_t = x_t + \mathbf{Lap}\left(\frac{1}{\varepsilon}\right)$
   — This would maintain privacy, but sum of $n$ noises accumulates to $\approx \sqrt{n}/\varepsilon$

2) Using composition: Every time $t \in [n]$ we release $\hat{c}_t = \left(\sum_{i=1}^{t} x_t\right) + \mathbf{Lap}(b)$
   — We would need $b \approx \sqrt{n}/\varepsilon$ due to composition

**<u>How can we do better?</u>**
- Observe: in solution (1) every user affects only one computation, so no need for composition, but the noises accumulate. In solution (2) we do not accumulate noises, but each one must be big to account for composition over $n$ computations
- We want something in between

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$

2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$

3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root

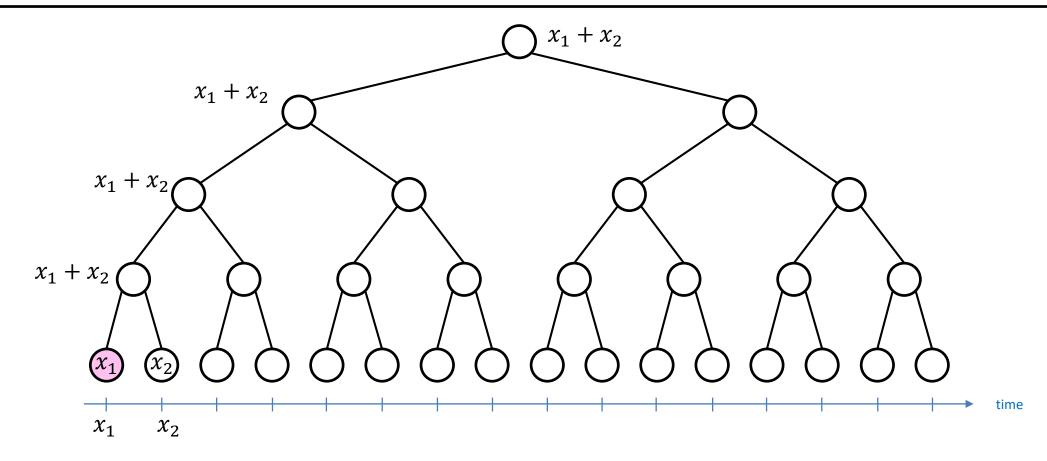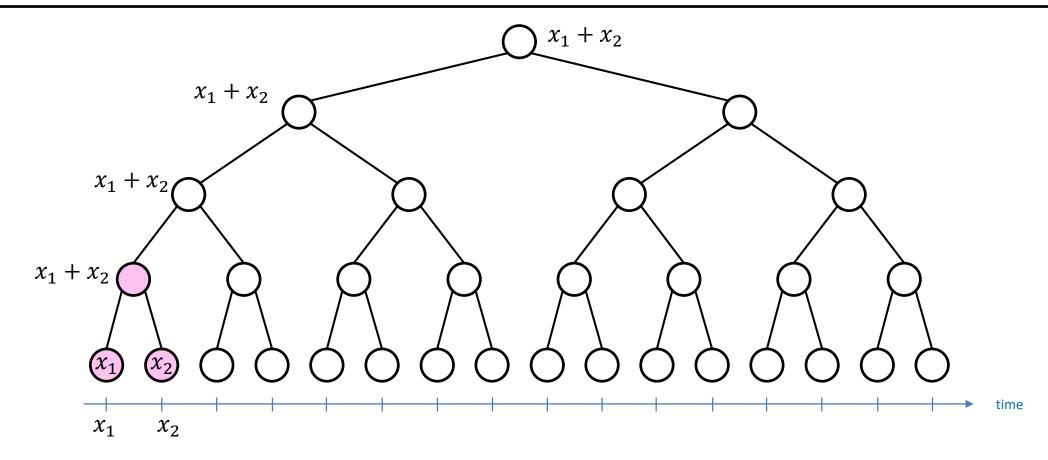4) When a subtree is "full" we release the content of its root

