# Generalizing Greenwald-Khanna Streaming Quantile Summaries for Weighted Inputs

March 13, 2023

Sepehr Assadi
Rutgers University

Nirmit Joshi
Northwestern University

Milind Prabhu
University of Michigan

Vihan Shah
Rutgers University
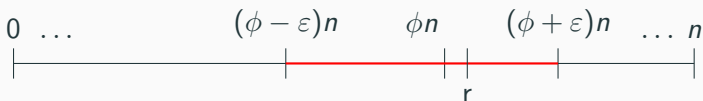
Rajiv Gandhi
Rutgers University–Camden

# Introduction

- **Input:**
    1. $S = \{x_1, \ldots, x_n\}$ of elements from an ordered universe (in the streaming fashion).
    2. Fixed approximation parameter $\varepsilon > 0$.
- **Goal:** At the end of the stream, for any $\phi \in (0, 1]$, we want to estimate $\phi$-quantile of $S$ up to an additive error of $\varepsilon$.
- On queried for any $\phi \in (0, 1]$, we want to be able to return $x \in S$ such that

$$(\phi - \varepsilon)n \leq rank(x, S) \leq (\phi + \varepsilon)n.$$



- Rank of an element:

$$rank(x, S) = |\{y \in S \mid y \leq x\}|$$

- **Input**:
  1. A weighted stream $S_w = \{(x_1, w_1), \ldots, (x_n, w_n)\}$.
  2. $w(x)$ is a positive integer.

$$W_n = \sum_{i=1}^{n} w(x_i)$$

  3. Fixed $\varepsilon > 0$

| 1 | 1 | 1 | 1 | 3 | 3 | 6 | 6 | 6 | 6 | 6 | 10 | 10 | 10 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|----|----|----|----|

# Weighted Generalized Problem
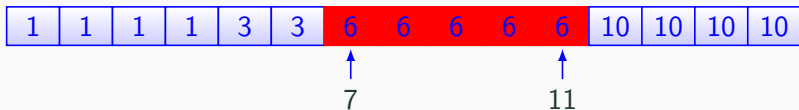
- **Input**:
  1. A weighted stream $S_w = \{(x_1, w_1), \ldots, (x_n, w_n)\}$.
  2. $w(x)$ is a positive integer.

  $$W_n = \sum_{i=1}^{n} w(x_i)$$

  3. Fixed $\varepsilon > 0$

| 1 | 1 | 1 | 1 | 3 | 3 | 6 | 6 | 6 | 6 | 6 | 10 | 10 | 10 | 10 |

- **Range of rank:**

| 1 | 1 | 1 | 1 | 3 | 3 | 6 | 6 | 6 | 6 | 6 | 10 | 10 | 10 | 10 |

7          11

- **Goal:** At the end, for any $\phi \in (0, 1]$, we want to return an element $x_j$ such that

$$(\text{Range of Ranks of } x_j) \cap \big[(\phi - \varepsilon)W_n, (\phi + \varepsilon)W_n\big] \neq \emptyset$$
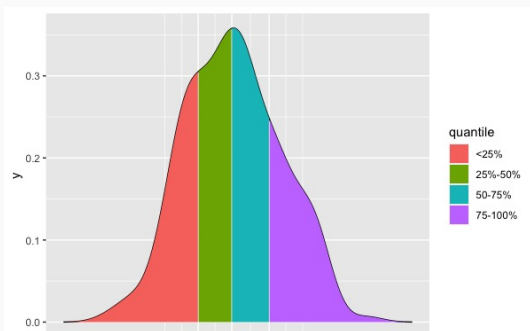


$$(\phi \pm \varepsilon)W_n$$

## Motivation

A fundamental problem in:

