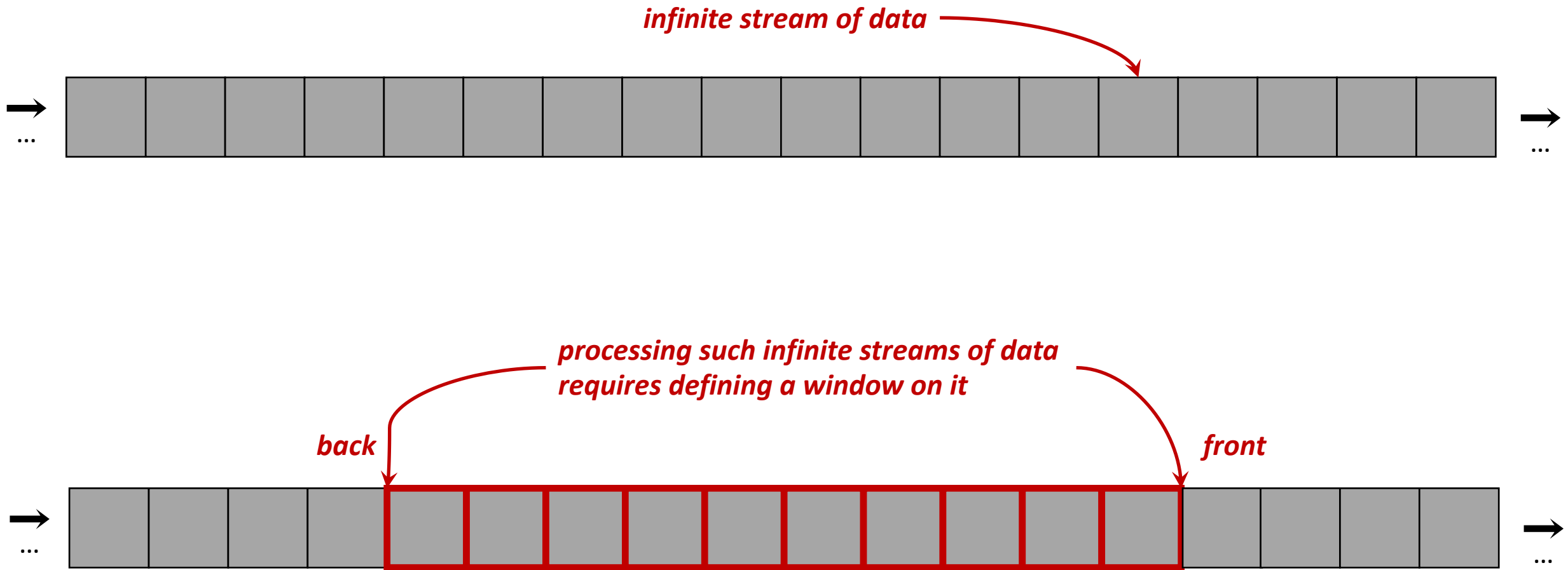# Optimal and General Out-of-Order Sliding Window Aggregation