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$
2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$
3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root
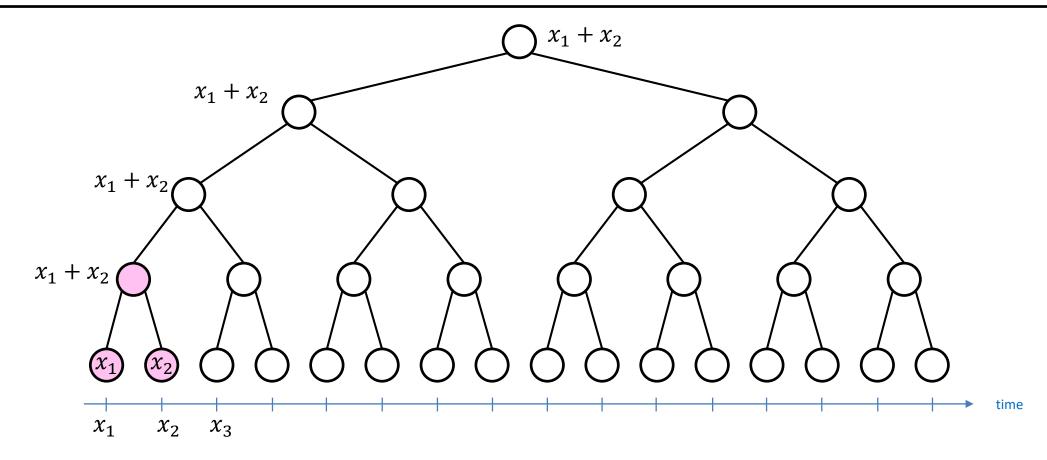4) When a subtree is "full" we release the content of its root



time

$x_1$

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$
2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$
3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root
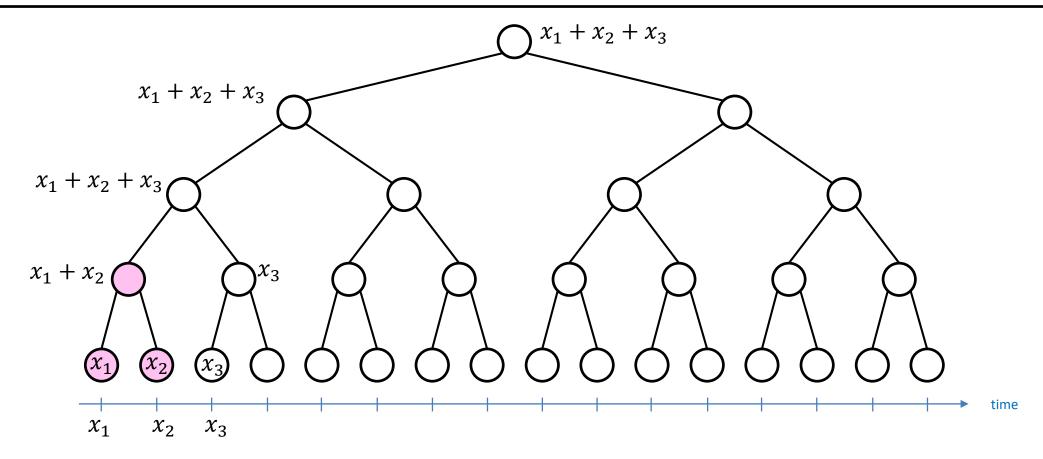4) When a subtree is "full" we release the content of its root

