

Differential Privacy on Fully Dynamic Streams

Yuan Qiu¹ and Ke Yi²

¹Southeast University

²Hong Kong University of Science and Technology



東南大學
SOUTHEAST UNIVERSITY



香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

Backgrounds – Fully Dynamic Streams

- A data stream is a sequence of pairs (s_t, x_t)
 - Insertion ($s_t = +$) $\Rightarrow D_t = D_{t-1} \cup \{x_t\}$
 - Deletion ($s_t = -$) $\Rightarrow D_t = D_{t-1} - \{x_t\}$
 - No-op ($s_t = \perp$) $\Rightarrow D_t = D_{t-1}$
- Notation
 - N_t : number of *updates* up to time t
 - n_t : *data size* at time t

$$N_t^+ + N_t^- = N_t \gg n_t = N_t^+ - N_t^-$$



Increases over time



Fluctuates and even hits zero




Timestamp	Update	Dataset
0	-	\emptyset
1	$+, a$	$\{a\}$
2	$+, b$	$\{a, b\}$
3	$+, c$	$\{a, b, c\}$
4	\perp	$\{a, b, c\}$
5	$-, b$	$\{a, c\}$
6	$+, d$	$\{a, c, d\}$
7	$-, c$	$\{a, d\}$
8	$-, a$	$\{d\}$
9	$-, d$	\emptyset
...

Backgrounds – Differential Privacy

- DP: for any neighboring instances $D \sim D'$ and any subset of outputs Y ,
$$\Pr[\mathcal{M}(D) \in Y] \leq e^\epsilon \cdot \Pr[\mathcal{M}(D') \in Y] + \delta.$$
- For streams, neighboring instances differ by **one update**.

Timestamp	Update	Dataset
0	-	\emptyset
1	$+, a$	$\{a\}$
2	$+, b$	$\{a, b\}$
3	$+, c$	$\{a, b, c\}$
4	\perp	$\{a, b, c\}$
5	$-, b$	$\{a, c\}$
...



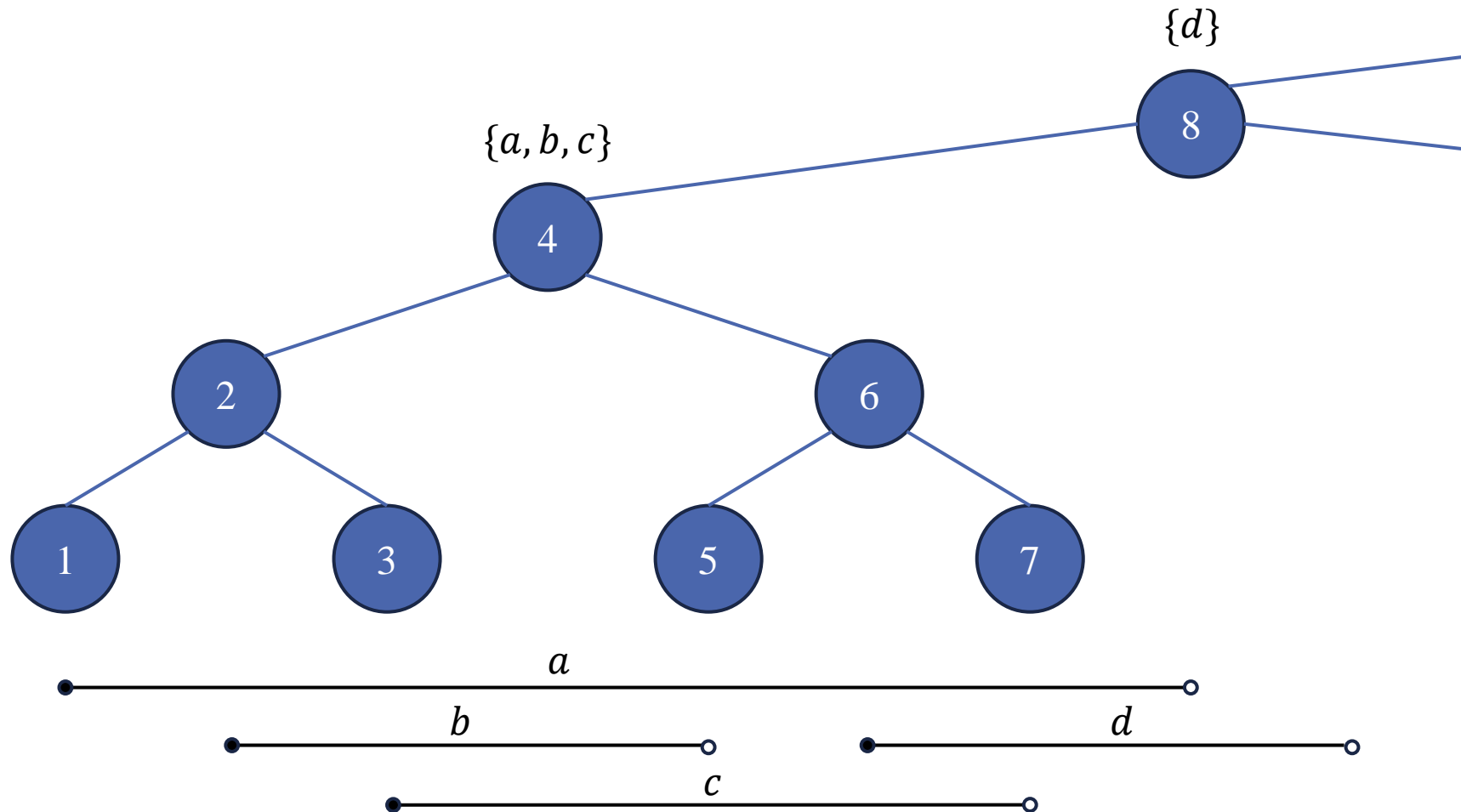
Timestamp	Update	Dataset
0	-	\emptyset
1	$+, a$	$\{a\}$
2	$+, b$	$\{a, b\}$
3	$+, c$	$\{a, b, c\}$
4	$+, e$	$\{a, b, c, e\}$
5	$-, b$	$\{a, c, e\}$
...