Kanat Tangwongsan[#], Martin Hirzel[+], **Scott Schneider**[+]
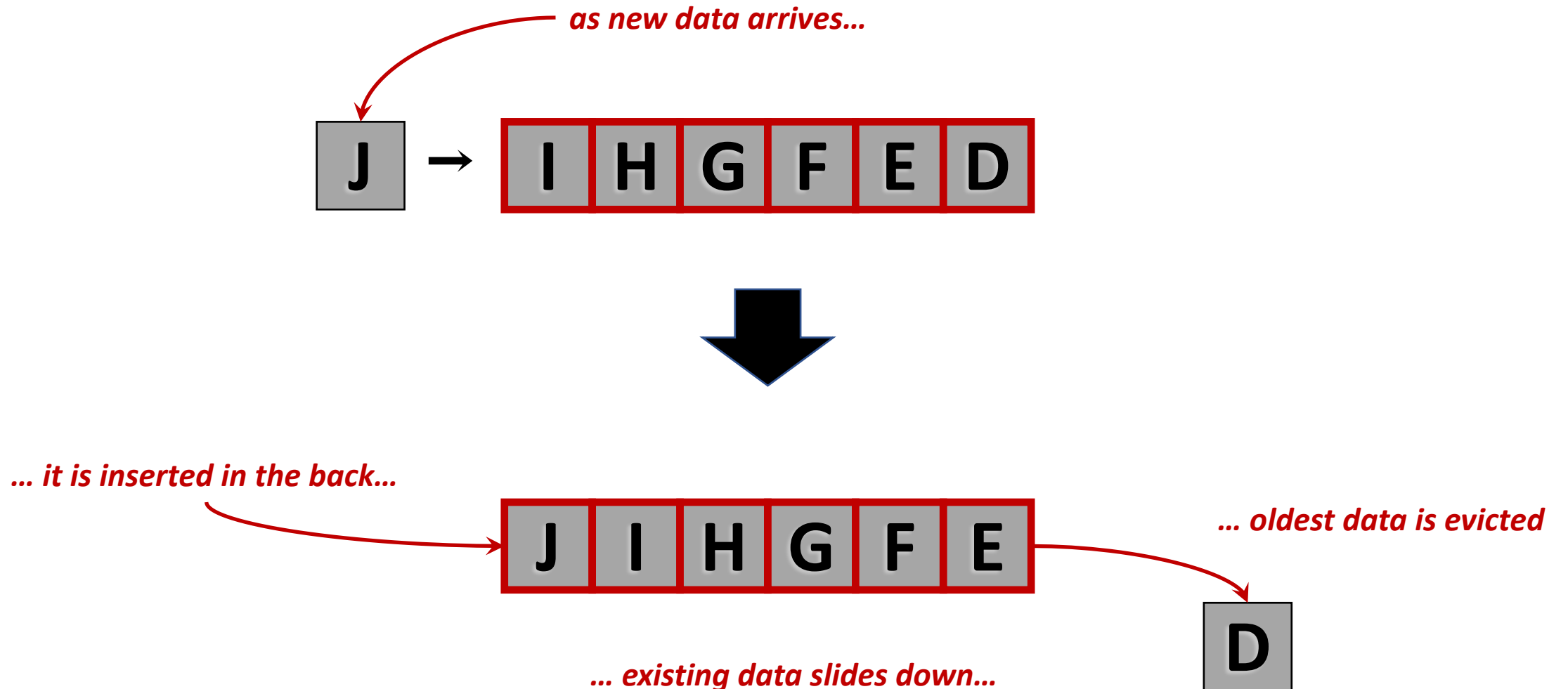
*[#]Mahidol University International College*
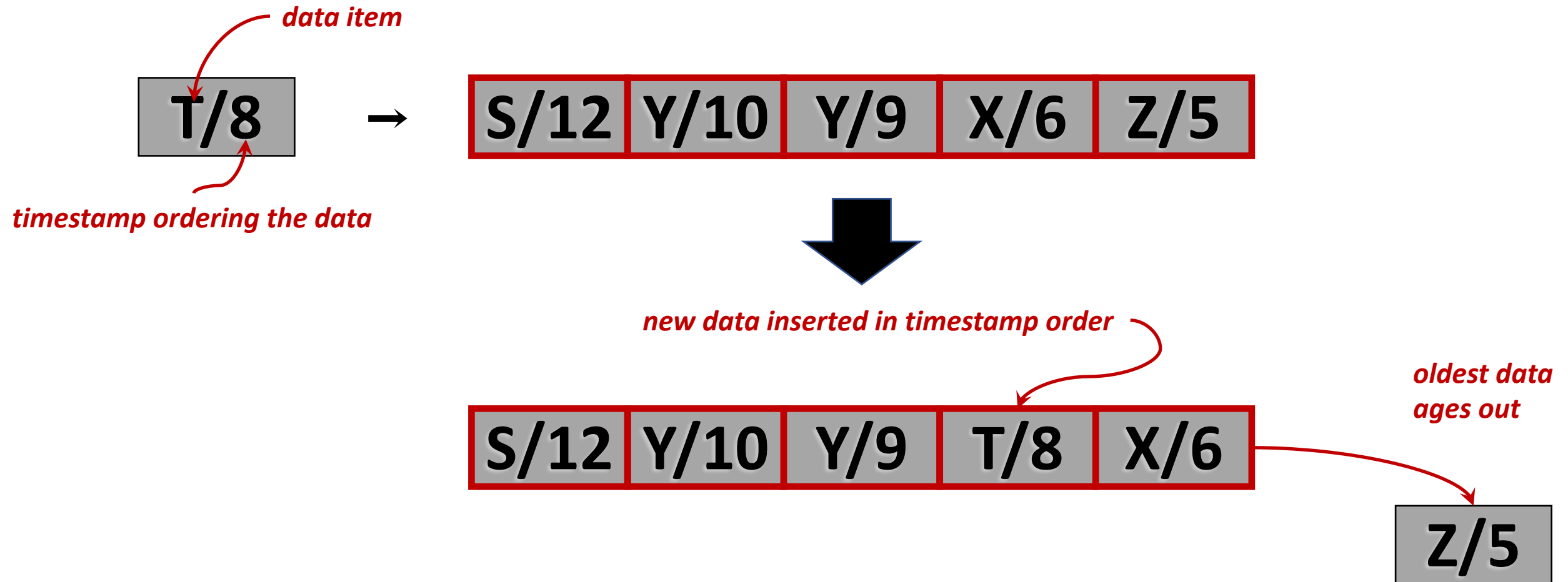*[+]IBM T. J. Watson Research Center*