- Data Mining and Data Science
- Machine Learning
- Computer Science

Quantiles provide concise information about the data distribution as they allow us to estimate the CDF of the underlying distribution.
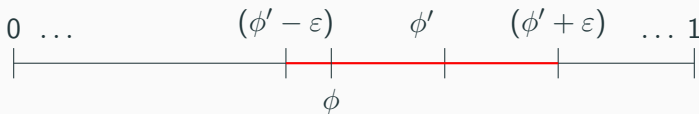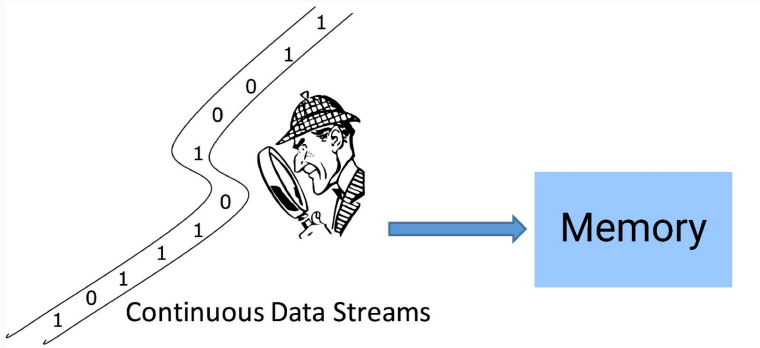
- $S = \{x_1, \ldots, x_n\}$ known apriori

- Which elements to store to approximately answer quantiles queries?

- Store $\varepsilon$-quantile, $3\varepsilon$-quantile, $5\varepsilon$-quantile, . . . . This requires only $O(1/\varepsilon)$ elements.

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ | $x_7$ | $x_8$ | $x_9$ | $x_{10}$ | $x_{11}$ | $x_{12}$ | $x_{13}$ | $x_{14}$ | $x_{15}$ |

- This is also necessary!! We must store $\Omega(1/\varepsilon)$ elements.

Continuous Data Streams

Memory $\lll$ Input Size

Elements $x_1, \ldots, x_n$ come one by one in any arbitrary order.

## Related Work

**Deterministic Algorithms:**

- Manku, Rajagopalan and Lindsay [MRL'98, SIGMOD] the MRL algorithm: uses $O(\frac{1}{\varepsilon} \log^2(\varepsilon n))$ space.

- Greenwald and Khanna [GK'01, SIGMOD] proposed the GK algorithm: uses $O(\frac{1}{\varepsilon} \log(\varepsilon n))$ space.

- The best-known 22-year-old deterministic quantile summary.

**Randomized Algorithm:**

- [KLL'16, FOCS]: answers with probabilistic guarantee $(1 - \delta)$ and achieves $O((\frac{1}{\varepsilon}) \log \log(1/\varepsilon\delta))$ space.

### Question 1

Can we improve this space-bound of $O(\frac{1}{\varepsilon} \log(\varepsilon n))$ bound? Or is this optimal?

### Answer

Cormode and Veselý [CV'20, PODS] recently resolved this question by proving $\Omega(\frac{1}{\varepsilon} \log(\varepsilon n))$ lower bound.

### Question 2

Can we simplify the GK Algorithm so that it allows for generalization to related problems, such as the weighted quantile problem?

### Answer

This paper!!

# Results

## Results for the Unweighted Setting

### Result 1

A simple and greedy algorithm that admits $O(\frac{1}{\varepsilon}\log^2(\varepsilon n))$ space guarantee.

- Similar to the "GK-Adaptive" [LWYC'16, VLDB] that has no theoretical guarantee.
- Leads to intuitions behind the counter-intuitive choices of the GK algorithm.

### Result 2

A new simpler description of the GK algorithm, which requires $O(\frac{1}{\varepsilon}\log(\varepsilon n))$ space.

- **Trivial way:** Feed multiple copies of the same element.
- Update time $O(\max_i w_i)$– prohibitively large.

> ### Result 3
>
> A non-trivial extension of the GK algorithm for weighted inputs that uses
>
> - $O(\frac{1}{\varepsilon} \log(\varepsilon n))$ space.
>
> - $O(\log(1/\varepsilon) + \log\log(\varepsilon n))$ update time per element.
>
> assuming weights are $\mathrm{poly}(n)$ and $\varepsilon \geq 1/n^{1-\delta}$ for any $\delta > 0$

- If $\varepsilon \approx 1/n$, even information-theoretically $\Omega(1/\varepsilon) = \Omega(n)$ elements needed.
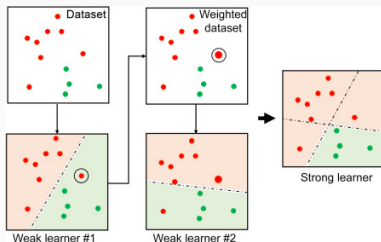
This matches the best unweighted case guarantees.