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$

2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$

3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root

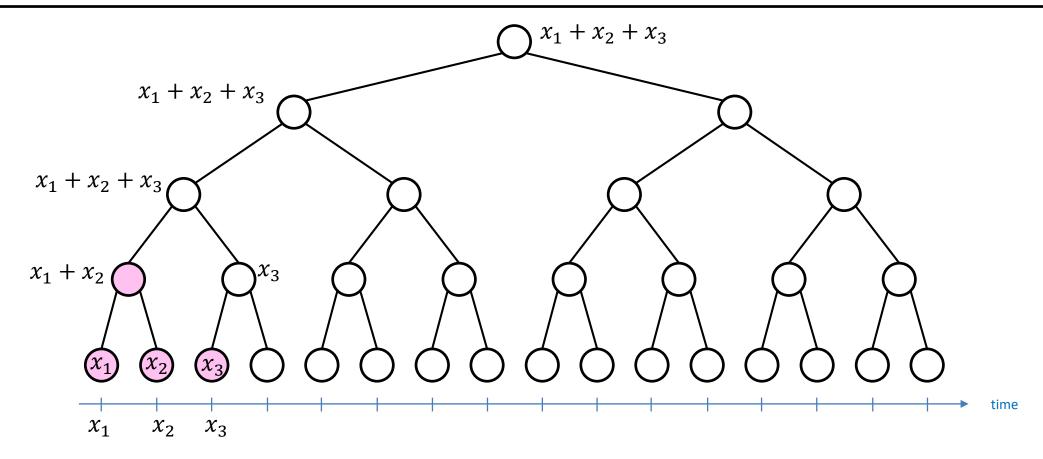4) When a subtree is "full" we release the content of its root

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$
2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$
3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root
4) When a subtree is "full" we release the content of its root

# Private counter under continual observation

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$
2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$
3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root
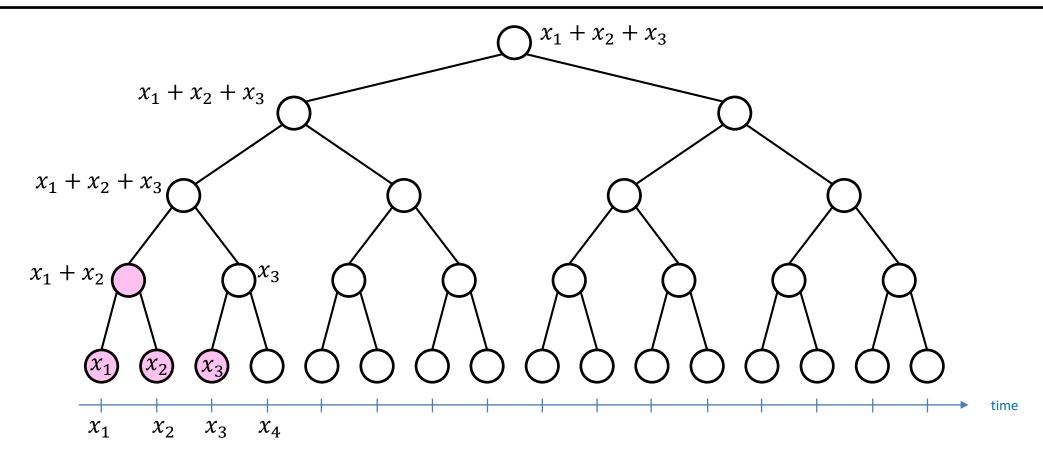4) When a subtree is "full" we release the content of its root

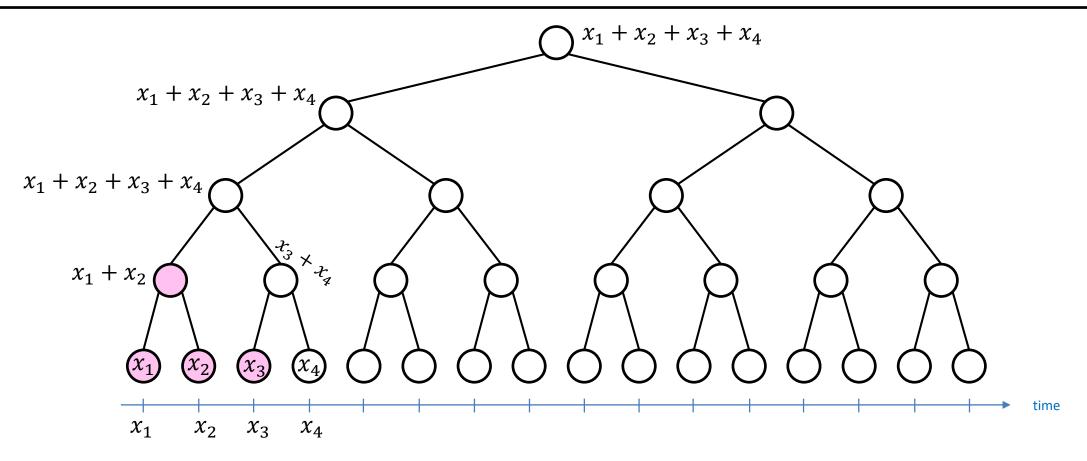# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \dots, n\}$
2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$
3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root
4) When a subtree is "full" we release the content of its root