Backgrounds – DP Linear Queries

- Linear queries: $f(D) = \sum_{x \in D} f(x)$
 - Counting: $f(x) \equiv 1 \implies f(D) = |D|$
 - Mean: $f(x) = \frac{x}{n} \implies f(D) = \frac{\sum_{x \in D} x}{n} = \mathbb{E}[x]$
 - Variance: $f(x) = \frac{x^2}{n} \implies f(D) = \frac{\sum_{x \in D} x^2}{n} = \mathbb{E}[x^2]$
- For *static data*: PMW mechanism has error $\tilde{O}(\sqrt{n})$
- For *insertion-only streams*: Binary Tree mechanism has error $\tilde{O}(\sqrt{N_t})$
- Direct extension: Separate the stream into D_t^+ and D_t^-
 - Problem: The error becomes $\tilde{O}(\sqrt{N_t^+} + \sqrt{N_t^-}) = \tilde{O}(\sqrt{N_t})$
 - Our target is $\tilde{O}(\sqrt{n_t})$

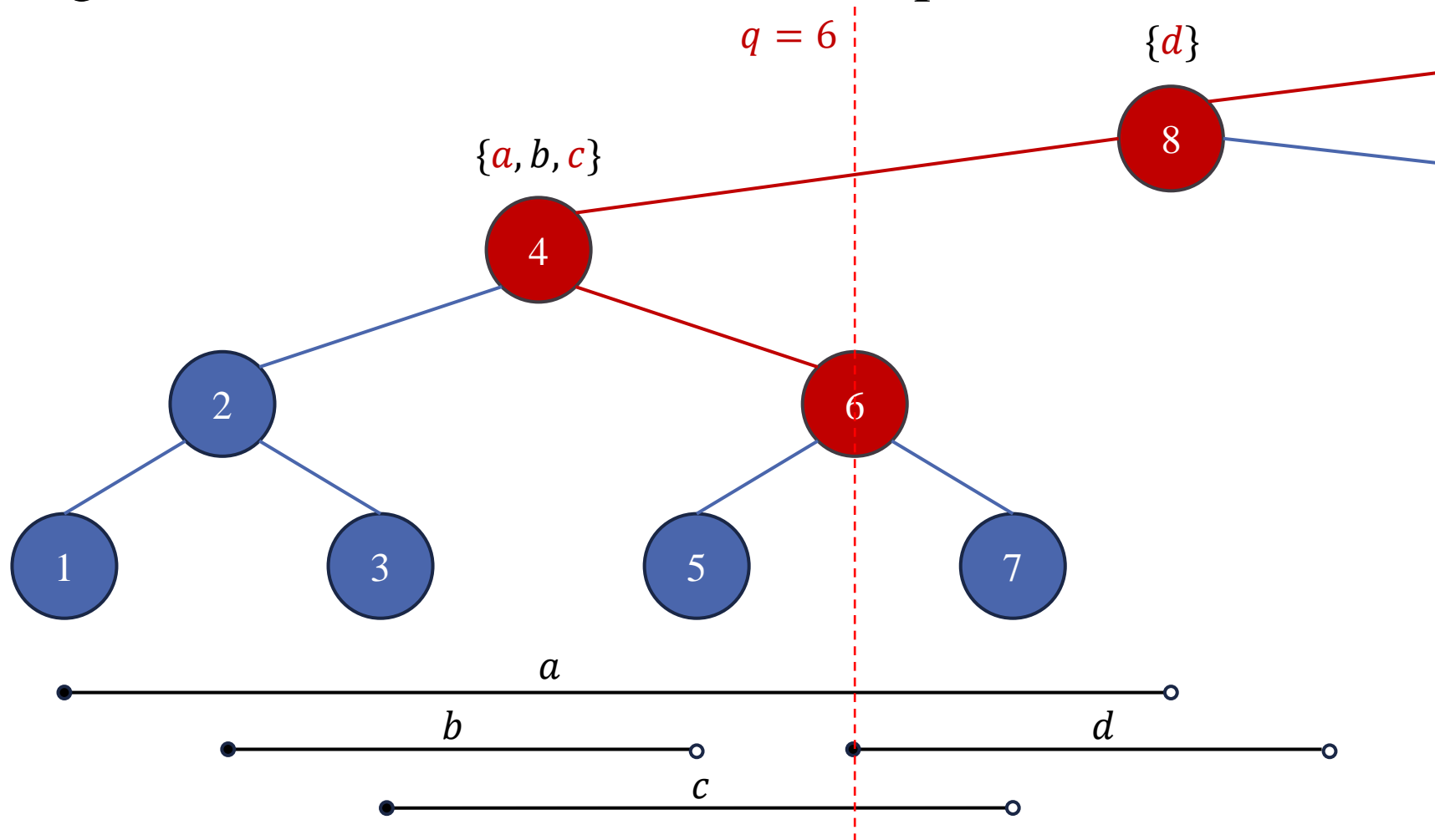
Interval Tree for Stabbing Queries

- Each (insertion, deletion) of an item can be considered an interval.