# Application

Weak
Learning

- Combine weak predictors to boost the confidence and accuracy



- Algorithm: XGBoost library by Nvidia
- Uses the weighted extension of the MRL algorithm: $O(\frac{1}{\varepsilon} \log^2(\varepsilon n))$ space
- Our GK extension can be used here....

# Basic Setup: Unweighted to Weighted Extension

## Unweighted Quantile Summary (QS)

- QS: A data structure that allows us to answer $\varepsilon$-approximate quantile queries
- It simply stores a subset of elements seen so far.

$$QS = \{e_1, \ldots, e_s\}$$

$$e_1 < e_2 < \cdots < e_s$$

- For each element $e \in QS$:

$$rmin(e_i) = \text{lower bound on the rank of } e_i$$

$$rmax(e_i) = \text{upper bound on the rank of } e_i$$



- Space Complexity=$|QS| = \#$ of elements stored at any time.

For each element $e_i \in \texttt{QS}$:

$$g_i := \text{rmin}(e_i) - \text{rmin}(e_{i-1})$$

$$\Delta_i = \text{rmax}(e_i) - \text{rmin}(e_i)$$



$$\boxed{(\text{rmin}, \text{rmax}) \iff (g, \Delta)}$$

$$\text{High } (g, \Delta) \iff \text{High}$$
$$\text{uncertainty in the Ranks}$$

# Insert and Delete

- We can define **Insert(x)** operation
- **Delete($e_i$)**: Just forget $e_i$ and keep rmin and rmax values unchanged.

### Delte($e_i$)

1. Delete $e_i$ from QS.

2. Keep rmin and rmax values changed

3. Update $g_{i+1} = g_{i+1} + g_i$.

- `WQS`: A data structure that allows us to answer $\varepsilon$-approximate weighted quantile queries

- It simply stores a subset of elements seen so far.

$$\text{WQS} = \{e_1, \ldots, e_s\}$$

$$e_1 < e_2 < \cdots < e_s$$

- Store $w(e_i)$ for each element

- For each element $e \in \text{WQS}$:

  $\text{rmin}(e_i) = \text{lower bound}$ on the rank of first copy $e_i$

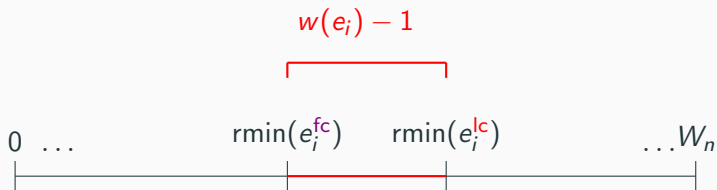  $\text{rmax}(e_i) = \text{upper bound}$ on the rank of first copy $e_i$

$$\text{rmin}(j\text{-th copy of } e_i) = \text{rmin}(e_i) + j - 1$$

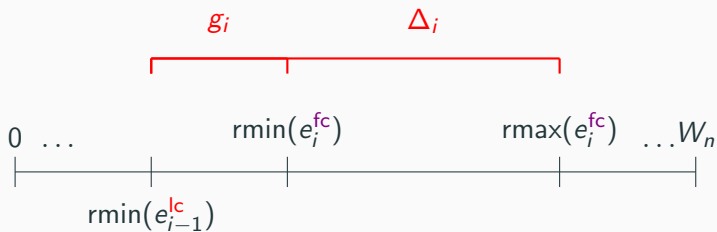$$\text{rmax}(j\text{-th copy of } e_i) = \text{rmax}(e_i) + j - 1$$