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$
2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$
3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root
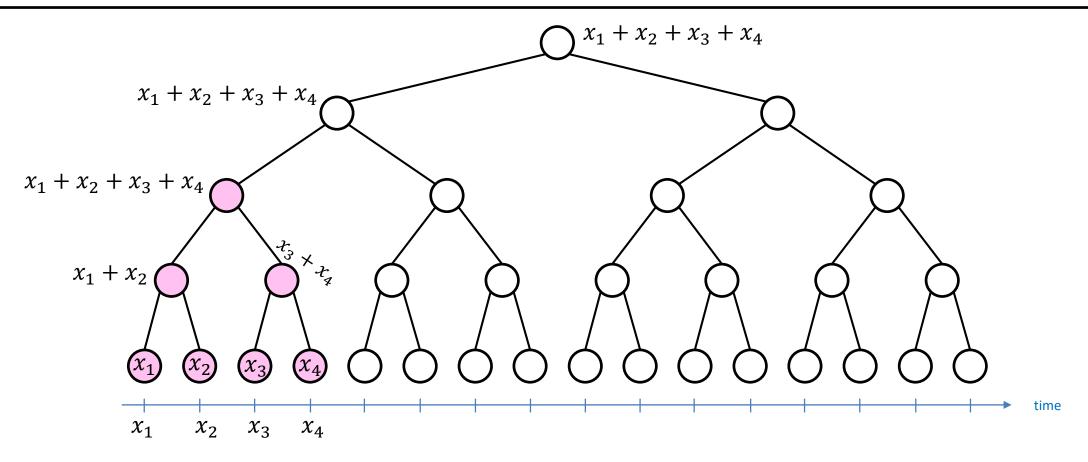4) When a subtree is "full" we release the content of its root

# Private counter under continual observation

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$
2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$
3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root
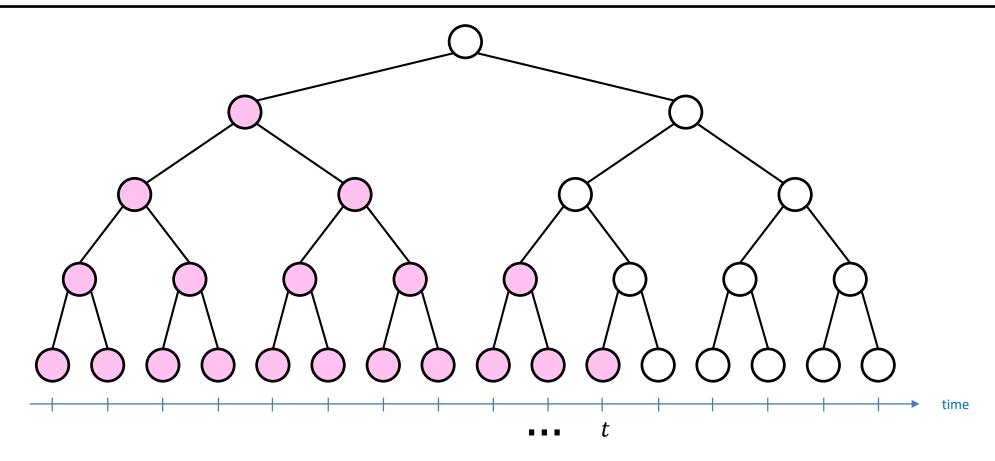4) When a subtree is "full" we release the content of its root

# Private counter under continual observation

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$
2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$
3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root
4) When a subtree is "full" we release the content of its root