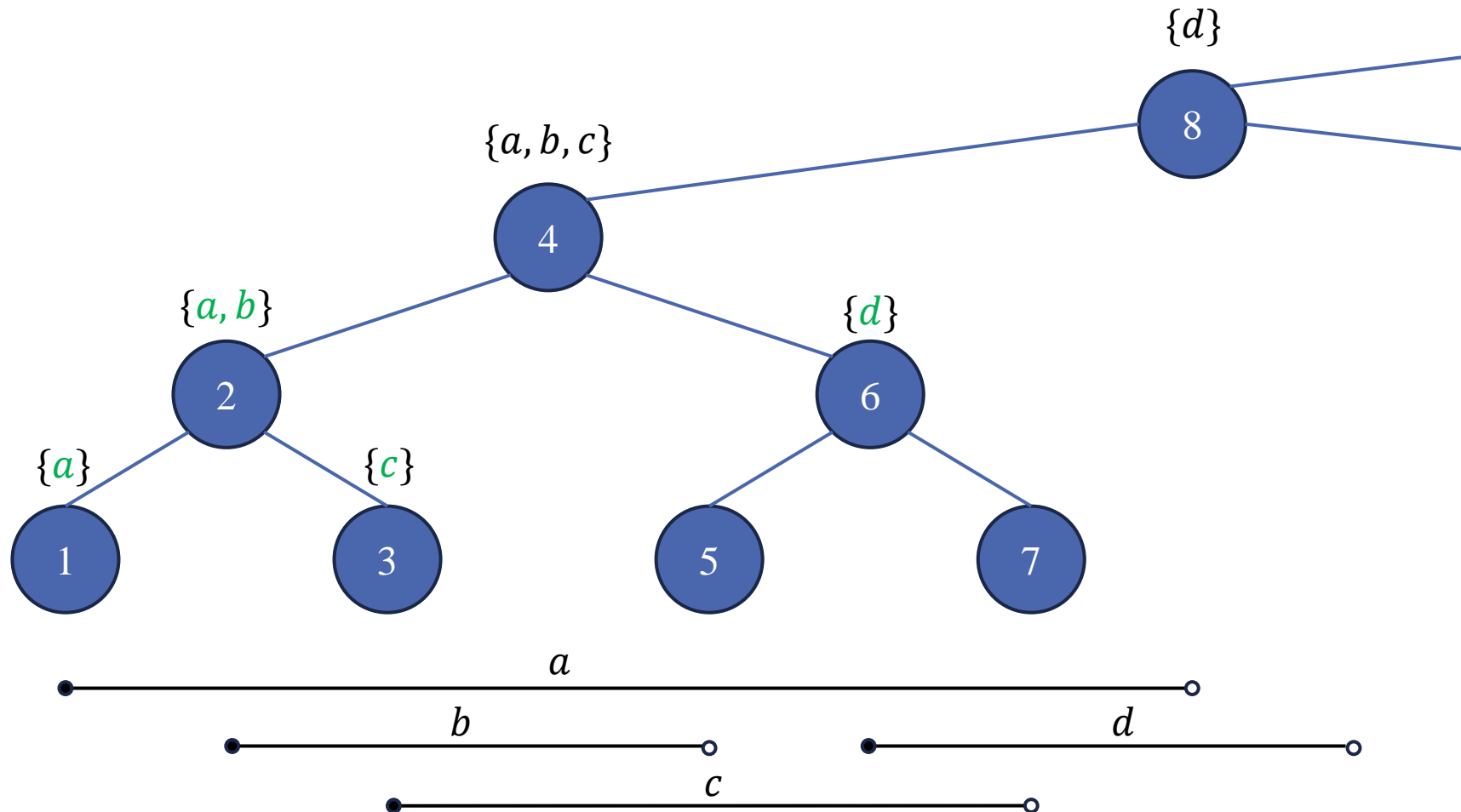
Interval Tree for Stabbing Queries

- Querying: visit nodes on the root-to-node path.