In particular,

$$\text{rmin}(e_i^{\text{lc}}) = \text{rmin}(e_i) + w(e_i) - 1$$

$$\text{rmax}(e_i^{\text{lc}}) = \text{rmax}(e_i) + w(e_i) - 1$$
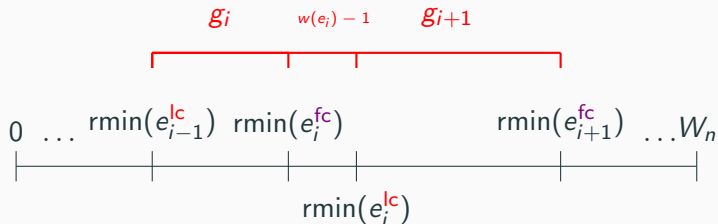
$$\Delta_i := \mathsf{rmax}(e_i) - \mathsf{rmin}(e_i)$$

$$g_i := \mathsf{rmin}(e_i^{\mathsf{fc}}) - \mathsf{rmin}(e_{i-1}^{\mathsf{lc}})$$
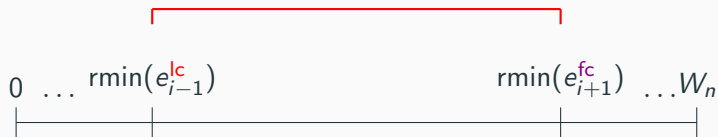
You can also define **Insert(x)** operation.



$$g_{i+1} := g_{i+1} + g_i + w(e_i) - 1$$

$$G_i = g_i + w(e_i) - 1$$

---

Delte($e_i$)

1. Delete $e_i$ from QS.
2. Update $g_{i+1} = g_{i+1} + g_i + w(e_i) - 1 = g_{i+1} + G_i$.

---

**Note:** $G_i = g_i$ in the unweighted case

### Algorithm Sketch

1. **Insertion Step:** Insert all arriving elements $x$ in the chunk using *Insert(x)*.

2. **Deletion Step:** Delete a few elements from QS, according to some rule.

**The only cleverness of the algorithm is in the deletion step!**

1. **To be able to answer the quantile queries**

2. **Minimize the space**

**Let's focus on the first...**

**Quantitatively, what $(g, \Delta) \rightarrow$ allows answering quantile queries??**

**Unweighted:** After $n$ insertions, for all elements $e_i \in \texttt{QS}$

$$g_i + \Delta_i \leq \varepsilon n$$

then we can answer any $\phi$-quantile query with $\varepsilon$-precision.



**Weighted:**

$$\boxed{g_i + \Delta_i \leq \varepsilon W_n}$$

Delete elements as long as $(g, \Delta)$ invariant holds....

Allows us to answer quantile queries... ✓

Space complexity ×

What is the magical GK deletion rule?

# Simplified GK for the Unweighted Case

**Bands ≈ Geometric grouping of elements**

- Band $\alpha$ contains $\approx 2^{\alpha}$ chunks of $1/\varepsilon$ elements
- After $n$ insertions we have, # of bands $= O(\log(\varepsilon n))$.

| Band number: | 1 | | 2 | | | | 3 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Chunk number: | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 |

### Segment

The **segment** of an element $e_i$ in QS, denoted by $seg(e_i)$, is defined as the maximal set of consecutive elements

$$e_j, e_{j+1}, \cdots, e_{i-1}$$

in QS with b-value strictly less than b-value($e_i$).

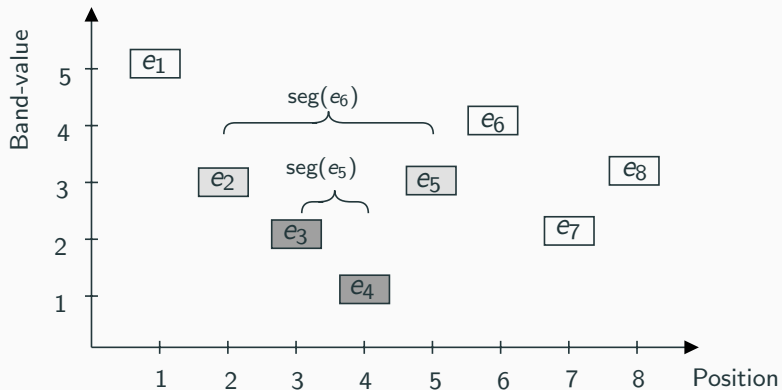Treat the element and its segment as one unit!