# Private counter under continual observation

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$
2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$
3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root
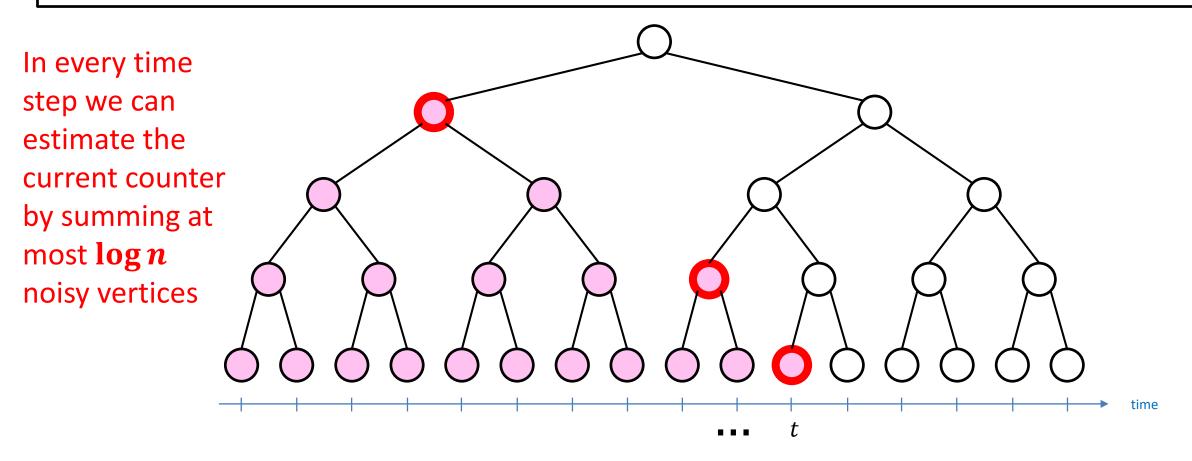4) When a subtree is "full" we release the content of its root

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$

2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$

3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root

4) When a subtree is "full" we release the content of its root

# Private counter under continual observation

[Dwork, Naor,  Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \dots, n\}$

2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$

3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root

4) When a subtree is "full" we release the content of its root



In every time step we can estimate the current counter by summing at most $\log n$ noisy vertices

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \dots, n\}$

2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$

3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root

4) When a subtree is "full" we release the content of its root

# Private counter under continual observation

[Dwork, Naor,  Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$
2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$
3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root
4) When a subtree is "full" we release the content of its root

**Privacy analysis:**

- Once a subtree is "full" then its root is never updated again
- Thus, we release the content of every node exactly once after adding Laplace noise
- Changing one $x_t$ affects only $\log m$ of these noise, so noise $\frac{\log n}{\varepsilon}$ suffices "by composition"

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

1) Define a binary tree whose leaves correspond to time steps $\{1, 2, 3, \ldots, n\}$
2) We initialize every node with independent random noise from $\mathbf{Lap}\left(\frac{\log n}{\varepsilon}\right)$
3) In time $t$ we get $x_t$ and add it to all the nodes along the path from leaf $t$ till the root
4) When a subtree is "full" we release the content of its root

**Privacy analysis:**
- Once a subtree is "full" then its root is never updated again
- Thus, we release the content of every node exactly once after adding Laplace noise
- Changing one $x_t$ affects only $\log m$ of these noise, so noise $\frac{\log n}{\varepsilon}$ suffices "by composition"

**Utility analysis:**
- At any time $t$ we can compute an estimated counter by summing at most $\log n$ nodes
- So we are only summing $\log n$ noises, each of magnitude $\approx \log n$
- Overall error is $\frac{polylog\, n}{\varepsilon}$

# Private counter under continual observation

[Dwork, Naor,  Pitassi, Rothblum]

**Thm:** Every $(1, 0)$-DP algorithm for this problem must have error $\Omega(\log n)$

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

**Thm:** Every $(\mathbf{1}, \mathbf{0})$-DP algorithm for this problem must have error $\Omega(\log n)$

- Suppose there is a private algorithm $\mathcal{A}$ such that w.p. $\mathbf{2/3}$ all of its estimates are accurate to within error $\frac{\log n}{16}$