# Windows on data streams

*infinite stream of data*

*processing such infinite streams of data requires defining a window on it*

*back*

*front*

# Sliding windows

*as new data arrives…*

J → | I | H | G | F | E | D |

*… it is inserted in the back…*

*… oldest data is evicted*

| J | I | H | G | F | E |

D

*… existing data slides down…*

# Real data tends to have *timestamps*



data item

timestamp ordering the data

S/12  Y/10  Y/9  X/6  Z/5

new data inserted in timestamp order

oldest data ages out

S/12  Y/10  Y/9  T/8  X/6

Z/5

# A query on the window is an *aggregation*

| S/12 | Y/10 | Y/9 | C/6 | Z/2 |
|------|------|-----|-----|-----|

min: **C**
mincount: **1**

| C/14 | S/12 | Y/10 | Y/9 | C/6 | Z/2 |
|------|------|------|-----|-----|-----|

min: **C**
mincount: **2**

| C/14 | S/12 | Y/10 | Y/9 | C/6 | A/4 | Z/2 |
|------|------|------|-----|-----|-----|-----|

min: **A**
mincount: **1**

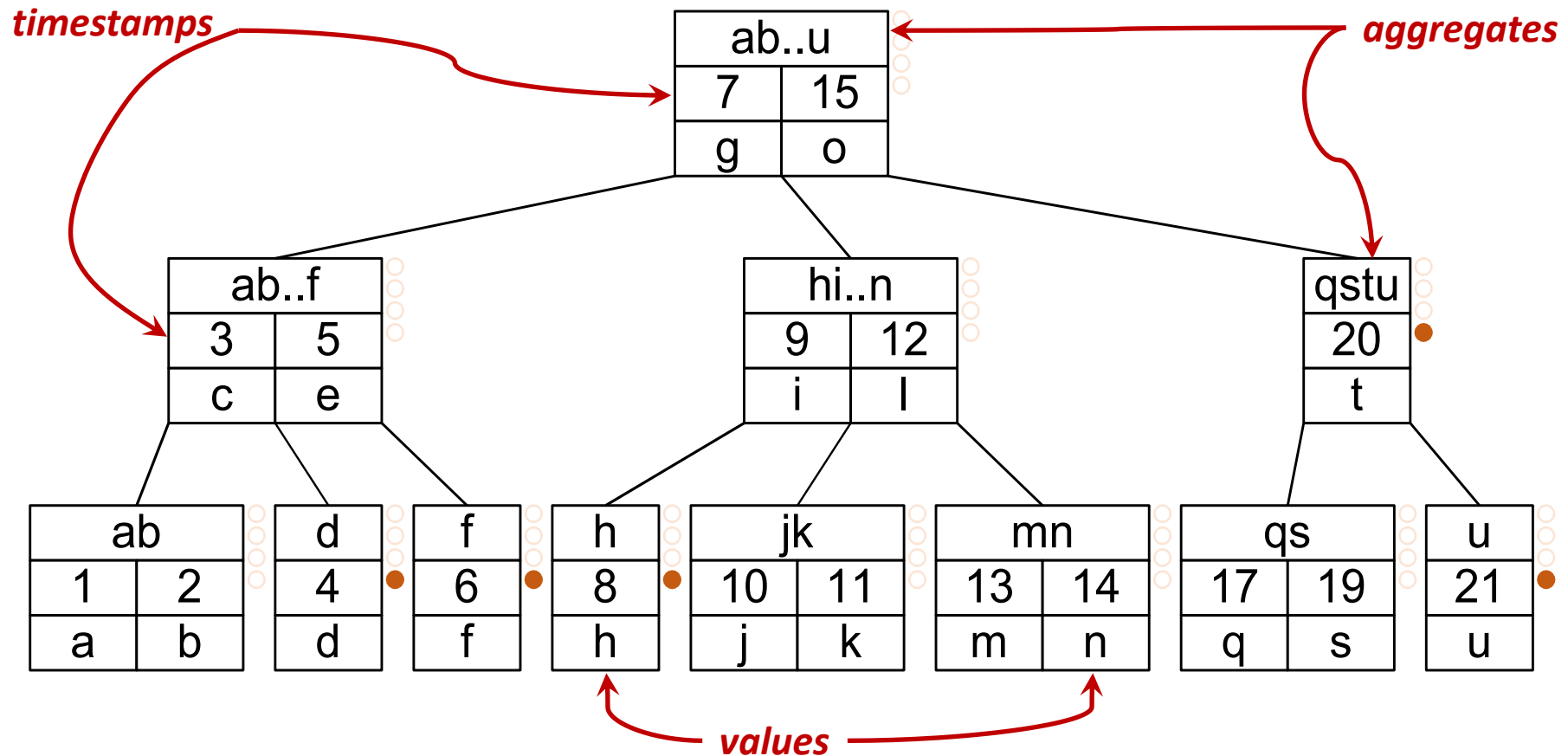| C/14 | S/12 | Y/10 | Y/9 | C/6 |
|------|------|------|-----|-----|