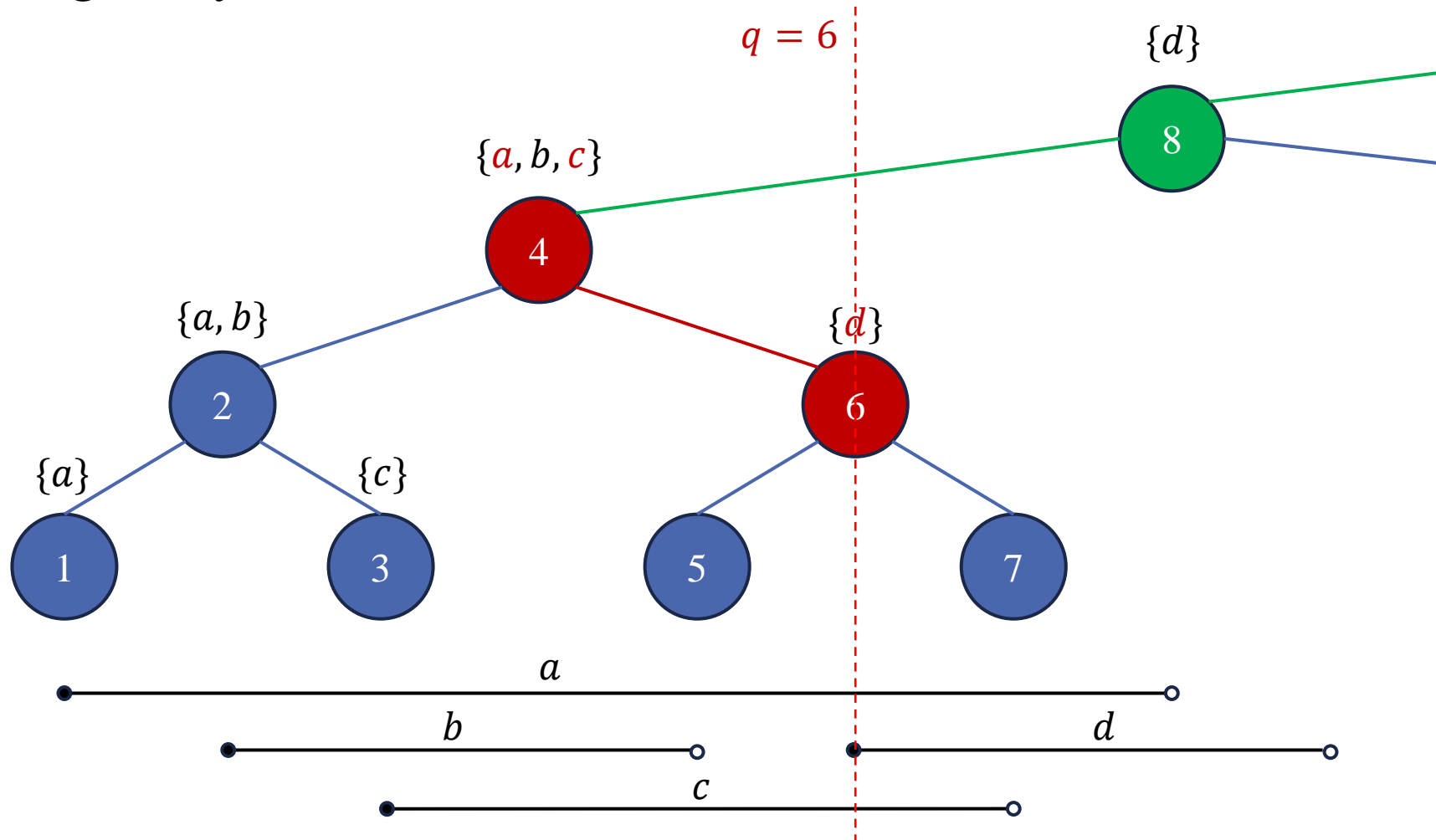
Contribution 1: Online Interval Tree

- Duplicate each node $O(\log t)$ times.



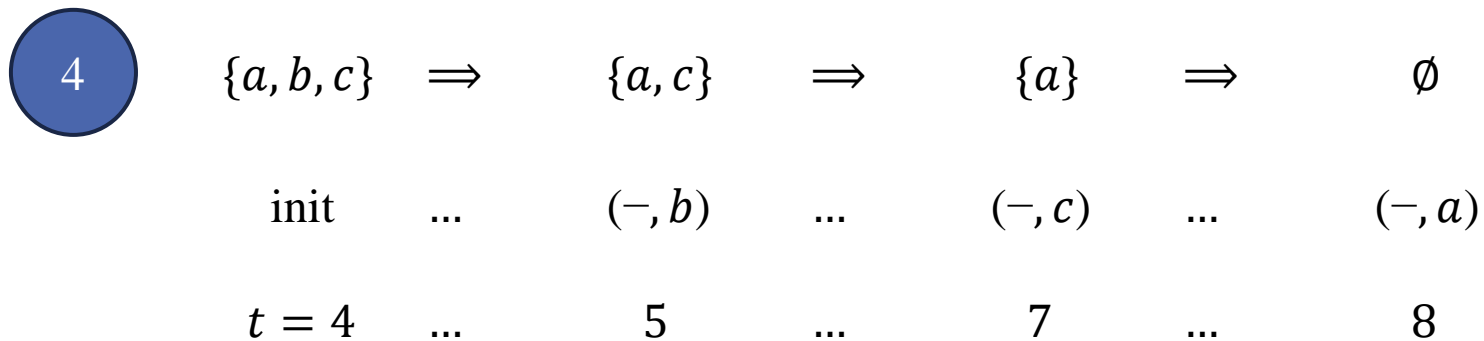
Contribution 1: Online Interval Tree

- Querying: only consider “left” nodes.



Contribution 2: Deletion-only Mechanism

- For each node, it is a deletion-only problem
 - Node i starts with n_i elements, but we can have $n_i \gg n_t$
- Our solution: Track the number of deletions
 - In each round, privately compute: $|D|$, $|D_t^-|$, $F(D)$, $F(D_t^-)$
 - If there are few deletions, then $n_i = O(n_t)$, we simply use $F(D) - F(D_t^-)$
 - If approximately more than $\frac{n_i}{2}$ items are deleted, we allocate new budgets
 - There can only be $O(\log N_t)$ restarts



Conclusion

- There is an (ε, δ) -DP mechanism for arbitrary linear queries on fully dynamic streams with error $\tilde{O}(\sqrt{n_t})$.
- \tilde{O} hides $\frac{1}{\varepsilon}$ and polylog factors in $\frac{1}{\delta}$, $\frac{1}{\beta}$, t , and N_t
- This improves over separating D_t^+ and D_t^- , whose error was $\tilde{O}(\sqrt{N_t})$