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

**Thm:** Every $(\mathbf{1}, \mathbf{0})$-DP algorithm for this problem must have error $\Omega(\log n)$

- Suppose there is a private algorithm $\mathcal{A}$ such that w.p. $\mathbf{2/3}$ all of its estimates are accurate to within error $\dfrac{\log n}{16}$

- Construct a collection $H = \left\{ \vec{x}^{(1)}, \vec{x}^{\left(\frac{\log n}{4}\right)}, \dots \right\}$ of input sequences where $\vec{x}^{(i)} = \Big(\mathbf{0}, \mathbf{0}, \dots, \mathbf{0}, \underbrace{\mathbf{1}, \mathbf{1}, \mathbf{1}, \dots, \mathbf{1}}_{\frac{\log n}{4} \text{ ones from time } i}, \mathbf{0}, \dots, \mathbf{0}\Big)$

# Private counter under continual observation

[Dwork, Naor, Pitassi, Rothblum]

**Thm:** Every $(1, 0)$-DP algorithm for this problem must have error $\Omega(\log n)$

- Suppose there is a private algorithm $\mathcal{A}$ such that w.p. $2/3$ all of its estimates are accurate to within error $\frac{\log n}{16}$

- Construct a collection $H = \left\{ \vec{x}^{(1)}, \vec{x}^{\left(\frac{\log n}{4}\right)}, \dots \right\}$ of input sequences where $\vec{x}^{(i)} = \left(0, 0, \dots, 0, \underbrace{1, 1, 1, \dots, 1}_{\frac{\log n}{4} \text{ ones from time } i}, 0, \dots, 0\right)$

- The algorithm has error at most $\frac{\log n}{16}$ so if we run it on input $\vec{x}^{(i)}$ then:

  — Before time $i$ the estimation must be at most $\frac{\log n}{16}$

  — After time $i + \frac{\log n}{4}$ the estimation must be at least $\frac{3 \log n}{16}$

- Hence, a sequence of answers cannot by good for more than one $\vec{x}^{(i)}$

# Private counter under continual observation

**Thm:** Every $(1, 0)$-DP algorithm for this problem must have error $\Omega(\log n)$

- Suppose there is a private algorithm $\mathcal{A}$ such that w.p. $2/3$ all of its estimates are accurate to within error $\frac{\log n}{16}$

- Construct a collection $H = \left\{ \vec{x}^{(1)}, \vec{x}^{\left(\frac{\log n}{4}\right)}, \dots \right\}$ of input sequences where $\vec{x}^{(i)} = \left(0, 0, \dots, 0, \underbrace{1, 1, 1, \dots, 1}_{\frac{\log n}{4} \text{ ones from time } i}, 0, \dots, 0\right)$

- The algorithm has error at most $\frac{\log n}{16}$ so if we run it on input $\vec{x}^{(i)}$ then:
  - Before time $i$ the estimation must be at most $\frac{\log n}{16}$
  - After time $i + \frac{\log n}{4}$ the estimation must be at least $\frac{3 \log n}{16}$

- Hence, a sequence of answers cannot by good for more than one $\vec{x}^{(i)}$ and for every $i \neq \ell$ we have

$$\frac{2}{3} \leq \Pr\begin{bmatrix} \mathcal{A}(\vec{x}^{(i)}) \text{ returns} \\ \text{good answers} \\ \text{for } \vec{x}^{(i)} \end{bmatrix} \leq e^{2\frac{\log n}{4}} \cdot \Pr\begin{bmatrix} \mathcal{A}(\vec{x}^{(\ell)}) \text{ returns} \\ \text{good answers} \\ \text{for } \vec{x}^{(i)} \end{bmatrix} = \sqrt{n} \cdot \Pr\begin{bmatrix} \mathcal{A}(\vec{x}^{(\ell)}) \text{ returns} \\ \text{good answers} \\ \text{for } \vec{x}^{(i)} \end{bmatrix}$$

# Private counter under continual observation

__Thm:__ Every $(1, 0)$-DP algorithm for this problem must have error $\Omega(\log n)$

- Suppose there is a private algorithm $\mathcal{A}$ such that w.p. $2/3$ all of its estimates are accurate to within error $\frac{\log n}{16}$