| A/4 | Z/2 |
|-----|-----|

min: **C**
mincount: **2**

# Problem statement

- We want a data structure for a sliding window that can:
  - **insert** data items with timestamps that are out-of-order by distance $d$ in amortized $O(\log d)$, reducing to $O(1)$ when $d=0$
  - **evict** data items based on timestamps in amortized $O(\log d)$, reducing to $O(1)$ when $d=0$
  - **query** the aggregations on the window in worst-case $O(1)$
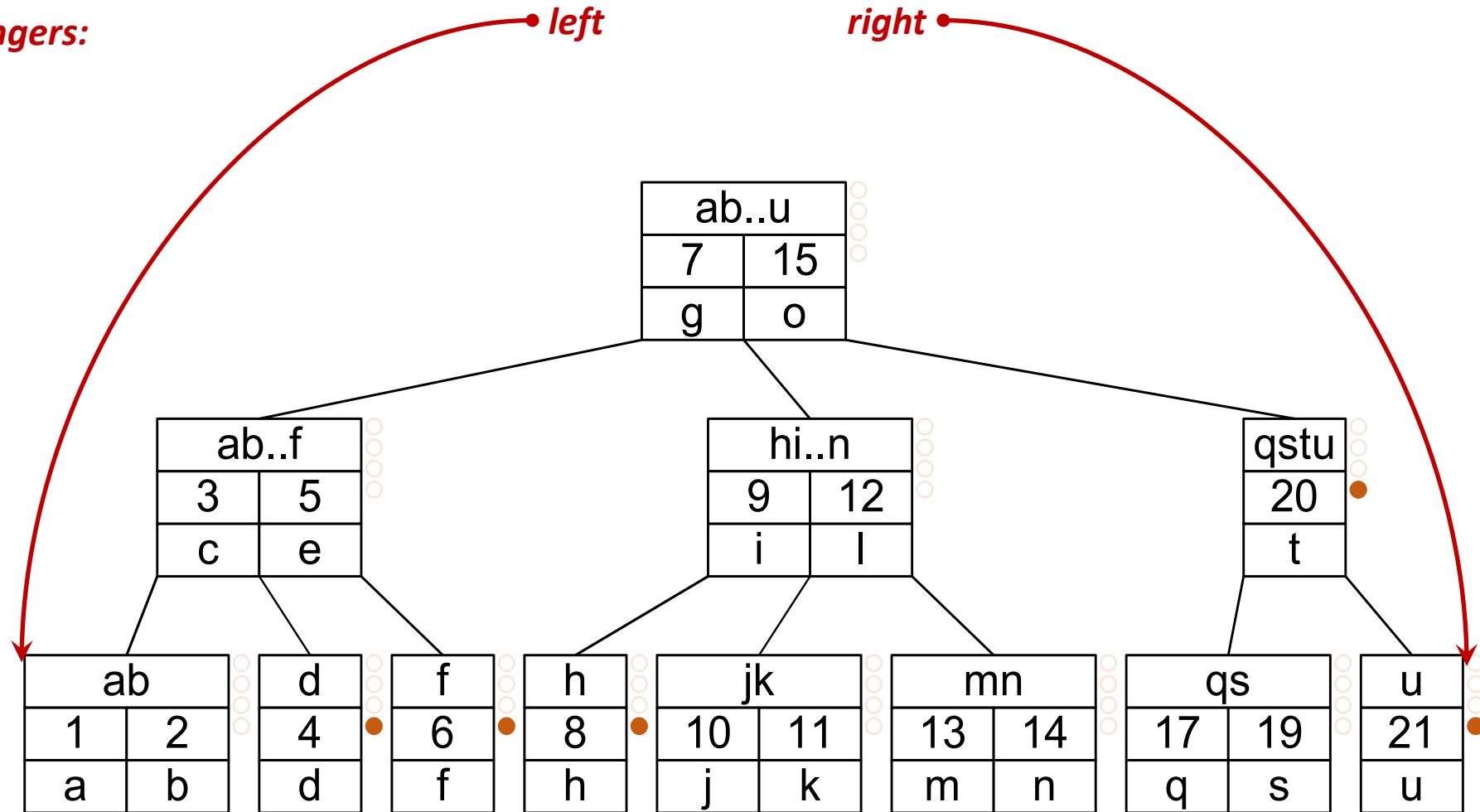- How?
  - B-Trees! (heavily modified)