---

### Simpler GK Algorithm

For arriving item $x_k$:

1. **Insertion Step:** Insert arriving element $x_k$ using *Insert($x_k$)*.

2. **Deletion Step:** For any $e_i \in \texttt{QS}$, delete $e_i$ along with seg($e_i$) if the following two conditions hold:

   b-value($e_i$) $\leq$ b-value($e_{i+1}$) and $g_i^* + g_{i+1} + \Delta_{i+1} \leq \varepsilon k$

---

$$g_i^* = g_i + \sum_{j \in \text{seg}(e_i)} g_j$$

$(g + \Delta) \leq \varepsilon n$ invariant holds after the deletion!

- Only $O(1/\varepsilon)$ elements per band:

$$\boxed{\text{Total elements} = O(\log(\varepsilon n)/\varepsilon).}$$

Delete element without segment?
$O(\frac{1}{\varepsilon} \log^2(\varepsilon n))$ space
Greedy!!

# Non-trivial extension of the GK for Weighted Inputs

$$\text{b-value}(x) = \text{b-value}(\text{first copy of } x)$$



b-value(1)  b-value(3) b-value(6)  b-value(10)

Different copies may have different b-values in the corresponding unweighted stream because of weights!!

$$\# \text{ bands} = O(\log \varepsilon W_n)$$

Treat the element and its segment as one unit!

---

### Weighted extension GK Algorithm

For any arriving item $(x_k, w(x_k))$:

1. **Insertion Step:** Run *Insert($x_k$)*.

2. **Deletion Step:** For any $e_i \in$ WQS, delete $e_i$ along with seg($e_i$) if the following two conditions hold:

   b-value($e_i$) $\leq$ b-value($e_{i+1}$) and $G_i^* + g_{i+1} + \Delta_{i+1} \leq \varepsilon W_k$

---

$$G_i^* = G_i + \sum_{j \in \text{seg}(e_i)} G_j$$

$$(G = g + w - 1)$$

- $W_n =$ Total weight of $n$ elements
- Using similar counting argument:

$$\text{Space Complexity} = O((1/\varepsilon)\log(\varepsilon W_n)).$$

Under assumption weights are $\mathrm{poly}(n)$:

$$\boxed{\text{Space Complexity} = O\left(\frac{1}{\varepsilon}\log(\varepsilon n)\right).}$$

- This is not just some "smart" implementation of the "trivial" GK extension!!

### Trivial extension GK Algorithm

For any arriving item $(x_k, w(x_k))$:

1. **Insertion Step:** Insert $w(x_k)$ copies of $x_k$ in QS.
2. **Deletion step:** Run unweighted GK deletion rule on QS.
3. At the end, collapse multiple copies of the remaining elements into one element.

**May delete a partial number of copies of one element into the other and then delete remaining copies later.....**

### Weighted extension GK Algorithm

For any arriving item $(x_k, w(x_k))$:

1. **Insertion Step:** Run *Insert($x_k$)*.

2. **Deletion Step:** For any $e_i \in$ WQS, delete $e_i$ along with seg($e_i$) if the following two conditions hold:

   b-value($e_i$) $\leq$ b-value($e_{i+1}$)  and  $G_i^* + g_{i+1} + \Delta_{i+1} \leq \varepsilon W_k$

**Runtime:** Already Update time doesn't depend on $\max_{i \in [n]} w_i$.

### Goal Accomplished!

## Faster Runtime

- We can get $O(\log |\text{WQS}|) = O(\log(1/\varepsilon) + \log\log(\varepsilon n))$ runtime.
- Store WQS as a balanced Binary Search Tree (BST).
- Insert and Delete takes $O(\log |\text{WQS}|)$ time.
- But still, deciding which elements to delete takes time linear in $|\text{WQS}|$......
- Perform deletion only after $|\text{WQS}|$ doubles by delaying deletions (the space increases only by a constant factor)
- Total time spent over $n$ insertion is

$$O\left(n \cdot (\log(1/\varepsilon) + \log\log(\varepsilon n))\right)$$

$O\left(\log(1/\varepsilon) + \log\log(\varepsilon n)\right)$ amortized update-time