- Construct a collection $H = \left\{ \vec{x}^{(1)}, \vec{x}^{\left(\frac{\log n}{4}\right)}, \dots \right\}$ of input sequences where $\vec{x}^{(i)} = \Big( 0, 0, \dots, 0, \underbrace{1, 1, 1, \dots, 1}_{\frac{\log n}{4} \text{ ones from time } i}, 0, \dots, 0 \Big)$

- The algorithm has error at most $\frac{\log n}{16}$ so if we run it on input $\vec{x}^{(i)}$ then:
  - Before time $i$ the estimation must be at most $\frac{\log n}{16}$
  - After time $i + \frac{\log n}{4}$ the estimation must be at least $\frac{3 \log n}{16}$

- Hence, a sequence of answers cannot by good for more than one $\vec{x}^{(i)}$ and for every $i \neq \ell$ we have

$$\frac{2}{3} \leq \Pr \begin{bmatrix} \mathcal{A}(\vec{x}^{(i)}) \text{ returns} \\ \text{good answers} \\ \text{for } \vec{x}^{(i)} \end{bmatrix} \leq e^{2\frac{\log n}{4}} \cdot \Pr \begin{bmatrix} \mathcal{A}(\vec{x}^{(\ell)}) \text{ returns} \\ \text{good answers} \\ \text{for } \vec{x}^{(i)} \end{bmatrix} = \sqrt{n} \cdot \Pr \begin{bmatrix} \mathcal{A}(\vec{x}^{(\ell)}) \text{ returns} \\ \text{good answers} \\ \text{for } \vec{x}^{(i)} \end{bmatrix}$$

- So,

$$\Pr \begin{bmatrix} \mathcal{A}(\vec{x}^{(\ell)}) \text{ returns} \\ \text{good answers} \\ \text{for } \vec{x}^{(i)} \end{bmatrix} \geq \frac{2}{3\sqrt{n}}$$

$$\Pr\left[\mathcal{A}\big(\vec{x}^{(\ell)}\big)\text{ fails}\right] \geq \Pr\left[\bigcup_{i\neq\ell}\left\{\begin{array}{c}\mathcal{A}\big(\vec{x}^{(\ell)}\big)\text{ returns}\\ \text{good answers}\\ \text{for }\vec{x}^{(i)}\end{array}\right\}\right] \underbrace{=}_{\substack{\text{disjoint}\\ \text{events}}} \sum_{i\neq\ell}\Pr\left[\begin{array}{c}\mathcal{A}\big(\vec{x}^{(\ell)}\big)\text{ returns}\\ \text{good answers}\\ \text{for }\vec{x}^{(i)}\end{array}\right] \gtrsim \frac{n}{\log n}\cdot\frac{1}{\sqrt{n}} > 1$$

- Suppose there is a private algorithm $\mathcal{A}$ such that w.p. **2/3** all of its estimates are accurate to within error $\frac{\log n}{16}$

- Construct a collection $\boldsymbol{H} = \left\{\vec{x}^{(1)}, \vec{x}^{\left(\frac{\log n}{4}\right)}, \dots\right\}$ of input sequences where $\vec{x}^{(i)} = \Big(\boldsymbol{0,0,\dots,0,} \underbrace{\boldsymbol{1,1,1,\dots,1}}_{\frac{\log n}{4}\text{ ones from time } i}, \boldsymbol{0,\dots,0}\Big)$

- The algorithm has error at most $\frac{\log n}{16}$ so if we run it on input $\vec{x}^{(i)}$ then:
  - Before time $\boldsymbol{i}$ the estimation must be at most $\frac{\log n}{16}$
  - After time $\boldsymbol{i} + \frac{\log n}{4}$ the estimation must be at least $\frac{3\log n}{16}$