# FiBA: Finger B-Tree Aggregator

*start with a B-Tree modified with aggregates:*
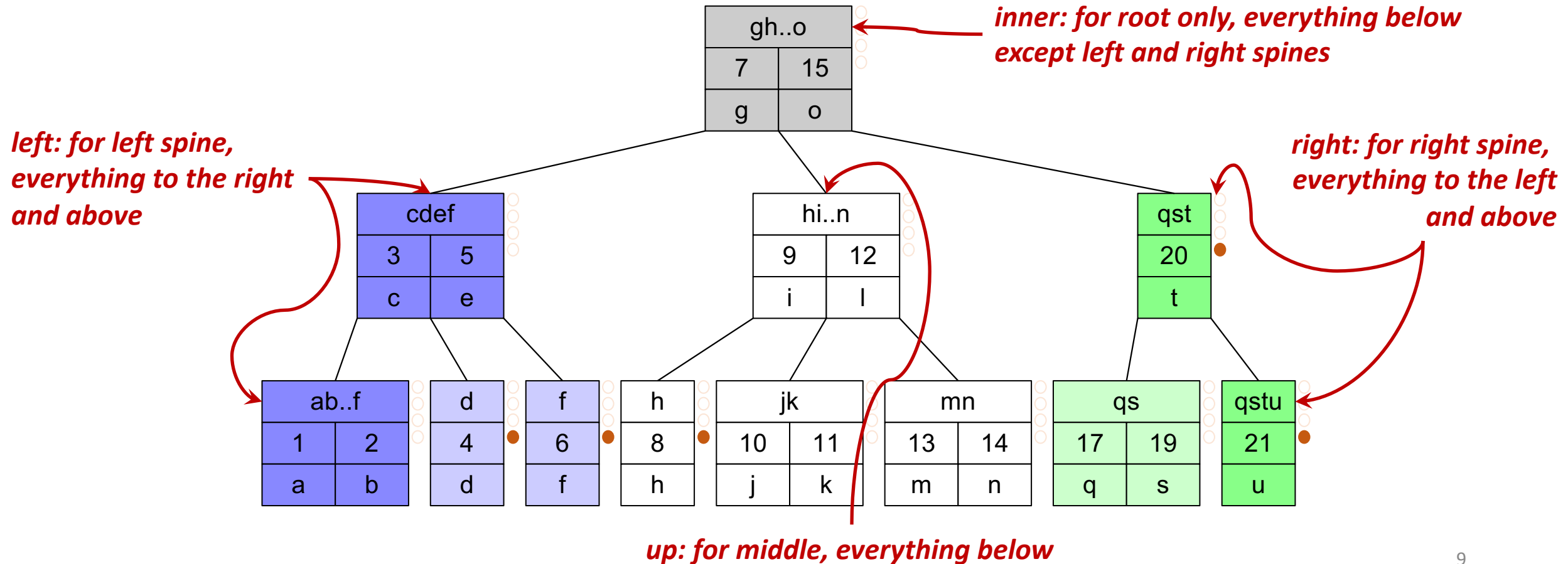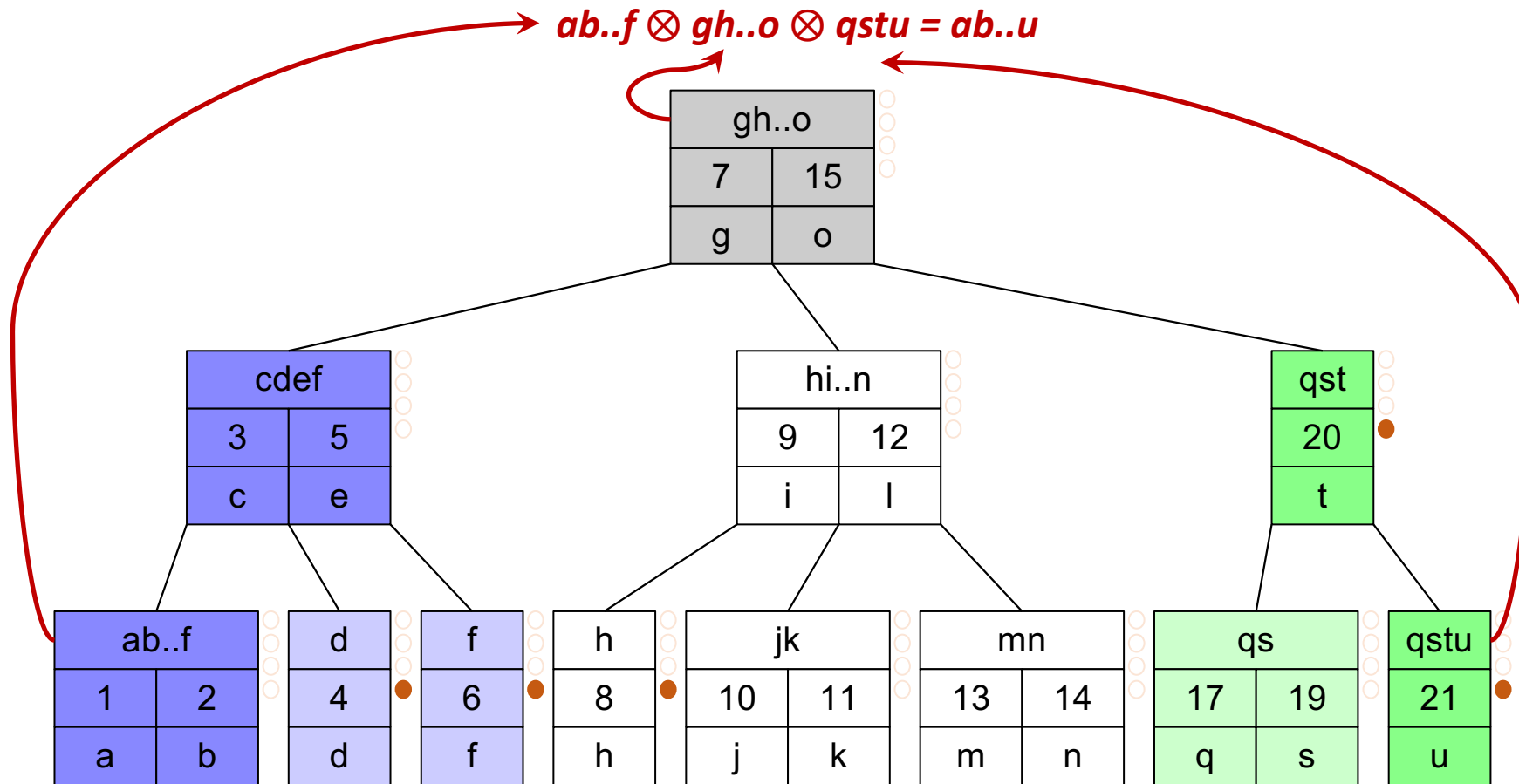
# FiBA: Finger B-Tree Aggregator