- Standard Techniques of delaying deletions
- Spread the time required for the deletion step which is linear in $|\texttt{WQS}|$
- Over the next few insertions.
- Interleave between Insertions and Deletions!!
- More details in the full version: arXiv 2303.06288

> **Main result**
>
> A non-trivial extension of the GK algorithm for weighted inputs, under mild assumptions:
>
> - $O(\frac{1}{\varepsilon} \log(\varepsilon n))$ space.
>
> - $O(\log(1/\varepsilon) + \log\log(\varepsilon n))$ update time per element.

## Open Question

- Stream $S_1$ of length $n_1$ and $S_2$ of length $n_2$
- $\mathtt{QS_1} = \mathcal{A}(S_1)$ and $\mathtt{QS_2} = \mathcal{A}(S_2)$ with size $f(n_1)$ and $f(n_2)$.

> ### Mergeable Summaris
>
> An algorithm $\mathcal{A}$ creates mergeable summaries if we can create
>
> $$\mathtt{QS} = \mathcal{A}(S_1 \cup S_2)$$
>
> just using $\mathtt{QS_1}$ and $\mathtt{QS_2}$. We then have $|\mathtt{QS}| = f(n_1 + n_2)$.
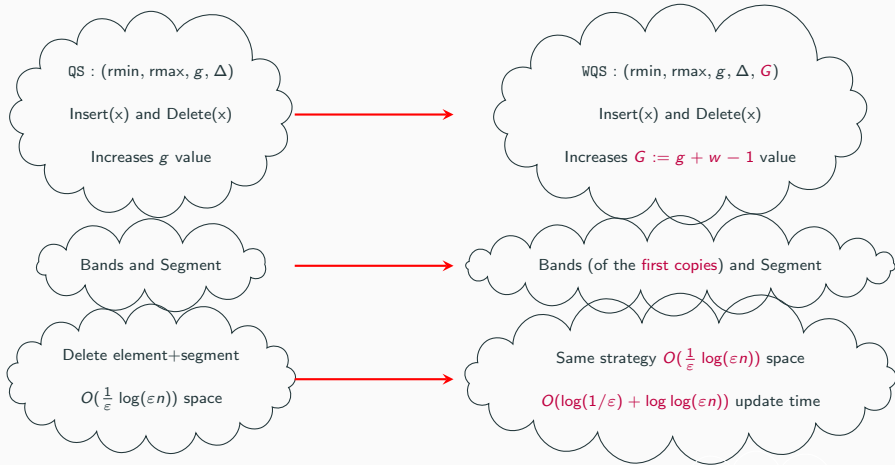
- The MRL summaries are mergeable: $O((1/\varepsilon) \log^2(\varepsilon n))$ space.
- Is the GK summary also mergeable?

  **Any algorithm that uses optimal space, i.e.**
  $f(n) = O(\frac{1}{\varepsilon} \log(\varepsilon n))$ **and produce mergeable summaries?**
- Important to parallelize the algorithm.

QS : (rmin, rmax, $g$, $\Delta$)

Insert(x) and Delete(x)

Increases $g$ value

WQS : (rmin, rmax, $g$, $\Delta$, $G$)

Insert(x) and Delete(x)

Increases $G := g + w - 1$ value

Bands and Segment

Bands (of the first copies) and Segment

Delete element+segment

$O(\frac{1}{\varepsilon} \log(\varepsilon n))$ space

Same strategy $O(\frac{1}{\varepsilon} \log(\varepsilon n))$ space

$O(\log(1/\varepsilon) + \log \log(\varepsilon n))$ update time

GK mergeable? Optimal mergeable summaries?

**Thank you!**
**Questions?**