- Hence, a sequence of answers cannot by good for more than one $\vec{x}^{(i)}$ and for every $\boldsymbol{i} \neq \ell$ we have

$$\frac{2}{3} \leq \Pr\left[\begin{array}{c}\mathcal{A}\big(\vec{x}^{(i)}\big)\text{ returns}\\ \text{good answers}\\ \text{for }\vec{x}^{(i)}\end{array}\right] \leq e^{2\frac{\log n}{4}}\cdot\Pr\left[\begin{array}{c}\mathcal{A}\big(\vec{x}^{(\ell)}\big)\text{ returns}\\ \text{good answers}\\ \text{for }\vec{x}^{(i)}\end{array}\right] = \sqrt{n}\cdot\Pr\left[\begin{array}{c}\mathcal{A}\big(\vec{x}^{(\ell)}\big)\text{ returns}\\ \text{good answers}\\ \text{for }\vec{x}^{(i)}\end{array}\right]$$

- So,

$$\Pr\left[\begin{array}{c}\mathcal{A}\big(\vec{x}^{(\ell)}\big)\text{ returns}\\ \text{good answers}\\ \text{for }\vec{x}^{(i)}\end{array}\right] \geq \frac{2}{3\sqrt{n}}$$

# Negative result for continual observation

**Recall the XOR-sum problem:**

- The input of every user $i$ is a pair $(j_i, b_i) \in \{1, 2, \ldots, n\} \times \{0, 1\}$
- The goal: estimate $\sum_{j=1}^{n} \bigoplus_{i : j_i = j} b_i$

# Negative result for continual observation

**Recall the XOR-sum problem:**
- The input of every user $i$ is a pair $(j_i, b_i) \in \{1, 2, \ldots, n\} \times \{0, 1\}$
- The goal: estimate $\sum_{j=1}^{n} \bigoplus_{i : j_i = j} b_i$

**As we mentioned:** Can be solved with error $\approx \frac{1}{\varepsilon}$ in the centralized model

**Informal theorem:** Any continual-DP algorithm for this problem must have error $\Omega(\text{poly } n)$

# Negative result for continual observation

---

**Recall the XOR-sum problem:**

- The input of every user $i$ is a pair $(j_i, b_i) \in \{1, 2, \dots, n\} \times \{0,1\}$
- The goal: estimate $\sum_{j=1}^{n} \bigoplus_{i:j_i=j} b_i$

---

**As we mentioned:** Can be solved with error $\approx \frac{1}{\varepsilon}$ in the centralized model

**Informal theorem:** Any continual-DP algorithm for this problem must have error $\Omega(\text{poly } n)$

**Proof:**

- An algorithm $\mathcal{A}$ for this problem can be used to answer **many** "Hamming distance queries"; contradiction...

# Negative result for continual observation

---

**Recall the XOR-sum problem:**
- The input of every user $i$ is a pair $(j_i, b_i) \in \{1, 2, \dots, n\} \times \{0, 1\}$
- The goal: estimate $\sum_{j=1}^{n} \bigoplus_{i:j_i=j} b_i$

---

**As we mentioned:** Can be solved with error $\approx \frac{1}{\varepsilon}$ in the centralized model

**Informal theorem:** Any continual-DP algorithm for this problem must have error $\Omega(\text{poly } n)$

**Proof:**
- An algorithm $\mathcal{A}$ for this problem can be used to answer **many** "Hamming distance queries"; contradiction…
- Specifically, given input dataset $X = \left(x_1, \dots, x_{\sqrt{n}}\right)$, feed $(1, x_1), \dots, \left(\sqrt{n}, x_{\sqrt{n}}\right)$ to algorithm $\mathcal{A}$
- Then, given a query $Y = \left(y_1, \dots, y_{\sqrt{n}}\right)$:
  - Feed $(1, y_1), \dots, \left(\sqrt{n}, y_{\sqrt{n}}\right)$ to algorithm $\mathcal{A}$ to obtain an answer $z$
  - Feed $(1, y_1), \dots, \left(\sqrt{n}, y_{\sqrt{n}}\right)$ to algorithm $\mathcal{A}$ again
- After $\approx \sqrt{n}$ queries (total input length $n$) we can reconstruct $X$ contradicting the privacy of $\mathcal{A}$

# Streaming/online settings
## Today's Outline

✓ 1. Private streaming algorithms

✓ 2. Privacy under continual observation