add fingers:                    • left          right •

```
                                    ab..u
                                  ┌──────┐
                                  │ 7 │15 │
                                  │ g │ o │
                                  └──────┘

        ab..f                     hi..n                        qstu
      ┌──────┐                  ┌──────┐                      ┌──────┐
      │ 3 │ 5 │                 │ 9 │12 │                     │    │20 │
      │ c │ e │                 │ i │ l │                     │    │ t │
      └──────┘                  └──────┘                      └──────┘

  ab      d      f      h      jk      mn       qs       u
┌────┐  ┌──┐  ┌──┐  ┌──┐  ┌────┐  ┌────┐   ┌────┐   ┌──┐
│1 │2 │  │4 │  │6 │  │8 │  │10│11│  │13│14│   │17│19│   │21│
│a │b │  │d │  │f │  │h │  │j │k │  │m │n │   │q │s │   │u │
└────┘  └──┘  └──┘  └──┘  └────┘  └────┘   └────┘   └──┘
```

8

# FiBA: Finger B-Tree Aggregator

*define position-aware aggregates:*

inner: for root only, everything below
except left and right spines

left: for left spine,
everything to the right
and above

right: for right spine,
everything to the left
and above

up: for middle, everything below

9

# FiBA: Finger B-Tree Aggregator

*answer queries by combining left finger, root and right finger:*

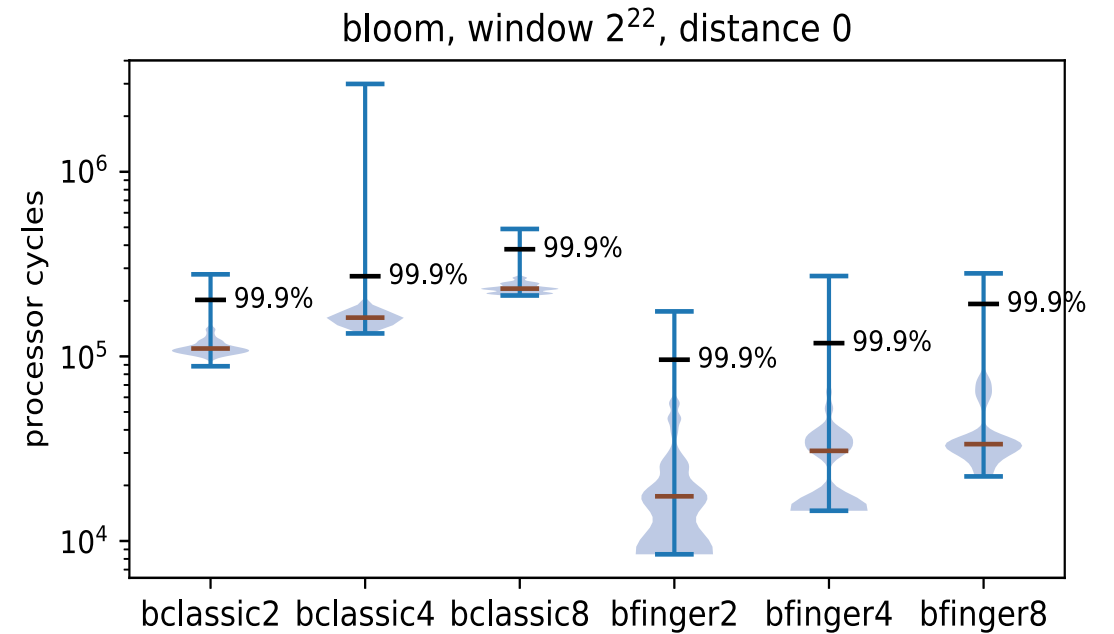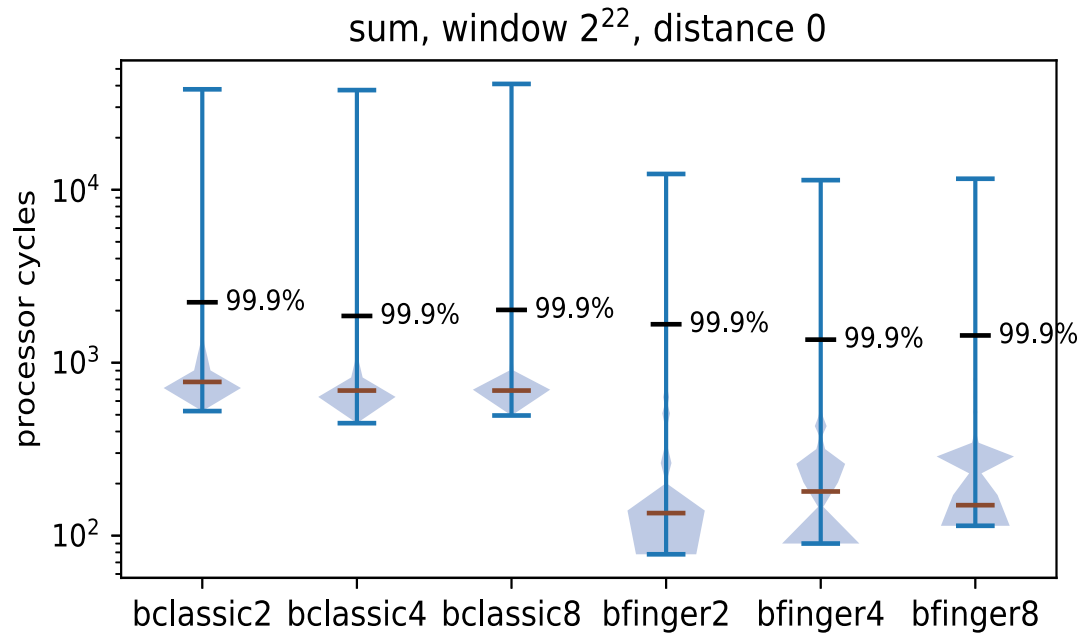$$ab..f \otimes gh..o \otimes qstu = ab..u$$

# FiBA: Intuition for why it works

- **Fingers** allow *O(1)* time access to oldest and youngest
  - without fingers, searching for a data item would be standard *O(log n)*
  - but we are dealing with a time-based window, where we are biased towards inserting at the young end and evicting from the old end
- Specially defined aggregates **shield** sections of the tree
  - updating a value in one section of the tree is unlikely to cause repairs to aggregates elsewhere
- Choice of min and max arity plus lazy splitting and merging **avoids** unnecessary tree rebalancing
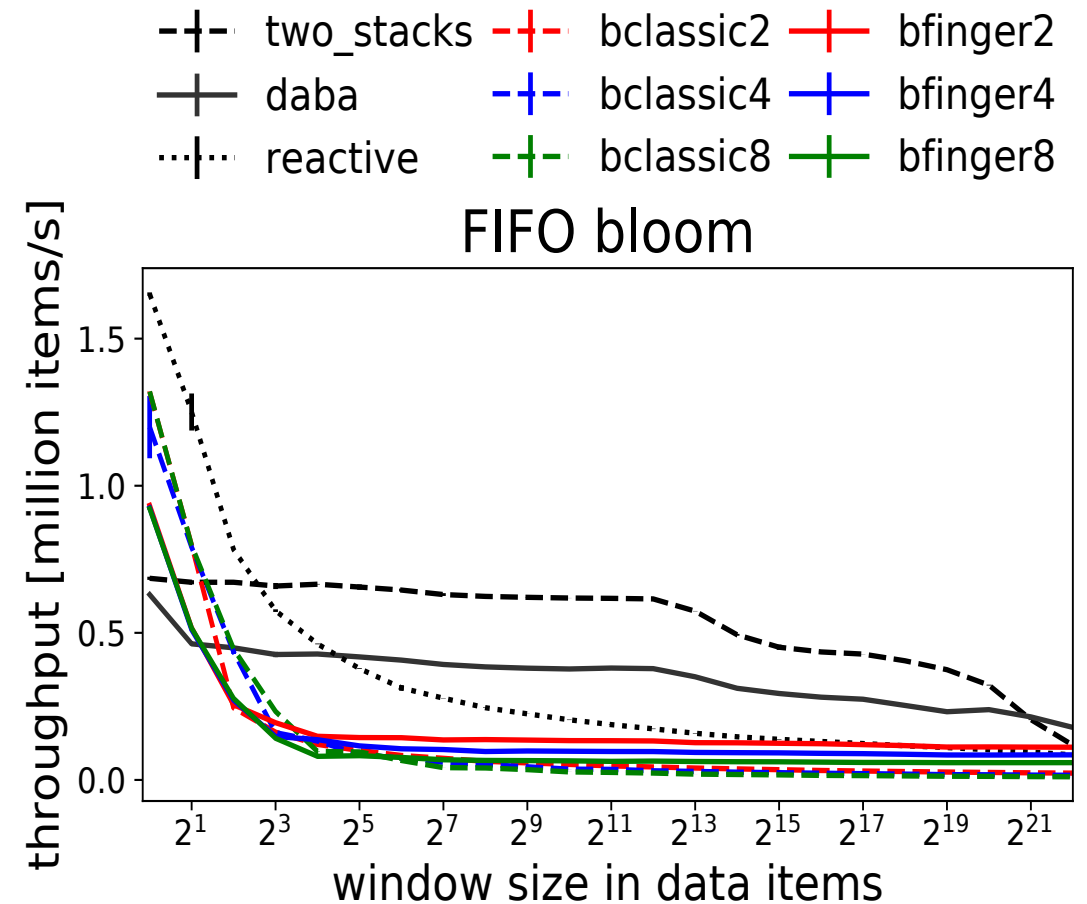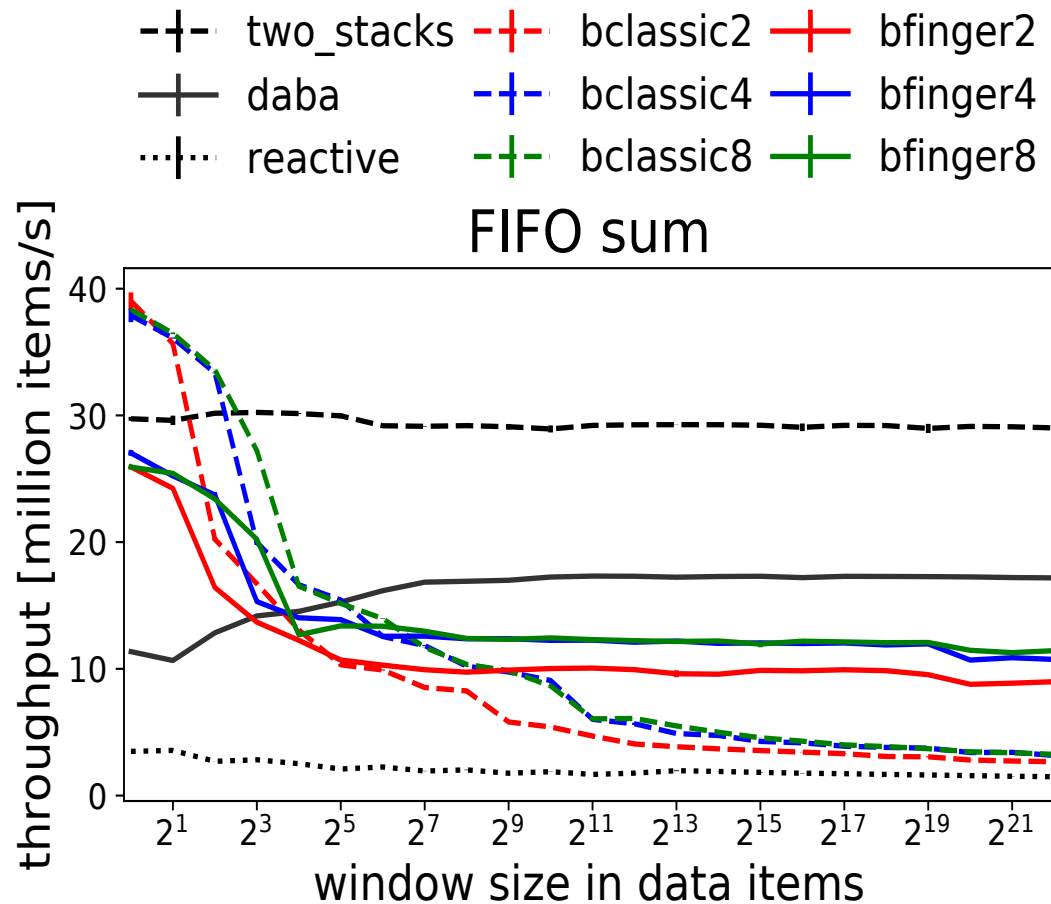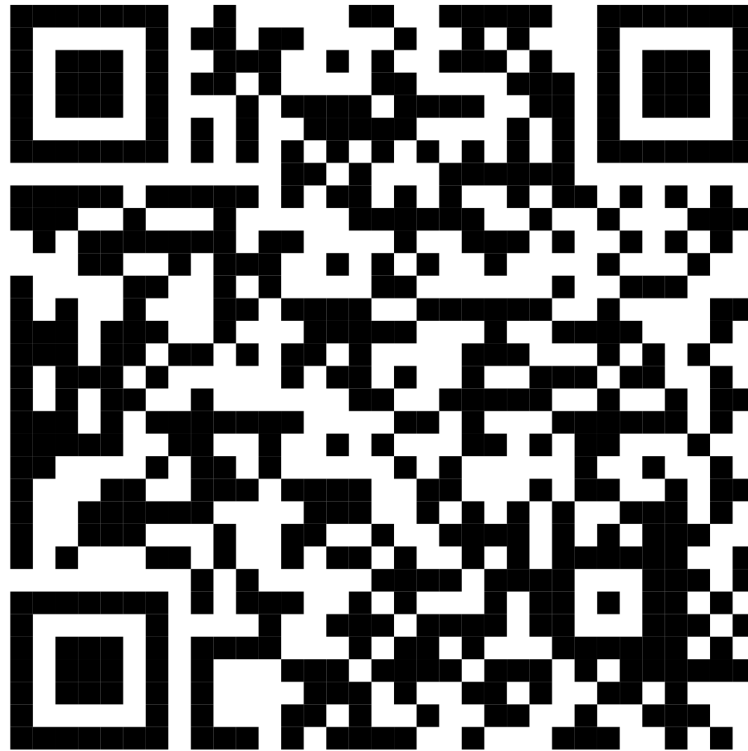
# Out-of-Order Throughput



OoO sum, window $2^{22}$

OoO bloom, window $2^{22}$

Legend: bclassic2, bclassic8, bfinger4, bclassic4, bfinger2, bfinger8

x-axis: out-of-order distance

y-axis: throughput [million items/s]

# In-Order Latency



sum, window $2^{22}$, distance 0

bloom, window $2^{22}$, distance 0

# In-Order Throughput

# Questions?

# Backup

# Out-of-Order Latency



sum, window $2^{22}$, distance $2^{20}$

bloom, window $2^{22}$, distance $2^{20}$